Bayesian statistical methods for realizing personalized medicine

Junwei Shen

Doctor of Philosophy

Department of Epidemiology, Biostatistics and Occupational Health

McGill University Montréal, Québec July 2024

A thesis submitted to McGill University in partial fulfillment of the requirements of the degree of Doctor of Philosophy © Copyright Junwei Shen, 2024

Dedication

This thesis is dedicated to my mother.

Acknowledgements

I am indebted to several individuals who have supported me in completing this work. First and foremost, I would like to express my deepest gratitude to my supervisors, Dr. Erica Moodie and Dr. Shirin Golchi, for their unwavering support and continuous encouragement throughout my PhD journey. I would like to thank Erica for her availability, prompt responses to my requests and questions, and her dedication to the students. Her insightful feedback and timely guidance have been crucial to my progress and success. I am grateful to Shirin for introducing me to the fascinating areas of Bayesian methods and clinical trials. Her expertise and enthusiasm in these fields have been instrumental in the completion of this work.

I sincerely thank Dr. David Benrimoh for providing information about the motivating trial for my first project and sharing precious experience about data requests for my third project. I acknowledge the financial support provided by Mitacs, by my supervisors, by the Department of Epidemiology and Biostatistics, and by the Fonds de recherche du Québec - Nature et technologies. I also acknowledge the computational support provided by Digital Research Alliance of Canada.

I would like to extend my appreciation to Dr. James Hanley, Dr. Alexandra M. Schmidt, Dr. Michal Abrahamowicz, Dr. Andrea Benedetti, and Dr. Qihuang Zhang for broadening my knowledge and generously sharing their experience. I thank all administrative staff in the Department of Epidemiology and Biostatistics, specially André Yves Gagnon and Katherine Hayden for their unhesitating help whenever I encountered administrative issues. I am grateful to my fellow students Haoyu Wu, James Willard, Vanessa Mcnealis, Marc Parsons, Ajmary Jaman, Xianglin Zhao, Mingchi Xu, Helen Bian, Peiyuan Huang, and Renjie Peng, for many interesting talks and making my life at McGill enjoyable.

Many thanks to my friends Zhenyu Ouyang for our wonderful travels which took me back

to our teenage years and cheered me up, and Gengyuan Zhang for the casual and friendly communications during difficult times.

My heartfelt appreciation goes to my mother for her endless love, support, and understanding. I could not have completed this without her.

Preface

This manuscript-based PhD thesis contains new research addressing challenges in design and analysis of studies aimed at personalized medicine. The thesis consists of six chapters: an introduction, a literature review, three chapters corresponding to three stand-alone manuscripts, and a conclusion. A complete bibliography is presented after the appendices. Chapter 1 gives a general overview of this thesis, and Chapter 2 provides a comprehensive review of the concepts and theories used in this thesis. These chapters were written entirely by Junwei Shen (JS) and corrected by Erica E. M. Moodie (EEMM) and Shirin Golchi (SG). Chapter 3 was conceptualized by JS, SG, EEMM, and David Benrimoh (DB). JS conducted the methodology, designed and conducted the simulation study, and wrote the manuscript under the guidance of SG and EEMM. DB introduced the motivating trial. SG, EEMM, and DB corrected and edited this chapter.

For Chapter 4, JS, EEMM, and SG contributed to the conceptualization of the research. JS carried out the methodological work, simulation study, real data analysis, and writing. SG and EEMM provided substantial guidance in developing the theory and analyzing the real data. The work was reviewed and edited by EEMM and SG.

The work in Chapter 5 was conceptualized by JS, EEMM, and SG. All methodological work, simulation study, real data analysis, and writing were conducted by JS with EEMM and SG as advisors and editors.

The conclusion and the discussion of limitations and future work in Chapter 6 were written by JS and edited by EEMM and SG.

Abstract

The personalized medicine (PM) approach provides patients with treatments that align most closely with their individual profiles. Challenges can arise when learning about individualized treatments, both in the design and the analysis of clinical trials aimed at learning about PM. In my thesis, I use Bayesian approaches to tackle problems in both of these domains, first focusing on the design of a randomized trial suitable for learning about proposed PM strategies, and next, focusing on estimation approaches that can be applied to multi-center studies or multiple randomized studies to leverage large datasets which may be subject to restrictions on data-sharing across sites or trials.

To integrate PM into routine clinical practice, it is imperative to evaluate the efficacy and safety of tools for individualized treatments in randomized controlled trials. Many such tools are best provided at a group level, making cluster-randomized trials (CRTs) an appealing option. However, CRTs are less efficient compared with individually randomized trials and flexible designs for CRTs are not common. In the first manuscript, motivated by a CRT designed to assess the effectiveness of a clinical decision support system for physicians, I develop two Bayesian group sequential designs for CRTs to allow for early stopping for efficacy at pre-planned interim analyses. One design sequentially enrolls the *clusters*, and individual participants for each cluster are recruited all at once. The other enrolls all clusters at one time, but the *individual participants* for each cluster are sequentially enrolled. I explore and compare the design operating characteristics of the two designs in simulations, and provide some practical recommendations.

The second challenge that I tackle arises in estimation of individualized treatment rules (ITRs) from multiple sources or data sites. ITR estimation can suffer from low power when attempting to detect the often subtle variability in treatment effect, making collaboration across sites or pooling of data across trials attractive. However, sensitive individual information may sometimes not be shared due to policy restrictions, motivating approaches that

avoid individual-level data sharing for ITR estimation. In my second manuscript, I adopt a two-stage Bayesian meta-analysis approach to estimate ITRs using multisite data without disclosing individual-level data beyond the sites. However, different data sites may recruit from different populations, making it infeasible to estimate identical models or all parameters of interest at all sites, and the number of non-zero parameters in the model for the treatment rule may be small. Simulations show that the proposed approach can provide consistent estimates of the parameters which characterize the optimal ITR. I apply Bayesian meta-analysis with shrinkage priors to estimate the optimal Warfarin dose strategy using the International Warfarin Pharmacogenetics Consortium data.

In the third manuscript, drawing from the network meta-analysis literature, I extend methods in the second manuscript to synthesize information across multiple, independent randomized trials where it is possible that not all trials included the same set of treatment randomization options. An application of the proposed method to data from three depression studies: Sequenced Treatment Alternatives to Relieve Depression (STAR*D), Establishing Moderators and Biosignatures of Antidepressant Response in Clinical Care (EMBARC), and Research Evaluating the Value of Augmenting Medication with Psychotherapy (REVAMP) is also presented.

My thesis contributes to statistical trial literature by proposing new efficient designs that may be appealing to clinical scientists, and to ITR literature by developing novel approaches that allow researchers to estimate optimal strategies using siloed data sources.

Abrégé

L'approche de médecine personnalisée (MP) offre aux patients des traitements qui correspondent le mieux à leurs profils individuels. Des défis peuvent survenir lors de l'apprentissage des traitements individualisés; optimaux tant dans la conception que dans l'analyse des essais cliniques visant à apprendre sur la MP. Dans ma thèse, j'utilise des approches bayésiennes pour aborder des problèmes dans ces deux domaines, en me concentrant d'abord sur la conception d'un essai randomisé adapté à l'apprentissage sur les stratégies proposées de MP, puis sur les approches d'estimation applicables aux études multicentriques ou à la combinaison d'études randomisées pour exploiter de grands ensembles de données qui peuvent être soumis à des restrictions de partage des données entre sites ou essais.

Pour intégrer la MP dans la pratique clinique courante, il est impératif d'évaluer l'efficacité et la sécurité des outils de traitements individualisés dans des essais contrôlés randomisés. Beaucoup de ces outils sont mieux fournis à un niveau de groupe, ce qui rend des essais randomisés par grappes (CRTs) attrayants. Cependant, les CRTs sont moins efficaces que les essais randomisés individuels et des devis flexibles pour des CRTs sont rares. Dans le premier manuscrit, motivé par un CRT conçu pour évaluer l'efficacité d'un système de soutien à la décision clinique pour des médecins, je développe deux plans séquentiels bayésiens de groupes pour des CRTs permettant un arrêt précoce pour efficacité lors des analyses intermédiaires planifiées. Une conception recrute séquentiellement des *regroupements*, et les participants individuels pour chaque groupe sont recrutés en une seule fois. L'autre recrute toutes les groupes en une fois, mais les *participants individuels* pour chaque groupe sont recrutés séquentiellement. J'explore et compare les caractéristiques de fonctionnement des deux conceptions dans des simulations et fournis quelques recommandations pratiques.

Le deuxième défi que j'aborde concerne l'estimation des règles de traitement individualisées (ITR) à partir de multiples sources ou sites de données. L'estimation des ITR peut souffrir d'une faible puissance lorsqu'on tente de détecter la variabilité souvent fiable de l'effet du traitement, ce qui rend la collaboration entre sites ou le regroupement de données entre essais attrayant. Cependant, des informations individuelles sensibles peuvent parfois ne pas être partagées en raison de restrictions administratives, motivant des approches qui évitent le partage des données individuelles pour l'estimation des ITR. Dans mon deuxième manuscrit, j'adopte une approche de méta-analyse bayésienne en deux étapes pour estimer les ITR en utilisant des données multisites sans divulguer les données individuelles au-delà des sites. Cependant, les différents sites de données peuvent recruter à partir de différentes populations, ce qui rend infaisable l'estimation de modèles identiques ou de tous les paramètres d'intérêt sur tous les sites. De plus, le nombre de paramètres non nuls dans le modèle de la règle de traitement peut être faible. Les simulations montrent que l'approche proposée peut fournir des estimations cohérentes des paramètres qui caractérisent l'ITR optimal. J'applique la méta-analyse bayésienne avec des méthodes de rétrécissement pour estimer la stratégie optimale de dose de Warfarine en utilisant les données du Consortium International de Pharmacogénétique de la Warfarine.

Dans le troisième manuscrit, en m'inspirant de la littérature sur la méta-analyse en réseau, j'étends je généralise les méthodes dans le deuxième manuscrit pour synthétiser l'information à travers plusieurs essais randomisés indépendants où il est possible que tous les essais n'aient pas inclus le même ensemble d'options de randomisation de traitement. Une application de la méthode proposée aux données de trois études sur la dépression: Sequenced Treatment Alternatives to Relieve Depression (STAR*D), Establishing Moderators and Biosignatures of Antidepressant Response in Clinical Care (EMBARC), et Research Evaluating the Value of Augmenting Medication with Psychotherapy (REVAMP) est présentée.

Ma thèse contribue à la littérature sur les essais statistiques en proposant de nouvelles conceptions efficaces qui peuvent attirer les scientifiques cliniques, et à la littérature sur les ITR en développant de nouvelles approches permettant aux chercheurs d'estimer des stratégies optimales en utilisant des sources de données cloisonnées.

Table of contents

1	Intr	oduction	1
2	Lite	rature review	5
	2.1	Cluster-randomized trials	5
	2.2	Group sequential designs	7
		2.2.1 Frequentist methods	8
		2.2.2 Bayesian methods	9
	2.3	Individualized treatment rules	10
	2.4	Individual participant data meta-analysis	17
		2.4.1 Two-stage approach	19
		2.4.2 One-stage approach	21
	2.5	Network meta-analysis	22
	2.6	Summary	26
3	Bay	esian group sequential designs for cluster-randomized trials	27
	3.1	Introduction	30
	3.2	Motivating setting	32
	3.3	Methods	33
		3.3.1 Two Bayesian group sequential designs for cluster-randomized trials .	33
		3.3.2 Early stopping at interim analysis	35

	3.4	Simulation studies	39
	3.5	Discussion	44
4	Spa	rse two-stage Bayesian meta-analysis for individualized treatments	49
	4.1	Introduction	52
	4.2	Methods	55
		4.2.1 Preliminaries	55
		4.2.2 Two-stage IPD meta-analysis	57
		4.2.3 Sparsity	59
	4.3	Simulation studies	61
		4.3.1 Overview	61
		4.3.2 Results	62
	4.4	Estimating an optimal Warfarin dose strategy	63
		4.4.1 Context and data source	63
		4.4.2 Analyses and results	67
	4.5	Discussion	68
5	Two	o-stage Bayesian network meta-analysis of individualized treatment rules	
	for	multiple treatments with siloed data	75
	5.1	Introduction	79
	5.2	Methods	81
		5.2.1 Preliminaries	81
		5.2.2 Two-stage Bayesian network meta-analysis	84
	5.3	Simulation studies	88
		5.3.1 Aims	88
		5.3.2 Data-generating mechanisms	89
		5.3.3 Estimands, methods, and performance metrics	93
		5.3.4 Results	95

	5.4	Estimating individualized depression treatment	96
	5.5	Discussion	100
6	Con	clusion	109
	6.1	Summary	109
	6.2	Limitations and Future Work	111
	6.3	Concluding Remarks	118
Aj	ppen	dices	119
A	Арр	pendix to Manuscript 1	120
	A.1	Simulation for binary outcomes	120
	A.2	Simulation results for continuous outcomes with cluster size m=16	127
в	Арр	pendix to Manuscript 2	129
	B.1	Link with a one-stage approach	129
	B.2	Data sparsity: A second toy example	134
	B.3	Simulation studies: ADEMP reporting	135
		B.3.1 Aims	135
		B.3.2 Data-generating mechanisms	135
		B.3.3 Estimands, methods, and performance metrics	140
	B.4	Model details in a sparse data setting in simulation	141
	B.5	Additional simulation results	143
	B.6	Analysis of Warfarin data	155
		B.6.1 Visual inspection of the overlap assumption	155
		B.6.2 Details of the models in the Warfarin analysis	156
С	App	pendix to Manuscript 3	163
_	_		

References

168

List of Tables

3.1	Simulation parameters for continuous and binary outcome	41
4.1	Simulation results for the many covariates setting.	74
5.1	Sites included in each network.	91
5.2	Coefficients in multinomial probabilities for treatment assignment	92
5.3	Selected simulation results when the data are generated under an unstructured	
	between-site heterogeneity model	96
5.4	Selected simulation results when the data are generated under the assumption	
	of common between-site heterogeneity.	97
5.5	Patient characteristics for STAR*D, EMBARC, and REVAMP studies	107
5.6	Blip parameter estimates (posterior medians) and the 95% posterior credible	
	intervals for the real data application.	108
B.1	Parameters in the propensity score model for binary treatment simulations in	
	different scenarios	136
B.2	Site-specific blip function parameter estimates $\hat{\xi}_i^{(1)}$ in the Warfarin analysis.	160
B.3	Site-specific blip function parameter estimates $\hat{\xi}_i^{(2)}$ in the Warfarin analysis.	161
C.1	Simulation results of \hat{psi}_{h0} , $h = 3,, 7$, for the setting where the data are	
	generated under an unstructured between-site heterogeneity model	164

- C.2 Simulation results of $\hat{\psi}_{h1}$, h = 3, ..., 7, for the setting where the data are generated under an unstructured between-site heterogeneity model. 165
- C.3 Simulation results of $\hat{\psi}_{h0}$, h = 3, ..., 7, for the setting where the data are generated under the assumption of common between-site heterogeneity. . . . 166

List of Figures

2.1	Illustration of the network structure for the hypothetical example	23
3.1	An illustrative example of the two designs.	35
3.2	Simulation results of false positive rate for a continuous outcome setting with	
	single interim analysis.	42
3.3	Simulation results of power for a continuous outcome setting with single in-	
	terim analysis	43
3.4	Simulation result of false positive rate for a continuous outcome setting with	
	multiple interim analyses and $m = 8. \ldots \ldots \ldots \ldots \ldots \ldots$	43
3.5	Selected simulation results of power for a continuous outcome setting with	
	multiple interim analyses and $m = 8$	48
4.1	Simulation results of $\hat{\psi}_0$ for the small sample size and the binary treatment	
	setting with a half-Cauchy (0,1) prior.	64
4.2	Simulation results of dVF for the small sample size and the binary treatment	
	setting with a half-Cauchy (0,1) prior.	65
4.3	Simulation results for the small sample size and the sparse data setting	66
5.1	Graphics of simulated networks.	90
5.2	Network structure of analysis of STARD, EMBARC, and REVAMP data	98

A.1	Simulation results of false positive rate for a binary outcome setting with	
	single interim analysis.	121
A.2	Simulation results of power for a binary outcome setting with single interim	
	analysis	122
A.3	Simulation results of false positive rate for a binary outcome setting with	
	multiple interim analysis and $m = 8$	123
A.4	Simulation results of power for a binary outcome setting with multiple interim	
	analyses and $m = 8$	124
A.5	Simulation results of false positive rate for a binary outcome setting with	
	multiple interim analyses and $m = 16.$	125
A.6	Simulation results of power for a binary outcome setting with multiple interim	
	analyses and $m = 16. \ldots$	126
A.7	Simulation results of false positive rate for a continuous outcome setting with	
	multiple interim analyses and $m = 16. \dots \dots \dots \dots \dots \dots \dots$	127
A.8	Simulation results of power for a continuous outcome setting with multiple	
	interim analyses and $m = 16$	128
B.1	Additional simulation results of $\hat{\psi}_0$ for the small sample size and the binary	
	treatment setting.	145
B.2	Additional simulation results of dVF for the small sample size and the binary	
	treatment setting.	146
B.3	Simulation results of $\hat{\psi}_1$ for the small sample size and the binary treatment	
	setting.	147
B.4	Simulation results of $\hat{\psi}_0$ for the large sample size and the binary treatment	
	setting	148
B.5	Simulation results of $\hat{\psi}_1$ for the large sample size and the binary treatment	
	setting	149
B.6	Simulation results of dVF for the large sample size and binary treatment setting	.150

B.7	Simulation results for the small sample size and the continuous treatment	
	setting	51
B.8	Simulation results for the large sample size and the continuous treatment setting.1	52
B.9	Additional simulation results for the small sample size and the sparse data	
	setting	.53
B.10	Simulation results for the large sample size and the sparse data setting 1	54
B.11	Distribution of the generalized propensity score by dose group in the Warfarin	
	analysis	55

Abbreviations

- ${\bf BUP}$ bup ropion
- ${\bf BUS}$ buspirone
- **CDSS** clinical decision support system
- \mathbf{CIT} citalopram
- \mathbf{CIT} + \mathbf{BUP} CIT augmented with BUP
- \mathbf{CRT} cluster-randomized trial
- **DTR** dynamic treatment regime
- \mathbf{dVF} difference between the value functions
- dWOLS dynamic weighted ordinary least squares
- **EMBARC** Establishing Moderators and Biosignatures of Antidepressant Response for Clinical Care
- $\mathbf{ESCIT} \ \mathrm{escitalopram}$
- HDRS-17 17-item Hamilton Depression Rating Scale
- $\mathbf{ICC}\,$ intra-cluster correlation coefficient
- \mathbf{INR} international normalized ratio

 ${\bf IPD}\,$ individual participant data

IPTW inverse probability of treatment weighted

 $\mathbf{ITR}\,$ individualized treatment rule

MCMC Markov Chain Monte Carlo

MDD major depressive disorder

 $\mathbf{MVN}\,$ multivariate normal distribution

 ${\bf PM}\,$ personalized medicine

QIDS-16 16-item Quick Inventory of Depressive Symptomatology

REVAMP Research Evaluating the Value of Augmenting Medication with Psychotherapy

 \mathbf{SER} sertraline

STAR*D Sequenced Treatment Alternatives to Relieve Depression

SUTVA stable unit treatment value assumption

 \mathbf{VEN} venlafaxine

Chapter 1

Introduction

Treatment responses can vary among patient subgroups, which are defined by combinations of factors such as demographics, genetic makeup, and clinical measures. To enhance health outcomes, personalized medicine (PM) aims to provide treatments that are closely aligned with individual patient profiles by leveraging this variability in treatment effects (Kosorok and Laber, 2019). Implementing individualized treatments in real-world settings requires statistical analysis of clinical trial data to establish rules for tailored treatment assignments, as well as clinical trial designs to evaluate the clinical utility and efficacy of tools that apply these rules in practice. In my thesis, I tackle problems in both domains. I first focus on the design of a randomized trial suitable for learning about proposed PM strategies. Next, I explore estimation approaches that can be applied to multi-center studies or multiple randomized studies to leverage large datasets, which may be subject to restrictions on data sharing across sites or trials.

Bayesian statistics offers a powerful tool for integrating prior knowledge, data, and uncertainty from heterogeneous sources. Statistical modelling within the Bayesian paradigm is increasingly gaining popularity in PM. Bayesian adaptive designs that incorporate individual characteristics in design components (e.g., the randomization ratio) or analyses have been proposed to identify subgroup treatment effect and more efficiently treat patients (e.g., Zhou et al., 2008; Zhang et al., 2021). Bayesian approaches to estimate personalized treatment strategies have also been proposed in the literature (e.g., Arjas and Saarela, 2010; Murray et al., 2018; Logan et al., 2019; Hahn et al., 2020). My thesis focuses on Bayesian approaches for the design and analysis of studies in PM.

With the increasing number of scientific findings on PM, information technology platforms and systems have been developed to deliver targeted information to guide clinical decisions. One example is the clinical decision support system (CDSS) designed by Aifred Health for physicians treating patients with depression (Benrimoh et al., 2018). The CDSS takes patient characteristics, such as sociodemographic and clinical information, as inputs and employs a deep-learning model to recommend a range of possible treatments based on their efficacy. Physicians can then determine how to utilize this information in their clinical decision-making processes for individual patients. However, the real-world efficacy of the CDSS remains to be evaluated, which motivates the research presented in Chapter 3.

Randomized controlled trials are the gold standard for evaluating efficacy. Since the CDSS is used at the physician level, a cluster-randomized trial (CRT) that randomizes physicians as "clusters" is appropriate. Compared to individually randomized trials with the same sample size, CRTs often suffer from low power to detect true differences between treatment arms due to the reduced variability of responses in clustered samples, resulting from the positive correlation between subjects within the same cluster (Donner and Klar, 2004). To overcome this and potentially reduce the required sample sizes, I propose two Bayesian group sequential designs for CRTs in Chapter 3. The philosophy behind group sequential designs is that groups of participants are sequentially enrolled, with pre-planned interim assessments of cumulative data over the course of the trial. These interim assessments provide opportunities for early trial termination (Pocock, 1977). Group sequential designs naturally fit into the Bayesian framework, as estimates are updated based on the accumulated information from interim data (Chevret, 2012). The two Bayesian group sequential designs I consider in Chapter 3 differ in their approaches to sequential recruitment: given the maximum number of clusters and maximum cluster size, either entire clusters or individual participants within the same cluster are sequentially recruited. Both continuous and binary outcomes are considered, and Bayesian models are used for interim analyses to estimate the posterior probability of efficacy, which can be used to inform the choice to terminate or continue the trial.

The CDSS implements an individualized treatment rule (ITR). Given patient-level information, candidate treatments ordered by their predicted efficacy are obtained from a deep learning model. Without considering practical factors such as tolerance or side effects, the optimal treatment can be the one with the highest predicted efficacy. Generally, an ITR is a decision rule that customizes treatment assignments based on individual patient characteristics at a single decision point. An optimal ITR optimizes expected patient outcomes. Estimation of the optimal ITRs is essential to the development of the CDSS and should be carefully considered before conducting efficacy assessment. In a regression-based approach such as Q-learning (Watkins, 1989; Sutton and Barto, 2018), G-estimation (Robins, 2004), or dynamic weighted ordinary least squares (dWOLS) (Wallace and Moodie, 2015), the expected outcome is modelled as a function of treatment, covariates, and their interactions. The optimal treatment for a given patient is the one that leads to the best estimated outcome. While treatment-covariate interactions are essential in such approaches, their estimation from a single study can suffer from low power (Greenland, 1983), potentially requiring multisite collaboration to increase sample sizes. Common regression-based approaches rely on individual-level data to estimate ITRs, but in multisite studies, sharing highly sensitive patient-level health data across sites may be constrained. This challenge is addressed in Chapter 4.

Chapter 4 describes a two-stage Bayesian meta-analysis approach for ITR estimation. In the first stage, site-specific analyses are conducted using patient-level data within each site. In the second stage, site-specific ITRs are pooled in a Bayesian hierarchical model, treating the parameter estimates characterizing these ITRs from the first stage as the data, rather than the individual-level data itself. This approach takes between-site heterogeneity into account. Sparsity in both the data and the model are also considered. I further apply the proposed approach to estimate an individualized Warfarin dosing strategy using data from the International Warfarin Pharmacogenetics Consortium (2009).

The method proposed in Chapter 4 applies to binary or continuous treatment settings, and assumes that all participating sites have the same treatment sets. This assumption may not always be met, as for some diseases the treatment landscape can be quite heterogeneous. Due to time or funding constraints, it may not be possible to include all candidate treatments in all sites, but an ITR of all available treatments is typically preferred. In Chapter 5, I extend the method proposed in Chapter 4 to a multiple treatment setting with each site encompassing varying sets of treatment options, using network meta-analysis techniques (Cipriani et al., 2013). The method is illustrated through an analysis of data from the Sequenced Treatment Alternatives to Relieve Depression (STAR*D) study (Rush et al., 2004), Establishing Moderators and Biosignatures of Antidepressant Response for Clinical Care (EMBARC) study (Trivedi et al., 2016), and Research Evaluating the Value of Augmenting Medication with Psychotherapy (REVAMP) study (Trivedi et al., 2008) to establish an ITR for the treatment of depression.

My PhD thesis is in a thesis-by-manuscript format: Chapters 3, 4, and 5 were originally written as stand-alone manuscripts and therefore, there is some overlap with Chapter 2, which provides a literature review of the core concepts and theories I use in this thesis. Chapter 3 has been published in *Stat.* Chapters 4 and 5 are under review in statistical journals. I conclude in Chapter 6 with an overview of the three manuscripts, a discussion of limitations, and some directions for future work.

Chapter 2

Literature review

This chapter provides a review of the core concepts and theories I use in this thesis and consists of six sections. Section 2.1 reviews CRTs. Section 2.2 describes group sequential methods of clinical trials in both frequentist and Bayesian frameworks. Section 2.3 summarizes the statistical framework of estimating ITRs including basic concepts, assumptions, and common estimation methods. An introduction to individual participant data (IPD) meta-analysis and network meta-analysis is provided in Sections 2.4 and 2.5, respectively. A brief summary is given in Section 2.6.

2.1 Cluster-randomized trials

CRTs differ from individually randomized trials in that clusters of participants, rather than individual participants, are randomly assigned to different treatment arms (Hales and Moulton, 2017). CRTs are particularly appealing when the intervention under study is naturally applied to groups of individuals. For example, the CDSS I consider in this thesis is designed for physicians, each with multiple patients, naturally forming a group of participants. Randomizing by clusters is logically sound in such cases. In contrast, randomizing by individuals in this setting may lead to contamination, as physicians might find it challenging to treat patients from different treatment arms strictly differently. This contamination effect can dilute outcome differences between treatment arms, bias the treatment effect estimates, and undermine the study's reliability and validity.

Even for interventions that can be delivered at the individual level, cluster randomization may sometimes offer better logistical and administrative convenience. For instance, Ross et al. (1993) studied the effect of Vitamin A supplementation on child mortality in north Ghana. In this study, child mortality is a rare outcome, requiring a large sample size to detect an assumed effect size. Additionally, the study area was rural and underserved by health services. It was practically difficult to randomize by individual children, as Vitamin A and placebo had to be delivered by field workers through monthly household visits and identifying the treatment group for each individual child would have been excessively burdensome (Hales and Moulton, 2017). Therefore, in this study, while Vitamin A can be assigned to individual children, the study area was divided into geographical clusters, which were then used for randomization.

A key feature of CRTs is that outcomes of subjects within the same cluster are usually more similar, and this positive correlation is measured by intra-cluster correlation coefficient (ICC). Standard methods for design and analysis of individually randomized trials assume independence between subject-level outcomes. When this assumption is violated, these methods may be inappropriate. Analysis of CRTs can be conducted at either cluster or individual level. In a cluster-level analysis, a summary measure is obtained for each cluster, and these measures are then compared between treatment groups by standard statistical methods, as cluster-level outcomes can be assumed to be independent in CRTs. For example, in a two-arm CRT with a continuous outcome, a cluster-level analysis might compute a mean outcome for each cluster, and these mean outcomes may be compared between the two groups by a two-sample t-test. While the cluster-level analysis is conceptually simple, it is inefficient due to information loss when summarizing individual outcomes within the same cluster into a single measure. Alternatively, regression analysis can be performed on individual-level data. In this case, regression models must account for the ICC. Common regression models used for analysis of CRTs include random-effects models and generalized estimation equations (Hales and Moulton, 2017).

In addition to the different statistical methods required to account for the ICC in the analysis of CRTs, the ICC also results in lower power of CRTs compared with individually randomized trials that randomize the same number of individuals, as variability of responses is reduced. Research on adaptive methods for improving trial efficiency and flexibility in CRTs is limited. Lake et al. (2002) and van Schie and Moerbeek (2014) consider interim re-estimation of the required number of clusters and the required number of individuals per cluster, respectively. Zou et al. (2005) introduce group sequential designs for CRTs with binary outcomes in the frequentist framework. In addition to the work focused on parallel two-arm CRTs, Grayling et al. (2017) and Grayling et al. (2022) explore group sequential designs and response-adaptive randomization, respectively, for stepped-wedge CRTs, where all clusters enter the trial in the control arm and are gradually moved to the intervention group over a number of time periods (Hales and Moulton, 2017). In this thesis, I focus on Bayesian group sequential designs. A brief review of group sequential designs is provided in the next section.

2.2 Group sequential designs

Group sequential designs are a family of designs that allow interim decision making according to pre-determined decision rules in clinical trials. Unlike fixed-sample designs, where the sample size is pre-determined to detect an assumed effect with sufficient power and a controlled false positive rate, and a final analysis is conducted once all data are collected, group sequential designs divide participant entry into multiple groups. Interim assessments of accumulating data are pre-planned, enabling researchers to make early decisions about continuing or terminating the trial (Pocock, 1977).

Group sequential designs are more efficient than fixed-sample designs. Throughout the trial, evidence of treatment efficacy or futility may emerge early. In such cases, early trial termination can save resources and accelerate the trial process. In addition, interim adjustment such as early stopping or randomization adjustment is ethically necessary to ensure participants are not assigned to ineffective or even unsafe treatments. Both frequentist and Bayesian approaches have been proposed for group sequential designs.

2.2.1 Frequentist methods

In frequentist group sequential designs, a hypothesis test of the treatment effect involves constructing test statistics, and the distributions of these statistics are derived under appropriate assumptions. Conclusions about statistical significance and trial termination decisions are made by comparing the observed statistics to a decision boundary or computing a p-value (Pocock, 1977).

Consider an example with two treatment arms, A and B, and up to K waves of patient recruitment, each with size n for both arms. Assume the outcome is normally distributed with known variance σ^2 and unknown means μ_A and μ_B for treatments A and B, respectively. Without loss of generality, the treatment with a larger mean outcome is superior. Suppose that the hypotheses are formulated as $H_0: \theta \leq 0$ against $H_A: \theta > 0$, where $\theta = \mu_A - \mu_B$ represents the treatment effect. At analysis point $k \leq K$, a sufficient statistics for θ can be constructed as $D_k = (nk)^{-1} (\sum_{i=1}^k \sum_{j=1}^n y_{A,ij} - \sum_{i=1}^k \sum_{j=1}^n y_{B,ij})$, where $y_{A,ij}$ and $y_{B,ij}$ are the observed outcomes for the *j*th participant in the *i*th sample group in treatment groups Aand B, respectively. The statistics D_k is also normally distributed with mean θ and variance $I_k^{-1} = 2\sigma^2/(nk)$. Then, given a preset value of c_k , a decision rule can be constructed at any interim point: if the standardized statistics $D_k\sqrt{I_k} > c_k$, one should reject the null hypothesis, stop the trial, and claim that treatment A is superior to treatment B; otherwise, the next group of samples are recruited and observed.

The decision boundaries c_k , k = 1, ..., K, should be specified such that the overall false positive rate is below a desired value, e.g., 0.05. Pocock (1977) proposes using constant values, resulting in equal significance levels at each analysis. O'Brien and Fleming (1979) propose $c_k \propto \sqrt{\frac{K}{k}}$ and the test significance levels at each analysis increase as the trial progresses. Both Pocock (1977) and O'Brien and Fleming (1979) provide the decision boundaries for certain values of K and the overall false positive rate. While these authors focus on two-sided tests, extensions to one-sided tests have also been studied (e.g., Demets and Ware, 1980; DeMets and Ware, 1982). Lan and DeMets (1983) propose an alpha-spending function, which characterizes the rate at which the overall false positive rate is spent. The decision boundaries are then determined by the alpha-spending function.

2.2.2 Bayesian methods

In contrast to frequentist methods, Bayesian group sequential designs employ Bayesian analysis and summaries of the posterior distribution to make interim decisions. Prior information can be incorporated into design and analysis, and inference as well as decision making is based on the posterior distributions given the data. Consider the same example as in Section 2.2.1. A Bayesian decision criterion could be constructed as follows: at analysis point k, if $P(\theta > 0 | y_{A,ij}, y_{B,ij}; i = 1, ..., k, j = 1, ..., n) > U_k$, the trial stops; otherwise, the next group of samples is recruited and observed. Alternative criteria based on predictive distributions, such as the posterior predictive probability of declaring trial success or statistical significance at the end of the study, can also be used (Dmitrienko and Wang, 2006; Gsponer et al., 2014; Saville et al., 2014). In most cases, the posterior probabilities in the decision criteria are not analytically tractable and can be evaluated using Markov chain Monte Carlo (MCMC) methods.

While false positive rate and power are frequentist concepts, they often need to be assessed

for Bayesian designs due to regulatory requirements (Food and Drug Administration, 2010). Similar to c_k in frequentist group sequential designs, the decision boundaries U_k are usually chosen to control the overall false positive rate. A common choice is to use a constant value, i.e., $U_1 = \ldots = U_K = U$, where U can be obtained through simulations (Berry et al., 2010). Methods to calculate varying U_k based on alpha-spending functions in Bayesian group sequential designs have also been proposed (e.g., Zhu and Yu, 2017; Shi and Yin, 2019).

Bayesian group sequential designs can also be approached from a decision-theoretic framework. A loss or utility function needs to be defined, and the optimal decision rule at each analysis point is obtained by optimizing the posterior expected loss or utility (Lewis and Berry, 1994; Zhou and Ji, 2023). Bayesian group sequential designs in the decision-theoretic framework can incorporate factors such as the costs of patient enrollment by specifying certain loss or utility functions. The control of false positive rate or other operating characteristics can be achieved by placing constraints on optimization of expected loss or utility (Ventz and Trippa, 2015).

In summary, Bayesian group sequential designs offer two main advantages over frequentist methods: first, they allow for the inclusion of prior information, which is particularly appealing in scenarios such as rare diseases, where the sample size is limited; second, they use interim stopping criteria based on posterior (predictive) probabilities, which align more closely with clinical decision making and are more intuitive for investigators (Gsponer et al., 2014). Though there are different types of Bayesian group sequential designs, in this thesis, I constrain my attention to stopping rules for efficacy based on posterior probabilities.

2.3 Individualized treatment rules

In this section, I introduce the general framework of ITRs, including the concepts, the assumptions, and the common estimation methods. Let Y denote a continuous outcome of interest, where larger values of Y are preferable. Let A denote the treatment received

by the patient. Let X be a vector of pre-treatment covariates. Uppercase, lowercase, and bold denote random variables, realizations of random variables, and vectors respectively. Throughout this section, let i index individual patients.

An ITR $d(\mathbf{X}) : \mathbf{X} \to A$ is a decision rule that utilizes patient-level information, \mathbf{X} , such as demographics, genetic makeup, or disease history, to customize treatment plans at a single decision point. Statistical estimation of optimal ITRs can be considered in the potential outcomes framework (Rubin, 1974; Splawa-Neyman et al., 1990). Let Y^a be the potential outcome a patient would experience if assigned treatment a. The axiom of consistency states that the potential outcome under the observed treatment and the observed outcome should agree: $Y = Y^a$, if A = a. Define a value function of the ITR $d(\mathbf{X})$, $\mathcal{V}(d)$, as the expected potential outcome if all patients in the population were treated according to d, i.e., $\mathcal{V}(d) = E(Y^{d(\mathbf{X})})$. The optimal ITR is defined as $d^{\text{opt}} = \operatorname{argmax}_d \mathcal{V}(d)$.

Statistical models for the optimal ITR from randomized trials and observational studies rely on several assumptions:

- 1. Stable unit treatment value assumption (SUTVA): a patient's outcome is not influenced by other patients' treatments (Rubin, 1980).
- 2. No unmeasured confounding (Robins, 1997): given X = x, the treatment assignment A is independent of potential outcomes Y^a for all possible a.
- 3. Positivity: there is a positive probability of receiving every possible treatment for every combination of covariate values that occur among individuals in the population (Cole and Hernán, 2008).

SUTVA assumes no interference between individual patients. This may be violated in certain cases, for example, in the case of a vaccination against an infectious disease. It is highly possible that the potential risk of the disease for an individual, regardless of his or her vaccination status, will be reduced if the surrounding people are vaccinated. The methodology developed in Chapter 4 and 5 is not applicable to such cases. The no unmeasured confounding assumption implies $E(Y^a \mid \mathbf{X}) = E(Y^a \mid A = a, \mathbf{X})$. This assumption is expected to be valid in ideal randomized trials by design. However, in observational studies, it is not testable, and expert knowledge is required to collect data on as many potential confounders as possible to make this assumption approximately true. The positivity assumption ensures that all possible treatments can be observed in every patient subgroup defined by the covariates. This assumption is often satisfied in randomized trials by design and can be empirically assessed in observational studies by checking the overlap region of the distribution of propensity score, i.e., the conditional probability of receiving the treatment of interest, across treatment groups. If the positivity assumption is violated, estimation of ITRs requires extrapolation, which may introduce bias.

There are two general approaches to estimate the optimal ITRs: value search approaches and regression-based approaches. For a value search method, a space of possible ITRs \mathcal{D} is first specified, and an appropriate method is used to estimate the value of each candidate ITR to find the best one. Inverse probability weighting (Robins, 2000) is the building block of many methods in this category. The value function can be expressed as

$$\mathcal{V}(d) = E(Y^d) = E\left[\frac{I\{A = d(\mathbf{X})\}}{p(A \mid \mathbf{X})}Y\right],$$

giving the inverse probability of treatment weighted (IPTW) estimator based on n samples

$$\widehat{\mathcal{V}}(d) = \frac{1}{n} \sum_{i=1}^{n} \frac{I\{A_i = d(\boldsymbol{x}_i)\}}{p(A_i = a_i \mid \boldsymbol{x}_i)} y_i.$$

The propensity score $p(A_i = a_i | \boldsymbol{x}_i)$ is either known in randomized trials or estimated from a propensity score model, e.g., logistic regression, in observational studies. The optimal ITR is obtained as $d^{\text{opt}} = \operatorname{argmax}_{d \in \mathcal{D}} \widehat{\mathcal{V}}(d)$. Finding the optimal ITR by maximizing $E[I\{A = d(\mathbf{X})\}\{p(A \mid \mathbf{X})\}^{-1}Y]$ is equivalent to identifying an ITR that minimizes $E[I\{A \neq d(\mathbf{X})\}\{p(A \mid \mathbf{X})\}^{-1}Y]$, which can be considered as a weighted misclassification error (Zhao et al., 2012). Therefore, the problem of estimating the optimal ITR can then be recast as a classification problem and machine learning methods can be used to estimate the optimal ITR. For example, Zhao et al. (2012) propose using a hinge loss function to approximate the non-smooth indicator function in $E[I\{A \neq d(\mathbf{X})\}\{p(A \mid \mathbf{X})\}^{-1}Y]$ and solving the optimization problem through support vector machine techniques. Classification-based methods have also been studied in Zhang et al. (2012); Zhou et al. (2017); Zhu et al. (2017); Liu et al. (2018); Zhang et al. (2020).

Value search methods directly estimate the value of ITRs and pick the one with the optimized value, often using non-parametric or semi-parametric models. Limited assumptions are required for the data in value search methods, giving more flexibility. However, the estimation of the value function can be unstable, resulting in high variability of the estimated optimal ITR (Chakraborty and Moodie, 2013). Furthermore, classification-based methods that adopt machine learning algorithms often generate uninterpretable ITRs, which can be difficult to use in clinical practice. In contrast, regression-based methods indirectly find the optimal ITR through a two-step procedure: first, the expected outcomes or contrasts of expected outcomes are modelled by regression models; then, the optimal ITR is identified to maximize the expected outcomes or contrasts of expected outcomes. Common regression-based methods include Q-learning (Watkins, 1989; Sutton and Barto, 2018), G-estimation (Robins, 2004), and dWOLS (Wallace and Moodie, 2015).

In regression-based approaches, the expected outcomes can be modelled by a treatment-free function f and a blip function γ (Robins, 2004) through:

$$E(Y \mid A = a, \boldsymbol{X} = \boldsymbol{x}) = f(\boldsymbol{x}^{(\boldsymbol{\beta})}; \boldsymbol{\beta}) + \gamma(a, \boldsymbol{x}^{(\boldsymbol{\psi})}; \boldsymbol{\psi}).$$

Here, both $x^{(\beta)}$ and $x^{(\psi)}$ are subvectors of x, and include predictive covariates and covariates

that interact with treatment (prescriptive variables), respectively. The prescriptive covariate vector $\boldsymbol{x}^{(\psi)}$ is a subvector of $\boldsymbol{x}^{(\beta)}$. The treatment-free function f represents the expected outcome at a reference treatment or a "zero" level of the treatment. Therefore, the function f is unrelated to treatment assignment. The blip function γ characterizes the difference in the expected potential outcome of patients between receiving treatment A = a and the reference A = 0, i.e., $\gamma(a, \boldsymbol{x}^{(\psi)}; \boldsymbol{\psi}) = E(Y^a - Y^0 \mid \boldsymbol{X} = \boldsymbol{x}^{(\psi)}; \boldsymbol{\psi})$. The functions f and γ are parameterized by vectors $\boldsymbol{\beta}$ and $\boldsymbol{\psi}$, respectively. The optimal ITR can then be identified as $d^{\text{opt}}(\boldsymbol{x}) = \operatorname{argmax}_a \gamma(a, \boldsymbol{x}^{(\psi)}; \boldsymbol{\psi})$. Its estimation relies on the estimation of $\boldsymbol{\psi}$, requiring correct specification of both f and γ in Q-learning.

Consider a binary treatment setting $A \in \{0, 1\}$. By definition, $\gamma(0, \boldsymbol{x}^{(\boldsymbol{\psi})}; \boldsymbol{\psi}) = 0$. A common choice for the blip function satisfying this condition is $\gamma(a, \boldsymbol{x}^{(\boldsymbol{\psi})}; \boldsymbol{\psi}) = ag(\boldsymbol{x}^{(\boldsymbol{\psi})}; \boldsymbol{\psi})$, where g is a function only depending on the prescriptive covariates $\boldsymbol{x}^{(\boldsymbol{\psi})}$ and is specified by the analysts. For illustration purposes, consider linear models for both f and g. In this setting, the outcome model becomes

$$E(Y \mid A = a, \boldsymbol{X} = \boldsymbol{x}) = \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\beta})} + a \boldsymbol{\psi}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\psi})},$$

where we now assume $x^{(\beta)}$ and $x^{(\psi)}$ both contain a leading column of ones, respectively corresponding to an intercept and a main effect of treatment. The parameters β and ψ are estimated based on *n* samples by solving

$$(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\psi}}) = \operatorname*{argmin}_{(\boldsymbol{\beta}, \boldsymbol{\psi})} \frac{1}{n} \sum_{i=1}^{n} \left(y_i - \boldsymbol{\beta}^{\top} \boldsymbol{x}_i^{(\boldsymbol{\beta})} - a_i \boldsymbol{\psi}^{\top} \boldsymbol{x}_i^{(\boldsymbol{\psi})} \right)^2.$$

Therefore, in a single stage setting with a continuous outcome, Q-learning can be implemented as a linear regression model. The estimated optimal ITR is $\hat{d}^{\text{opt}}(\boldsymbol{x}) = I(\hat{\boldsymbol{\psi}}^T \boldsymbol{x}^{(\boldsymbol{\psi})} > 0)$. Similarly, for a binary or an unbounded discrete-valued outcome, Q-learning can be implemented as logistic regression and Poisson regression (Moodie et al., 2014), respectively. While Q-learning is straightforward, the outcome regression models must be correctly specified to yield consistent estimators of model parameters. In reality, the true relationship between treatment, covariates, and outcome may be complex and may not be captured by simple linear models. In this case, flexible specifications in either frequentist or Bayesian frameworks can be considered (e.g., Moodie et al., 2014; Wager and Athey, 2018; Murray et al., 2018; Logan et al., 2019; Hahn et al., 2020). For example, Moodie et al. (2014) use generalized additive models for Q-learning. A generalized additive model is still additive but the relationship between the covariates and the outcome can be modelled by smooth functions, e.g., the outcome model may take the form

$$E(Y \mid A = a, \mathbf{X} = \mathbf{x}) = \boldsymbol{\beta}^{\top} \mathbf{x}^{(\boldsymbol{\beta})} + a \boldsymbol{\psi}^{\top} \mathbf{x}^{(\boldsymbol{\psi})} + f_1(x_1) + f_2(x_2) + f_3(x_3) + \cdots,$$

where x_1, x_2, x_3, \ldots are different specific covariates contained in $\mathbf{x}^{(\beta)}$, and f_1, f_2, f_3, \cdots are smooth functions, e.g., penalized regression splines. Murray et al. (2018) provide a general framework for Bayesian machine learning approaches to Q-learning. While they focus on a dynamic treatment regime (DTR) setting, which is an extension of ITRs to multiple stages, they argue that in a single stage setting, standard Bayesian regression models would be sufficient. Specifically, they describe the use of Bayesian additive regression trees, which in an ITR setting, is similar to the work of Logan et al. (2019). The latter propose that the outcome could be modelled by an ensemble of regression trees in an additive fashion:

$$E(Y \mid A = a, \boldsymbol{X} = \boldsymbol{x}) = \sum_{s} \phi(\boldsymbol{x}, a; T_s, M_s),$$

where $\phi(\boldsymbol{x}, a; T_s, M_s)$ represents a binary tree function, T_s represents a specific tree structure, including interior and terminal nodes as well as decision rules in the interior nodes, and M_s includes the function values at the terminal nodes. Then, priors are assigned to both T_s and M_s . For details of prior specification, see Murray et al. (2018) and Logan et al. (2019). In such a model, posterior predictive samples of the outcome for an individual with $\boldsymbol{X} = \boldsymbol{x}$ and under treatment A = a can be obtained from MCMC. The treatment that maximizes $E(Y | \mathbf{X} = \mathbf{x}, A = a)$, approximated by the mean of the posterior predictive samples, is identified as the optimal treatment. While these flexible models provide a powerful way to model the complex relationship, the estimated optimal ITRs are usually hard to interpret.

Semi-parametric approaches that are more robust to model misspecification, such as Gestimation (Robins, 2004) and dWOLS (Wallace and Moodie, 2015), have also been developed. In G-estimation, the blip parameter vector $\boldsymbol{\psi}$ is estimated by solving the estimating equation

$$\sum_{i=1}^{n} \left[S(A_i) - E\{S(A_i) \mid \boldsymbol{x_i}, \widehat{\boldsymbol{\alpha}}\} \right] \left(G_i(\boldsymbol{\psi}) - E[G_i(\boldsymbol{\psi}) \mid \boldsymbol{x_i}, \widehat{\boldsymbol{\beta}}(\boldsymbol{\psi})] \right) = 0,$$

where S(A) is a vector-valued function chosen by the analysts to contain prescriptive covariates; a common choice is $S(A) = \partial \gamma(a, \boldsymbol{x}^{(\psi)}; \boldsymbol{\psi}) / \partial \boldsymbol{\psi}$, ensuring that the estimating function is linked to the blip function, and is of the same dimensionality of the parameters $\boldsymbol{\psi}$ of γ . The parameter estimate $\hat{\boldsymbol{\alpha}}$ is obtained from a treatment model, e.g., a logistic regression for a binary treatment assignment. A model is posited for the treatment-free component $G(\boldsymbol{\psi}) = Y - \gamma(a, \boldsymbol{x}^{(\boldsymbol{\psi})}; \boldsymbol{\psi})$, and the parameters are estimated in terms of $\boldsymbol{\psi}$, i.e., $\hat{\boldsymbol{\beta}}(\boldsymbol{\psi})$.

The dWOLS method estimates parameters β and ψ by minimizing a weighted least squares. With linear models for both f and γ ,

$$(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\psi}}) = \operatorname*{argmin}_{(\boldsymbol{\beta}, \boldsymbol{\psi})} \frac{1}{n} \sum_{i=1}^{n} w_i \left(y_i - \boldsymbol{\beta}^{\top} \boldsymbol{x}_i^{(\boldsymbol{\beta})} - a_i \boldsymbol{\psi}^{\top} \boldsymbol{x}_i^{(\boldsymbol{\psi})} \right)^2.$$

Here, the weight w_i is chosen to satisfy a balancing property and may depend on $E(A_i | \mathbf{x}_i, \hat{\boldsymbol{\alpha}})$, e.g., $w_i = |a_i - E(A_i | \mathbf{x}_i, \hat{\boldsymbol{\alpha}})|$. Both G-estimation and dWOLS satisfy the doublyrobust property: only one of the two models, a treatment model or an outcome model, needs to be correctly specified to obtain a consistent estimator of $\boldsymbol{\psi}$. However, dWOLS is slightly more restrictive, in that the treatment-free model must contain all tailoring variables whereas in G-estimation, the treatment-free model could follow any specification, even a null (intercept-only) model.

2.4 Individual participant data meta-analysis

Clinical decision making should be guided by evidence from high-quality studies. However, a single study may not be sufficient to draw definitive conclusions, as different studies can report varying or even conflicting results. In this case, it is often hard to know which result is more reliable, and combining evidence from multiple studies is necessary to enhance reliability, confidence, and generalizability.

Meta-analysis is a widely-used quantitative approach for evidence synthesis (Borenstein et al., 2021). Traditional meta-analysis relies on aggregate data extracted from study publications, which refers to the information averaged or estimated across all participants in the study. For example, to obtain an overall measure of the effect of a particular treatment with evidence from multiple studies, treatment effect estimates and their uncertainty (e.g., standard errors or confidence intervals) can be obtained from individual studies and synthesized into a weighted average. Aggregate data meta-analysis is subject to limitations. First, even if individual studies target the same research question, they may differ in study designs, populations, data measurements, and analysis methods. These differences can result in substantial heterogeneity in estimates, and a weighted average may be less informative and meaningful. Second, in the PM paradigm, it is important to identify individual-level treatment-covariate interactions, but aggregate data meta-analysis such as meta-regression is prone to ecological bias (Berlin et al., 2002; Simmonds and Higgins, 2007). The relationship identified based on the aggregate data may not reflect the relationship at the individual level. In addition, identification of treatment-covariate interactions based on aggregate data may suffer from low power, as typically there are only a small number of studies that cannot provide sufficient variability in study-level aggregated covariates.
An alternative approach that can address these limitations is individual participant data (IPD) meta-analysis. IPD refers to the raw information for each participant in a study, such as baseline characteristics, treatment assignment, and outcome (Riley and Fisher, 2021). Having access to IPD provides an opportunity to standardize eligibility criteria, outcome or covariate definitions, and analysis methods. It also allows for more flexible modelling. For example, covariate adjustment can be incorporated at the individual level, and missing data mechanisms, if applicable, can also be modelled consistently. In addition, IPD from unpublished studies may be included to reduce publication bias (Riley et al., 2010). While IPD meta-analysis is increasingly in demand, its implementation can be resource-intensive. It may take a long time to obtain IPD, as study investigators may not be contactable or they may not be willing to share their IPD. Even if IPD are available, more advanced statistical expertise may be required for flexible modelling, and there may be practical and ethical concerns about using or managing patient-level data.

One may approach the IPD meta-analysis using a two-stage method, which resembles the aggregate data meta-analysis in that the IPD for each study are analyzed separately to produce study-level aggregate data (e.g., treatment effect estimates), which are then combined to obtain an overall summary result. Alternatively, a one-stage analysis pools the IPD from different sources/studies, and analyzes them in a single statistical model.

The two-stage approach is conceptually straightforward. Standard methods for aggregate data meta-analysis can be used in the second stage, making it more accessible to researchers, especially non-statisticians, who often are already familiar with aggregate data meta-analysis. In contrast, the one-stage approach provides a flexible framework to incorporate data at both individual and study levels. Moreover, the one-stage approach avoids the need to synthesize study-level estimates which are typically assumed to follow some probability distributions (e.g., a normal distribution) in the two-stage approach. Instead, it directly models the outcome distributions, which is advantageous when included studies have small sample

sizes.

Researchers have also explored the differences in the results obtained from the one- and twostage approaches either theoretically, empirically, or via simulations (Smith and Williamson, 2007; Koopman et al., 2008; Bowden et al., 2011; Morris et al., 2018; Kontopantelis, 2018, e.g.,). When the same modelling assumptions and estimation methods are used, differences in the results between one- and two-stage approaches tend to be small. However, the twostage approach is often computationally faster and naturally separates within-study and across-study information, which is important when estimating treatment-covariate interactions (Riley and Fisher, 2021). In addition, even if access to IPD is feasible for individual studies, a two-stage approach may be necessary if IPD from different studies cannot be released and merged into a single dataset directly, which is the case considered in Chapters 4 and 5. In the next section, I use a toy example for illustration.

2.4.1 Two-stage approach

Throughout this section, assume a continuous outcome Y, a binary treatment $A \in \{0, 1\}$ and a single covariate X. Therefore, the IPD in this example include information on treatment assignment, covariate X, and outcome Y for each individual. Suppose that, in addition to the treatment effect, the treatment-covariate interaction is also of interest. In the first stage, the IPD will be analyzed separately for each study. Various modelling options can be considered, however, I use a linear regression model for illustration purposes, i.e.,

$$y_{ij} = \alpha_i + \beta_i x_{ij} + \theta_i a_{ij} + \psi_i a_{ij} x_{ij} + \epsilon_{ij},$$

$$\epsilon_{ij} \sim N(0, \sigma_i^2).$$
(2.1)

Here, *i* and *j* index the study and individual participants for each single study, respectively. The parameters α_i , β_i , $\theta_i \psi_i$, and σ_i^2 represent separate intercept, covariate effect, treatment effect, treatment-covariate interaction effect, and error variance for study *i*, respectively. Model (2.1) is only fitted among observations in study *i*, using only within-trial information to estimate the parameters, which automatically avoids study-level confounding and ecological bias.

Treatment effect estimates $\hat{\theta}_i$, treatment-covariate interaction estimates $\hat{\psi}_i$, and their variances var $(\hat{\theta}_i)$, var $(\hat{\psi}_i)$ can be obtained from model (2.1), and are then synthesized in the second stage to produce an overall estimate under either a common effect or a random effects assumption. Since the second-stage analysis can be applied analogously to both treatment effect estimates and treatment-covariate interactions, I only describe the analysis for treatment effect here.

Under the common effect assumption, the $\hat{\theta}_i$ s are estimates of a common treatment effect θ . That is, the true treatment effect is assumed to be the same for all studies, and the differences in the observed study estimates are purely due to chance. The meta-analysis model in the second stage can be written as

$$\hat{\theta}_i \sim N(\theta, \operatorname{var}(\hat{\theta}_i)).$$
 (2.2)

In model (2.2), the treatment effect estimates $\hat{\theta}_i$ are assumed to be normally distributed, which is plausible for a sufficiently large sample size. In addition, the variance of $\hat{\theta}_i$ is assumed to be estimated with no error. Estimate of the overall treatment effect θ can be analytically obtained as a weighted average: $\hat{\theta} = (\sum_{i=1}^{K} \hat{\theta}_i w_i) (\sum_{i=1}^{K} w_i)^{-1}$, and the variance $\operatorname{var}(\hat{\theta}) = (\sum_{i=1}^{K} w_i)^{-1}$, with $w_i = \operatorname{var}^{-1}(\hat{\theta}_i)$.

The common effect assumption may not be appropriate when differences exist in the true treatment effects across studies, referred to as between-study heterogeneity in meta-analysis.

To allow for between-study heterogeneity, the model (2.2) is extended to

$$\hat{\theta}_i \sim N\{\theta_i, \operatorname{var}(\hat{\theta}_i)\},$$

$$\theta_i \sim N(\theta, \tau^2).$$
(2.3)

Here, the different true treatment effects for individual studies are assumed to come from a normal distribution with mean θ and between-study variance τ^2 . Setting $\tau^2 = 0$ will reduce model (2.3) to model (2.2). Given an estimate of τ^2 , the maximum likelihood estimator of θ can also be expressed as a weighted average: $\hat{\theta} = (\sum_{i=1}^{K} \hat{\theta}_i w_i^*) (\sum_{i=1}^{K} w_i^*)^{-1}$ with $w_i^* =$ $(\operatorname{var}(\hat{\theta}_i) + \hat{\tau}^2)^{-1}$. There are various approaches for estimating τ^2 ; see, e.g., Cochran (1954); Paule and Mandel (1982); DerSimonian and Laird (1986); Hardy and Thompson (1996); DerSimonian and Kacker (2007). Bayesian approaches can also be used in the second stage. A likelihood is defined by model (2.2) or model (2.3). Then, prior distributions are required for the unknown model parameters, such as θ and τ .

2.4.2 One-stage approach

The one-stage approach analyzes IPD from all studies together in a single statistical model, which typically is hierarchical to account for clustering of samples from the same study. Consider the same example described in Section 2.4.1. A one-stage model without treatmentcovariate interactions could be

$$y_{ij} = \alpha_i + \beta_i x_{ij} + \theta_i a_{ij} + \epsilon_{ij},$$

$$\epsilon_{ij} \sim N(0, \sigma_i^2).$$

This model is fitted on all IPD. The parameters α_i , β_i , and θ_i are assumed to be either common (e.g., $\theta_1 = \ldots = \theta_K$), stratified (e.g., $\theta_1, \ldots, \theta_K$ are different fixed effects), or random (e.g., $\theta_1, \ldots, \theta_K$ are different but come from a common distribution). The residual variances can also be stratified or common. When treatment-covariate interaction is incorporated, Riley and Fisher (2021) recommend that within-study and across-study information should be disentangled to avoid ecological bias. They propose two approaches: (1) stratify all parameters outside the interaction terms by study, or (2) center the covariate about its study-specific mean and include another term for study-specific covariate mean to explain the between-study heterogeneity in the treatment effect.

2.5 Network meta-analysis

Traditional meta-analysis approaches typically compare two treatments at a time. However, in real-world scenarios, multiple treatments are often available for a single disease, and a single study will rarely compare all of these treatments simultaneously. Instead, only a subset of the treatments is included in individual studies. Network meta-analysis assesses the comparative effectiveness of multiple treatments by synthesizing evidence from a network of studies (Salanti, 2012; Cipriani et al., 2013). It relies on simultaneous analysis of both direct and indirect evidence. Here, direct evidence of a treatment comparison refers to the information from studies comparing the two treatments of interest, while indirect evidence refers to the information from studies comparing the treatments of interest with one or more common comparators.

Consider an example with three treatments, d_1 , d_2 , and d_3 . Suppose some studies compare d_1 and d_2 , while others compare d_1 and d_3 . In network meta-analysis, the network structure of studies is usually represented by a graph where a node represents a treatment and a line (or an edge) connecting two treatments indicates that direct comparisons exist for the two treatments (see Figure 2.1 for the network structure of the hypothetical example). For real studies, it is also common that the sizes of the nodes are proportional to the number of participants in the corresponding treatment groups, and the widths of the lines are also proportional to the number of studies including the two treatments connected by the lines.



Figure 2.1: Illustration of the network structure for the hypothetical example.

While no direct comparisons are possible between d_2 and d_3 , indirect comparison can be made through the common comparator d_1 . The indirect estimate of the effect of d_2 relative to d_3 , $\hat{\theta}_{d_2d_3}^{\text{indirect}}$, can be obtained by subtracting the pairwise meta-analytic estimate of studies of d_3 and d_1 , $\hat{\theta}_{d_3d_1}$, from the estimate of studies of d_2 and d_1 , $\hat{\theta}_{d_2d_1}$, i.e., $\hat{\theta}_{d_2d_3}^{\text{indirect}} = \hat{\theta}_{d_2d_1} - \hat{\theta}_{d_3d_1}$. When the direct estimate $\hat{\theta}_{d_2d_3}^{\text{direct}}$ is also available, it can be combined with the indirect estimate to obtain a mixed estimate, which is potentially more precise than a single direct estimate or an indirect estimate (Cooper et al., 2011; Caldwell et al., 2015). The validity of indirect and mixed comparisons relies on the transitivity assumption that the studies comparing different sets of treatments should be sufficiently similar with respect to all important factors other than the treatments being compared (Cipriani et al., 2013). This assumption is often evaluated qualitatively, e.g., by comparing the distributions of effect modifiers across studies.

The consistency assumption, which states that indirect and direct evidence for the same treatment comparison are in agreement, is the statistical manifestation of transitivity (Cipriani et al., 2013). For the example provided, the consistency assumption suggests a consistency equation: $\theta_{d_2d_3} = \theta_{d_2d_1} - \theta_{d_3d_1}$, where $\theta_{d_2d_1}$, $\theta_{d_3d_1}$, and $\theta_{d_2d_3}$ are the true overall effects of d_2 versus d_1 , d_3 versus d_1 , and d_2 versus d_3 , respectively. Unlike transitivity, consistency can be quantitatively assessed when both indirect and direct estimates are available. Meth-

ods of assessing consistency in network meta-analysis can be broadly categorized as local or global approaches. Local approaches evaluate consistency for particular comparisons in the network. For example, for a particular treatment comparison, the node-splitting method (Dias et al., 2010b) calculates a direct estimate from studies including both the treatments being compared and an indirect estimate from the remaining studies, and the difference between the two estimates are examined. Global approaches assess consistency across the entire network. For example, models can be constructed by relaxing the consistency assumptions (Lu and Ades, 2006), e.g., modifying the consistency equation to add an inconsistency factor $\delta_{d_1d_2d_3}$: $\theta_{d_2d_3} = \theta_{d_2d_1} - \theta_{d_3d_1} + \delta_{d_1d_2d_3}$. Consistency can be assessed by examining the inconsistency factors or comparing models under consistency and inconsistency.

To evaluate the comparative effectiveness of multiple treatments, a common reference treatment should be identified, but this treatment does not have to be present in every study. Typically, the effects of the remaining treatments relative to this common reference treatment are of particular interest and can be used to derive the effects of other treatment comparisons under the consistency assumption. Except for consistency, the statistical framework of network meta-analysis is similar to that of traditional meta-analysis. The study-level estimates can be combined under a common effect or a random effects assumption to produce overall estimates. For example, assume that d_3 is the common reference treatment. Given studylevel treatment effect estimates $\hat{\theta}_{i,d_2d_1}$, $\hat{\theta}_{i,d_3d_1}$ and their variances $\operatorname{var}(\hat{\theta}_{i,d_2d_1})$, $\operatorname{var}(\hat{\theta}_{i,d_3d_1})$, the network meta-analysis model assuming random effects and consistency could be

$$\begin{aligned} \widehat{\theta}_{i,d_2d_1} &\sim N(\theta_{i,d_2d_1}, \operatorname{var}(\widehat{\theta}_{i,d_2d_1})), \\ \theta_{i,d_2d_1} &\sim N(\theta_{d_2d_3} + \theta_{d_3d_1}, \tau^2_{d_2d_1}), \\ \widehat{\theta}_{i,d_3d_1} &\sim N(\theta_{i,d_3d_1}, \operatorname{var}(\widehat{\theta}_{i,d_3d_1})), \\ \theta_{i,d_3d_1} &\sim N(\theta_{d_3d_1}, \tau^2_{d_3d_1}). \end{aligned}$$

Here, different between-study variances $(\tau_{d_2d_1}^2, \tau_{d_3d_1}^2)$ are assumed for different treatment

comparisons. This may be feasible in this particular example, as only two different treatment comparisons in two-arm studies are considered. For multi-arm studies, vectors of treatment effect estimates as well as variance-covariance matrices of the estimate vector can be obtained. For these studies, a model under the random effects assumption will require a betweenstudy variance-covariance matrix. For example, if there also exist studies including d_1 , d_2 , and d_3 , under the random effects assumption, treatment effect estimate vectors $\hat{\theta}_{i,d_1d_2d_3} =$ $(\hat{\theta}_{i,d_1d_3}, \hat{\theta}_{i,d_2d_3})$ and their variance-covariance matrices $\hat{\Sigma}(\hat{\theta}_{i,d_1d_2d_3})$ can be modelled by

$$\widehat{\boldsymbol{\theta}}_{i,d_1d_2d_3} \sim \text{MVN}(\boldsymbol{\theta}_{i,d_1d_2d_3}, \widehat{\boldsymbol{\Sigma}}(\widehat{\boldsymbol{\theta}}_{i,d_1d_2d_3})),$$

$$\boldsymbol{\theta}_{i,d_1d_2d_3} \sim \text{MVN}(\boldsymbol{\theta}_{d_1d_2d_3}, \boldsymbol{\Sigma}_{d_1d_2d_3}),$$
(2.4)

where MVN represents a multivariate normal distribution. The parameters $\theta_{i,d_1d_2d_3}$ and $\theta_{d_1d_2d_3}$ are the study-specific and overall true treatment effect parameter vector, and $\Sigma_{d_1d_2d_3}$ is a 2 × 2 variance-covariance matrix for the between-study heterogeneity. The between-study variance-covariance matrix could be unstructured in theory, but this will induce a large number of parameters when many treatments and multi-arm studies are available. In this case, parameter estimators obtained from a limited number of studies may suffer from low precision. Alternatively, a common specification that can reduce the number of model parameters to be estimated is that the between-study variance-covariance matrix has diagonals of a between-study variance that is common to all treatment comparisons, and off-diagonals of half of the common between-study variances (White et al., 2012; Riley and Fisher, 2021).

Parameters in network meta-analysis models are usually estimated in the Bayesian framework. Bayesian approaches naturally provide a convenient way to estimate ranking probabilities, i.e., the probability of each treatment to be the best, the second best, etc. In addition, posterior probability statements are more interpretable and intuitive for decision making. Frequentist estimation is also possible by formulating the model as a multivariate meta-regression and using data augmentation techniques (White et al., 2012). Network meta-analysis based on IPD has the same advantages as the traditional IPD meta-analysis. An additional value of IPD for network meta-analysis is the detection and reduction of inconsistencies. For example, distributions of covariates can be compared across studies. If the covariate distributions in studies that provide direct evidence are different from those in studies giving indirect evidence, the transitivity and also the consistency assumption may not be reliable. In this case, including patient-level covariates in the model may help reduce both inconsistency and between-study heterogeneity.

2.6 Summary

The literature review introduced important concepts and theories I use in this thesis. I first discussed CRTs and group sequential designs, which are related to the research work presented in Chapter 3. Then, I described the assumptions and some estimation methods for ITRs. An overview of IPD meta-analysis and network meta-analysis were also presented for Chapters 4 and 5.

Chapter 3

Bayesian group sequential designs for cluster-randomized trials

Preamble to Manuscript 1. Clinical decision support systems (CDSSs) are crucial for integrating PM into routine clinical practice by providing physicians with tailored treatment recommendation information for individual patients. Therefore, evaluating their efficacy and safety through randomized controlled trials is essential. This chapter is motivated by a real-world CRT of a machine learning-based CDSS designed for physicians treating patients with depression. CRTs often have low power due to the positively correlated responses of subjects within the same cluster. Group sequential designs can enhance trial efficiency by recruiting participants sequentially and incorporating pre-planned interim analyses to decide on trial continuation or early termination. However, most group sequential designs, whether frequentist or Bayesian, focus on individually randomized trials. This manuscript develops two Bayesian group sequential designs for CRTs which differ in how participants are sequentially recruited. The power and false positive rate of the two designs are evaluated and compared across two outcome types and a wide range of scenarios in a simulation study. The corresponding manuscript was published in *Stat* (Shen et al., 2022).

Bayesian group sequential designs for cluster-randomized trials

Junwei Shen¹, Shirin Golchi¹, Erica E.M. Moodie¹, David Benrimoh^{2,3}.

¹Department of Epidemiology, Biostatistics, and Occupational Health, McGill University, Quebec, Canada ²Aifred Health, Quebec, Canada ³Department of Psychiatry, McGill University, Quebec, Canada

This thesis contains the accepted version of the corresponding paper published in Stat (Shen et al., 2022).

Abstract

Flexible approaches have been proposed for individually randomized trials to save time or reduce sample size. However, flexible designs for cluster-randomized trials in which groups of participants rather than individuals are randomized to treatment arms are less common. Motivated by a cluster-randomized trial designed to assess the effectiveness of a machine-learning based clinical decision support system for physicians treating patients with depression, two Bayesian group sequential designs for cluster-randomized trials are proposed to allow for early stopping for efficacy at pre-planned interim analyses. The difference between the two designs lies in the way that participants are sequentially recruited. Given a maximum number of clusters as well as maximum cluster size allowed in the trial, one design sequentially recruits clusters with the given maximum cluster size, while the other recruits all clusters at the beginning of the trial but sequentially enrolls individual participants until the trial is stopped early for efficacy or the final analysis has been reached. The design operating characteristics are explored via simulations for a variety of scenarios and two outcome types for the two designs. We make recommendations for Bayesian group sequential designs of cluster-randomized trials based on the simulation results.

3.1 Introduction

Randomized controlled trials, which can ensure that subjects assigned to each treatment group are comparable with respect to all characteristics of interest to draw a causal conclusion, have played an essential role in evaluating the effectiveness of interventions (Friedman et al., 2015). In most trials designed to assess the effect of a drug or a treatment, individual participants are randomly allocated to each treatment arm. However, some interventions naturally operate at a group level, or target either a social network or physical environment. For example, in a trial assessing the clinical utility, safety, and potential effectiveness of a machine-learning based clinical decision support system (CDSS) developed by Aifred Health (Benrimoh et al., 2018; Popescu et al., 2021; Benrimoh et al., 2021; Tanguay-Sela et al., 2022), the decision support tool is designed for physicians, naturally forming clusters of individual patients being treated by that physician (Mehltretter et al., 2020b,a). For these types of interventions, a population-level effect is of more interests to researchers.

In addition, randomizing by individuals in this setting may lead to a contamination effect between trial arms, as physicians might struggle to treat patients from different treatment arms strictly differently. Contamination can cause dilution bias and affect the reliability and validity of the study. One way to reduce the possibility of contamination is to randomize by physicians who act as 'clusters' (Torgerson, 2001; Puffer et al., 2005). Thus, it may not be advisable to randomize individual participants to different treatment arms, and groups of subjects being treated by their clinicians can instead be randomly assigned to the treatment arms in what is known as a *cluster-randomized trial* (Donner and Klar, 2000; Puffer et al., 2005; Hales and Moulton, 2017; Turner et al., 2017). In other cluster-randomized trials, clusters may be formed by subjects sharing other common features such as geographical areas, communities, or worksites (Hales and Moulton, 2017).

Due to the difference in randomization unit, the design and analysis of cluster-randomized trials differ from individually randomized trials (Klar and Donner, 2001; Campbell et al.,

2007, 2001). The independence assumption between participants' outcomes is violated in a cluster-randomized trial since subjects from the same cluster tend to have more similar responses than subjects from a different cluster. Positive correlation between subjects within the same cluster, as measured by intra-cluster correlation coefficient (ICC), should be considered carefully when analyzing data from cluster randomized trials. The positive correlation reduces the variability of responses in a clustered sample and thus reduces the statistical power to detect true differences between treatment arms relative to trials that randomize the same number of individuals (Killip et al., 2004; Donner and Klar, 2004).

A natural way to improve trial efficiency without undermining the validity and integrity of the trial is through the use of group sequential methods which divide participant entry into a number of groups and allow for planned assessments of the evidence in the cumulative data that incorporate decision-making regarding trial continuation (Chow et al., 2005; Pocock, 1977). In addition to flexibility and efficiency, group sequential designs are attractive to clinical scientists because they may reflect medical practice in the real world and they are ethical with respect to the need to determine and monitor efficacy as well as safety of the treatment (Chow and Chang, 2008).

Group sequential designs naturally fit into the Bayesian framework as results or estimates are continually updated based on the accumulated information from interim data (Chevret, 2012). Group sequential designs based on Bayesian approaches have been extensively studied in recent years (Berry et al., 2010; Gsponer et al., 2014; Yin et al., 2017; Freedman et al., 1994; Lewis and Berry, 1994; Freedman and Spiegelhalter, 1989). However, most group sequential designs, regardless of whether they are frequentist or Bayesian, focus on individually randomized trials and are categorized under the broad family of adaptive trial designs. Adaptive designs for cluster-randomized trials are less common and the incorporation of adaptive features poses significant statistical challenges. Some specific adaptive features such as sample size re-estimation and frequentist group sequential analyses have been proposed in combination with cluster-randomized trials (Lake et al., 2002; van Schie and Moerbeek, 2014; Grayling et al., 2017; Zou et al., 2005). However, no formal statistical design of Bayesian adaptive cluster-randomized trials has been developed. We address this gap by proposing two Bayesian group sequential designs for cluster-randomized trials.

The organization of this paper is as follows. The motivating trial is briefly described in section 3.2, followed by the development of two Bayesian group sequential designs and models for continuous and binary outcomes in section 3.3. Simulation studies are carried out in section 3.4, assessing the performance of the two proposed designs across a range of scenarios. The paper concludes in section 3.5.

3.2 Motivating setting

Aifred Health has designed a machine-learning based clinical decision support system (CDSS) for physicians treating patients with depression (Benrimoh et al., 2018; Popescu et al., 2021; Benrimoh et al., 2021; Tanguay-Sela et al., 2022). Patient characteristics of interests such as sociodemographic information, clinical information and medical history are input into the CDSS. The CDSS then, using a deep learning model, outputs the predicted efficacy for a number of possible treatments for that patient (Mehltretter et al., 2020b,a). Treatments are ordered by efficacy and presented to the physician when they reach the treatment selection step of a clinical algorithm based on best practice guidelines (Popescu et al., 2021; Kennedy et al., 2016). Physicians with the CDSS can, on an individual patient basis, decide whether or not to use the information presented by CDSS as part of their medical decision-making.

An important step is to establish the clinical utility, safety, and potential effectiveness of a tool such as the CDSS. Practically, an intervention such as the CDSS must be delivered at the physician level, such that a cluster-randomized design would be appealing for the reasons described above. Participating physicians could be randomized, with patients recruited from physicians' usual practices in order to approximate the real-world clinical conditions and populations as closely as possible.

Typical outcomes in a study of depression in such a trial could be a continuous measure of symptoms (e.g., a visual analog scale, the Quick Inventory of Depressive Symptomatology (Trivedi et al., 2004), the 9-question depression scale from the Patient Health Questionnaire (Kroenke and Spitzer, 2002), the World Health Organization Disability Assessment Schedule 2.0 (Gold, 2014), etc.) or a binary measure of treatment response, minimum clinically significant change in symptoms, or remission, perhaps defined by a dichotomization of a standard depression score (Beck and Alford, 2014; Smarr and Keefer, 2020; McGlothlin and Lewis, 2014). The importance of mental health treatment and the often relatively slow rate of patient accrual motivate the use of a Bayesian group sequential design to ensure adequate sample size and the possibility of early termination due to treatment effectiveness.

3.3 Methods

3.3.1 Two Bayesian group sequential designs for cluster-randomized trials

Suppose that the maximum number of clusters and maximum cluster size are equal across the two treatment arms. Let K, n, m be the number of interim analyses (not including the final analysis), the maximum number of clusters for each treatment arm, and the maximum cluster size, respectively. For simplicity of exposition, we will assume that all clusters enroll the same number of participants. Two Bayesian designs, design 1 and design 2 in the remainder of this paper, are developed to sequentially enroll participants and analyze interim data in different ways.

In design 1, [n/(K + 1)] clusters enter the trial at the start where [x] denotes the largest integer not exceeding x, and m individual participants will be enrolled for each cluster. At the subsequent analysis point, if the accumulated information up to the current analysis is sufficient to conclude the efficacy of the treatment or the final analysis has been reached, the trial will be terminated. Otherwise, another [n/(K+1)] new clusters will be enrolled, and m individual participants will be recruited for each new cluster. The trial then proceeds to the next analysis with new samples. This procedure is repeated until termination.

In design 2, at the beginning of the trial, all n clusters enter the trial, but only [m/(K+1)]individual participants are enrolled for each cluster. At the subsequent analysis point, if the trial is not to be terminated, another [m/(K+1)] individual participants are recruited for the same n clusters. The trial then proceeds to the next analysis. This procedure is repeated until termination.

The fundamental difference between the two designs lies in the way that participants are sequentially recruited. In design 1, the *clusters* are sequentially enrolled, and individual participants for each cluster are recruited all at once. In design 2, all clusters are enrolled at one time, but the *individual participants* for each cluster are sequentially enrolled. For illustration, consider an example for one treatment arm with K = 1, n = 4 and m = 6. That is, there is only one interim analysis planned (thus two analyses in total, including the final analysis). The maximum number of clusters is 4, and the maximum cluster size is 6. Figure 3.1 gives a graphical illustration. The circles represent different individual participants for the corresponding clusters. In design 1, clusters 1 and 2 will be first enrolled, and for each, six individual participants will be recruited. The interim analysis is based on data from clusters 1 and 2. If we decide to continue the trial, then we will further recruit clusters 3 and 4 and their respective six individual participants for the final analysis. The final analysis is based on data from all the four clusters. In design 2, all four clusters will be recruited at the start, but for each cluster, only three participants will be enrolled. If evidence based on the 12 individuals from 4 clusters is unable to conclude the efficacy of the intervention at the interim analysis, then an additional three individual participants for each cluster will be recruited, and the trial proceeds to the final analysis.



Figure 3.1: An illustration of the two designs for one treatment arm when there is only one pre-planned interim analysis, and at most four clusters with a total of six individual participants within each cluster over the course of the trial. The black labelled circles represent different individual participants for each cluster.

3.3.2 Early stopping at interim analysis

At each interim look, one should determine whether to stop the trial early or continue based on the interim result. The goal is to evaluate the efficacy of the treatment by testing the hypothesis

$$H_0: \theta \leq \delta \quad vs \quad H_A: \theta > \delta,$$

where θ is the mean difference for continuous outcome and risk difference for binary outcome, and δ is the minimal important difference. The efficacy of the treatment can be concluded and the trial can be stopped early for efficacy at the *k*th interim analysis if

$$P(\theta > \delta \mid D_k) > U, \tag{3.1}$$

where D_k is all the available data up to the kth analysis, and U is the decision boundary.

Continuous outcomes

Based on the context of the motivating trial and without loss of generality, assume that a smaller value of the continuous outcome is preferred (as most depression and disability rating scales associate worse symptoms with larger total scores). Then, $\theta = \mu_c - \mu_t$ where μ_c , μ_t are the mean outcome for control and treatment groups, respectively. We further assume that at the kth analysis, $k = 1, \ldots, K+1$, there are n_k clusters with m_k observations in each cluster. Let Y_{ij} be the continuous outcome of the *i*th subject in the *j*th cluster at current analysis point, $j = 1, \ldots, n_k$, and $i = 1, \ldots, m_k$. The normality assumption in cluster-randomized trials states

$$\mu_j \sim N(\mu, \sigma_B^2)$$
 and $Y_{ij} \mid \mu_j \sim N(\mu_j, \sigma_W^2)$,

where μ_j is the cluster-specific mean, μ is the population mean, σ_W^2 and σ_B^2 are within- and between-cluster variances respectively, and they can be related via the ICC, $\rho = \sigma_B^2/(\sigma_W^2 + \sigma_B^2)$. It can be shown that the marginal distribution of Y_{ij} is normal with mean μ and variance $\sigma_B^2 + \sigma_W^2$. Let $\mathbf{Y} = (\mathbf{Y}_1^\top, \dots, \mathbf{Y}_{n_k}^\top)^\top$ be the response vector where $\mathbf{Y}_j = (Y_{1j}, \dots, Y_{m_kj})^\top$, for $j = 1, \dots, n_k$. The covariance structure satisfies

$$\operatorname{Var}(Y_{ij}) = \sigma_B^2 + \sigma_W^2$$
$$\operatorname{Cov}(Y_{ij}, Y_{sj}) = \sigma_B^2, \quad \forall i \neq s$$
$$\operatorname{Cov}(Y_{ij}, Y_{tl}) = 0, \quad \forall j \neq l$$
$$(3.2)$$

so that $Y \sim \text{MVN}(\boldsymbol{\mu}, \Sigma)$ where $\boldsymbol{\mu}$ is a vector of $n_k \times m_k$ with all elements equal to $\boldsymbol{\mu}$ and Σ is a block matrix of the form

$$egin{pmatrix} \Sigma_{Y_1} & 0 & \cdots & 0 \ 0 & \Sigma_{Y_2} & \cdots & 0 \ dots & dots & \ddots & 0 \ 0 & 0 & \cdots & \Sigma_{Y_{n_k}} \end{pmatrix}$$

and $\Sigma_{Y_j} = \text{cov}(Y_j)$ specified in equation (3.2), and MVN indicates a multivariate normal distribution.

In the Bayesian framework, we assume a normal prior for μ ,

$$\mu \sim N(a, b^2)$$

where a is the prior mean and b^2 is the prior variance. Then, using Bayes' theorem, the posterior distribution for μ is:

$$\mu \mid Y \sim N\left(\frac{b^2 \mathbf{1}^\top \Sigma^{-1} Y + a}{b^2 \mathbf{1}^\top \Sigma^{-1} \mathbf{1} + 1}, \frac{b^2}{b^2 \mathbf{1}^\top \Sigma^{-1} \mathbf{1} + 1}\right)$$

The above result applies to both μ_t and μ_c . The quantity of interest $P(\mu_c - \mu_t > \delta \mid Y)$ can be obtained analytically from the posterior distributions of μ_t and μ_c , as the difference between two normal random variables is still normal.

At any interim analysis, if $P(\mu_c - \mu_t > \delta \mid Y) > U$, the trial is stopped early for efficacy. Otherwise, the trial continues to enroll clusters/participants and proceeds to the next analysis. These steps are repeated until either the trial is stopped early for efficacy or reaches the final analysis.

Binary outcomes

For binary outcomes, assume that a larger proportion is preferred (e.g., a larger proportion of patients meeting the criteria for treatment response). Let $\theta = \pi_t - \pi_c$ where π_c , π_t are the population proportion for the control and treatment groups, respectively. For binary outcomes, the joint distribution of all observations cannot be obtained analytically and thus a tractable form of the posterior distribution of the parameters of interest is not available. Therefore, we propose the following hierarchical model

$$\pi_j \sim \text{Beta}(\alpha, \beta)$$

 $r_j \mid \pi_j \sim \text{Binomial}(m_k, \pi_j)$

where r_j is the number of events in the *j*th cluster and π_j is the cluster-specific proportion, $j = 1, \ldots, n_k$. Then, under the model, the cluster-specific proportion π_j have the meanvariance relationship

$$\operatorname{Var}(\pi_j) = E(\pi_j) \{ 1 - E(\pi_j) \} \frac{1}{\alpha + \beta + 1}.$$
(3.3)

However, in cluster-randomized trials with binary outcomes, it is assumed that

$$E(\pi_j) = \pi,$$

$$Var(\pi_j) = \rho \pi (1 - \pi),$$
(3.4)

where π is the population proportion and ρ is the ICC.

To make (3.3) and (3.4) consistent, define two transformed parameters

$$\pi = \frac{\alpha}{\alpha + \beta},$$
$$v = \alpha + \beta,$$

where π is exactly the mean of the Beta distribution and v measures the information in the corresponding Beta distribution. Also, due to the consistency of (3.3) and (3.4), once ρ is fixed or estimated, v can also be determined through $v = (1 - \rho)/\rho$. Thus, the only free parameter for the hierarchical model is π .

In the Bayesian framework, assume a uniform prior for π ,

$$\pi \sim \text{Unif}[0, 1].$$

Then the posterior probability $P(\pi_t - \pi_c > \delta \mid R)$, where $R = (r_{1c}, \ldots, r_{n_k,c}, r_{1t}, \ldots, r_{n_k,t})$ includes the number of events per cluster within each treatment group, can be approximated by

$$\hat{\pi} = \frac{1}{M} \sum_{i=1}^{M} I\{\pi_{ci}^{post} - \pi_{ti}^{post} > \delta\}$$

where π_{ci}^{post} , π_{ti}^{post} are sampled from the corresponding posterior distributions of π_c and π_t , and M is the number of Monte Carlo samples drawn from the posterior. The posterior distributions of π_c and π_t are not analytically available; posterior samples can be obtained using Markov Chain Monte Carlo (MCMC) implemented in any Bayesian software. In this paper, we use RStan (Stan Development Team, 2021, 2020).

3.4 Simulation studies

The false positive rate and power cannot be obtained analytically in Bayesian adaptive trials since the sampling distributions of the test statistics (i.e., posterior probability statements (3.1) in section 3.3.2) are not known. Therefore simulation studies are required to specify the decision boundaries and other design characteristics (Berry et al., 2010). For both outcomes, a single interim analysis was explored first; multiple interim analyses were then investigated. In the presence of cluster-level covariates, stratified randomization could be considered to

help reduce covariate imbalance between treatment arms. In our simulation, covariates were not considered, and thus stratified randomization was not required; permuted block randomization with block size 4 was implemented for treatment allocation. The minimal important difference was set as 0 in the simulation but it is straightforward to generalize to other values. For each scenario, 3000 simulation replications were performed. The performance of designs was compared based on overall false positive rate and power, accounting for the rejection of the null at either interim or final analysis. The false positive rate is estimated as

False positive rate =
$$\frac{\text{Number of times the null hypothesis is falsely rejected}}{\text{Number of simulation runs}}$$
,

where falsely rejecting the null hypothesis means that for some $k \leq K + 1$,

$$P_{\theta=0}(\hat{\theta}_k > \delta \mid D_k) > U$$

for $\hat{\theta}_k$ the estimated mean difference for continuous outcome or estimated risk difference for binary outcome at the *k*th analysis, θ is the true mean or risk difference, δ is minimal important difference, and $\delta = 0$ in our simulation. D_k is all the available data up to the *k*th analysis, and *U* is the decision boundary as described in section 3.3.2.

Power is estimated as

$$Power = \frac{Number of times the null hypothesis is correctly rejected}{Number of simulation runs}$$

where correctly detecting the difference means that for some $k \leq K + 1$,

$$P_{\theta=\theta_A}(\hat{\theta}_k > \delta \mid D_k) > U$$

for θ_A the value of θ under the alternative hypothesis.

For continuous outcomes, the prior mean and variance for the population mean for both

	Continuous		Binary	
Parameters	Single interim analysis	Multiple interim analyses	Single interim analysis	Multiple interim analyses
Population mean for control [†] (μ_c)	0	0	\	\
Within-cluster variance (σ_W^2)	1	1	\setminus	\backslash
Baseline risk‡ (π_c)	\setminus	\setminus	0.25, 0.35, 0.45	0.25,0.35,0.45
Number of clusters per group (n)	20, 40, 60	20, 40, 60	20, 40, 60	20, 40, 60
True treatment effect (θ)	$0, 0.1, \ldots, 0.9$	$0, 0.1, \ldots, 0.9$	0, 0.1, 0.2, 0.3	0, 0.1, 0.2
Decision boundary (U)	0.95, 0.98	0.95, 0.98	0.95, 0.98	0.95, 0.98
Intra-cluster correlation coefficient (ρ)	0.05, 0.1, 0.15	0.05, 0.1, 0.15	0.05, 0.1	0.05, 0.1
Cluster size (m)	8	8, 16	8	8, 16
Number of interim looks (K)	1	1, 2, 3	1	1, 3

Table 3.1: Simulation parameters for continuous and binary outcome

[†] Not applicable to binary outcomes.

‡ Not applicable to continuous outcomes.

groups are fixed at 0 and 1000, respectively. In real-world cluster-randomized trials, the ICC tends to be small (e.g. around 0.05 in primary care trials (Campbell, 2000)). Similar values of ICC were explored in simulation and the between-cluster variance σ_B^2 was determined through $\sigma_B^2 = \sigma_W^2 \rho / (1 - \rho)$. To generate clustered continuous data, we first generate the n cluster-specific means from normal distributions with mean μ_c (control group) and μ_t (treatment group) and variance σ_B^2 . Then, within each cluster, m samples are drawn from a normal distribution with mean equal to the cluster-specific mean and variance σ_W^2 . The resulting $m \times n$ samples are expected to satisfy the preset correlation structure. For binary outcomes, clustered binary data are generated via Beta and binomial distributions; see Appendix A.1 for details. All simulation parameters are summarized in Table 3.1.

Only results for continuous outcomes are presented here. The results for binary outcomes can be found in the Appendix A.1. Figures 3.2 and 3.3 show the false positive rate and power when a single interim analysis is planned. In general, the two designs have comparable power but design 1 has higher false positive rates. The accrual of information (measured by Fisher information) for design 1 is linear with increasing number of analyses, as the newly recruited clusters are independent from already-recruited clusters, while diminishing return in the accrual of information for design 2 can be expected due to the correlation between responses on individuals from the same cluster. The information at final analysis, if reached, is the same for both designs. However, at any interim analysis, design 2 contains more information than design 1, leading to smaller false positive rates.

The effect of ICC on false positive rate is quite small whereas the effect on power is more evident. Under the same conditions, if the underlying ICC is higher, power is lower. In conventional cluster-randomized trials, the sample size needed to maintain a sufficient power is jointly determined by ICC, effect size, cluster size as well as outcome variance. In the Bayesian group sequential design, the decision boundary should also be determined to achieve desired operating characteristics and sample sizes. In all scenarios, with U = 0.95, the false positive rates for both designs are well above 0.05. With U = 0.98, both false positive rates and powers are reduced. In practice, the value of U is specified to achieve a small false positive rate (e.g. 0.05), while maintaining a required level of power. A larger decision boundary corresponds to a more conservative test since more evidence is required to conclude the treatment efficacy. However, an over-conservative boundary may increase the cost as more samples may be needed to maintain sufficient power.



Figure 3.2: Plot of false positive rate versus ICC when n = 20, 40, 60 for (a) U = 0.95 and (b) U = 0.98 with single interim analysis planned. The dashed lines show the false positive rate of 0.05.

Figures 3.4 and 3.5 show the false positive rate and power when multiple interim looks are built into the study design. Design 1 still has higher false positive rate and similar



Figure 3.3: Plot of power versus true treatment effect when $n = 20, 40, 60, \rho = 0.05, 0.1, 0.15$ for (a) U = 0.95 and (b) U = 0.98 with single interim analysis planned. The dashed lines show the power of 0.8.

power compared with design 2. As was observed in Figures 3.2 and 3.3, a larger decision boundary can reduce false positive rate and the resulting reduction in power can be remedied by recruiting a larger number of clusters. Adding more interim analyses may increase the false positive rate and power. Given a fixed decision boundary, with more interim analyses planned, there is higher chance of rejecting the null hypothesis. The effect of cluster size on both false positive rate and power is very modest. The results for a larger cluster size can be found in the Appendix A.2.



Figure 3.4: Plot of false positive rate versus number of interim looks for n = 20, 40, 60, $\rho = 0.05, 0.1, 0.15, m = 8$ for (a) U = 0.95 and (b) U = 0.98. The dashed lines show the false positive rate of 0.05.

To conclude, both designs perform better in terms of false positive rate with U = 0.98 as compared to U = 0.95, and design 1 has a higher false positive rate as compared to design 2 in almost all scenarios. The choice between the two designs as well as the decision boundaries depends on the research question, phase of the study, and the information we have about the treatment. For example, in some early phase clinical trials, if a relatively high false positive rate is acceptable, then a smaller decision boundary (e.g. U = 0.95) may be recommended. However, in late phase clinical trials, controlling false positive rate may be more important, and so design 2 with a larger decision boundary may be preferred. For the situation we explored, a single interim analysis is preferred no matter which design or decision boundary is chosen, as the false positive rate with single interim analysis is lowest and an increase in power by incorporating multiple interim analyses is not needed. General recommendations require further exploration regarding the trade-off between false positive rate and power as well as other practical consideration.

3.5 Discussion

Motivated by a potential real-world cluster-randomized trial, we explore the statistical properties of Bayesian group sequential designs for cluster-randomized trial. We explore stopping rules for efficacy in a cluster-randomized trial. Interim analyses may be planned over the course of the trial and at each interim analysis the trial may be stopped early if sufficient evidence is established to conclude efficacy otherwise the trial proceeds to the next analysis. We proposed two designs which sequentially enroll participants in different ways. The first design sequentially enrolls clusters, and individual participants for each cluster are recruited all together. The second recruits all clusters at the start of the trial, then sequentially enrolls batches of participants. The difference between the designs lies in the sequential enrollment of participants and thus the information accrual process. In design 1, as the clusters recruited for each analysis are independent of each other, the information is accrued linearly with increasing number of analysis stages. However, in design 2, more information may be accumulated at the start of the trial, and due to the positive ICC there are diminishing returns in information accumulation with more participants enrolled for the same clusters.

Interim analyses of both continuous and binary outcomes are performed based on common models. For continuous outcomes, on the basis of the normality assumption, we obtained the analytical form of the posterior distribution of the population mean. Based on Monte Carlo simulation, the posterior probability of efficacy can be easily estimated by drawing random samples directly from the two posterior distributions calculated from the accumulated data. one for each treatment group. For design 1, due to the independence between the clusters enrolled at each analysis, it is also possible to update prior distribution at each analysis by the posterior distribution obtained from the previous analysis and base each analysis on new samples only. This computation allows us to use only interim analysis estimates rather than the original data for later analysis. However, for design 2, analysis should be performed on all accumulated data in order to avoid loss of information regarding the correlation between the responses on individuals within the same clusters and thus avoid inflation of false positive rate. For binary outcomes, due to the complex correlation structure, the posterior distribution of population proportion cannot be obtained analytically. Instead, a hierarchical model was used and the posterior probability of efficacy is estimated via MCMC.

Given that our motivating trial will be the first clinical trial of CDSS, accurate estimates for design parameters such as effect size are not known with much certainty. The simulation parameters in Table 3.1 were informed by rough estimates provided by Aifred Health but covered a wider range of scenarios. Simulation results showed that design 2 with a single interim analysis would be recommended for our motivating trial based on design operating characteristics, as it has smaller false positive rates and comparable powers to design 1 or multiple interim analyses. However, in practice the feasibility of the two designs is also a determinant factor. For example, if the recruitment process of physicians is very slow, it may not be possible to recruit all physicians simultaneously or in a short time. In this case, design 1 may be preferred. The company would need to consider these issues carefully before implementing a real trial. A more general recommendation for other cluster-randomized trials requires further exploration, as different cluster-randomized trials may have different research goals or different design parameters that we did not explore in our simulation. In addition, increasing cluster size may not necessarily bring much improvement in design performance. Keeping a small cluster size may not hurt design operating characteristics very much but in some instances it may save time and cost. In our motivating trial, a smaller cluster size is also more realistic, as it reduces the burden on the individual clinician to find suitable patients from their practice.

There are some limitations to our work. First, we only considered implementing a stopping rule for efficacy. There are many other adaptive features that we have not considered. One example is adaptive randomization in cluster-randomized trials. Second, we only focus on two-arm trials. Extension to multi-arm trials which may involve arm dropping can be considered in the future. Third, we only considered continuous and binary outcomes. The work can be extended to other outcome types such as time-to-event outcomes. An extension to survival outcomes based on survival models can also improve the framework of Bayesian adaptive cluster-randomized trial. In addition, the effect of unequal cluster sizes has not been explored for our designs. However, the effect of varying cluster sizes on design operating characteristics in standard cluster-randomized trials has been extensively discussed (Campbell et al., 2007; Eldridge et al., 2006; Guittet et al., 2006; Kerry and Martin Bland, 2001; Manatunga et al., 2001). Typically, designs are most efficient for equal cluster sizes. Inflation of false positive rate will occur for imbalanced studies, and with the same sample size, imbalanced trials may be underpowered compared with their balanced counterparts. Similar impacts of unequal cluster sizes on power or false positive rate may also exist for our Bayesian adaptive cluster-randomized trials, and this would need to be determined in future work. Also, practical issues in planning a Bayesian group sequential design for a clusterrandomized trial are not taken into account in our paper. For example, sometimes it may not be realistic to plan the preferred design recommended by design operating characteristics. In this case, feasibility may be the main reason for design choice. These context-specific issues require more exploration in the future.

Declaration of competing interest

D.B. is a director, shareholder, and employee of Aifred Health.

Acknowledgements

This work was supported by funding from MITACS Accelerate Grant #ACC IT18791.

Supporting information

The following supporting information is available as part of the online article:

Appendix A.1. Simulation for binary outcomes.

Appendix A.2. Simulation results for continuous outcomes with cluster size m=16.



Figure 3.5: Plot of power versus number of interim looks for $n = 20, 40, 60, \theta = 0.2, 0.5, 0.8, m = 8$ with the subpanels (a)-(f) indicating all possible combinations of $\rho = 0.05, 0.1, 0.15$ and U = 0.95, 0.98. The dashed lines show the power of 0.8.

Chapter 4

Sparse two-stage Bayesian meta-analysis for individualized treatments

Preamble to Manuscript 2. In Manuscript 1, I described two statistical designs to address the low power of CRTs when evaluating the efficacy of the CDSS. Since such a CDSS relies on a particular ITR, it is crucial to accurately estimate the ITR before assessing the CDSS's efficacy. One way to increase the power and generalizability of estimating ITRs in regression-based approaches is to combine data from multiple studies and sites. While individual-level data are thought to be essential for estimating ITRs in regression-based approaches, sharing these highly sensitive health data can be restricted by policies, creating a methodological gap. Manuscript 2 addresses this challenge with a two-stage Bayesian metaanalysis approach. The proposed two-stage approach is compared to a one-stage approach, where individual-level data from different sites are pooled together and analyzed in a single statistical model to estimate an ITR. Practical solutions to sparsity in both the data and the model are also discussed. An application of the proposed method to the International Warfarin Pharmacogenetics Consortium data (2009) is presented. This manuscript has been submitted for peer review.

Sparse two-stage Bayesian meta-analysis for individualized treatments

Junwei Shen¹, Erica E.M. Moodie¹, Shirin Golchi¹.

¹Department of Epidemiology, Biostatistics, and Occupational Health, McGill University

Abstract

Individualized treatment rules tailor treatments to patients based on clinical, demographic, and other characteristics. Estimation of individualized treatment rules requires the identification of individuals who benefit most from the particular treatments and thus the detection of variability in treatment effects. To develop an effective individualized treatment rule, data from multisite studies may be required due to the low power provided by smaller datasets for detecting the often small treatment-covariate interactions. However, sharing of individuallevel data is sometimes constrained. Furthermore, sparsity may arise in two senses: different data sites may recruit from different populations, making it infeasible to estimate identical models or all parameters of interest at all sites, and the number of non-zero parameters in the model for the treatment rule may be small. To address these issues, we adopt a two-stage Bayesian meta-analysis approach to estimate individualized treatment rules which optimize expected patient outcomes using multisite data without disclosing individual-level data beyond the sites. Simulation results demonstrate that our approach can provide consistent estimates of the parameters which fully characterize the optimal individualized treatment rule. We estimate the optimal Warfarin dose strategy using data from the International Warfarin Pharmacogenetics Consortium, where data sparsity and small treatment-covariate interaction effects pose additional statistical challenges.

4.1 Introduction

In the traditional one-size-fits-all approach, patients with the same disease receive the same treatment regardless of their individual characteristics. This strategy can be suboptimal, as the best treatment for one patient might be ineffective for another due to treatment effect heterogeneity among patient subgroups. The personalized medicine approach has emerged as an alternative, leveraging this treatment effect heterogeneity in clinical decision making to recommend the most appropriate treatment to individual patients (Chakraborty and Moodie, 2013). Dynamic treatment regimes (DTRs), which consist of a sequence of decision rules, formalize the statistical framework of personalized medicine in the setting of multiple treatment stages (Chakraborty and Moodie, 2013; Chakraborty and Murphy, 2014; Laber et al., 2014). When only a single treatment stage is considered, a DTR reduces to an individualized treatment rule (ITR), which is the focus of our work.

Several methods have been proposed to estimate optimal ITRs (i.e., rules to optimize expected patient outcomes) or their multi-stage counterparts for various outcome or exposure types based on individual-level data. Among the many alternatives, regression-based approaches indirectly estimate the optimal ITR by first modelling the expected outcome as a function of treatment, covariates, and their interactions and then determining the treatment that, for each covariate combination, will optimize the estimated expected outcome. Common regression-based approaches include Q-learning (Watkins, 1989; Sutton and Barto, 2018), G-estimation (Robins, 2004), and dynamic weighted ordinary least squares (dWOLS) (Wallace and Moodie, 2015). One practical challenge of ITR estimation is the low power for detecting treatment-covariate interactions (Greenland, 1983). To increase the sample size and generalizability of the findings, multisite data are attractive. Ideally, in a multisite study one could pool the individual-level data from different sites together and analyze the pooled data in a common analysis center. However, this may be infeasible due to, e.g., institutional policies which restrict the sharing of individual-level information. There-

fore, it is desirable to develop valid methods for ITR estimation that avoid sharing individual records.

Different strategies have been proposed for analyses under constrained data sharing (Rassen et al., 2013), but few have focused specifically on restrictions on data sharing in the context of ITR estimation (Danieli and Moodie, 2022; Moodie et al., 2022) and none account for the possibility that parameters of interest may vary across centers or sites. Meta-analysis, a widely-used approach for evidence synthesis, can avoid releasing individual-level records when analyzing multisite data (Rassen et al., 2013). However, classic meta-analysis techniques based on aggregate data may be unsuitable for identifying treatment effects within subgroups under heterogeneity (Berlin et al., 2002; Simmonds and Higgins, 2007), which is the target of personalized medicine, and so it is unclear whether a meta-analysis approach to ITR estimation is feasible. While using individual participant data (IPD) for metaanalysis (Riley and Fisher, 2021) appears promising for ITR analysis, a one-stage IPD metaanalysis requires combining all individual-level data into a single dataset, and thus cannot be used in settings where individual records cannot be shared and thus a two-stage approach is required. In this work, we adopt a Bayesian two-stage IPD meta-analysis approach to estimate the optimal ITR using multisite data without the need for sharing individual-level information across sites. We note that treatment-covariate interactions sometimes are described as "treatment effect heterogeneity" in the literature. However, in the meta-analysis literature, heterogeneity typically refers to the variability across sites. To avoid confusion, from now on, we adopt the meta-analysis tradition and reserve the word *heterogeneity* for variability across sites.

Warfarin is a widely-used oral anticoagulant for thrombosis and thromboembolism treatment and prevention (Rettie and Tai, 2006). Establishing an optimal Warfarin dose strategy is of vital importance due to the narrow therapeutic window and the large interindividual variability in patients' response to the drug (Rettie and Tai, 2006; International Warfarin
Pharmacogenetics Consortium, 2009). The data from the International Warfarin Pharmacogenetics Consortium was collected in nine countries from four continents. These data have suggested factors that are possibly associated with Warfarin dosing (International Warfarin Pharmacogenetics Consortium, 2009). However, statistical challenges arise in the context of a two-stage Bayesian meta-analysis of the optimal Warfarin dose strategy.

One challenge is data sparsity, a term we use to refer to the phenomenon of not observing a sufficient number of patients with a given set of characteristics. For example, VKORC1 genotype (AA, AG or GG) is potentially a tailoring variable for Warfarin dosage. In certain sites all patients fall under only one or two categories of VKORC1 genotype and therefore inference about the interaction between VKROC1 genotype and Warfarin dose cannot be made based on the data from these sites. A naïve approach would remove sites with sparse data, leading to a significant loss of information.

A second challenge is model sparsity, i.e., extremely small (practically zero) treatmentcovariate interactions resulting from variables that are irrelevant for the dosing decision. The small effect estimates can lead to invalid dose recommendations (e.g., outside of the appropriate range) or the estimated optimal dose for an individual being highly volatile or sensitive to the parameter estimates. Including covariates with no tailoring effect also makes the estimated optimal dosing strategy unnecessarily complex.

We address data sparsity by reparametrizing the likelihood such that the site-specific estimates are linked to the correct set of parameters. To address the second challenge, we use shrinkage priors (van Erp et al., 2019), opting for a horseshoe prior due to its proven advantages in maintaining large effects and efficiently handling sparsity (Carvalho et al., 2010). Methods including the notation and assumptions are described in section 4.2. In section 4.3, we explore the performance of the proposed method across a range of scenarios via simulation studies. In section 4.4, we estimate an optimal individualized Warfarin dose strategy without the need for sharing patient-level data. The paper concludes with a discussion.

4.2 Methods

4.2.1 Preliminaries

Let Y denote a continuous outcome of interest, where larger values of Y are preferable. Let A denote the binary or continuous treatment received by the patient, and X be a vector of pre-treatment covariates. Let Y^a be the potential outcome a patient would experience if assigned treatment a. Uppercase, lowercase, and bold denote random variables, realizations of random variables, and vectors respectively.

An ITR $d(\mathbf{X}) : \mathbf{X} \to A$ tailors treatment to patients based on individual characteristics, \mathbf{X} . An optimal ITR $d^{\text{opt}}(\mathbf{X})$ maximizes the value function under $d^{\text{opt}}(\mathbf{X})$, that is, the expected potential outcome $E(Y^{d(\mathbf{X})})$ if all patients in a population are treated according to $d(\mathbf{X})$. Identification of an optimal ITR relies on several assumptions: (i) the stable unit treatment value assumption (SUTVA): a patient's outcome is not influenced by other patients' treatment (Rubin, 1980); (ii) no unmeasured confounding (Robins, 1997); (iii) positivity: $p(A = a \mid \mathbf{X} = \mathbf{x}) > 0$ almost surely for all possible \mathbf{x} and a (Cole and Hernán, 2008). Additional modelling assumptions will also be required for the (regression-based) meta-analytic approach that we pursue.

To identify the optimal ITR, an outcome model could be specified and decomposed into two components: $E(Y \mid A = a, \mathbf{X} = \mathbf{x}) = f(\mathbf{x}^{(\beta)}; \boldsymbol{\beta}) + \gamma(a, \mathbf{x}^{(\psi)}; \boldsymbol{\psi})$, where both $\mathbf{x}^{(\beta)}$ and $\mathbf{x}^{(\psi)}$ are subvectors of \mathbf{x} , and include predictive covariates and covariates that interact with treatment (prescriptive variables), respectively. The prescriptive covariate vector $\mathbf{x}^{(\psi)}$ is a subvector of $\mathbf{x}^{(\beta)}$. The parameter vectors in the treatment-free function $f(\mathbf{x}^{(\beta)}; \boldsymbol{\beta})$ and the blip function $\gamma(a, \mathbf{x}^{(\psi)}; \boldsymbol{\psi})$ (Robins, 2004) are denoted by $\boldsymbol{\beta}$ and $\boldsymbol{\psi}$, respectively. The treatment-free function depends on covariates $\mathbf{x}^{(\beta)}$ but not treatment, and thus is not relevant to treatment or dosing decisions. The difference in the expected potential outcome of patients receiving treatment A = a and reference A = 0, with the same prescriptive covariates $\mathbf{X} = \mathbf{x}^{(\psi)}$, is represented by the blip function: $\gamma(a, \mathbf{x}^{(\psi)}; \psi) = E(Y^a - Y^0 | \mathbf{X} = \mathbf{x}^{(\psi)}; \psi)$. The blip function satisfies $\gamma(0, \mathbf{x}^{(\psi)}; \psi) = 0$ and the optimal ITR is then defined as $d^{\text{opt}}(\mathbf{x}) = \arg \max_a \gamma(a, \mathbf{x}^{(\psi)}; \psi)$, since the treatment assignment influences the expected outcome only through the blip function. Therefore, estimation of the optimal ITR requires correct specification of the blip function and estimation of the blip parameter ψ . A common form for the blip function in the binary treatment setting $A \in \{0,1\}$ is $\gamma(a, \mathbf{x}^{(\psi)}; \psi) = ag(\mathbf{x}^{(\psi)}; \psi)$, and linear models could be assumed for both f and g. Under a linearity assumption, the outcome model becomes

$$E(Y \mid A = a, \boldsymbol{X} = \boldsymbol{x}) = \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\beta})} + a \boldsymbol{\psi}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\psi})}, \qquad (4.1)$$

where, in this form, we assume that $\boldsymbol{x}^{(\beta)}$ and $\boldsymbol{x}^{(\psi)}$ have been augmented by a column of ones to ensure an intercept and main effect of treatment, respectively. In this setting, the optimal ITR is given by $d^{\text{opt}}(\boldsymbol{x}) = I(\boldsymbol{\psi}^{\top}\boldsymbol{x}^{(\psi)} > 0)$. For a continuous treatment (e.g., a dose of a drug), a blip function should be specified so that the optimal treatment can be an interior point of the set of possible treatments. For example, a quadratic or higher order term for the treatment could be included in the blip function: $\gamma(a, \boldsymbol{x}^{(\psi)}; \boldsymbol{\psi}) = (\boldsymbol{\psi}^{(1)\top}, \boldsymbol{\psi}^{(2)\top})(a\boldsymbol{x}^{(\boldsymbol{\psi}^{(1)})}, a^2 \boldsymbol{x}^{(\boldsymbol{\psi}^{(2)})})^{\top}.$

Various approaches are available for unbiased and consistent estimation of the blip parameter ψ . For example, basic Q-learning can fit a standard linear regression to the model in equation (4.1) or indeed a more flexible model; consistency of the estimation is guaranteed under correct model specification (Chakraborty and Moodie, 2013). Other regression-based methods such as G-estimation and dWOLS, both of which are doubly robust, could also be employed (Robins, 2004; Wallace and Moodie, 2015). Bayesian approaches have also been proposed, such as Bayesian G-computation (Arjas and Saarela, 2010), a Bayesian machine learning approach to Q-learning (Murray et al., 2018), Bayesian additive regression trees (Logan et al., 2019), and Bayesian causal forest (Hahn et al., 2020).

4.2.2 Two-stage IPD meta-analysis

We now describe a two-stage IPD meta-analysis approach in the Bayesian framework to estimate the optimal ITR when multisite data are available but data sharing across sites is not allowed. Let K be the number of sites. For illustration purposes, we assume linear models for both treatment-free and blip functions. The extension to other parametric outcome models is straightforward. The site-specific outcome model can be written as: $E(Y_{ij} | \mathbf{X} = \mathbf{x}_{ij}, A =$ $a_{ij}) = \boldsymbol{\beta}_i^{\top} \mathbf{x}_{ij}^{(\beta)} + a_{ij} \boldsymbol{\psi}_i^{\top} \mathbf{x}_{ij}^{(\psi)}$, where $i \in \{1, \ldots, K\}$ and $j \in \{1, \ldots, n_i\}$ index the site and individual patient in a given site respectively, and n_i is the number of patients in site *i*. The predictive and prescriptive covariate vectors are denoted by $\mathbf{x}_{ij}^{(\beta)}$ and $\mathbf{x}_{ij}^{(\psi)}$, respectively. The *p*-dimensional site-specific treatment-free parameter $\boldsymbol{\beta}_i = (\beta_{i0}, \ldots, \beta_{i,p-1})$ and *q*-dimensional blip parameter $\boldsymbol{\psi}_i = (\psi_{i0}, \ldots, \psi_{i,q-1})$ have similar interpretations to those in equation (4.1), except that the target here is the site-specific ITR. The site-specific variables included in the vectors $\mathbf{x}_{ij}^{(\beta)}, \mathbf{x}_{ij}^{(\psi)}$ are identical across sites; we later relax that assumption.

Suppose that our interest is not in the site-specific optimal ITRs but a common optimal ITR that could be applied to all sites and, more generally, to future patients at comparable sites that may not have contributed data to the estimation. When the site-specific parameters are not identical across sites, it may be reasonable to assume that a common distribution exists for the varying site-specific blip parameters:

$$\boldsymbol{\psi}_{\boldsymbol{i}} \sim \mathrm{MVN}(\boldsymbol{\psi}, \boldsymbol{\Sigma}_{\boldsymbol{\psi}}),$$
 (4.2)

where MVN represents the multivariate normal distribution, and $\psi = (\psi_0, \dots, \psi_{q-1})$ and Σ_{ψ} are the common mean vector and variance-covariance matrix, respectively.

The two-stage IPD meta-analysis approach estimates a common optimal ITR by first conducting separate analyses of site-specific optimal ITRs and then combining the site-specific optimal ITRs via a hierarchical model. Specifically, at the first stage, each site obtains estimates of blip parameters $\hat{\psi}_{it}$ and the associated standard deviations $\mathrm{sd}(\hat{\psi}_{it})$, for $t = 0, \ldots, q - 1$, As mentioned in the last section, $\hat{\psi}_{it}$ and $\mathrm{sd}(\hat{\psi}_{it})$ can be obtained from various approaches such as Q-learning or dWOLS, using only site-specific data. At the second stage, only those site-specific estimates will be transferred to a common analysis center (and thus the individual-level data are preserved), and combined in a Bayesian hierarchical model:

$$\hat{\psi}_{it} \sim N(\psi_{it}, \operatorname{sd}(\hat{\psi}_{it})^2), \quad \psi_{it} \sim N(\psi_t, \sigma_{\psi_t}^2),
\psi_t \sim p_{\psi_t}(\psi_t), \quad \sigma_{\psi_t}^2 \sim p_{\sigma_{\psi_t}^2}(\sigma_{\psi_t}^2).$$
(4.3)

Here, ψ_{it} and ψ_t are the (t + 1)th elements of the site-specific and common blip parameter vectors. The between-site heterogeneity associated with ψ_{it} is denoted by $\sigma_{\psi_t}^2$. Prior distributions p_{ψ_t} and $p_{\sigma_{\psi_t}^2}$ can be assigned for the unknown parameters ψ_t and $\sigma_{\psi_t}^2$. Popular prior choices include a normal prior with large variance for the mean parameter ψ_t and a half-Cauchy prior for the variance component parameter σ_{ψ_t} (Gelman, 2006; Gelman et al., 2013). We use these prior choices in our simulation studies and the optimal Warfarin dosing analysis that follows. The Bayesian hierarchical model can be easily fitted in any Bayesian software. In this paper, we use RStan (Stan Development Team, 2020, 2021).

The two-stage IPD meta-analysis approach requires each site to provide estimates of blip parameters and the associated standard deviations only, which avoid sharing the individuallevel data. At the first stage, the treatment-free parameters are estimated together with the blip parameters. However, they are irrelevant for the optimal treatment decision and will not be used in the second stage. We make no common distribution assumptions of sitespecific treatment-free parameters in (4.2). In reality, if the blip parameters come from a common distribution, then it is highly likely that the treatment-free parameters are also from a common distribution. However, this is not required for our approach. The model at the second stage depends on unbiasedness, consistency, and normality of site-specific estimates from the first stage. Therefore, alternative stage-one models could be considered for different site-specific ITRs, as long as unbiasedness, consistency, and normality are assured. In the following sections, we assume no model misspecification exists. We establish the link between the proposed two-stage approach and a one-stage approach based on full individual-level data in Appendix B.1 of the Supplementary Materials.

4.2.3 Sparsity

So far, we have assumed that the specific variables included in the vectors $x_{ij}^{(\beta)}$, $x_{ij}^{(\psi)}$ are identical across sites. However, estimation of identical models at all sites may be infeasible due to heterogeneity of patient populations across sites (data sparsity). In addition, when a large number of covariates are available, we may believe that the number of non-zero parameters in the model is small (model sparsity). We now describe our approach to sparsity.

Data sparsity

Data sparsity occurs when insufficiently many patients with a given set of characteristics are represented in the samples at all sites. Then, not all site-specific parameters can be estimated for sites with sparse data. We employ a Bayesian hierarchical model which borrows information across sites; this requires modification of likelihood contribution for sites with data sparsity.

To illustrate this, consider a toy example, with a binary covariate $X \in \{0, 1\}$ (i.e., p = q = 2), and the following true outcome model for an individual at site i: $E(Y | X) = \beta_{i0} + \beta_{i1}X + A(\psi_{i0} + \psi_{i1}X)$. Therefore, ψ_{i0} is the difference in the mean outcome between A = 1 and A = 0 when X = 0 in site i; ψ_{i1} is the difference in the treatment effect between patients with X = 1 and X = 0 in site i; $\psi_{i0} + \psi_{i1}$ is the treatment effect for patients with X = 1 in site i. Then we consider two scenarios for data sparsity.

In the first scenario, X = 0 for all patients in site *i*. In this case, the main effect of X, β_{i1} , and the treatment-covariate interaction ψ_{i1} cannot be estimated. The site-specific outcome model reduces to $E(Y \mid X) = \gamma_{i0} + A\xi_{i0}$. Since this model is fitted among patients with $X = 0, \xi_{i0}$ is the treatment effect for patients with X = 0. That is, $\xi_{i0} = \psi_{i0}$. The likelihood contribution of site *i* is then $\hat{\xi}_{i0} \sim N(\psi_{i0}, \operatorname{sd}(\hat{\xi}_{i0})^2)$. In the second scenario, X = 1 for all patients in site *i*. The same reduced outcome model is fitted among patients with X = 1. Therefore, ξ_{i0} is the treatment effect for patients with X = 1 in site *i*, i.e., $\xi_{i0} = \psi_{i0} + \psi_{i1}$. In this case, without examining the data one might naïvely link the estimate $\hat{\xi}_0$ to ψ_0 via $\hat{\xi}_0 \sim N(\psi_0, \operatorname{sd}(\hat{\xi}_0)^2)$, but the correct specification is $\hat{\xi}_{i0} \sim N(\psi_{i0} + \psi_{i1}, \operatorname{sd}(\hat{\xi}_{i0})^2)$. A second toy example is provided in Appendix B.2 of the Supplementary Materials.

Model sparsity

Often, many potential tailoring variables are available but only a few are truly predictive of patient response to treatment. Including all available covariates in the ITR estimation will result in near-zero treatment-covariate interactions, which may result in invalid treatment recommendations and an uninterpretable optimal ITR. To address this, shrinkage priors which aim to shrink small effects towards zero are considered. We use horseshoe priors, though alternatives exist (see van Erp et al., 2019, for a review). Specifically, for treatment-covariate interactions, we assume $\psi_t \sim N(0, \tau^2 \lambda_t^2)$, where λ_t and τ are local and global shrinkage parameters, respectively, and $\lambda_t, \tau \sim$ Half-Cauchy(0, 1). The shrinkage prior is not placed on the main effect of treatment, ψ_0 , as we assume the treatment under consideration has at least some effect on the outcome of interest. We select the treatment-covariate interactions based on posterior credible intervals (van Erp et al., 2019): if $[\psi_{t,0.025}, \psi_{t,0.975}]$ does not include zero, the treatment-covariate interaction corresponding to ψ_t will be selected, where $\psi_{t,0.025}$ and $\psi_{t,0.975}$ are the 2.5th and 97.5th percentiles of the posterior distribution of ψ_t , respectively.

4.3 Simulation studies

The simulation study is reported following the scheme proposed in Morris et al. (2019). A brief summary is given below, with details provided in Appendix B.3 of the Supplementary Materials.

4.3.1 Overview

The simulation study aims to evaluate ITR estimation for a continuous outcome when the individual-level data from multisite studies are not shared across sites, varying: the confounding across sites (both variable set and strength), the degree of heterogeneity across sites, and the choice of prior distribution used in the analysis. Simulations consider both the case of binary treatments and, inspired by our motivating example of the International Warfarin Pharmacogenetics Consortium, a continuous dose. We include a sparse data setting (again, mimicking an aspect of the Warfarin data) and explore the use of shrinkage priors in a setting where many covariates are available, but most are not relevant for optimal treatment decisions. See Appendices B.3 and B.4 for details.

Performance is measured via the bias of estimators of the blip parameters relative to their true values, the standard deviation of the estimators, the difference between the value function (dVF) under the true optimal ITR and the value function under the estimated optimal ITR, and the standard deviation of dVF when the estimated treatment rule was applied to the same population. For the many covariates setting, these measures are assessed over: (1) a full set of 2000 iterations, and (2) a subset of iterations where the non-zero treatment-covariate interactions are correctly selected. Results are compared to those obtained from a one-stage analysis based on the full individual-level data.

4.3.2 Results

Here, we present the simulation results for the small sample size and binary treatment under different confounding scenarios, heterogeneity levels, and half-Cauchy (0,1) prior. For brevity, only results for ψ_0 and dVF are presented here; we also present the estimates of blip parameters in the sparse data setting with small sample size, different heterogeneity levels and half-Cauchy (0,1) prior, and results for ψ_1, ψ_2 and dVF in the many covariates setting; all other results are presented in the Supplementary Materials, Appendix B.5.

Estimates of ψ_0 are presented in Figure 4.1. Relative bias is typically less than 1%, and neither relative bias nor standard deviation vary much across different confounding scenarios. When the confounding effect is larger, the relative bias is slightly larger. Relative bias is similar for different heterogeneity levels, but the variability of the estimators increases with the heterogeneity level. The variation in $\hat{\psi}_0$ in the common rule setting is greater than that in the varying effects setting, a consequence of the data-generating mechanism was such that heterogeneity in ψ_{i0} is greater in the common rule setting than in the varying effects setting.

The dVF is shown in Figure 4.2. A smaller dVF corresponds to better ITR estimation. The values of the estimated optimal ITR are comparable across different confounding scenarios. When the confounding effect is larger, the dVF is slightly larger. In the common rule and common effect settings, the dVF is near 0 and varies little. The dVF is larger with increasing heterogeneity. This is unsurprising, since this is a scenario where a single ITR will not provide the optimal treatment for all individuals; rather, the truly optimal treatment is site-specific. However, implementing site-specific rules in a real-world setting is impractical.

Blip parameter estimates in the sparse data setting are shown in Figure 4.3. The estimators are unbiased and variability increases with heterogeneity. In the sparse data setting, the dVF (not shown) is zero for all but one or two simulation runs, and it does not change much with different heterogeneity levels.

Choice of prior made little difference to performance. The model results also do not differ much between the one- and two-stage approaches, but the variability in the two-stage approach is slightly larger than in the one-stage approach when the heterogeneity is small since the two-stage approach cannot share information across sites.

The relative bias and standard deviation of $\hat{\psi}_1$ and $\hat{\psi}_2$, the proportion of selection and the dVF between the true and estimated optimal ITR for the many covariates setting are reported in Table 4.1. It is more difficult to correctly detect the treatment-covariate interaction when the effect is small or a large number of noisy variables is present. When only 10 candidate tailoring variables are considered, non-zero ψ_1 is detected in the one- and two-stage models across 75.4 % and 71.8% simulation runs, respectively; these numbers drop to 49.2% and 46.5% when 20 covariates are considered. In all cases, the larger ψ_2 is correctly identified for all simulated datasets. Due to the poorer performance in detecting the small non-zero ψ_1 , the relative bias of $\hat{\psi}_1$ assessed over all simulation runs is large; among the simulation iterations when ψ_1 is correctly identified, the relative bias is 1 - 2%. Consistent observations can be found for dVF. The horseshoe prior and the selection criterion we adopt may not accurately detect smaller treatment-covariate interactions. However, if the non-zero effect has been identified, we achieve low bias for parameter estimation and small dVF, corresponding to good ITR estimation.

4.4 Estimating an optimal Warfarin dose strategy

4.4.1 Context and data source

Warfarin is a widely-used oral anticoagulant for thrombosis and thromboembolism treatment and prevention. Establishing an optimal Warfarin dose strategy is vital due to the narrow therapeutic window and the large interindividual variability in patients' response to the drug (International Warfarin Pharmacogenetics Consortium, 2009). The international normalized ratio (INR) is a measure of the time needed for the blood to clot. It should be closely moni-



Model 🖻 One-stage 🖻 Two-stage

Figure 4.1: Simulation results for the small sample size and the binary treatment setting. Performance of the methods is assessed over 2000 iterations. Estimates (posterior means) of ψ_0 are shown under different confounding scenarios, heterogeneity levels ($I^2 = 0.1, 0.2, 0.3$), and half-Cauchy (0,1) prior. The triangles represent the mean of the estimates in each case. The dashed line shows the true value of 2.5.

tored for the safety and effectiveness of Warfarin dosing. Many methods have been proposed for finding the optimal dose rule, and clinical factors, demographics, and genetic variability may play an essential role in interindividual variations in the required dose of Warfarin. The International Warfarin Pharmacogenetics Consortium(2009) compared several Warfarin dose algorithms and concluded that a pharmacogenetic algorithm in which both genetic and



Model 🖻 One-stage 🖻 Two-stage

Figure 4.2: Simulation results for the small sample size and the binary treatment setting. Performance of the methods is assessed over 2000 iterations. The difference in the value function (dVF) between the true and estimated optimal ITR is shown under different confounding scenarios, heterogeneity levels ($I^2 = 0.1, 0.2, 0.3$), and half-Cauchy (0,1) prior. The triangles represent the mean of the estimates in each case. A smaller dVF corresponds to a better optimal ITR estimation.

clinical variables are used to inform the appropriate Warfarin dose performs best.

We use the International Warfarin Pharmacogenetics Consortium data and the proposed approaches to data and model sparsity to estimate an optimal individualized Warfarin dose strategy without sharing individual-level data. Following the work of Schulz and Moodie (2021); Danieli and Moodie (2022), observations with missing values are removed, leading to a sample size of n = 2853 from 14 different sites. One of the 14 sites was removed as it only includes seven patients, too few to allow a center-specific analysis. Therefore, the



Model 🖻 One-stage 🖻 Two-stage

Figure 4.3: Simulation results for the small sample size and the sparse data setting. Performance of the methods is assessed over 2000 iterations. Estimates (posterior means) of ψ_0 , ψ_1 , ψ_2 , ψ_3 are shown under different heterogeneity levels ($I^2 = 0.1, 0.2, 0.3$) and half-Cauchy (0,1) prior. The triangles represent the mean of the estimates in each case. The dashed lines show the true values of ψ_0 , ψ_1 , ψ_2 , ψ_3 .

sample size is reduced to n = 2846. The final dataset used in the analyses includes several variables: patient age (binned into 9 groups), sex, race, weight and height centered by the site mean, an indicator for taking amiodarone (an important interacting drug of Warfarin), and VKORC1 and CYP2C9 genotypes, where the latter two variables are genes that may be associated with the interindividual variation in Warfarin dose requirement. Information on these variables is collectively denoted by the vector \boldsymbol{x} including the leading constant term of one. The stable Warfarin dose and the corresponding INR for each patient are also available in the data.

4.4.2 Analyses and results

In accordance with previous work (Schulz and Moodie, 2021; Danieli and Moodie, 2022), the outcome variable is defined as $Y = -\sqrt{|2.5 - \text{INR}|}$ such that larger values of Y are clinically preferable. When Y is closer to zero, the INR is closer to the midpoint of the therapeutic window. The observed distribution of the outcomes is roughly symmetric, with values varying from -1.38 to 0. We define a treatment-free function including the main effects of all available variables for the stage-one analysis of the site-specific data. The blip function is assumed to be quadratic in Warfarin dose, including the main effects of dose, squareddose, and their interactions with the available variables not including height, weight, and race $(\boldsymbol{x}^{(\psi_i^{(1)})} = \boldsymbol{x}^{(\psi_i^{(2)})} = \boldsymbol{x}/\{\text{weight, height, race}\})$. Therefore, the outcome model for site *i* is $E(Y_i \mid \boldsymbol{X} = \boldsymbol{x}, A = a) = \beta_i^{\top} \boldsymbol{x} + a \psi_i^{(1) \top} \boldsymbol{x}^{(\psi_i^{(1)})} + a^2 \psi_i^{(2) \top} \boldsymbol{x}^{(\psi_i^{(2)})}$, where *a* is the Warfarin dose, and β_i represents the main effects of the available variables on the outcome in site *i* through the treatment-free function. The site-specific blip parameter vectors $\boldsymbol{\psi}_i^{(1)} = (\psi_{i0}^{(1)}, \ldots, \psi_{i7}^{(1)})$ and $\boldsymbol{\psi}_i^{(2)} = (\psi_{i0}^{(2)}, \ldots, \psi_{i7}^{(2)})$ include main effects of dose and squared-dose (i.e., $\psi_{i1}^{(1)}$, and $\psi_{i0}^{(2)}$), and their interactions with all predictors contained in $\boldsymbol{x}^{(\psi_i^{(1)})} = \boldsymbol{x}^{(\psi_i^{(2)})}$ (i.e., $\psi_{i1}^{(1)}, \ldots, \psi_{i7}^{(1)}$

To assess the impact of attempting to preserve individual-level data, the proposed two-stage IPD meta-analysis is implemented. A frequentist linear regression is used as the stage-one model to obtain the site-specific blip parameter estimates $\hat{\psi}_i^{(1)}$ and $\hat{\psi}_i^{(2)}$. The site-specific blip parameter estimates were found to be small in magnitude (i.e., close to zero), raising the question of whether the available covariates have important tailoring effects on the optimal Warfarin dosing in the dataset. Additionally, in some sites there are not enough patients to estimate the site-specific amiodarane effect, VKORC1 or CYP2C9 genetic effects. Therefore, the Bayesian hierarchical model in the second stage has to be adapted for data sparsity to estimate the common blip parameters $\psi^{(1)} = (\psi_0^{(1)}, \ldots, \psi_7^{(1)})$ and $\psi^{(2)} = (\psi_0^{(2)}, \ldots, \psi_7^{(2)})$. Horseshoe priors (Carvalho et al., 2010) are assumed for all treatment-covariate interactions

to select variables that truly influence Warfarin dosing by shrinking small effects to zero. Normal priors with mean zero and variance 10,000 are used for the main effects of dose and squared-dose, as Warfarin dose does have effects on the outcome and thus there is no reason to shrink its effects towards zero. The details of the model is described in the Supplementary Materials, Appendix B.6.

The posterior distribution of the optimal dose for an individual with vector \boldsymbol{x} is then approximated by substitution of posterior samples of $\psi_t^{(u)}$, $t = 0, \ldots, 7$, u = 1, 2 in the maximizer to the common blip function, i.e., $a^{\text{opt}} = -(\boldsymbol{\psi}^{(1)T}\boldsymbol{x}\boldsymbol{\psi}_i^{(1)})/(2\boldsymbol{\psi}^{(2)T}\boldsymbol{x}\boldsymbol{\psi}_i^{(2)})$. No significant treatment-covariate interactions are selected. Therefore, the optimal dose is fully determined by $\psi_0^{(1)}$ and $\psi_0^{(2)}$, leading to the same recommended dose distribution to all patients and the optimal dose is 40.12 mg/week (posterior median). Not knowing the true optimal dose in this real-data analysis, we also compare the results of the two-stage approach that avoids disclosing site-specific individual-level data to a one-stage approach that requires all individual-level data to be used at once (and thus requires sharing of data outside the sites at which they were collected). A consistent conclusion is obtained in the one-stage approach that none of the covariates under consideration have significant tailoring effect on the optimal Warfarin dose, and the posterior median estimate of the common recommended dose is 38.80 mg/week.

4.5 Discussion

Estimation of optimal ITRs requires the identification of individuals who benefit most from particular treatments and thus the detection of treatment-covariate interactions. Due to the low power associated with detecting the interactions, large datasets may be required to develop an effective optimal ITR, which motivates collaboration across different sites. In multisite studies, the sharing of individual-level data is sometimes constrained, which poses statistical challenges for estimating ITR. In this paper, we adopt a Bayesian two-stage IPD meta-analysis approach to estimate the optimal ITR using multisite data without sharing individual-level data.

In the presence of treatment-covariate interactions, traditional meta-regression based on aggregate data are prone to ecological bias and may not reflect the individual-level interactions. However, a two-stage IPD meta-analysis avoids such bias by estimating treatment-covariate interactions within each site separately at the first stage and then synthesising these estimates at the second stage (Berlin et al., 2002; Simmonds and Higgins, 2007; Fisher et al., 2011). In our Bayesian two-stage IPD meta-analysis approach, we estimate site-specific optimal ITRs using the regression-based method of Q-learning implemented via linear regression at the first stage, however alternatives such as dWOLS or G-estimation could also have been employed. At the second stage, the site-specific blip parameter estimates are shared to the common analysis center where a Bayesian hierarchical model is used to combine the estimates. We also consider sparsity: estimation of identical models or all parameters of interest at all sites may be infeasible due to heterogeneity of patient populations across sites, and the number of non-zero parameters is often small. We address data sparsity by reparametrizing the likelihood in the second stage such that the site-specific estimates are linked to the correct set of parameters, and use a horseshoe prior and a credible interval selection criterion to select significant treatment-covariate interactions to account for model sparsity. Simulations demonstrate that our approach gives consistent estimation of the common blip parameters which fully characterize the optimal ITR. When the site-specific optimal ITRs are not very heterogeneous, the value function of the estimated optimal ITR is also close to that of the true optimal ITR.

We estimate an optimal Warfarin dosing strategy using data from the International Pharmacogenetics Warfarin Consortium. Shrinkage priors are used to select the covariates that truly have effect on the optimal Warfarin dosing. We compare the results obtained from the two-stage approach with the results obtained from a one-stage approach using the full individual-level data. We found that both approaches are consistent in the sense that none of the covariates are selected for the Warfarin dosing strategy, and both provide very close dose recommendations (40.12 vs. 38.80 mg/week). We also emphasize that although our results may provide some guidance for the establishment of the optimal Warfarin dosing strategy, it is unlikely to provide the whole picture. Several important predictors are not included in the dataset such as alcohol consumption or Vitamin K intake (International Warfarin Pharmacogenetics Consortium, 2009), and thus could not be considered in the analysis.

The simulation results and Warfarin analysis illustrate our approach's potential to avoid individual-level data sharing when using multisite data to estimate the optimal ITR. Our approach allows for heterogeneity across sites, which has not been explored in the previous work of Danieli and Moodie (2022); Moodie et al. (2022) but may be a more reasonable and realistic assumption. In addition, our approach is quite flexible in the sense that at the first stage we could use any regression-based method to estimate the site-specific blip parameters. as long as the unbiasedness, normality and consistency of the estimators are guaranteed. We focus on simple and interpretable regression, but this approach does require correct specification of the regression model. If the regression model is misspecified, the site-specific estimators obtained in the first stage could be biased for the true site-specific parameters, leading to biased common parameter estimators (and thus a suboptimal ITR) in the second stage. Semi-parametric alternatives that offer double robustness (Robins, 2004; Wallace and Moodie, 2015) could also be considered for stage-one models, allowing for the use of more flexible specifications of the covariate effects (i.e., the treatment-free model, which is essentially a nuisance model). Additionally, the use of a Bayesian approach in the second stage allows for the seamless incorporation of the accumulating or external information regarding the optimal ITR from various sources. We may also incorporate variable selection without much additional efforts by using shrinkage priors. This is useful, especially in observational studies where many variables that are irrelevant for the treatment decision may also be collected. In this paper, we use a horseshoe prior and select the treatment-covariate interactions if the 95% posterior credible intervals do not include zero. The horseshoe prior and credible interval selection criterion are standard in Bayesian variable selection. Simulation shows that our choice may not be powerful enough to detect very small effects. Other shrinkage priors and alternative selection criteria can also be considered (Bondell and Reich, 2012; Hahn and Carvalho, 2015; Li and Pati, 2017; van Erp et al., 2019).

The proposed approach can be easily adapted for data sparsity, as demonstrated both in simulation and with the Warfarin data. The data sparsity within individual sites considered in this paper only occurs in covariates. The sites in the Warfarin study, as well as in our simulated examples, recruited from different (but not entirely distinct) populations, but all treatment choices under investigation were available in all sites. Sparsity in covariates may restrict the generalization of the estimated site-specific ITRs to a broader population. For example, without additional assumptions, the estimated site-specific ITRs for sites with only White patients in the samples may not generalize well to the non-White population or a target population including non-White people. However, this is not, in general, of concern since our estimated of interest is a common rather than site-specific optimal ITR. However, a related concern is that the positivity (or overlap) assumption, often made in the context of causal inference, is violated. There was, in fact, a moderate lack of overlap in the Warfarin data (visual inspection of the overlap can be found in Figure B.11 in the Appendix B.6 of the Supplementary Materials). In the presence of the non-overlap, estimation of parameters of interest requires extrapolation which may introduce bias. However, the induced bias might be negligible, if the model is correctly specified and the relationship between treatment, covariates and outcome is consistent across the covariate space, which is possible in our case. A more concerning issue resulting from the lack of overlap is the increased variability of parameter estimators, especially when the estimated effects are all very small (as in this analysis) leading to a more variable and less reliable estimate of the optimal ITR.

Although this paper focuses on a continuous outcome, it could be easily extended to outcomes of other types by directly using the existing model of the particular outcome type as the stage-one model (Tchetgen Tchetgen et al., 2010; Linn et al., 2017; Kidwell et al., 2018; Simoneau et al., 2020). Besides the concerns regarding the individual-level data sharing, the two-stage approach is also computationally efficient compared with the one-stage approach using the full individual-level data (Burke et al., 2017). In our Warfarin analysis, consistent conclusions are obtained in both approaches. However, the one-stage approach requires evaluating the likelihood with a large amount of IPD and estimating a large number of parameters simultaneously in a single model, which increases the computational burden.

Our approach also has some limitations. First, each site is required to have sufficient statistical knowledge of ITR estimation, as site-specific optimal ITRs are estimated separately at each site. This might increase the funding burden for each site in practice. Second, our simulations show that in comparison with the estimated optimal ITR obtained under the scenario of homogeneous optimal ITRs across sites, the estimated optimal ITR obtained in the presence of heterogeneity in site-specific optimal ITRs performs worse in terms of the value function. A larger number of sites might be needed to account for the heterogeneity in the site-specific optimal ITRs. However, it may not be realistic to have a large number of sites. Heterogeneity could arise from the diversity in various aspects across sites such as patient characteristics, treatment delivery, measurements, and study designs (Higgins et al., 2019). When planning a multisite trial for optimal ITR estimation, one may consider standardizing the research protocol (e.g., eligibility criteria, sampling, and study designs) to reduce extraneous heterogeneity. Standardizing research staff training might also be necessary to ensure that each site follows the common protocol strictly and that the data are measured and collected uniformly (Weinberger et al., 2001; Noda et al., 2006).

Another approach to aggregate heterogeneous linear ITRs in the binary treatment setting across sites could be a maximin projection learning method proposed by Shi et al. (2018).

However, this method assumes a common main effect of treatment for which the estimation requires the pooling of all individual-level data (Shi et al., 2018) and thus cannot be applied in our setting, where we aim to avoid sharing any individual-level records. The extension of the maximin projection learning method to varying main treatment effects across sites and more generalized settings (e.g., continuous treatment) may be of interest to consider in future as an alternative approach to the privacy-preserving estimation of ITRs.

In addition to the pooling and other approaches to ITR estimation without sharing individual data noted in the introduction, a more general approach with stronger guarantees on privacy that has been pursued in other non-ITR contexts is the use of differentially private algorithms (Dwork, 2008). Standard linear regression is not differentially private, therefore, our two-stage approach may not be differentially private, although it does offer some degree since the treatment-free parameter estimates never need to be published or shared under our proposed approach. An interesting avenue of future work is to investigate whether the sharing of only the blip model estimates would pose a violation of differential privacy.

Acknowledgements

This work is supported by an award from the Canadian Institutes of Health Research CIHR FDN-16726. EEMM is a Canada Research Chair (Tier 1) in Statistical Methods for Precision Medicine and acknowledges the support of a chercheur de mérite career award from the Fonds de Recherche du Québec, Santé. SG acknowledges support from the Natural Sciences and Engineering Research Council of Canada (NSERC), Canadian Institute for Statistical Sciences (CANSSI) and Fonds de recherche du Québec - Sante (FRQS).

Supporting Information

Appendices B.1 - B.6, referenced in Sections 4.2 - 4.4, are available with this paper.

Number of covariates	Model	Set	ψ_1		ψ_2		dVF (SD)
			RB (SD)	% Selection	RB (SD)	% Selection	
	One-stage	Full	-25.703 (0.257)	75 400	-1.503(0.228)	100 000	$0.135\ (0.178)$
10		Selected	-1.463(0.174)		-1.503(0.228)		$0.037\ (0.046)$
10	Two-stage	Full	-28.748(0.265)	71 800	$-1.549\ (0.230)$	100 000	$0.153\ (0.188)$
		Selected	-0.764(0.177)		$-1.549\ (0.230)$		$0.039\ (0.048)$
	One-stage	Full	$-49.967\ (0.275)$	49.222	-1.261(0.237)	100 000	$0.241 \ (0.210)$
00		Selected	$1.647\ (0.171)$		-1.261(0.237)		$0.032\ (0.042)$
07	Two-stage	Full	-52.554 (0.273)	46 475	-1.332(0.242)	100 000	$0.257\ (0.211)$
		Selected	$2.089\ (0.171)$		-1.332(0.242)		$0.034\ (0.045)$

Chapter 5

Two-stage Bayesian network meta-analysis of individualized treatment rules for multiple treatments with siloed data

Preamble to Manuscript 3. In Manuscript 2, I considered a binary or continuous treatment setting and assume the same treatment sets across all sites. However, this assumption may not hold when there are numerous candidate treatments, but only a subset of them are available at individual sites due to time or funding constraints. In Manuscript 3, I extend the two-stage Bayesian meta-analysis approach presented in Manuscript 2 to accommodate multiple treatment settings with varying treatment sets across sites, drawing on techniques from network meta-analysis literature. A simulation study is conducted to investigate the finite sample behavior of the proposed method. Specifically, the feasibility of the common between-site heterogeneity assumption, which is recommended in network meta-analysis to reduce model complexity, is explored. The proposed method is used to establish an ITR for the treatment of depression, recommending from among six therapeutic options. This manuscript has been submitted for peer review.

Two-stage Bayesian network meta-analysis of individualized treatment rules for multiple treatments with siloed data

Junwei Shen¹, Erica E.M. Moodie¹, Shirin Golchi¹.

¹Department of Epidemiology, Biostatistics, and Occupational Health, McGill University

Abstract

Individualized treatment rules leverage patient-level information to tailor treatments for individuals. Estimating these rules, with the goal of optimizing expected patient outcomes, typically relies on individual-level data to identify the variability in treatment effects across patient subgroups defined by different covariate combinations. To increase the statistical power for detecting treatment-covariate interactions and the generalizability of the findings, data from multisite studies are often used. However, sharing sensitive patient-level health data is sometimes restricted. Additionally, due to funding or time constraints, only a subset of available treatments can be included at each site, but an individualized treatment rule considering all treatments is desired. In this work, we adopt a two-stage Bayesian network meta-analysis approach to estimate individualized treatment rules for multiple treatments using multisite data without disclosing individual-level data beyond the sites. Simulation results demonstrate that our approach can provide consistent estimates of the parameters which fully characterize the optimal individualized treatment rule. We illustrate the method's application through an analysis of data from the Sequenced Treatment Alternatives to Relieve Depression (STAR*D) study, the Establishing Moderators and Biosignatures of Antidepressant Response for Clinical Care (EMBARC) study, and the Research Evaluating the Value of Augmenting Medication with Psychotherapy (REVAMP) study.

5.1 Introduction

It is widely recognized in medical research that treatment responses often exhibit variability among different patient subgroups. Personalized medicine leverages this heterogeneity in treatment effect to enhance healthcare service quality by delivering tailored treatments to individual patients (Kravitz et al., 2004; Chakraborty and Moodie, 2013; Kosorok and Laber, 2019). An individualized treatment rule (ITR) is a decision rule that utilizes patientlevel information, such as demographics, genetic makeup, or disease history, to customize treatment plans at a single decision point. An optimal ITR guides treatment selection for individual patients with the goal of optimizing patient outcomes. Estimating the optimal ITR is essential for the practice of personalized medicine and has attracted significant research focus.

Regression-based approaches are commonly employed to indirectly estimate the optimal ITR. These methods model the expected patient outcome as a function of treatment, covariates, and treatment-covariate interactions. Then, the optimal treatment is determined as the one that leads to the best estimated outcome for any given covariate profile. Q-learning (Watkins, 1989; Sutton and Barto, 2018), G-estimation (Robins, 2004), and dynamic weighted ordinary least squares (dWOLS) (Wallace and Moodie, 2015) are three popular regression-based approaches. Estimating treatment-covariate interactions based on individual-level data is essential in these approaches. With advancements in technologies, the availability of large collections of health data from multiple data sites has facilitated the identification of factors that contribute to differential treatment responses, provided a higher statistical power of treatment-covariate interaction estimation that cannot be offered by a single site (Greenland, 1983), and improved the generalizability of the findings. However, patient-level health data are typically highly sensitive, and their disclosure could cause a violation to data sharing agreements or policies, presenting a challenge for ITR estimation with multisite data. Therefore, valid approaches to ITR estimation without releasing patient-level information

are desired.

Several approaches have been proposed to avoid individual-level data sharing for ITR estimation. Spicker et al. (2024) investigate differential privacy (Dwork, 2006) in the context of dynamic treatment regimes, which is an extension of ITRs to multiple treatment decision points. Instead of regression-based approaches, they focus on an outcome weighted learning method (Zhao et al., 2012), which frames the estimation of ITRs as a classification loss minimization problem and identifies the optimal treatment through support vector machine classifiers. Danieli and Moodie (2022) study the use of data pooling (Saha-Chaudhuri and Weinberg, 2017) and distributed regression (Rassen et al., 2013) to protect individual-level data from release in multisite studies in the context of ITR estimation with generalized dWOLS for continuous outcomes. In their approach, estimators characterizing the optimal ITR are computed using data summaries (e.g., pooled data or matrix products) shared by each single site, rather than individual-level data. Moodie et al. (2022) also explore distributed regression in dynamic weighted survival modelling, a generalization of dWOLS to survival outcomes (Schulz and Moodie, 2021). One limitation of both approaches is that they typically assume parameters of interest are fixed and common to all sites. To overcome this limitation, our recent work (Shen et al., 2024) adopted a two-stage Bayesian meta-analysis approach, which requires only site-specific analyses of individual-level data within each site and sharing site-specific estimates, as summary data, to construct a common optimal ITR in settings where all sites assigned the same treatment options. Conventional meta-analysis approaches typically assume that the treatment is binary and that each site consists of the same treatment comparison, limiting their applicability in a wide range of diseases where the treatment landscape can be quite heterogeneous, as is the case with conditions like depression. In such cases, due to funding or time constraints, only a subset of available treatment options can be delivered in each site, and yet establishing an optimal ITR that considers all treatments is often desired. Analogously, in our motivating example, we wish to draw inferences using randomized trial data from trials whose randomization groups are overlapping but not identical. In this manuscript, we consider ITR estimation in multisite studies without sharing individual-level data, when more than two treatments are available and each site may encompass different sets of treatment assignment options.

An extension of classic meta-analysis to multiple treatments is network meta-analysis (Salanti, 2012; Cipriani et al., 2013). Network meta-analysis compares multiple treatments within a network of studies, involving the simultaneous analysis of direct evidence obtained from head-to-head trials and indirect evidence from studies including the treatments of interest and one or more common comparator treatments, when comparing any two treatments in the network. This drives the extension of the two-stage Bayesian meta-analysis approach proposed in our previous work (Shen et al., 2024) to the current setting where treatments are not common across all sites or studies, which is the objective of this work.

The remainder of this paper is organized as follows: Section 5.2 describes the proposed method including the notations and assumptions. A simulation study is presented in Section 5.3 to explore the performance of ITR estimation using the proposed method. Section 5.4 demonstrates the application of the proposed method via an analysis of real data from three randomized clinical trials for the treatment of depression. The paper concludes with a discussion in Section 5.5.

5.2 Methods

5.2.1 Preliminaries

Consider the data (\mathbf{X}, A, Y) , where \mathbf{X} includes pre-treatment covariates, $A \in \mathcal{A} = \{d_1, \ldots, d_G\}$ represents the treatment received by individual patients with G unique options. Without loss of generality, we assume d_1 is the reference treatment, and $\tilde{\mathbf{A}} = (I(A = d_2), \ldots, I(A = d_G))^{\top}$ codes the treatment assignment in a vector of dummy variables. We denote Y to be the continuous outcome of interests, with larger values preferred. We use uppercase, lowercase, and bold letters to denote random variables, their observed values, and vectors, respectively.

We make the following assumptions: (i) the stable unit treatment value assumption (SUTVA): a patient's outcome is not influenced by other patients' treatment (Rubin, 1980); (ii) no unmeasured confounding (Robins, 1997); (iii) positivity: there is a positive probability of receiving every possible treatment for every combination of covariate values that occur among individuals in the population (Cole and Hernán, 2008).

Define a treatment-free function $f(\mathbf{x}) = E(Y \mid A = d_1, \mathbf{X} = \mathbf{x})$, which represents the expected outcome at the reference treatment d_1 for patients with covariates $\mathbf{X} = \mathbf{x}$. A blip function $\gamma(a, \mathbf{x})$ (Robins, 2004) is defined such that $\gamma(d_h, \mathbf{x}) = E(Y \mid A = d_h, \mathbf{X} = \mathbf{x}) - E(Y \mid A = d_1, \mathbf{X} = \mathbf{x})$ for $h \neq 1$, and $\gamma(d_1, \mathbf{x}) = 0$. Therefore, $\gamma(d_h, \mathbf{x})$ is the expected difference in the outcomes between receiving treatment d_h and the reference treatment d_1 for patients with covariates $\mathbf{X} = \mathbf{x}$. For example, it can be the main effect of d_h and interaction effects between d_h and covariates \mathbf{x} . With f and γ , the outcome can be decomposed:

$$E(Y \mid A = a, \boldsymbol{X} = \boldsymbol{x}) = f(\boldsymbol{x}) + \gamma(a, \boldsymbol{x}).$$

We aim to identify the optimal ITR, i.e., a decision rule that, given individual characteristics, outputs a tailored treatment which can maximize the expected outcome. The treatment-free function f is not related to any terms of treatments d_2, \ldots, d_G . Therefore, the optimal ITR $d^{\text{opt}}(\boldsymbol{x})$ only depends on γ , i.e., $d^{\text{opt}}(\boldsymbol{x}) = \arg \max_{a \in \{d_1, \ldots, d_G\}} \gamma(a, \boldsymbol{x})$. However, the estimation of the optimal ITR requires model specifications for both f and γ . For example, we can posit functional forms: $f(\boldsymbol{x}) = \boldsymbol{\beta} w(\boldsymbol{x})$ and $\gamma(a, \boldsymbol{x}) = z(\tilde{\boldsymbol{a}}) \boldsymbol{\psi} l(\boldsymbol{x})$, where w, z, and l are multivariate functions specified by analysts, with $z(\tilde{\boldsymbol{a}}) = 0$ for $a = d_1$ to ensure the condition $\gamma(d_1, \boldsymbol{x}) = 0$ is met for every possible \boldsymbol{x} . The dimensions of w, z, l and parameters $\boldsymbol{\beta}$ and $\boldsymbol{\psi}$ should be compatible to guarantee both f and γ output scalar values. In this manuscript, for illustration purposes, we assume $w(\boldsymbol{x}) = \boldsymbol{x}^{(\boldsymbol{\beta})}, z(\tilde{\boldsymbol{a}}) = \tilde{\boldsymbol{a}}$ and $l(\boldsymbol{x}) = \boldsymbol{x}^{(\boldsymbol{\psi})}$. We use $\boldsymbol{x}^{(\boldsymbol{\beta})}$ and $\boldsymbol{x}^{(\boldsymbol{\psi})}$ to indicate that not all collected variables in \boldsymbol{x} , but those related to patient outcomes or treatment selection are included in f and γ . A total of p covariates that contribute to the outcome (predictive variables) are included in $\boldsymbol{x}^{(\beta)}$, among which q have tailoring effects on treatment assignment (prescriptive variables) and are included in $\boldsymbol{x}^{(\psi)}$. Both $\boldsymbol{x}^{(\beta)}$ and $\boldsymbol{x}^{(\psi)}$ are augmented with an intercept term and are subvectors of \boldsymbol{x} , and $\boldsymbol{x}^{(\psi)}$ is also contained in $\boldsymbol{x}^{(\beta)}$, i.e., $\boldsymbol{x}^{(\psi)} = (1, x_1, \dots, x_q)^{\top}$ and $\boldsymbol{x}^{(\beta)} = (1, x_1, \dots, x_q, x_{q+1}, \dots, x_p)^{\top}$ with $q \leq p$. Alternative parametric choices such as nonlinear models can also be considered for w and l. Given a linear specification of w, z and l, for example, the outcome model becomes

$$E(Y \mid A = a, \mathbf{X} = \mathbf{x}) = \underbrace{\boldsymbol{\beta}^{\top} \mathbf{x}^{(\boldsymbol{\beta})}}_{\text{treatment-free function}} + \underbrace{\tilde{\boldsymbol{a}}^{\top} \boldsymbol{\psi} \mathbf{x}^{(\boldsymbol{\psi})}}_{\text{blip function}}, \tag{5.1}$$

where the treatment-free function f and blip function γ are parameterized by a (p + 1)dimensional vector $\boldsymbol{\beta}$ and a $(G - 1) \times (q + 1)$ matrix $\boldsymbol{\psi}$, i.e.,

$$\boldsymbol{\psi} = \begin{pmatrix} \psi_{20} & \psi_{21} & \cdots & \psi_{2q} \\ \psi_{30} & \psi_{31} & \cdots & \psi_{3q} \\ \vdots & \vdots & \ddots & \vdots \\ \psi_{G0} & \psi_{G1} & \cdots & \psi_{Gq} \end{pmatrix},$$

respectively, and ψ_{ht} is the main effect of treatment d_h (t = 0) or the interaction effect between x_t and $I(a = d_h)$ (t = 1, ..., q). We use $\psi_{\cdot t} = (\psi_{2t}, ..., \psi_{Gt})^{\top}$ and $\psi_{h\cdot} = (\psi_{h0}, ..., \psi_{hq})^{\top}$ to represent the (t + 1)th column and the *h*th row of ψ , containing all blip parameters related to a given covariate x_t and a given treatment d_h , respectively. Then, the outcome model can also be written as

$$E(Y \mid A = a, \mathbf{X} = \mathbf{x}) = \boldsymbol{\beta}^{\mathsf{T}} \mathbf{x}^{(\boldsymbol{\beta})} + \sum_{t=0}^{q} \boldsymbol{\psi}_{\cdot t}^{\mathsf{T}} \tilde{\mathbf{a}} x_{t},$$

$$= \boldsymbol{\beta}^{\mathsf{T}} \mathbf{x}^{(\boldsymbol{\beta})} + \sum_{h=2}^{G} I(a = d_{h}) \boldsymbol{\psi}_{h\cdot}^{\mathsf{T}} \mathbf{x}^{(\boldsymbol{\psi})}.$$
 (5.2)

Since the treatment assignment only influences the outcome through the blip model γ , the parameter $\boldsymbol{\psi}$ will solely determine the optimal ITR. Given the parameter $\boldsymbol{\psi}$ in the model (5.2), the optimal ITR $d^{\text{opt}}(\boldsymbol{x}) = \arg \max_{a \in \{d_1, \dots, d_G\}} \gamma(a, \boldsymbol{x})$ can be written as

$$d^{\text{opt}}(\boldsymbol{x}) = \begin{cases} d_{h_0}, & h_0 = \arg \max_{h \in \{2, \dots, T\}} \boldsymbol{\psi}_{\boldsymbol{h}}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\psi})} & \text{and} & \boldsymbol{\psi}_{\boldsymbol{h}_0}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\psi})} > 0 \\ \\ d_1, & \boldsymbol{\psi}_{\boldsymbol{h}}^{\mathsf{T}} \boldsymbol{x}^{(\boldsymbol{\psi})} \le 0 & \forall h = 2, \dots, G \end{cases}$$

The parameter ψ can be estimated by different approaches. In the absence of model misspecification, consistent and unbiased estimators of ψ can be obtained from a Q-learning method (Chakraborty and Moodie, 2013), which, in our setting, reduces to a standard linear regression for the model (5.1). Doubly robust alternatives such as dWOLS or G-estimation can also be employed. Bayesian approaches have also been proposed for ITR estimation, such as Bayesian G-computation (Arjas and Saarela, 2010), a Bayesian machine learning approach to Q-learning (Murray et al., 2018), Bayesian additive regression trees (Logan et al., 2019), and Bayesian causal forest (Hahn et al., 2020).

5.2.2 Two-stage Bayesian network meta-analysis

In this section, we describe the use of a two-stage Bayesian network meta-analysis approach to avoid disclosing individual-level data, when estimating the optimal ITR for multiple treatments using multisite data. We first describe the model when all G treatments are present in all sites and then explain the extension to varying sets of treatments across sites. Suppose we have K sites. For site $i \in \{1, \ldots, K\}$, the outcome model can be expressed as

$$E(Y_{ij} \mid A = a_{ij}, \boldsymbol{X} = \boldsymbol{x}_{ij}) = \boldsymbol{\beta}_{i}^{\top} \boldsymbol{x}_{ij}^{(\boldsymbol{\beta})} + \sum_{t=0}^{q} \boldsymbol{\psi}_{i\cdot t}^{\top} \tilde{\boldsymbol{a}}_{ij} \boldsymbol{x}_{ijt},$$

where $j \in \{1, ..., n_i\}$ indexes individual patients within each site and n_i is the number of patients in site *i*. We include index *i* in β_i and $\psi_{i\cdot t}$ to indicate that these parameters are site-specific and can vary across sites. The varying site-specific blip parameters $\psi_{i\cdot t}$ are

assumed to come from a common distribution, i.e.,

$$\boldsymbol{\psi}_{i\cdot t} = (\psi_{i2t}, \dots, \psi_{iGt})^{\top} \sim \text{MVN}(\boldsymbol{\psi}_{\cdot t}, \boldsymbol{\Sigma}_{t}), \qquad (5.3)$$

where $\text{MVN}(\boldsymbol{\psi}_{\cdot t}, \boldsymbol{\Sigma}_{t})$ represents a multivariate normal distribution with mean $\boldsymbol{\psi}_{\cdot t}$ and variancecovariance matrix $\boldsymbol{\Sigma}_{t}$. The common parameters $\boldsymbol{\psi}_{\cdot t}$, $t = 0, \ldots, q$, fully characterize a common optimal ITR applicable to a broader population, encompassing subpopulations from various sites in the dataset, and potentially for future patients at comparable sites. Estimation of ψ_{ht} , $h = 2, \ldots, G$, $t = 0, \ldots, q$, is of primary interest.

In the two-stage Bayesian meta-analysis approach, we first obtain a set of estimates for the site-specific parameters by conducting analyses on data from each single site, and then combine these estimates in a Bayesian hierarchical model to obtain estimates of common parameters $\psi_{\cdot t}$ and thus a common optimal ITR. That is, in the first stage, estimates for $\psi_{i\cdot t}$, i.e., $\hat{\psi}_{i\cdot t}$, and the corresponding $(G-1) \times (G-1)$ variance-covariance matrix $\hat{\Sigma}(\hat{\psi}_{i\cdot t})$ can be acquired from approaches mentioned in Section 5.2.1, based on solely site-specific data. In the second stage, these site-level estimates, rather than individual records, are shared to a central analysis site and combined in a Bayesian hierarchical model:

$$\hat{\psi}_{i\cdot t} \sim \text{MVN}(\psi_{i\cdot t}, \widehat{\Sigma}(\hat{\psi}_{i\cdot t})),$$

$$\psi_{i\cdot t} \sim \text{MVN}(\psi_{\cdot t}, \Sigma_{t}),$$

$$\psi_{ht} \sim p_{\psi_{ht}}(\psi_{ht}),$$

$$\Sigma_{t} \sim p_{\Sigma_{t}}(\Sigma_{t}).$$
(5.4)

Here, prior distributions $p_{\psi_{ht}}$ and $p_{\Sigma_t}(\Sigma_t)$ can be assigned for the unknown parameters ψ_{ht} and Σ_t . A popular prior choice for ψ_{ht} could be a normal prior with large variance (Gelman et al., 2013). The between-site heterogeneity matrix Σ_t could be structured under the assumption that the between-site heterogeneity is the same across different treatment comparisons (White et al., 2012; Riley and Fisher, 2021). In this case, a common specification in the network meta-analysis literature (White et al., 2012; Riley and Fisher, 2021) is that Σ_t has diagonal elements σ_t^2 and off-diagonal elements $0.5\sigma_t^2$ (i.e., the correlation between any two treatment contrasts is 0.5), where σ_t^2 is the between-site variance associated with ψ_{ht} for all $h = 2, \ldots, G$. By assuming the same variance for all ψ_{ht} with a given t and fixing the correlation between ψ_{h_1t} and ψ_{h_2t} , $h_1 \neq h_2$, $h_1, h_2 = 2, \ldots, G$, at 0.5, we have that the variance of the contrast $\psi_{h_1t} - \psi_{h_2t}$ is given by : $\operatorname{var}(\psi_{h_1t} - \psi_{h_2t}) = \operatorname{var}(\psi_{h_1t}) + \operatorname{var}(\psi_{h_2t}) - 2\operatorname{cov}(\psi_{h_1t}, \psi_{h_2t}) = \sigma_t^2$. This is referred to as the *common between-site heterogeneity* assumption, and the contrast $\psi_{h_1t} - \psi_{h_2t}$ will be needed in a consistency equation (5.6) described later. This specification also reduces the number of parameters to be estimated in Σ_t , and can also improve model convergence. With the structured heterogeneity, a prior is needed only for σ_t or σ_t^2 . In this manuscript, we use a half-Cauchy prior for σ_t , however alternatives may also be employed (Gelman, 2006). If the variance-covariance matrix Σ_t is deemed to be unstructured, i.e., a separate between-site heterogeneity is to be estimated for each different treatment comparison, we can have the decomposition $\Sigma_t = UVU$, where U is a diagonal matrix of between-site standard deviations and V is an unknown correlation matrix. Then, priors can be assigned to those between-site standard deviations and also correlation matrix V. We still use a half-Cauchy prior for standard deviation parameters and a Lewandowski-Kurowicka-Joe (LKJ) prior for correlation matrix V (Lewandowski et al., 2009; Stan Development Team, 2021). We consider both heterogeneity structures in the simulation studies. The Bayesian hierarchical model is implemented in RStan (Stan Development Team, 2020, 2021).

Model (5.4) requires all G candidate treatments under consideration to be observed at all sites. However, in reality, some treatments are not administered in specific sites possibly due to the insufficient sample size and funding to implement a large number of treatments. When the set of treatments differs across sites, we can still implement a two-stage approach. However, in this setting, not all ψ_{iht} in the site-specific outcome models are estimable and modifications based on the network meta-analysis approach are made to (5.4). To proceed, the treatment set in site *i* is denoted by $\mathcal{A}_i = \{d_{a_i^{(1)}}, \ldots, d_{a_i^{(\nu_i)}}\}$, where ν_i is the number of treatments in site *i*, $\nu_i < G$ and $1 \leq a_i^{(1)} < a_i^{(2)} < \cdots < a_i^{(\nu_i)} \leq G$. Without loss of generality, we assume $d_{a_i^{(1)}}$ is the reference treatment for site *i*. When d_1 is available in site *i*, we have $d_{a_i^{(1)}} = d_1$ (i.e., the site-specific reference treatment is the common reference treatment). Otherwise, $d_{a_i^{(1)}} \neq d_1$. Then, with treatment set \mathcal{A}_i , we can fit a site-specific outcome model

$$E(Y_{ij} \mid A = a_{ij}, \boldsymbol{X} = \boldsymbol{x}_{ij}) = \boldsymbol{\beta}_{i}^{\top} \boldsymbol{x}_{ij}^{(\boldsymbol{\beta})} + \sum_{t=0}^{q} \tilde{\boldsymbol{\psi}}_{i\cdot t}^{\top} \tilde{\boldsymbol{a}}_{ij}^{(2)} \boldsymbol{x}_{ijt},$$

where $\tilde{a}_{ij}^{(2)} = (I(a = d_{a_i^{(1)}}), \dots, I(a = d_{a_i^{(\nu_i)}}))^{\top}$ is a subvector of \tilde{a}_{ij} , and the vector $\tilde{\psi}_{i,i} = (\tilde{\psi}_{i,a_i^{(2)}a_i^{(1)},t}, \dots, \tilde{\psi}_{i,a_i^{(\nu_i)}a_i^{(1)},t})^{\top}$ includes the estimable site-specific blip parameters, such that $\tilde{\psi}_{i,a_i^{(\tilde{h})}a_i^{(1)},t}$, $\tilde{h} = 2, \dots, \nu_i$, are the main effect (t = 0) or the treatment-covariate interaction $(t = 1, \dots, q)$ of treatment $d_{a_i^{(\tilde{h})}}$ relative to the site-specific reference treatment $d_{a_i^{(1)}}$. The estimates and common means of $\tilde{\psi}_{i,a_i^{(\tilde{h})}a_i^{(1)},t}$ are denoted by $\hat{\tilde{\psi}}_{i,a_i^{(\tilde{h})}a_i^{(1)},t}$ and $\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)},t}$, respectively. However, the parameters of interest are still ψ_{ht} , $h = 2, \dots, G$, $t = 0, \dots, q$, which characterize a common optimal ITR.

With the site-specific estimates $\hat{\tilde{\psi}}_{i\cdot t}$ and the associated $(\nu_i - 1) \times (\nu_i - 1)$ variance-covariance matrix $\hat{\Sigma}(\hat{\tilde{\psi}}_{i\cdot t})$, the first two levels of model (5.4) will be modified to

$$\begin{split} &\tilde{\boldsymbol{\psi}}_{\boldsymbol{i}\cdot\boldsymbol{t}} \sim \text{MVN}(\tilde{\boldsymbol{\psi}}_{\boldsymbol{i}\cdot\boldsymbol{t}}, \widehat{\boldsymbol{\Sigma}}(\hat{\boldsymbol{\psi}}_{\boldsymbol{i}\cdot\boldsymbol{t}})), \\ &\tilde{\boldsymbol{\psi}}_{\boldsymbol{i}\cdot\boldsymbol{t}} \sim \text{MVN}(\tilde{\boldsymbol{\psi}}_{\boldsymbol{i}\cdot\boldsymbol{t}}^{(2)}, \widetilde{\boldsymbol{\Sigma}}_{\boldsymbol{i}t}), \end{split}$$
(5.5)

where $\tilde{\psi}_{i\cdot t}^{(2)} = (\tilde{\psi}_{a_i^{(2)}a_i^{(1)},t}, \dots, \tilde{\psi}_{a_i^{(\nu_i)}a_i^{(1)},t})^{\top}$ is a vector of length $\nu_i - 1$, and $\widetilde{\Sigma}_t$ is a $(\nu_i - 1) \times (\nu_i - 1)$ variance-covariance matrix reflecting between-site heterogeneity. We include index i in both $\tilde{\psi}_{i\cdot t}^{(2)}$ and $\widetilde{\Sigma}_{it}$ to indicate their dependence on the treatment set \mathcal{A}_i , and the vector $\tilde{\psi}_{i\cdot t}^{(2)}$ includes common rather than site-specific blip parameters.

When $d_{a_i^{(1)}} = d_1$, we have $\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)},t} = \psi_{a_i^{(\tilde{h})}t}$, and the vector $\tilde{\psi}_{i\cdot t}^{(2)}$ is a subvector of $\psi_{\cdot t}$:

 $\tilde{\psi}_{i\cdot t}^{(2)} = (\tilde{\psi}_{a_i^{(2)}a_i^{(1)},t}, \dots, \tilde{\psi}_{a_i^{(\nu_i)}a_i^{(1)},t})^{\top} = (\psi_{a_i^{(2)}t}, \dots, \psi_{a_i^{(\nu_i)}t})^{\top}$. Therefore, ψ_{ht} can still be estimated through borrowing information across sites but with lower precision, as not all $\hat{\psi}_{i,t}$ include information relevant to estimating ψ_{ht} . When $d_{a_i^{(1)}} \neq d_1$, all ψ_{iht} are not estimable and $\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)},t} \neq \psi_{a_i^{(\tilde{h})}t}$. We make the consistency assumption as common in the network meta-analysis literature. In network meta-analysis, two treatments d_1 and d_2 can be either (1) directly compared in head-to-head studies, referred to as direct evidence, or (2) indirectly compared via studies comparing d_1 or d_2 with one or more common comparator treatments (i.e., indirect evidence). Then, the consistency assumption states that the indirect and direct estimates are in agreement (Salanti, 2012). In our setting, the consistency assumption ensures that we can link $\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)},t}$, $\tilde{h} = 2, \dots, \nu_i$, to ψ_{ht} , $h = 2, \dots, G$, through the equation:

$$\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)},t} = \psi_{a_i^{(\tilde{h})}t} - \psi_{a_i^{(1)}t}.$$
(5.6)

Then, priors can be assigned to all between-site variance and common mean parameters as in model (5.4).

5.3 Simulation studies

The simulation study is reported following the ADEMP (aims, data-generating mechanisms, estimands, methods, and performance measures) scheme proposed in the work of Morris et al. (2019).

5.3.1 Aims

We aim to evaluate ITR estimation for a continuous outcome and multiple treatments when individual-level data are protected from disclosure via a two-stage Bayesian network metaanalysis approach, under assumptions regarding (1) network sizes, (2) network shapes, (3) the true between-site heterogeneity, and (4) the assumed between-site heterogeneity in the Bayesian hierarchical model. Points (1) - (3) concern the data-generating mechanisms, while (4) concerns the analysis model. We have explored ITR estimation with the two-stage Bayesian pairwise meta-analysis under different confounding scenarios, between-site heterogeneity levels, prior choices and sample sizes in the previous work (Shen et al., 2024); note that in this previous work, we did not consider settings where the assigned treatments varied across site. While that previous work did not include settings where treatments offered varied by site, we expect similar results can be obtained when we have a network of studies and thus do not consider those particular features in the simulations here.

5.3.2 Data-generating mechanisms

The network structures considered in the simulation are shown in Figure 5.1. Sites included in each network structure with their site-specific treatment set \mathcal{A}_i are summarized in Table 5.1. A network can comprise sites with different treatment arms. Both networks (a) and (b) depicted in Figure 5.1 include three treatments d_1 , d_2 , and d_3 . However, in network (a), each site only includes two treatments: either comparing d_1 and d_2 or comparing d_1 and d_3 , while in network (b), we also have a third site only including d_2 and d_3 , forming a loop. Similarly, for networks (c) and (d), a larger treatment set is considered and the networks may or may not include loops. Network (e) reflects the network structure for the real data application to three trials considering six treatments of depression described in Section 5.4. With a given site-specific treatment set \mathcal{A}_i , the number of sites could be 1 or 3. That is, for any particular pair of site-specific treatments, there is either 1 or 3 sites that considered that set of treatment options. For each site, the sample size is fixed at 300.


Figure 5.1: Graphics of simulated networks (a) – (e). Network (e) reflects the network structure for the real data application to three trials described in Section 5.4. Connecting lines indicate the two treatments can be directly compared. The treatment d_1 is considered as the common reference treatment in each network.

For site *i*, we first generate a random number s_i uniformly from $\{0, 1, 2\}$. Then, a continuous covariate X_1 and a binary covariate X_2 are generated from the following distributions:

$$X_{1} \sim \begin{cases} N(5,1) & s_{i} = 0, \\ 6\text{Beta}(4,4) + 2 & s_{i} = 1, \\ U[2,8] & s_{i} = 2, \end{cases} \quad \text{Bernoulli}(0.5) \quad s_{i} = 0, \\ \text{Bernoulli}(0.3) & s_{i} = 1, \\ \text{Bernoulli}(0.7) & s_{i} = 2, \end{cases}$$

Given the site-specific treatment set $\mathcal{A}_i = \{d_{a_i^{(1)}}, \ldots, d_{a_i^{(\nu_i)}}\}$, the treatment assignment A follows a multinomial distribution with probabilities determined by

$$P(A = d_h \mid x_1, x_2, \mathcal{A}_i) = \frac{\exp(\alpha_{0,h} + \alpha_{1,h}x_1 + \alpha_{2,h}x_2)}{\sum_{d_h \in \mathcal{A}_i} \exp(\alpha_{0,h} + \alpha_{1,h}x_1 + \alpha_{2,h}x_2)},$$

where the coefficients $\alpha_{0,h}$, $\alpha_{1,h}$, and $\alpha_{2,h}$ are shown in Table 5.2. That is, although our realdata analysis focuses on randomized trial data, we perform our simulations under a more general setting where treatment allocation may depend on covariates.

Network	Treatment set \mathcal{A}	Number of arms per site	Site-specific treatment set \mathcal{A}_i
(a)	$\{d_1,d_2,d_3\}$	2	$\{d_1, d_2\}, \{d_1, d_3\}$
(b)	$\{d_1,d_2,d_3\}$	2	$\{d_1, d_2\}, \{d_1, d_3\}, \{d_2, d_3\}$
(c)	$\{d_1, d_2, d_3, d_4, d_5\}$	2 3	$ \{ d_1, d_2 \}, \{ d_1, d_3 \}, \{ d_1, d_4 \}, \{ d_1, d_5 \} \{ d_2, d_3, d_5 \}, \{ d_3, d_4, d_5 \} $
(d)	$\{d_1, d_2, d_3, d_4, d_5, d_6, d_7\}$	2	$ \{ d_1, d_2 \}, \{ d_1, d_3 \}, \{ d_1, d_4 \}, \{ d_1, d_5 \}, \\ \{ d_3, d_7 \}, \{ d_5, d_6 \} $
(e)	$\{d_1, d_2, d_3, d_4, d_5, d_6\}$	2 5 4	$\{d_1,d_2\}\ \{d_1,d_2,d_3,d_4,d_5\}\ \{d_1,d_2,d_3,d_4,d_5\}$

Table 5.1: Sites included in each network.

Treatment	$\alpha_{0,h}$	$\alpha_{1,h}$	$\alpha_{2,h}$
d_1	0	0	0
d_2	0	0.03	0.08
d_3	0.01	0.09	0.03
d_4	0.05	0.02	0.09
d_5	0.08	0.02	0.04
d_6	0.01	0.01	0.08
d_7	0.08	0.09	0.03

Table 5.2: Coefficients in multinomial probabilities for treatment assignment.

Suppressing the individual-specific subscript, the continuous outcome for an individual at site i is generated by

$$Y_i = \beta_{i0} + \beta_{i1}x_1 + \beta_{i2}x_2 + (\boldsymbol{\psi}_{i\cdot \mathbf{0}}^\top + \boldsymbol{\psi}_{i\cdot \mathbf{1}}^\top x_1)\tilde{\boldsymbol{a}} + \epsilon,$$

where $\tilde{\boldsymbol{a}} = (I(a = d_2), \dots, I(a = d_7))^{\top}$, $\boldsymbol{\psi}_{i\cdot t} = (\psi_{i2t}, \dots, \psi_{i7t})^{\top}$, $t = 0, 1, \beta_{i0} + \beta_{i1}x_1 + \beta_{i2}x_2$ is the site-specific treatment-free function, and $(\boldsymbol{\psi}_{i\cdot 0}^{\top} + \boldsymbol{\psi}_{i\cdot 1}^{\top}x_1)\tilde{\boldsymbol{a}}$ is the site-specific blip function. The random error ϵ follows a normal distribution with mean zero and residual variance $\sigma_{\epsilon}^2 = 0.25$. We note that in the above outcome generation model, seven treatments are assumed, whereas in all scenarios shown in Table 5.1, all sites employ fewer than seven treatments. This common form of data generation is appropriate. When a treatment d_h is not present in a given site *i*, the blip parameters related to d_h , i.e., ψ_{ih0} and ψ_{ih1} , will not contribute to the outcome as $I(a = d_h) = 0$. This outcome generation model is different from the outcome model we fit in the first stage:

$$E(Y_i \mid a, x_1 x_2) = \beta_{i0} + \beta_{i1} x_1 + \beta_{i2} x_2 + (\tilde{\boldsymbol{\psi}}_{i \cdot \mathbf{0}}^\top + \tilde{\boldsymbol{\psi}}_{i \cdot \mathbf{1}}^\top x_1) \tilde{\boldsymbol{a}}^{(2)},$$

where a general definition of \tilde{a} and $\tilde{\psi}_{i,t}$ have been provided in Section 5.2.1, and the notations should be accordingly adapted here.

The site-specific parameters $\boldsymbol{\theta}_{i} = (\beta_{i0}, \beta_{i1}, \beta_{i2}, \boldsymbol{\psi}_{i\cdot 0}, \boldsymbol{\psi}_{i\cdot 1})$ in the outcome generation model are simulated by: $\beta_{is} \sim N(\beta_{s}, \sigma_{B}^{2})$, s = 0, 1, 2, and $\boldsymbol{\psi}_{i\cdot t} \sim \text{MVN}(\boldsymbol{\psi}_{\cdot t}, \boldsymbol{\Sigma}_{t})$, t = 0, 1. For the 6×6 variance-covariance matrix $\boldsymbol{\Sigma}_{t}$, we consider two scenarios:

- common between-site heterogeneity: Σ_t has diagonals σ_B^2 and off-diagonals $0.5\sigma_B^2$;
- varying between-site heterogeneity: Σ_t has diagonals $(0.7, 1, 1.3, 0.7, 1, 1.3)\sigma_B^2$ and offdiagonals $0.5\sigma_B^2$,

where the between-study variance σ_B^2 is derived from heterogeneity level $I^2 = \sigma_B^2/(\sigma_B^2 + \sigma_\epsilon^2) = 0.1$. We note that for each distinct t, the variance-covariance matrix Σ_t can be different. Fundamentally, under the common between-site heterogeneity mechanism, for a given t, the between-site variance for ψ_{ht} is the same regardless of h (i.e., the diagonals of Σ_t are equal with a given t), but the between-site variance in different Σ_t can be different. However, the between-site variance is not our primary interest. Therefore, for simplicity, we assume a single between-site variance parameter σ_B^2 in all Σ_t , as well as for β_s . The common treatment-free function parameters are $\beta_0 = 4$, $\beta_1 = 1$, $\beta_2 = 1$, and the common blip function parameters are $\psi_{\cdot 0} = (5, 8, 4, 6, 2, 3)$, and $\psi_{\cdot 1} = (-0.9, -1.6, -1.3, -1.5, -0.8, -1.1)$. Let $\omega_{d_h}(\mathbf{x}) = \psi_{h0} + \psi_{h1}x_1$, $d_h \in \mathcal{A}/\{d_1\}$ and $\omega_{d_1}(\mathbf{x}) = 0$. The common optimal ITR is given by $d^{\text{opt}}(\mathbf{x}) = \arg \max_{d_h \in \mathcal{A}} \omega_{d_h}(\mathbf{x})$, which, in all networks, can be reduced to

$$d^{\text{opt}}(\boldsymbol{x}) = \begin{cases} d_1 & x_1 > \frac{50}{9}, \\ d_2 & \frac{30}{7} < x_1 < \frac{50}{9}, \\ d_3 & x_1 < \frac{30}{7}. \end{cases}$$

5.3.3 Estimands, methods, and performance metrics

The estimands of interest are the common blip function parameters ψ_{ht} , $d_h \in \mathcal{A}$, t = 0, 1, which fully characterize the optimal ITR in each network. We implement a two-stage Bayesian network meta-analysis approach, using linear regression in the first stage and a Bayesian hierarchical model for the second stage. For the mean parameters, we use a normal prior with mean 0 and variance 10,000. Regarding variance-covariance matrix $\tilde{\Sigma}_{it}$ in (5.5), we consider two scenarios:

- 1. When only a single site exists for each different site-specific treatment set \mathcal{A}_i in the network, we lack sufficient data to estimate the between-site heterogeneity, and thus $\widetilde{\Sigma}_{it} = 0.$
- 2. When we have three sites for each unique sites-specific treatment set, priors will be assigned under different modelling assumptions:
 - Under common between-site heterogeneity assumption, $\widetilde{\Sigma}_{it}$ has diagonal entries σ_t^2 and off-diagonal entries $0.5\sigma_t^2$, and a half-Cauchy (0,1) prior is assigned to σ_t .
 - For unstructured $\widetilde{\Sigma}_t$ under varying between-site heterogeneity assumption, we have decomposition $\widetilde{\Sigma}_t = UVU$, where U is a diagonal matrix with diagonals $\sigma_{a_i^{(\tilde{h})}, a_i^{(1)}, t}$ and V is an unknown correlation matrix. Then, a half-Cauchy (0,1) prior and a LKJ (1) prior are assumed for $\sigma_{a_i^{(\tilde{h})}, a_i^{(1)}, t}$ in U and V, respectively.

We are unaware of any existing methods that can estimate common ITRs across multiple studies with differing treatment values at different sites without individual-level data being shared. Thus, our analyses consist of comparing different model specifications within our Bayesian network meta-analysis approach. As no alternatives were found in the literature, no competing methods were included. We assess: (i) the relative bias of blip parameter estimators, which can be calculated by the difference between the mean of the estimates and the true value, divided by the latter, (ii) the standard deviation of the estimates, (iii) the difference in the value function (dVF) under the true and estimated optimal ITR, where the value function with respect to an ITR is approximated by the expected outcome if all patients in a new cohort of size 100,000 were treated according to the ITR, and (iv) the empirical standard deviation of the dVF when the estimated treatment rule was applied to the same population.

5.3.4 Results

In this section, for the sake of space, only simulation results for $\hat{\psi}_{20}$, $\hat{\psi}_{21}$, and dVF are presented. All results related to other blip parameters are presented in the Supplementary Materials. Tables 5.3 and 5.4 show simulation results when the true site-specific blip parameters are generated under varying and common between-site heterogeneity assumptions, respectively. Across all scenarios, the relative bias remains below 1%. When multiple sites exist for each \mathcal{A}_i , in the Bayesian hierarchical model, we can assume either common betweenstudy heterogeneity or an unstructured form for $\tilde{\Sigma}_t$. Irrespective of the true generation mechanisms of site-specific blip parameters, the relative bias of blip parameter estimators is similar in the two specifications. However, the estimators have greater variability when varying between-study heterogeneity is assumed, which is reasonable due to the increased number of parameters to be estimated. The ITRs estimated under the common betweenstudy heterogeneity also have slightly higher values. Overall, in the explored scenarios, the results are insensitive to the heterogeneity assumptions in the model, but assuming the common between-site heterogeneity will result in fewer parameters and a more practically feasible model.

With more sites contributing to the common ITR estimation, we have more precise blip parameter estimators and a smaller dVF, corresponding to a better ITR estimation. When there is only one site for each \mathcal{A}_i , due to the limited information, the variability is higher even without considering any between-site heterogeneity, and dVF is also larger. Networks (a) and (b) only differ in the additional sites for d_2 and d_3 , which will not provide direct information for the parameters of interest (the main effects or treatment-covariate interactions relative to the common reference treatment d_1). However, we still have more precise parameter estimates and ITRs with higher values in network (b). This suggests indirect evidence can also help with the common ITR estimation. Networks (c), (d), and (e) present more complex structures with relatively limited data information. No obvious difference in results are observed, indicating that the complexity of the network structure may not significantly impact the ITR estimation as long as the model is adapted accordingly based on the consistency equation (5.6).

Table 5.3: Simulation results when the data are generated under an unstructured betweensite heterogeneity model. Relative Bias (%, denoted by RB) and standard deviations (SD) of $\hat{\psi}_{20}$ and $\hat{\psi}_{21}$, and difference in value function (dVF) between true and estimated ITR and its standard deviation are reported across different networks, numbers of sites and heterogeneity assumptions in the Bayesian hierarchical model based on 2000 simulation runs.

Network	Number of sites	Heterogeneity	$\hat{\psi}_{20}$	$\hat{\psi}_{21}$	dVF (SD)
			nd (SD)	nd (SD)	
	1	0	-0.062(0.417)	$0.311 \ (0.182)$	$0.162 \ (0.125)$
a	3	Common	$0.007 \ (0.325)$	-0.004 (0.204)	$0.087\ (0.075)$
	·	Varying	$0.016\ (0.471)$	-0.017(0.331)	$0.089\ (0.076)$
	1	0	$0.023\ (0.330)$	$0.341 \ (0.158)$	$0.126\ (0.101)$
b	3	Common	-0.103(0.235)	0.062(0.143)	$0.058\ (0.052)$
		Varying	-0.087 (0.276)	$0.137\ (0.169)$	$0.063\ (0.056)$
	1	0	-0.124(0.315)	-0.236(0.138)	$0.107\ (0.083)$
С	3	Common	$0.042 \ (0.197)$	0.177(0.112)	$0.051 \ (0.046)$
	5	Varying	$0.061 \ (0.238)$	$0.219\ (0.135)$	$0.058\ (0.052)$
	1	0	0.134(0.411)	-0.759(0.183)	$0.219\ (0.299)$
d	3	Common	-0.089(0.268)	-0.221 (0.155)	$0.091\ (0.079)$
	5	Varying	-0.088(0.470)	-0.267(0.326)	$0.095\ (0.094)$
	1	0	$0.107 \ (0.290)$	0.125(0.124)	0.109(0.111)
e	3	Common	$0.039\ (0.180)$	-0.205(0.092)	$0.042 \ (0.044)$
	0	Varying	$0.042 \ (0.188)$	-0.190(0.097)	$0.044\ (0.060)$

5.4 Estimating individualized depression treatment

In this section, we apply the proposed method to estimate an ITR for patients with major depressive disorder (MDD) using data from three studies: The Sequenced Treatment Alternatives to Relieve Depression (STAR*D) study (Rush et al., 2004), Establishing Moderators

Table 5.4: Simulation results when the data are generated under the assumption of common between-site heterogeneity. Relative Bias (%, denoted by RB) and standard deviations (SD) of $\hat{\psi}_{20}$ and $\hat{\psi}_{21}$, and difference in value function (dVF) between true and estimated ITR and its standard deviation are reported across different networks, numbers of sites and heterogeneity assumptions in the Bayesian hierarchical model based on 2000 iterations.

Network	Number of sites	Heterogeneity	$\hat{\psi}_{20}$ RB (SD)	$\hat{\psi}_{21}$ RB (SD)	dVF (SD)
	1	0	-0.064 (0.416)	0.315 (0.182)	0.148 (0.109)
a	9	Common	$0.017 \ (0.322)$	$0.023 \ (0.195)$	$0.081 \ (0.068)$
	3	Varying	-0.004(0.461)	$0.032 \ (0.327)$	$0.082 \ (0.068)$
	1	0	$0.024 \ (0.327)$	$0.321 \ (0.152)$	0.113(0.088)
b	3	Common	-0.104(0.230)	$0.061 \ (0.134)$	$0.052 \ (0.047)$
	0	Varying	-0.089(0.272)	0.149(0.162)	$0.056\ (0.051)$
	1	0	-0.123(0.315)	-0.232(0.135)	$0.099\ (0.074)$
с	3	Common	$0.042 \ (0.196)$	$0.177 \ (0.110)$	$0.047\ (0.044)$
	5	Varying	$0.057 \ (0.237)$	0.187(0.133)	$0.054\ (0.049)$
	1	0	0.135(0.411)	-0.760(0.183)	$0.194\ (0.259)$
d	3	Common	-0.090 (0.267)	-0.215(0.153)	$0.084\ (0.069)$
	5	Varying	-0.082(0.463)	-0.244(0.330)	$0.089\ (0.086)$
	1	0	$0.083 \ (0.299)$	0.168(0.145)	$0.190\ (0.278)$
е	3	Common	-0.030(0.197)	$0.284\ (0.115)$	$0.065\ (0.064)$
	5	Varying	$-0.041 \ (0.206)$	$0.297\ (0.116)$	$0.071 \ (0.126)$

and Biosignatures of Antidepressant Response for Clinical Care (EMBARC) study (Trivedi et al., 2016), and Research Evaluating the Value of Augmenting Medication with Psychotherapy (REVAMP) study (Trivedi et al., 2008).

All the three studies are multistage randomized trials, with details of their designs described elsewhere (Rush et al., 2004; Trivedi et al., 2016, 2008). STAR*D include four stages. Due to single treatment assignment in the first stage and limited sample size in stages 3 and 4, we use data from stage 2 where patients without a satisfactory clinical outcome to citalopram (CIT) in the first stage were randomized to seven treatments. Among these,



Figure 5.2: Network structure of analysis of STARD, EMBARC, and REVAMP data. The size of each node (in red) is proportional to the total sample size in the corresponding treatment group, and the width of the connecting line (in gray) between any two treatments is proportional to the number of studies that directly compared the two treatments.

our focus on medications only: venlafaxine (VEN), sertraline (SER), bupropion (BUP), CIT augmented with BUP (CIT + BUP) or buspirone (BUS). In the case of EMBARC, we focus on SER and BUP in the second stage, as patients received only one active treatment SER and placebo in the first stage. For REVAMP, data from the first stage is used where a medication algorithm was implemented for treatment assignment, and SER, BUP, VEN, and escitalopram (ESCIT) are included. Therefore, in total, 6 treatments are identified: $\mathcal{A} = \{\text{SER}, \text{BUP}, \text{VEN}, \text{CIT} + \text{BUP}, \text{CIT} + \text{BUS}, \text{ESCIT}\}$, forming a network structure as shown in Figure 5.2, which corresponds to the scenario 5.1e in the simulation. The common reference treatment is SER, as it was included in all three studies and is often considered as the front-line treatment of MDD (Cipriani et al., 2008). In this case, study-specific and common reference treatments were considered the same.

Depression severity is measured by the 17-item Hamilton Depression Rating Scale (HDRS-17), with a larger value corresponding to more severe symptoms. In our analysis, we consider the negative of HDRS-17 as the outcome. We choose covariates based on meta-reviews of antidepressant treatment outcome predictors and modifiers (Kessler et al., 2017; Perlman et al., 2019). The following covariates that are common in the three studies and potentially related to differential treatment effects were identified and included in the model at the first stage: (1) socio-demographic variables: age (in years), sex (male, female), race (White, Non-White), marital status (single, married, divorced/widowed), number of years in formal education, employment status (employed, unemployed), number of people in household; (2) clinical variables: age at onset of first MDD (in years), number of depressive episodes, chronicity of current episode, baseline HDRS-17, and baseline 16-item Quick Inventory of Depressive Symptomatology (QIDS-16) before receiving the treatments. Among these variables, race was only used as an adjustment variable rather than a tailoring variable for treatment assignment, as basing treatment decisions on racial or ethnic groups can lead to healthcare disparities and inequities (Vyas et al., 2020). Additionally, while the number of depressive episodes was small for many patients, there were also several large values (e.g., 120), making it unsuitable to include this variable as a continuous linear term in the model. Therefore, the number of depressive episodes was dichotomized using a cutoff point of four; that is, a binary variable was created based on whether the number of episodes is greater than or equal to four. This threshold value of four was chosen based on the data to ensure sufficient sample sizes for estimating the parameters associated with the dichotomized variable. Since the inclusion criteria differ among the three studies, patients who had their first MDD after the age of 30 in STAR*D and REVAMP studies were excluded from the analysis to make the study populations more similar, thereby making the positivity assumption more plausible. Information on the included variables is collectively denoted by the vector \boldsymbol{x} . Records with missing values are removed. Finally, in our analysis dataset, we have 407, 87, and 308 samples from STAR*D, EMBARC, and REVAMP studies, respectively. The distributions of covariates were summarized in Table 5.5.

Linear regression models with above mentioned covariates and their interactions with treatments were used to obtain site-specific blip parameter estimates and the corresponding variance-covariance matrix. Since for most treatment comparisons only one or two studyspecific estimates are available, a Bayesian hierarchical model with $\tilde{\Sigma}_{it} = 0$ was used to obtain the common blip parameter estimates. The estimated ITR can thus be expressed as $d^{\text{opt}}(\boldsymbol{x}) = \arg \max_{d_h \in \mathcal{A}} \hat{\omega}_{d_h}(\boldsymbol{x})$, where $\hat{\omega}_{\text{SER}}(\boldsymbol{x}) = 0$, and for $d_h \neq \text{SER}$,

$$\begin{split} \hat{\omega}_{d_h}(\boldsymbol{x}) = & \hat{\psi}_{h0} + \hat{\psi}_{h1} \text{Age} + \hat{\psi}_{h2} \text{Male} + \hat{\psi}_{h3} \text{Single} + \hat{\psi}_{h4} \text{Divorced/Widowed} \\ & + \hat{\psi}_{h5} \text{Years of Education} + \hat{\psi}_{h6} \text{Unemployed} + \hat{\psi}_{h7} \text{Number of People in Household} \\ & + \hat{\psi}_{h8} \text{Age of First MDD} + \hat{\psi}_{h9} \text{Number of Episodes} + \hat{\psi}_{h,10} \text{Non-chronic} \\ & + \hat{\psi}_{h,11} \text{Baseline HDRS-17} + \hat{\psi}_{h,12} \text{Baseline QIDS-16}. \end{split}$$

with parameter estimates $\hat{\psi}_{ht}$ and the corresponding 95% posterior credible intervals shown in Table 5.6.

All estimated effects including the main treatment effects have wide credible intervals that include zero. This is not surprising, given the limited number of studies available for this analysis. Additionally, most variables in the analysis are binary, providing less information than continuous variables. Applying the estimated ITR to the 802 patients in the analysis dataset, we found that the optimal treatment recommendation was SER, BUP, VEN, CIT + BUP, CIT + BUS, and ESCIT for 27, 142, 292, 163, 93, and 85 patients, respectively.

5.5 Discussion

An optimal ITR can be estimated in a regression-based approach by including predefined treatment-covariate interactions. To increase the power for detecting differential treatment effects by covariates, large collections of datasets from multiple sites or studies are often needed. Different sites or studies may have varying treatment sets, however, an ITR analysis of all available treatments is desired. This presents a methodological gap which, to our knowledge, has not previously been considered. To address this gap, we adopt a two-stage Bayesian meta-analysis approach: at the first stage, study-specific analyses are conducted on the single study data only; at the second stage, summary measures including blip parameter estimates and variance or covariance estimates are shared and combined in a Bayesian hierarchical model to estimate a common optimal ITR.

The conventional pairwise meta-analysis approach focuses on binary treatments and assumes that all studies consist of the same treatment comparisons. In this manuscript, we consider multiple treatments and different studies may encompass different sets of treatments. With different treatment sets across studies, the estimated ITRs using study-specific data will only include a subset of the available treatments, and treatments that are not present in the same study cannot be simultaneously considered in an estimated ITR. To address this issue and construct an ITR for all treatments, we employ a network meta-analysis approach, and construct the Bayesian hierarchical model at the second stage based on the consistency equation (5.6). The consistency assumption relates our parameters of interest and parameters that can be estimated with direct evidence, and is essential to ensure the validity of the Bayesian hierarchical model at the second stage. However, this assumption can be violated in certain cases. For example, a treatment in a given site may be inappropriate for patients in another site due to the heterogeneity in population. In this case, not only will the consistency assumption be violated, but so too will the positivity assumption if specific covariate combinations preclude receipt of particular treatments. Therefore, the proposed method will be inappropriate due to the violation of the assumptions in both the network meta-analysis and causal inference aspects of the analysis. The consistency assumption can be statistically assessed when both direct and indirect evidence are available. A discussion of these approaches can be found elsewhere (Lu and Ades, 2006; Dias et al., 2010a,b; White et al., 2012; Piepho et al., 2012; Donegan et al., 2013). We can reduce the possibility of inconsistency in both the design and analysis stages. When designing a multisite trial, it is recommended to standardize the study protocol to guarantee the same or similar populations, treatment delivery, and assessment of the outcomes and covariates (Weinberger et al., 2001; Noda et al., 2006). Before the two-stage analysis of the data, we may also pre-screen participants based on some analyst-defined harmonized eligibility criteria to ensure the samples in the analysis are similar across sites. If the consistency assumption is deemed unfeasible, incorporation of inconsistency may be considered (Lu and Ades, 2006). For example, an inconsistency factor $\delta_{a_i^{\tilde{h}}, a_i^{(1)}, 1}$ may be added to the consistency equation (5.6), i.e., $\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)}, t} = \psi_{a_i^{(\tilde{h})}t} - \psi_{a_i^{(1)}t} + \delta_{a_i^{\tilde{h}}, a_i^{(1)}, 1}$, and a model or a distribution can be posited for $\delta_{a_i^{\tilde{h}}, a_i^{(1)}, 1}$. However, this adaptation in our setting requires further investigation and the implications for positivity violations should also be considered; this is beyond the scope of this manuscript.

In the consistency equation (5.6), to identify parameters $\psi_{a_i^{(\tilde{h})}t}$ and $\psi_{a_i^{(1)}t}$ from $\tilde{\psi}_{a_i^{(\tilde{h})}a_i^{(1)},t}$ data information should be available for at least one of the two parameters $\psi_{a_i^{(\tilde{h})}t}$ and $\psi_{a_i^{(1)}t}$ This requires a connected network, i.e., every two treatments in the network can be either directly or indirectly compared. A disconnected network complicates the model (Goring et al., 2016; Stevens et al., 2018; Schmitz et al., 2018; Béliveau and Gustafson, 2021). It would be of interest to explore whether proposed methods for meta-analysis in a disconnected network setting can be adapted in our context. For example, in a random baseline model (Béliveau et al., 2017), a disconnected network will be connected by assuming the baseline effects are exchangeable across studies. In our context, this could be translated into the exchangeability of the site-specific treatment-free parameters and including their estimates in the Bayesian hierarchical model. One criticism of the random baseline model is that it breaks the randomization by assuming randomization not only within studies but also across studies. This limitation might be less of a concern in our case. Although we only assume a common distribution for the blip parameters in (5.3), in reality, it is also likely that the treatment-free parameters have a common distribution, as the populations in each site should be similar or be a subset of a larger target population. If the population (or the site-specific ITR) is totally different or unrelated across sites, estimating a common optimal ITR is not meaningful.

The estimated common optimal ITR is not necessarily better than the site-specific optimal ITRs at individual sites; however, it is applicable to a larger target population. Schnitzer et al. (2016) defined this large target population to be the union of populations represented by the individual sites, referring to it as a metapopulation. They also nonparametrically established target parameters in a network meta-analysis that can be causally interpreted on the metapopulation. Since their definition is model-independent, it may also help with the causal interpretation of the common optimal ITR in our work.

In this manuscript, we focus on a continuous outcome, and a linear regression model is used at the first stage. The implementation at the first stage assumes linearity and thus the results are sensitive to model specification. In practical cases, the linear relationship may not fully capture the true dynamics between covariates, treatments, treatment-covariate interactions, and the outcome. When the outcome model is misspecified in the Q-learning approach, the estimator in the first stage loses its consistency and unbiasedness, leading to biased estimation of the common blip parameters. However, linear models offer the advantage of interpretability, which is crucial in treatment decision-making. Alternatives such as dWOLS and G-estimation, which provide both double robustness and interpretability, can also be considered in the first stage; these options are particularly attractive in the context of randomized trials, where the treatment allocation is known by design. Extension to other outcome types is also straightforward; we can apply existing ITR estimation methods for a specific outcome type at the first stage (Tchetgen Tchetgen et al., 2010; Linn et al., 2017; Kidwell et al., 2018; Simoneau et al., 2020).

We evaluate the ITR estimation through simulations. In all scenarios explored, the proposed method yielded consistent estimation. Additionally, simulation results support the feasibility of assuming common between-site heterogeneity when specifying the structure of variancecovariance matrix $\tilde{\Sigma}_t$, regardless of the true underlying structure. The application of the proposed method is illustrated through an analysis of real data from the STAR*D, EMBARC, and REVAMP studies. Common covariates in the three studies that are considered to be related to the depression outcome or treatment response in the literature were selected and included in the model.

Although statistical associations between these covariates and the depression outcome have been established in the literature, in clinical settings some variables such as marital status and the number of people in the household are seldom considered by physicians when prescribing depression medications. We considered these covariates in the analysis as a proxy for social support, but whether and how to deploy these in clinical contexts requires careful consideration. We dichotomized the number of depressive episodes, which will lead to information loss. The cutoff point was determined solely based on the data and thus lacks clinical interpretations. In clinical settings, patients with two or more episodes are considered to have recurrent depressive disorders. However, most patients in the EMBARC study had two or more episodes, making a threshold value of two less appealing and feasible. There would be little information available to learn about treatment tailoring by episodes if the standard threshold is used. Moreover, only patients who had their first MDD before the age of 30 could be included in EMBARC study. Patients not satisfying this condition in the STAR*D and REVAMP studies were excluded from the analysis. While this reduced the sample sizes in individual studies, it could be practically recommended when populations from different studies are quite different as a means of ensuring the positivity assumption is met.

All estimated effects were relatively small in magnitude compared to their wide credible intervals, indicating a lack of strong evidence for the need to tailor treatment assignments based on the covariates considered in the analysis. This aligns with findings in the literature that baseline anxiety level and common socio-demographic variables, such as age, marital/employment status, or education level, do not contribute to the differential treatment effects (Rush et al., 2001; Archer et al., 2024). However, Noma et al. (2019) found that several variables including age, age at onset, and HDRS subscales could be potential effect modifiers for response to depression treatments through a meta-analysis. We also observed that the main effects of CIT+ BUP, and ESCIT are negative. However, combination therapy is generally expected to be superior to monotherapy (Henssler et al., 2016) and ESCIT is known to be more or at least similarly effective compared to a range of antidepressants including CIT, SER, VEN, and BUP (Kennedy et al., 2009; Kirino, 2012). The discrepancy could arise from the limited evidence available in our analysis, as both CIT + BUP and ES-CIT were only represented in a single study. Additionally, only 35 patients in the REVAMP study received ESCIT, resulting in estimates with low precision. No modelling assumptions were found to be heavily violated for the linear regressions at the first stage, but the R^2 was relatively low (about 30% to 40%). The covariates included in the analysis only reflect the socio-demographic and symptom information. Some important features that are predictive of the outcome, such as those related to genetic information or comorbidities (Perlman et al., 2019), are missing. These limitations highlight the need for further research to identify reliable effect modifiers and predictive covariates, as well as to obtain more data from additional samples or studies. Our analysis therefore provides an important proof-of-concept, and important additional refinements would be required before deploying findings from such an analysis in a clinical setting. Nevertheless, the results showcase a promising approach to leveraging multiple data-sources to learn about effect modification of a potentially large number of treatment options by important patient characteristics, leveraging those covariates to better allocate treatment for improved patient care.

Acknowledgements

This work is supported by an award from the Canadian Institutes of Health Research CIHR FDN-16726. EEMM is a Canada Research Chair (Tier 1) in Statistical Methods for Precision Medicine and acknowledges the support of a chercheur de mérite career award from the Fonds de Recherche du Québec, Santé. SG acknowledges support from the Natural Sciences and Engineering Research Council of Canada (NSERC), Candian Institute for Statistical Sciences (CANSSI) and Fonds de recherche du Québec - Sante (FRQS).

Variables	STAR*D	EMBARC	REVAMP
	(n = 407)	(n = 87)	(n = 308)
Age	38.41 (12.57)	38.70 (13.13)	40.06 (12.29)
Sex			
Female	251 (61.7)	53~(60.9)	191 (62.0)
Male	156 (38.3)	34(39.1)	117 (38.0)
Race			
White	348 (85.5)	$60 \ (69.0)$	259(84.1)
Non-White	59(14.5)	27 (31.0)	49(15.9)
Marital Status			
Married	165 (40.5)	23(26.4)	111 (36.0)
Single	136(33.4)	50(57.5)	113 (36.7)
Divorced/Widowed	106 (26.0)	14(16.1)	84(27.3)
Years of Education	13.78(2.81)	15.24(2.39)	14.72(2.62)
Employment Status			
Employed	243(59.7)	44 (50.6)	195~(63.3)
Unemployed/Retired	164 (40.3)	43 (49.4)	113 (36.7)
Number of People in Household	2.67(1.48)	2.62(2.31)	2.53(1.57)
Age of First MDD	$16.81 \ (6.11)$	16.51(5.41)	17.47(5.85)
Number of Episodes			
< 4	194 (47.7)	39(44.8)	245(79.5)
≥ 4	213(52.3)	48(55.2)	63(20.5)
Chronicity			
Chronic	95(23.3)	37(42.5)	294 (95.5)
Non chronic	312(76.7)	50(57.5)	14 (4.5)
Baseline HRSD-17	$16.91 \ (7.12)$	15.69(5.47)	20.85(4.30)
Baseline QIDS-16	11.86(4.73)	17.85(2.84)	15.13(3.55)

Table 5.5: Patient characteristics for STAR*D, EMBARC, and REVAMP studies.

	L.				
ψ_{ht}	BUP	VEN	CIT + BUP	CIT + BUS	ESCIT
Main treatment effect $(\hat{\psi}_{h0})$	$5.24 \ (-10.8, 21.48)$	-0.61(-21.4,21.1)	-2.87 $(-24.97, 18.97)$	$21.32\ (2.6, 39.7)$	-15.86(-46.17, 13.35)
Age $(\hat{\psi}_{h_1})$	$0.03 \ (-0.15, 0.2)$	-0.14 (-0.35, 0.08)	-0.08 $(-0.29, 0.12)$	-0.13 (-0.32, 0.07)	0.07 (-0.24, 0.38)
Male $(\hat{\psi}_{h_2})$	-0.49 $(-4.44, 3.49)$	-1.04(-5.8,3.91)	1 (-3.54, 5.61)	-2.09 $(-6.53, 2.1)$	1.72 (-3.43, 7.05)
Single $(\hat{\psi}_{h3})$	-0.66(-5.96,4.77)	$-0.42 \ (-7.3, 6.3)$	$1.57 \ (-4.53, 7.59)$	-3.79 (-9.74, 2.15)	2.26(-6.55,10.78)
Divorced/Widowed $(\hat{\psi}_{h4})$	0.95 (-3.94, 6.02)	5.75(0.1,11.38)	3.43 (-2.12, 8.96)	2.81 (-2.56, 8.29)	-0.47 ($-8.72,7.81$)
Years of Education $(\hat{\psi}_{h5})$	$-0.34 \ (-1.05, 0.38)$	0.08 (-0.78, 0.95)	$-0.32 \ (-1.24, 0.6)$	-0.87 (-1.65,-0.1)	-0.1 (-0.89, 0.73)
Unemployed/Retired $(\hat{\psi}_{h6})$	2.78 (-0.95, 6.54)	4.35 $(-0.76, 9.22)$	2.23(-2.39,6.9)	3.67 (-0.9, 7.84)	-2.56 $(-8.35, 3.22)$
Number of People in Household $(\hat{\psi}_{h\tau})$	-0.48 $(-1.62, 0.65)$	$0.19 \ (-1.94, 2.3)$	$0.03 \ (-1.71, 1.74)$	-1.38 (-2.92,0.17)	-0.05(-2.31, 2.24)
Age of First MDD $(\hat{\psi}_{n8})$	$-0.04 \ (-0.33, 0.24)$	-0.08(-0.48,0.31)	$0.16 \ (-0.24, 0.55)$	-0.03 (-0.37, 0.31)	0.23 $(-0.36, 0.83)$
Number of Episodes $(\hat{\psi}_{h9}) \ge 4$	-0.78 $(-4.96, 3.34)$	-0.4 (-5.67, 4.85)	1.18 (-3.7, 6.04)	-0.04 $(-4.54, 4.5)$	-0.96(-7.27, 5.39)
Non-chronic $(\hat{\psi}_{h,10})$	0.4 (-3.82, 4.77)	2.84(-2.96, 8.4)	1.68 (-3.65, 6.76)	0.35(-4.44,5.21)	1.53 (-8.93, 12)
Baseline HRSD-17 $(\hat{\psi}_{h,11})$	$0.01 \ (-0.34, 0.37)$	$0.3 \ (-0.15, 0.76)$	$0.29 \ (-0.13, 0.73)$	0.38 (-0.02,0.78)	0.26 (-0.54, 1.01)
Baseline QIDS-16 $(\hat{\psi}_{h,12})$	$0.1 \ (-0.49, 0.65)$	-0.04 (-0.74,0.67)	0.09 (-0.57,0.76)	-0.47 (-1.1, 0.15)	0.33 (-0.49,1.12)

Table 5.6: Blip parameter estimates (posterior medians) and the 95% posterior credible intervals for the real data application.

Chapter 6

Conclusion

6.1 Summary

The three manuscripts (Chapters 3, 4 and 5) presented in this thesis address practical challenges in PM and contribute novel insights to the literature on statistical trials and ITRs.

In Chapter 3, motivated by a potential real-world CRT of a CDSS, I develop two Bayesian group sequential designs for CRTs to improve trial efficiency. These designs aim to enhance trial efficiency by dividing participant recruitment into groups and allowing for early efficacy stopping. The two designs differ in participant recruitment methods: one sequentially recruits clusters with a preset maximum cluster size, while the other recruits all clusters at once but sequentially enrolls individual participants within each cluster until early trial termination for efficacy or final analysis. As literature on Bayesian adaptive CRTs is scarce, the proposed designs and the investigation of their operating characteristics serves as an important contribution to this area. Through simulation studies, I evaluate the power and false positive rates of both designs across various scenarios and two outcome types. Practical design recommendations are provided based on the simulation results. Chapter 4 focuses on addressing the challenge of estimating ITRs using multisite data under constrained individual-level data sharing. I adopt a two-stage Bayesian meta-analysis approach. In the first stage, site-specific analyses are conducted based on individual-level data within each site. In the second stage, site-specific estimates rather than individual-level data are shared and modelled in a Bayesian hierarchical model to obtain a common ITR. This approach allows for heterogeneity across sites, a potentially more realistic scenario that has not been thoroughly explored in previous work (Danieli and Moodie, 2022; Moodie et al., 2022). Challenging aspects that might occur in real-life data, in particular, sparsity in both the data and the model, are also discussed. In the presence of data sparsity, it is crucial to examine how the lack of variability in some covariates at some sites will change the interpretation of the site-specific estimates and link these estimates to correct common parameters. For model sparsity, I propose to use shrinkage priors and the 95% credible interval criterion to select the variables that truly have tailoring effects on treatment assignment. The simulation results suggest that the proposed approach can provide consistent estimation of the parameters that fully characterize the optimal ITR. The estimation results are stable across different confounding scenarios and prior choices explored in the simulation study. However, with a larger degree of heterogeneity across sites, the estimated optimal ITR will have a lower value, requiring more sites to obtain accurate estimation. The Warfarin analysis suggests none of the covariates under consideration have tailoring effects on Warfarin dosing and thus a common dose is recommended. Both the simulation study and the Warfarin analysis demonstrate that the proposed method produces similar estimates to a one-stage approach where individual-level data are pooled together and analyzed as a single dataset, illustrating its potential to estimate the optimal ITR using multisite data without the need for sharing individual-level data beyond the sites.

Chapter 4 considers a conventional pairwise meta-analysis approach, assuming all treatments are available at all sites. Building upon this, Chapter 5 extends the methodology from Chapter 4 to settings involving multiple treatments and varying treatment sets across sites. In this case, site-specific ITRs may not encompass all available treatments, and an ITR of all treatments needs to be derived from a network of evidence across multiple sites. To handle the network structure, a consistency assumption is made and the consistency equation (5.6) is incorporated in the Bayesian hierarchical model in the second stage. The simulation study suggests that the proposed two-stage Bayesian network meta-analysis approach can yield consistent estimators for the parameters characterizing the optimal ITR. It also supports the feasibility of assuming common between-site heterogeneity across treatment comparisons, a pivotal consideration in network meta-analysis for simplifying model complexity. This approach is illustrated using the data from STAR*D, EMBARC, and REVAMP studies to establish an ITR for the treatment of depression.

6.2 Limitations and Future Work

While the two designs proposed in Chapter 3 are straightforward, they are restricted with respect to stopping criteria, designs, and outcome types. Particularly I only considered stopping for efficacy, a two-arm design and binary and continuous outcomes. Extending the two designs to incorporate other decision criteria is feasible. For example, given the efficacy and futility decision boundaries U_e and U_f , a decision rule incorporating both efficacy and futility criteria can be as follows: at interim analysis k, if $P(\theta > \delta | D_k) > U_e$, the trial stops early for efficacy; if $P(\theta > \delta | D_k) < U_f$, the trial stops early for futility; otherwise, the trial will continue to recruit the next group of samples.

Extension to multi-arm CRTs is also an interesting direction for future work. In individually randomized trials, with multiple treatments, instead of a simple continue-or-terminate decision at interim analyses, researchers can choose to drop ineffective treatments or select effective ones. Among the multiple treatments, one will be identified as the control. Arm dropping or selection at each interim point can then be facilitated by assessing a series of decision criteria. For example, for each active treatment, the posterior probability that the treatment effect relative to the control arm is greater than a preset minimal important difference can be evaluated. If this probability is smaller than a pre-specified decision boundary, this specific treatment can be dropped. Extending this procedure to the proposed two designs can create new challenges. For example, design 2 recruits all clusters and randomizes them into different treatment arms at the start of the trial. If, at an interim point, one arm is dropped, it is unclear whether the clusters in the dropped arm should be completely removed from the trial or be reassigned to the remaining arms. Future work can compare how these choices will influence the analysis and design performance.

The proposed designs assume that randomization probabilities are fixed over the course of the trial. Adaptive randomization may also be incorporated for design 1. For instance, after assessing interim results, randomization probabilities can be adjusted in proportion to the posterior probability that a treatment arm yields better outcomes or is the best if multiple treatments are available. This response-adaptive randomization approach (Rosenberger et al., 2012) can improve resource allocation and also has ethical implications, as more clusters along with their individual participants are assigned to the better treatment. However, this adaptive feature may not be feasible for design 2, as all clusters are assigned to different treatment arms from the start.

Bayesian group sequential designs with survival outcomes have been less studied compared with continuous or binary outcomes. The Cox proportional hazards model is standard in analysis of survival outcomes. It is a semi-parametric model where the baseline hazard is unspecified and the parameters can be estimated by maximizing a partial likelihood function in the frequentist framework. Its Bayesian version, however, requires specification of the full likelihood. One can model the baseline hazard function by splines or directly assign priors to the baseline hazard function, such as the Gamma process (Ibrahim et al., 2001). To account for the correlation resulting from clustering, a frailty term may also be included (Hougaard, 1995). Then, decision criteria can be constructed similarly based on posterior statements of hazard ratio. A more recent approach for Bayesian group sequential designs for individually randomized trials with survival outcomes has been studied by Zhu et al. (2019). The number of events at each unique event time is modelled by Bernoulli distributions and the stopping rule is constructed based on the Bayes factor (Kass and Raftery, 1995): given the hypotheses H_0 and H_A , at each interim point, if the Bayes factor in favor of H_0 is smaller than a specified decision boundary, then a recommendation is made to stop the trial early, where the decision boundary is chosen to maintain the overall false positive rate at a desired level through an alpha-spending function. Adapting such approaches for CRTs with survival outcomes may be considered in future work.

As described in Section 2.2.2, Bayesian group sequential designs can also be developed based on decision theory, where loss or utility functions can be defined to incorporate factors such as different costs of recruiting clusters and individual patients. By optimizing expected loss or utility, trials can be designed to balance practical considerations, while maintaining statistical performance. Extending the framework of Bayesian adaptive CRTs to include the additional features and methodologies described above can significantly improve their utility and applicability. These extensions will provide researchers with more tools and flexibility in designing and conducting clinical trials, leading to more efficient and ethical studies.

Chapters 4 and 5 explore the two-stage Bayesian meta-analysis approach for ITR estimation in different settings. Both focus solely on continuous outcomes, using Q-learning implemented via linear regressions for the first stage. However, Q-learning does require correct specification of the outcome regression model. Model misspecification can lead to biased blip parameter estimators and thus suboptimal ITRs. To address this limitation, exploring alternative stage-one models such as G-estimation (Robins, 2004) and dWOLS (Wallace and Moodie, 2015) might be of interest. Both G-estimation and dWOLS offer double robustness, which means they can provide consistent estimators of blip parameters even if one of the treatment assignment and outcome models is misspecified. Such doubly robust models are particularly appealing when the data come from randomized trials in which the treatment assignment model is already known by the design.

Extensions to other outcome types, such as binary or survival outcomes, are also feasible given the existing methods for ITR estimation for these outcome types (Tchetgen Tchetgen et al., 2010; Linn et al., 2017; Kidwell et al., 2018; Simoneau et al., 2020). For example, Simoneau et al. (2020) extend dWOLS to survival outcomes and use accelerated failure time models. That is, the logarithm of survival time under a specific treatment can be decomposed into a treatment-free function, a blip function and some random error, where the treatment-free and blip functions are defined similarly as those in this thesis. Consistency and asymptotic normality of the blip parameter estimators have also been established. In this case, only the stage-one model which provides the estimates used in the meta-analysis procedure will be influenced and the Bayesian hierarchical model in the second stage can remain the same.

In this thesis, the Q-learning approach in the first stage is considered within the frequentist framework, while the model in the second stage is Bayesian, incorporating prior information for the common blip parameters through prior distributions. This choice is partly due to the fact that frequentist models for ITR estimation are well-established in the literature, although there is growing interest in Bayesian approaches; see, for example, Murray et al. (2018); Logan et al. (2019); Rodriguez Duque et al. (2022). However, there are instances where individual sites possess additional information about the site-specific ITRs. In such cases, Bayesian models may be considered for the first stage as well. This might allow for a more nuanced integration of site-specific information and potentially yield more accurate estimates of the optimal ITRs.

Although Chapters 4 and 5 focus on ITRs, healthcare services and treatments can span longterm periods. Some patient characteristics and outcomes can change over time, requiring treatment to be adapted over time for it to remain optimal. Consequently, a single ITR for the entire treatment period may not be feasible. DTRs extend ITRs to multiple treatment stages, and the regression-based approaches mentioned in this thesis are applicable to DTRs using backward induction. To extend the proposed methods to DTRs, a natural choice is that individual sites estimate their own DTRs and the two-stage model will be repeated for the multiple stages. However, such an approach would require significant further investigation as combining DTRs across sites where there is model selection at each stage may be significantly more challenging than estimating and combining single-stage ITRs.

Another approach that allows a quantification of whether data can be analyzed without releasing individual-level data is differential privacy. Differentially private algorithms to estimate the optimal DTR have been studied by Spicker et al. (2024), focusing on outcome weighted learning, which is a classification-based method. Standard linear regression, which I employed for Q-learning in the first-stage estimation, is not differentially private, but in the proposed two-stage approach, only blip model estimates from the first stage are shared. Exploring whether this poses a violation of differential privacy would be an interesting area of research.

A horseshoe prior and a 95% credible interval selection criterion are used in both the simulation and the Warfarin analysis in Chapter 4 to select variables that have tailoring effects on the optimal treatment. While these choices are common in Bayesian variable selection, simulation shows that they may not be powerful enough to detect very small effects. In real-world scenarios, without knowing the true effects, it is recommended to explore a series of shrinkage priors and selection criteria (Bondell and Reich, 2012; Hahn and Carvalho, 2015; Li and Pati, 2017; van Erp et al., 2019) and compare their results to obtain potentially more robust conclusions.

The method described in Chapter 5 relies on the consistency assumption in network metaanalysis literature, which guarantees that the indirect and direct evidence are in agreement and further leads to the consistency equation (5.6). Before applying the method to real-

world data, it is essential to carefully evaluate the validity of the consistency assumption using methods summarized in Section 2.5. However, for the real data application in Chapter 5, consistency assessment is not required. The consistency assumption is only necessary when both direct and indirect evidence exist. In the analysis of STAR*D, EMBARC, and REVAMP data, a common reference treatment, SER, is identified, and parameters of interest are all defined with respect to SER. Therefore, there is no indirect evidence related to parameters of interest in that real data analysis. In general cases, if the consistency assumption is not satisfied, it is necessary to incorporate the inconsistency into the Bayesian hierarchical model at the second stage (Lu and Ades, 2006). One possible approach is to add an inconsis- $\text{tency factor } \delta_{a_i^{\tilde{h}}, a_i^{(1)}, 1} \text{ to the consistency equation (5.6), i.e., } \tilde{\psi}_{a_i^{(\tilde{h})} a_i^{(1)}, t} = \psi_{a_i^{(\tilde{h})} t} - \psi_{a_i^{(1)} t} + \delta_{a_i^{\tilde{h}}, a_i^{(1)}, 1},$ and posit a model, such as a normal model, for the inconsistency factor. However, in my simulation study, I only consider consistent scenarios. Future work can focus on the performance of ITR estimation under inconsistency to accommodate more complex scenarios. This could involve exploring various types of inconsistency, for example, the loop and design inconsistencies defined in Higgins et al. (2012), and their impact on the estimation of ITRs.

Identifying the parameters of interest from the consistency equation (5.6) requires a connected network, where each pair of treatments can be compared either directly or indirectly. If the network is disconnected, it is challenging to estimate these parameters accurately due to identifiability issues. In such cases, exploring whether existing methods for metaanalysis in disconnected networks can be adapted would be of interest. For instance, one possible approach is the random baseline model (Béliveau et al., 2017), which connects a disconnected network by assuming the exchangeability of the baseline effects. To adapt this method for the proposed two-stage approach, in addition to blip parameter estimates, site-specific treatment-free parameter estimates will also be incorporated in the Bayesian hierarchical model. A common distribution can be assumed for the site-specific treatment-free parameters which are defined with respect to a common reference treatment. While this approach can address the disconnectedness, its use in randomized trials is criticized since the assumed common distribution for the baseline effects may interfere with randomization. This may not be a significant problem in the case considered in this thesis. Although I only assume a common distribution for the blip parameters, in reality, it is also likely that the treatment-free parameters have a common distribution, as the populations in each site should be similar or be a subset of a larger target population. If the populations (or the sitespecific ITRs) are totally different or unrelated across sites, estimating a common optimal ITR is not meaningful.

The real data analyses in Chapters 4 and 5 may offer some insights into individualized Warfarin dosing and individualized treatment of depression. However, caution is needed when using this information to support clinical decision making in practice. One significant limitation is that some important predictors might not be included in the datasets and, consequently, are not considered in the analyses. This can lead to incomplete or less accurate models. For instance, in the real data analysis presented in Chapter 5, while no modeling assumptions appear to be violated in the first stage, the low R^2 value indicates that the model has limited explanatory power. Therefore, the real data analyses should be viewed as illustrative examples of the proposed methods rather than guides for clinical practice. Further model refinement and validation with comprehensive datasets, including all relevant variables, will be necessary before these findings can be reliably applied in clinical settings.

Practical issues should also be carefully considered when implementing the methods presented in this thesis. For example, in Chapter 3, the two designs are compared primarily in terms of power and false positive rate in the simulation study. Depending on the contexts, other characteristics may also be considered, such as probability of early stopping or expected sample size. These characteristics evaluate the designs from a statistical perspective. However, feasibility is also a critical factor for design decisions. For instance, one design may exhibit better operating characteristics, but its implementation could be more time-consuming or expensive. In such cases, the balance between improved design performance and additional costs resulting from choosing this design must be carefully considered, requiring comprehensive investigation into specific contexts.

For the two-stage approach considered in Chapters 4 and 5, sufficient statistical expertise at individual sites is necessary for site-specific analyses. Consistency is crucial for the validity of the proposed two-stage Bayesian network meta-analysis approach, and simulation in Chapter 4 also shows the estimated optimal ITR has a lower value for a larger between-site heterogeneity. Therefore, when planning a multisite trial to estimate ITRs, it is important to reduce extraneous heterogeneity or inconsistency by uniform staff training, treatment delivery, data collection, and measurement.

6.3 Concluding Remarks

PM is of growing importance in healthcare. This thesis addresses practical challenges in PM by proposing new, efficient designs that are appealing to clinical scientists and developing novel approaches that allow researchers to estimate optimal treatment strategies using siloed data sources. Appendices

APPENDIX A

Appendix to Manuscript 1

A.1 Simulation for binary outcomes

To generate clustered binary data, we first generate n cluster-specific proportions from Beta distributions with Beta parameters determined by the preset population mean and ICC so that all cluster-specific proportions are strictly between 0 and 1. Then, within each cluster, the number of events were generated from a binomial distribution with the cluster size at current stage and the cluster-specific proportions. Then the resulting $m \times n$ binary variables satisfy the predetermined correlation structure.

The performance of the two designs in terms of false positive rate and power, when only one interim analysis is planned, is displayed in Figures A.1 and A.2. It can be observed that design 1 has higher false positive rate and both designs perform almost equally well in terms of power over the parameter space we explored. Also, an ad hoc decision boundary of U = 0.98 helps control false positive rate compared with U = 0.95 while the decrease in power coming with a larger decision boundary is acceptable. However, for design 2, even with a smaller decision boundary, the false positive rates for most cases are acceptable considering that they are fluctuating around 0.05. In addition, as in the continuous case, when the underlying ICC is larger, power is relatively small.



Figure A.1: Plot of false positive rate versus baseline risk when $n = 20, 40, 60, \rho = 0.05, 0.1$ for (a) U = 0.95 and (b) U = 0.98 with single interim analysis planned. The dashed lines show the false positive rate of 0.05.

The design performance for multiple interim analyses is shown in Figures A.3 and A.4. Design 1 still has higher false positive rates and the powers for the two designs are similar over the parameter space we explored. With multiple interim analyses, both false positive rate and power will increase for most cases. However, similar to continuous outcomes, when the available sample size is large, only a slight increase in power can be expected with more interim analyses planned. With multiple interim analyses a larger boundary value can evidently help reduce false positive rate, especially for design 1. The resulting decrease in power is acceptable for a sufficient sample size. However, with a small sample size, a larger decision boundary may lead to insufficient power. The results for a larger cluster size are also displayed in Figures A.5 and A.6.

Therefore, for binary outcomes, for the parameters explored in the simulation, design 2 may be recommended based on design operating characteristics since it has smaller false positive rate and performs comparably with design 1 in terms of power. Additionally, a single interim analysis is sufficient to control false positive rate while maintaining satisfactory power for the scenarios we explored. Regarding the decision boundary, if design 2 with single interim analysis is planned, for the settings in our simulation, the smaller decision boundary (e.g. U = 0.95) may be conservative enough for obtaining a satisfactory false positive rate as well



Figure A.2: Plot of power versus true treatment effect for $n = 20, 40, 60, \pi_c = 0.25, 0.35, 0.45$, with the subpanels (a)-(d) indicating all possible combinations of $\rho = 0.05, 0.1$ and U = 0.95, 0.98 for single interim analysis. The dashed lines show the power of 0.85.

as power. However, if due to feasibility, design 1 is preferred or multiple interim analyses are of interest a larger decision boundary (e.g. U = 0.98) may be considered. Developing a more general set of recommendations requires further exploration.



Figure A.3: Plot of false positive rate versus number of interim looks for n = 20, 40, 60, $\pi_c = 0.25, 0.35, 0.45, m = 8$ with subpanels (a)-(d) indicating all possible combinations of $\rho = 0.05, 0.1$ and U = 0.95, 0.98. The dashed lines show the false positive rate of 0.05.



Figure A.4: Plot of power versus number of interim looks for $n = 20, 40, 60, \pi_c = 0.25, 0.35, 0.45, m = 8$ with subpanels (a)-(h) indicating all possible combinations of $\rho = 0.05, 0.1, U = 0.95, 0.98$ and $\theta = 0.1, 0.2$. The dashed lines show the power of 0.8.



Figure A.5: Plot of false positive rate versus number of interim looks for n = 20, 40, 60, $\pi_c = 0.25, 0.35, 0.45, m = 16$ with subpanels (a)-(d) indicating all possible combinations of $\rho = 0.05, 0.1$ and U = 0.95, 0.98. The dashed lines show the false positive rate of 0.05.


Figure A.6: Plot of power versus number of interim looks for $n = 20, 40, 60, \pi_c = 0.25, 0.35, 0.45, m = 16$ with subpanels (a)-(h) indicating all possible combinations of $\rho = 0.05, 0.1, U = 0.95, 0.98$ and $\theta = 0.1, 0.2$. The dashed lines show the power of 0.8.

A.2 Simulation results for continuous outcomes with clus-



ter size m=16

Figure A.7: Plot of false positive rate versus number of interim looks for n = 20, 40, 60, $\rho = 0.05, 0.1, 0.15$, m = 16 for (a) U = 0.95 and (b) U = 0.98. The dashed lines show the false positive rate of 0.05.



Figure A.8: Plot of power versus number of interim looks for $n = 20, 40, 60, \theta = 0.2, 0.5, 0.8, m = 16$ with subpanels (a)-(f) indicating all possible combinations of $\rho = 0.05, 0.1, 0.15$ and U = 0.95, 0.98. The dashed lines show the power of 0.8.

APPENDIX B

Appendix to Manuscript 2

B.1 Link with a one-stage approach

In this section, we illustrate that, under certain assumptions, similar estimates of the blip function parameters ψ can be obtained in the proposed two-stage approach and a one-stage approach based on the full individual-level data. As mentioned in the main manuscript, in site *i*, we have the site-specific outcome model:

$$E(Y_{ij} \mid \boldsymbol{X} = \boldsymbol{x}_{ij}, A = a_{ij}) = \boldsymbol{\beta}_i^{\top} \boldsymbol{x}_{ij}^{(\boldsymbol{\beta})} + a_{ij} \boldsymbol{\psi}_i^{\top} \boldsymbol{x}_{ij}^{(\boldsymbol{\psi})},$$

where $i \in \{1, \ldots, K\}$ and $j \in \{1, \ldots, n_i\}$ index the site and individual patient in a given site respectively, and n_i is the number of patients in site i. The predictive and prescriptive covariate vectors are denoted by $\boldsymbol{x}_{ij}^{(\beta)}$ and $\boldsymbol{x}_{ij}^{(\psi)}$, respectively. The *p*-dimensional site-specific treatment-free function parameter and *q*-dimensional blip function parameter are denoted by $\boldsymbol{\beta}_i = (\beta_{i0}, \ldots, \beta_{i,p-1})$ and $\boldsymbol{\psi}_i = (\psi_{i0}, \ldots, \psi_{i,q-1})$, respectively. Then, with site-specific estimates $\hat{\psi}_{it}$ and the associated standard deviations $\mathrm{sd}(\hat{\psi}_{it})$, for $t = 0, \ldots, q - 1$, obtained from the stage-one models, a Bayesian hierarchical model is implemented in the second stage:

$$\begin{split} \hat{\psi}_{it} &\sim N(\psi_{it}, \mathrm{sd}(\hat{\psi}_{it})^2), \\ \psi_{it} &\sim N(\psi_t, \sigma_{\psi_t}^2), \\ \psi_t &\sim p_{\psi_t}(\psi_t), \\ \sigma_{\psi_t}^2 &\sim p_{\sigma_{\psi_t}^2}(\sigma_{\psi_t}^2), \end{split}$$

where ψ_{it} and ψ_t are the (t + 1)-th elements of the site-specific and common blip function parameter vectors. The between-site heterogeneity associated with ψ_{it} is denoted by $\sigma_{\psi_t}^2$. Prior distributions p_{ψ_t} and $p_{\sigma_{\psi_t}^2}$ can be assigned for the unknown parameters ψ_t and $\sigma_{\psi_t}^2$. The joint posterior distribution for the two-stage approach is then

$$p(\boldsymbol{\psi}, \boldsymbol{\psi}_{1}, \dots, \boldsymbol{\psi}_{K}, \boldsymbol{\sigma}_{\boldsymbol{\psi}}^{2} \mid \hat{\boldsymbol{\psi}}_{i}, \operatorname{var}(\hat{\boldsymbol{\psi}}_{i})) \propto \underbrace{\prod_{i=1}^{K} \prod_{t=0}^{q-1} p(\hat{\boldsymbol{\psi}}_{it} \mid \boldsymbol{\psi}_{it}, \operatorname{var}(\hat{\boldsymbol{\psi}}_{it}))}_{\text{Likelihood}} \times \underbrace{\prod_{i=1}^{K} \prod_{t=0}^{q-1} p(\boldsymbol{\psi}_{it} \mid \boldsymbol{\psi}_{t}, \boldsymbol{\sigma}_{\boldsymbol{\psi}_{t}}^{2}) p(\boldsymbol{\psi}, \boldsymbol{\sigma}_{\boldsymbol{\psi}}^{2})}_{\text{Prior}},$$

where $\hat{\boldsymbol{\psi}}_{\boldsymbol{i}} = (\hat{\psi}_{i0}, \dots, \hat{\psi}_{i,q-1}), \operatorname{var}(\hat{\boldsymbol{\psi}}_{\boldsymbol{i}}) = (\operatorname{var}(\hat{\psi}_{i0}), \dots, \operatorname{var}(\hat{\psi}_{i,q-1})), \text{ and }$

$$p(\boldsymbol{\psi}, \boldsymbol{\sigma}_{\boldsymbol{\psi}}^{2}) = \prod_{t=0}^{q-1} p(\psi_{t}) \prod_{t=0}^{q-1} p(\sigma_{\psi_{t}}^{2})$$

With the full individual-level data, a one-stage model can be implemented:

$$Y_{ij} = \boldsymbol{\beta}_{i}^{\top} \boldsymbol{x}_{ij}^{(\boldsymbol{\beta})} + a_{ij} \boldsymbol{\psi}_{i}^{\top} \boldsymbol{x}_{ij}^{(\boldsymbol{\psi})} + \epsilon_{ij},$$
$$= \sum_{s=0}^{p-1} \beta_{is} x_{ijs}^{(\boldsymbol{\beta})} + a_{ij} \sum_{t=0}^{q-1} \psi_{it} x_{ijt}^{(\boldsymbol{\psi})} + \epsilon_{ij},$$

where the residual error ϵ_{ij} follows a normal distribution with mean 0 and within-site residual

variance σ_i^2 . The site-specific parameters β_{is} , ψ_{it} for $i = 1, \ldots, K$, $s = 0, \ldots, p - 1$, $t = 0, \ldots, q - 1$ satisfy

$$\beta_{is} \sim N(\beta_s, \sigma_{\beta_s}^2),$$

$$\psi_{it} \sim N(\psi_t, \sigma_{\psi_t}^2),$$
(B.1)

where $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})$ and $\boldsymbol{\psi} = (\psi_0, \dots, \psi_{q-1})$ are the common treatment-free and blip function parameters, respectively. We note that the distributional assumption (B.1) is slightly different from the distributional assumption

$$\boldsymbol{\psi}_{i} \sim \mathrm{MVN}(\boldsymbol{\psi}, \boldsymbol{\Sigma}_{\boldsymbol{\psi}}).$$
 (B.2)

As discussed in the main manuscript, in the two-stage approach, the site-specific treatmentfree function parameter estimates are ignored in the second stage. Therefore, only assumption (B.2) is required to pool the blip function parameter estimates, although assumption (B.1) is also reasonable. In the Bayesian framework, priors will be assigned to the unknown parameters β_s , ψ_t , σ_i^2 , $\sigma_{\beta_s}^2$, $\sigma_{\psi_t}^2$, $i = 1, \ldots, K, s = 0, \ldots, p - 1, t = 0, \ldots, q - 1$:

$$\begin{split} \beta_s &\sim p_{\beta_s}(\beta_s), \qquad \qquad \psi_t \sim p_{\psi_t}(\psi_t), \\ \sigma_i^2 &\sim p_{\sigma_i^2}(\sigma_i^2), \qquad \qquad \sigma_{\beta_s}^2 \sim p_{\sigma_{\beta_s}^2}(\sigma_{\beta_s}^2), \qquad \qquad \sigma_{\psi_t}^2 \sim p_{\sigma_{\psi_t}^2}(\sigma_{\psi_t}^2). \end{split}$$

Therefore, the joint posterior distribution for the one-stage approach is

$$p(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\beta}_{1}, \dots, \boldsymbol{\beta}_{K}, \boldsymbol{\psi}_{1}, \dots, \boldsymbol{\psi}_{K}, \boldsymbol{\sigma}^{2}, \boldsymbol{\sigma}^{2}_{\boldsymbol{\beta}}, \boldsymbol{\sigma}^{2}_{\boldsymbol{\psi}} \mid \boldsymbol{Y}_{1}, \dots, \boldsymbol{Y}_{K})$$

$$\propto \underbrace{\prod_{i=1}^{K} \prod_{j=1}^{n_{i}} p(Y_{ij} \mid \boldsymbol{\beta}_{i}, \boldsymbol{\psi}_{i}, \boldsymbol{\sigma}^{2}_{i})}_{\text{Likelihood}} \times \underbrace{\prod_{i=1}^{K} \prod_{s=0}^{p-1} p(\beta_{is} \mid \beta_{s}, \boldsymbol{\sigma}^{2}_{\boldsymbol{\beta}_{s}}) \prod_{i=1}^{K} \prod_{t=0}^{q-1} p(\psi_{it} \mid \boldsymbol{\psi}_{t}, \boldsymbol{\sigma}^{2}_{\boldsymbol{\psi}_{t}}) p(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\sigma}^{2}, \boldsymbol{\sigma}^{2}_{\boldsymbol{\beta}}, \boldsymbol{\sigma}^{2}_{\boldsymbol{\psi}})}_{\text{Prior}},$$

where $\mathbf{Y}_{i} = (Y_{i1}, \ldots, Y_{i,n_{i}}), \, \boldsymbol{\sigma}^{2} = (\sigma_{1}^{2}, \ldots, \sigma_{K}^{2}), \, \boldsymbol{\sigma}_{\beta}^{2} = (\sigma_{\beta_{0}}^{2}, \ldots, \sigma_{\beta_{p-1}}^{2}), \, \boldsymbol{\sigma}_{\psi}^{2} = (\sigma_{\psi_{0}}^{2}, \ldots, \sigma_{\psi_{q-1}}^{2}),$ and independent priors can be assigned to $\beta_{s}, \psi_{t}, \sigma_{i}^{2}, \sigma_{\beta_{s}}^{2}, \sigma_{\psi_{t}}^{2}$ such that $p(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\sigma}^{2}, \boldsymbol{\sigma}_{\beta}^{2}, \boldsymbol{\sigma}_{\psi}^{2}) = \prod_{s=0}^{p-1} p(\beta_{s}) \prod_{t=0}^{q-1} p(\psi_{t}) \prod_{i=1}^{K} p(\sigma_{i}^{2}) \prod_{s=0}^{p-1} p(\sigma_{\beta_{s}}^{2}) \prod_{t=0}^{q-1} p(\sigma_{\psi_{t}}^{2}).$ Thus, all parameters are estimated at once in the one-stage approach, while only blip function parameters and their related between-site variances are estimated separately in the two-stage approach. To see the similarity between the two approaches, we show that, under certain assumptions, $\prod_{j=1}^{n_{i}} p(Y_{ij} \mid \boldsymbol{\beta}_{i}, \boldsymbol{\psi}_{i}, \sigma_{i}^{2})$ and $p(\hat{\psi}_{it} \mid \psi_{it}, \operatorname{var}(\hat{\psi}_{it}))$ carry the same information of ψ_{it} . Define $Y_{ijt} = Y_{ij} - \sum_{s=0}^{p-1} \beta_{is} x_{ijs}^{(\beta)} - a_{ij} \sum_{t' \neq t} \psi_{it'} x_{ijt'}^{(\psi)}$ and $\tilde{Y}_{ijt} = Y_{ij} - \sum_{s=0}^{p-1} \hat{\beta}_{is} x_{ijs}^{(\beta)} - a_{ij} \sum_{t' \neq t} \hat{\psi}_{it'} x_{ijt'}^{(\psi)}$. Without loss of generality, assume that the focus now is only on $\psi_{t_{0}}, \psi_{it_{0}}, \text{ and } \sigma_{\psi_{t_{0}}}^{2}$ for some $t_{0} \in \{0, \ldots, q-1\}$, and other parameters (e.g., $\sigma_{i}^{2}, \beta_{is}, \psi_{it}, t \neq t_{0}$) are nuisance parameters. The likelihood in the one-stage approach is

$$\prod_{j=1}^{n_i} p(Y_{ij} \mid \boldsymbol{\beta}_i, \boldsymbol{\psi}_i, \sigma_i^2) \propto \exp\left\{-\frac{1}{2\sigma_i^2} \sum_{j=1}^{n_i} \left(Y_{ijt_0} - a_{ij}\psi_{it_0}x_{ijt_0}^{(\psi)}\right)^2\right\}$$
$$\propto \exp\left\{-\frac{1}{2\sigma_i^2} \left(\psi_{it_0}^2 \sum_{j=1}^{n_i} a_{ij}^2 (x_{ijt_0}^{(\psi)})^2 - 2\psi_{it_0} \sum_{j=1}^{n_i} a_{ij}x_{ijt_0}^{(\psi)}Y_{ijt_0}\right)\right\}.$$

The likelihood in the two-stage approach is

$$\begin{split} p(\hat{\psi}_{it_0} \mid \psi_{it_0}, \operatorname{var}(\hat{\psi}_{it_0})) &\propto \exp\left\{-\frac{(\hat{\psi}_{it_0} - \psi_{it_0})^2}{2\operatorname{var}(\hat{\psi}_{it_0})}\right\} \\ &\propto \exp\left\{-\frac{(\psi_{it_0}^2 - 2\hat{\psi}_{it_0}\psi_{it_0})}{2\operatorname{var}(\hat{\psi}_{it_0})}\right\} \\ &\propto \exp\left\{-\frac{\psi_{it_0}^2\sum_{j=1}^{n_i}a_{ij}^2(x_{ijt_0}^{(\psi)})^2 - 2\psi_{it_0}\sum_{j=1}^{n_i}a_{ij}x_{ijt_0}^{(\psi)}\tilde{Y}_{ijt_0}}{2\operatorname{var}(\hat{\psi}_{it_0})\sum_{j=1}^{n_i}a_{ij}^2(x_{ijt_0}^{(\psi)})^2}\right\} \end{split}$$

since the ordinary least squares (OLS) estimator is

$$\hat{\psi}_{it_0} = \frac{\sum_{j=1}^{n_i} \tilde{Y}_{ijt_0} a_{ij} x_{ijt_0}^{(\psi)}}{\sum_{j=1}^{n_i} a_{ij}^2 (x_{ijt_0}^{(\psi)})^2}$$

When $\hat{\beta}_{is} = \beta_{is}$, and $\hat{\psi}_{it} = \psi_{it}$ for $s = 0, \dots, p-1, t \neq t_0$, that is, β_{is} and ψ_{it} are esti-

mated with negligible error in the site-specific linear regression model, then $\tilde{Y}_{ijt_0} \approx Y_{ijt_0}$ and $\operatorname{var}(\hat{\psi}_{it_0}) = \sigma_i^2 / \{\sum_{j=1}^{n_i} a_{ij}^2 (x_{ijt_0}^{(\psi)})^2\}$. Thus, we have

$$p(\hat{\psi}_{it_0} \mid \psi_{it_0}, \operatorname{var}(\hat{\psi}_{it_0})) \propto \exp\left\{-\frac{1}{2\sigma_i^2} \left(\psi_{it_0}^2 \sum_{j=1}^{n_i} a_{ij}^2 (x_{ijt_0}^{(\psi)})^2 - 2\psi_{it_0} \sum_{j=1}^{n_i} a_{ij} x_{ijt_0}^{(\psi)} Y_{ijt_0}\right)\right\},$$

and $p(\hat{\psi}_{it_0} \mid \psi_{it_0}, \operatorname{var}(\hat{\psi}_{it_0}))$ contains the same information of ψ_{it_0} as $\prod_{j=1}^{n_i} p(Y_{ij} \mid \beta_i, \psi_i, \sigma_i^2)$. This applies to all sites under the assumption that β_{is} and ψ_{it} are estimated with negligible error (i.e., $\hat{\beta}_{is} = \beta_{is}, \hat{\psi}_{it} = \psi_{it}$) in the stage-one linear regression models. This assumption is plausible and approximately true for a moderate to large sample size, given the unbiasedness and consistency of the OLS estimators, if the model is correctly specified. Then, with the same common distribution for $\psi_{it_0}, i = 1, \ldots, K$, and the same priors for ψ_{t_0} and $\sigma_{\psi_{t_0}}^2$, the posterior distribution of $\psi_{t_0}, \psi_{(t_0)} = (\psi_{1,t_0}, \ldots, \psi_{K,t_0}), \sigma_{\psi_{t_0}}^2$ conditional on Y_i, β_i, σ^2 , $\psi_{i(-t_0)} = \psi_i / \{\psi_{i,t_0}\}$, for $i = 1, \ldots, K$, in the one-stage approach

$$p(\psi_{t_0}, \boldsymbol{\psi}_{(t_0)}, \sigma_{\psi_{t_0}}^2 \mid \boldsymbol{Y}_i, \boldsymbol{\beta}_i, \boldsymbol{\psi}_{i(-t_0)}, \boldsymbol{\sigma}^2)$$

$$\propto \prod_{i=1}^K \prod_{j=1}^{n_i} p(Y_{ij} \mid \boldsymbol{\beta}_i, \boldsymbol{\psi}_i, \sigma_i^2) \prod_{i=1}^K p(\psi_{it_0} \mid \psi_{t_0}, \sigma_{\psi_{t_0}}^2) p(\psi_{t_0}) p(\sigma_{\psi_{t_0}}^2),$$

is equivalent to the joint posterior distribution of $\psi_{t_0}, \psi_{(t_0)}, \sigma^2_{\psi_{t_0}}$ given $\hat{\psi}_{(t_0)} = (\hat{\psi}_{1,t_0}, \dots, \hat{\psi}_{K,t_0})$ and $\operatorname{var}(\hat{\psi}_{(t_0)}) = (\operatorname{var}(\hat{\psi}_{1,t_0}), \dots, \operatorname{var}(\hat{\psi}_{K,t_0}))$ in the two-stage approach:

$$p(\psi_{t_0}, \psi_{(t_0)}, \sigma_{\psi_{t_0}}^2 \mid \hat{\psi}_{(t_0)}, \operatorname{var}(\hat{\psi}_{(t_0)}))$$

$$\propto \prod_{i=1}^{K} p(\hat{\psi}_{it_0} \mid \psi_{it_0}, \operatorname{var}(\hat{\psi}_{it_0})) \prod_{i=1}^{K} p(\psi_{it_0} \mid \psi_{t_0}, \sigma_{\psi_{t_0}}^2) p(\psi_{t_0}) p(\sigma_{\psi_{t_0}}^2),$$

which leads to similar estimates in the two approaches.

B.2 Data sparsity: A second toy example

A simply toy example is provided in the main text. Here, in a second example, we assume one categorical covariate X consisting of three levels (i.e., p = q = 3). The true outcome model for an individual at site *i* is $E(Y \mid X) = \beta_{i0} + \beta_{i1}X_2 + \beta_{i2}X_3 + A(\psi_{i0} + \psi_{i1}X_2 + \psi_{i2}X_3).$ We choose the first category as the reference, and two indicators X_2 , X_3 are created for the second and third categories. Therefore, ψ_{i0} is the treatment effect for patients in the first category in site i; $\psi_{i0} + \psi_{i1}$ is the treatment effect for patients in the second category in site i; and $\psi_{10} + \psi_{i1}$ is the treatment effect for patients in the third category in site i. When (i) all patients in site i have a covariate value that is in the same category, or (ii) none take the second (or the third) category, but there are patients in the first and the third (or the second) categories, the situations are similar to the first example, and we do not duplicate the discussion. We consider a different scenario where none lie in the reference category, but both the second and third categories are represented in the samples. In this case, one of the last two categories will automatically become the "new" reference. Without loss of generality, assume the second category as the new reference. The site-specific outcome model then becomes $E(Y \mid X) = \gamma_{i0} + \gamma_{i2}X_3 + A(\xi_{i0} + \xi_{i2}X_3)$, where ξ_{i0} is the treatment effect for patients in the second category in site *i*, i.e., $\xi_{i0} = \psi_{i0} + \psi_{i1}$; ξ_{i2} is the difference in treatment effects for patients between the third and the second categories in site i, i.e., $\xi_{i2} = \psi_{i2} - \psi_{i1}$. Then the likelihood contribution of site *i* becomes $\hat{\gamma}_{i0} \sim N(\psi_{i0} + \psi_{i1}, \mathrm{sd}(\hat{\gamma}_{i0})^2)$, $\hat{\gamma}_{i2} \sim N(\psi_{i2} - \psi_{i1}, \mathrm{sd}(\hat{\gamma}_{i2})^2)$. Therefore, it is essential to examine the data in each site to detect any cases of sparsity in variable levels before incorporating the site-specific estimates into the model. For each site with sparse data, we update the likelihood contribution in the Bayesian hierarchical model based on the impact of data sparsity on the model parameter interpretation. Then, priors can be assigned to the common mean parameters and variance component parameters as is shown in the main text.

B.3 Simulation studies: ADEMP reporting

B.3.1 Aims

The aim of the simulation study is to evaluate ITR estimation for a continuous outcome when the individual-level data from multisite studies is protected from release via a twostage IPD meta-analysis, under assumptions concerning (1) the confounder sets across sites, (2) the strength of confounding, (3) the degree of heterogeneity across sites, and (4) the choice of prior distribution. Points (1) - (3) concerns the data-generating mechanisms, while (4) concerns the analysis model.

B.3.2 Data-generating mechanisms

In the simulation, we primarily focus on the binary treatment setting but include a reduced set of scenarios for the continuous treatment to illustrate the use of the proposed approach in a dosing setting. We also include a sparse data setting which mimics a particular, challenging feature of the International Warfarin Pharmacogenetics Consortium data: not all parameters can be estimated at all sites due to differences in populations across sites. Additionally, a small simulation is conducted to explore the use of shrinkage priors when a number of covariates are available, but only some are truly relevant for optimal treatment decisions. For all settings (except for the simulations with many covariates), K = 10 sites with an average sample size of n = 50 (small sample size) or 200 (large sample size) are assumed for all scenarios. The site-specific sample sizes vary between 0.6n and 1.4n. For simulations with shrinkage priors in the many covariates setting, only a large sample size is assumed.

Binary treatment setting

Two covariates X_1, X_2 are considered and their distributions vary across sites: for sites 3, 6, and 9, $X_1 \sim N(5, 1), X_2 \sim \text{Bernoulli}(0.5)$; for sites 1, 4, 7, and 10, $X_1 \sim 6\text{Beta}(4, 4) + 2, X_2 \sim$ Bernoulli(0.3); for sites 2, 5, and 8, $X_1 \sim U[2, 8], X_2 \sim \text{Bernoulli}(0.7)$. The treatment assignment A follows a Bernoulli distribution with the propensity score $P_i(\boldsymbol{x})$ at site *i* determined by $P_i(A = 1 | \boldsymbol{X} = \boldsymbol{x}) = \left[1 + e^{-(\alpha_{i0} + \alpha_{i1}x_1 + \alpha_{i2}x_2)}\right]^{-1}$, where $(\alpha_{i0}, \alpha_{i1}, \alpha_{i2})$ for different confounding scenarios are given in Table B.1. In scenarios 1 and 2, propensity score models are identical across sites, and the confounding effect can be either large (scenario 1) or small (scenario 2). In scenarios 3 and 4, site-specific propensity score models with the same set of confounders are assumed for each site, and two different confounding effects are also assumed. In scenarios 5 and 6, both propensity score model parameters and confounder sets are different across sites.

	Scenario 1	Scenario 2		
α_{i0}	0.1	0.01		
α_{i1}	0.1	0.01		
α_{i2}	0.1	0.01		
	Scenario 3	Scenario 4		
α_{i0}	U[0.06, 0.14]	U[0.006, 0.014]		
α_{i1}	U[0.06, 0.14]	U[0.006, 0.014]		
α_{i2}	U[0.06, 0.14]	U[0.006, 0.014]		
	Scenario 5	Scenario 6		
α_{i0}	U[0.3,0.7]	U[0.03, 0.07]		
α_{i1}	$\begin{cases} 0 & i = 1, 3, 5, 7, 9 \\ U[0, 0, 6, 0, 1, 4] & i = 2, 4, 6, 8, 10 \end{cases}$	$\begin{cases} 0 & i = 1, 3, 5, 7, 9 \\ U[0, 006, 0, 014] & i = 2, 4, 6, 8, 10 \end{cases}$		
α_{i2}	$\begin{cases} U[0.3, 0.7] & i = 2, 4, 6, 8, 10 \\ U[0.3, 0.7] & i = 1, 3, 5, 7, 9 \\ 0 & i = 2, 4, 6, 8, 10 \end{cases}$	$\begin{cases} U[0.005, 0.014] & i = 2, 4, 6, 8, 10 \\ U[0.03, 0.07] & i = 1, 3, 5, 7, 9 \\ 0 & i = 2, 4, 6, 8, 10 \end{cases}$		

Table B.1: Parameters in the propensity score model for binary treatment simulations in different scenarios

Suppressing the individual-specific subscript, the continuous outcome for an individual at site *i* is generated by $Y_i = \beta_{i0} + \beta_{i1}x_1 + \beta_{i2}x_2 + a(\psi_{i0} + \psi_{i1}x_1) + \epsilon$, where the random error ϵ follows a normal distribution with mean zero and residual variance $\sigma_{\epsilon}^2 = 0.25$. For the site-specific parameters $\boldsymbol{\theta}_i = (\beta_{i0}, \beta_{i1}, \beta_{i2}, \psi_{i0}, \psi_{i1})$, we consider three different scenarios:

- common effect: $\theta_1 = \theta_2 = \ldots = \theta_{10} = \theta$ and $\theta = (\beta_0, \beta_1, \beta_2, \psi_0, \psi_1)$ is the common

population parameter;

- common rule: $\beta_{it} \sim N(\beta_t, \sigma_B^2), \ \psi_{i1} \sim N(\psi_1, \sigma_B^2), \text{ for } t = 0, 1, 2, \ i = 1, \dots, 10 \text{ and}$ $\psi_{10}/\psi_{11} = \psi_{20}/\psi_{21} = \dots = \psi_{10,0}/\psi_{10,1} = -5, \text{ where the between-site variance } \sigma_B^2 \text{ is}$ derived from heterogeneity level $I^2 = \frac{\sigma_B^2}{\sigma_B^2 + \sigma_\epsilon^2} \in \{0.1, 0.2, 0.3\};$
- varying effects: $\theta_i \sim \text{MVN}(\theta, \Sigma_{\theta})$, where Σ_{θ} is a 5 × 5 diagonal matrix where the between-site variance is derived from I^2 as in the common rule setting.

In all three scenarios, the common treatment-free parameters are $\beta_0 = 4$, $\beta_1 = 1$, $\beta_2 = 1$ and the common blip parameters are $\psi_0 = 2.5$, $\psi_1 = -0.5$ such that the common optimal ITR is given by $d^{\text{opt}}(\boldsymbol{x}) = I(\psi_0 + \psi_1 x_1 > 0) = I(x_1 < 5)$. The common effect setting assumes that all site-specific parameters are equal to the common population parameters as in the simulation studies in Danieli and Moodie (2022); Moodie et al. (2022). No heterogeneity exists in the site-specific blip parameters ψ_{i0} and ψ_{i1} and the site-specific optimal ITRs are $d_i^{\text{opt}}(\boldsymbol{x}) = I(\psi_{i0} + \psi_{i1}x_1 > 0)$. The varying effects setting assumes a common multivariate normal distribution for the site-specific parameters. The two blip parameters ψ_{i0} and ψ_{i1} are freely varying across sites. Therefore, heterogeneity exists in (ψ_{i0}, ψ_{i1}) and $d_i^{\text{opt}}(\boldsymbol{x})$. The common rule setting considers heterogeneity scenarios that can be viewed as intermediate between common effect and varying effects; the blip parameters ψ_{i0}, ψ_{i1} are varying across sites, however, the site-specific optimal ITRs $d_i^{\text{opt}}(\boldsymbol{x})$ are fixed by restricting the ratio ψ_{i0}/ψ_{i1} to be identical across sites. In this setting, heterogeneity only exists in the blip parameters but not the site-specific optimal ITRs.

Continuous treatment setting

For the continuous treatment setting, the same covariates X_1 , X_2 are generated in the same way as the binary treatment setting. The treatment $A \sim N(X_1, 1)$. The outcome for an individual at site *i* is generated by $Y_i = \beta_{i0} + \beta_{i1}x_1 + \beta_{i2}x_2 + a(\psi_{i0} + \psi_{i1}a + \psi_{i2}x_1) + \epsilon$, where the random error ϵ follows a normal distribution with mean zero and residual variance $\sigma_{\epsilon}^2 = 0.25$. Two different settings are considered for the site-specific parameters $\boldsymbol{\theta}_{i} = (\beta_{i0}, \beta_{i1}, \beta_{i2}, \psi_{i0}, \psi_{i1}, \psi_{i2})$:

- common effect: $\theta_1 = \theta_2 = \ldots = \theta_{10} = \theta$ and $\theta = (\beta_0, \beta_1, \beta_2, \psi_0, \psi_1, \psi_2)$ is the common population parameter;
- varying effects: $\theta_i \sim \text{MVN}(\theta, \Sigma_{\theta})$, where Σ_{θ} is a 6 × 6 diagonal matrix where the between-site variance is obtained from the heterogeneity level $I^2 = \frac{\sigma_B^2}{\sigma_B^2 + \sigma_\epsilon^2} \in \{0.1, 0.2, 0.3\}.$

In both settings, the common treatment-free parameters are $\beta_0 = 4$, $\beta_1 = 1$, $\beta_2 = 1$, and the common blip parameters are $\psi_0 = 1$, $\psi_1 = -2$, $\psi_2 = 1$. The common optimal ITR is $d^{\text{opt}}(\boldsymbol{x}) = \operatorname{argmax}_a(-2a^2 + a + ax_1) = (1 + x_1)/4.$

Sparse data setting

As discussed, it is possible in multisite studies that the site-specific parameters cannot be estimated due to an insufficient number of patients with a given set of characteristics. To show how the proposed method deals with this scenario, a small simulation focusing on a sparse data setting is performed. For simplicity, a binary treatment $A \sim \text{Bernoulli}(0.5)$ is considered. A binary covariate X_1 and a categorical covariate X_2 consisting of three levels are assumed and their distributions vary across sites: for sites 3, 6, and 9, $X_1 = 1, X_2 \sim$ Multinomial(1; 0, 0.5, 0.5); for sites 1, 4, 7, and 10, $X_1 = 0, X_2 \sim \text{Multinomial}(1; 0.5, 0, 0.5);$ for sites 2, 5, and 8, $X_1 \sim \text{Bernoulli}(0.5), X_2 \sim \text{Multinomial}(1; 1/3, 1/3)$. The continuous outcome for an individual at site *i* is generated by $Y_i = \beta_{i0} + \beta_{i1}x_1 + \beta_{i2}x_{2,2} + \beta_{i3}x_{2,3} + a(\psi_{i0} + \psi_{i1}x_1 + \psi_{i2}x_{2,2} + \psi_{i3}x_{2,3}) + \epsilon$, where the random error ϵ follows a normal distribution with mean 0 and residual variance $\sigma_{\epsilon}^2 = 0.25$. For X_2 , the first category is assumed as the reference level, and two binary indicators $X_{2,2}$ and $X_{2,3}$ are created for the second and third categories of X_2 . Two different settings are considered for the site-specific parameters $\boldsymbol{\theta}_i = (\beta_{i0}, \beta_{i1}, \beta_{i2}, \beta_{i3}, \psi_{i0}, \psi_{i1}, \psi_{i2}, \psi_{i3})$:

- common effect: $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \ldots = \boldsymbol{\theta}_{10} = \boldsymbol{\theta}$ and $\boldsymbol{\theta} = (\beta_0, \beta_1, \beta_2, \beta_3, \psi_0, \psi_1, \psi_2, \psi_3)$ is the common population parameter;
- varying effects: $\theta_i \sim \text{MVN}(\theta, \Sigma_{\theta})$, where Σ_{θ} is a 8 × 8 diagonal matrix where the between-site variance is obtained from the heterogeneity level $I^2 = \frac{\sigma_B^2}{\sigma_B^2 + \sigma_{\epsilon}^2} = 0.1, 0.2, 0.3.$

In both settings, the common treatment-free parameters are $\beta_0 = 4$, $\beta_1 = 1$, $\beta_2 = 1$, $\beta_3 = -1$, and the common blip parameters are $\psi_0 = 1$, $\psi_1 = 1$, $\psi_2 = -2.5$, $\psi_3 = 2$. The common optimal ITR is $d^{\text{opt}}(\boldsymbol{x}) = I(\psi_0 + \psi_1 x_1 + \psi_2 x_{2,2} + \psi_3 x_{2,3} > 0)$. The model details in a sparse data setting is described in Appendix B.4.

Many covariates setting

We consider two scenarios: a total of either 10 or 20 covariates is collected, but only three covariates, X_1 , X_2 , and X_3 are related to optimal treatment assignment. The covariates X_1 and X_2 are generated in the same way as in the binary treatment setting. We generate X_3 by the following distribution: for sites 3, 6, and 9, $X_3 \sim \text{Exponential}(1)$; for sites 1, 4, 7, and 10, $X_3 \sim \text{Exponential}(1.7)$; for sites 2, 5, and 7, $X_3 \sim \text{Exponential}(0.7)$. Other covariates are generated by $X_j \sim N(0,1)$, for $j \geq 4$. For simplicity, we assume a binary treatment $A \sim \text{Bernoulli}(0.5)$. The continuous outcome for an individual at site *i* is generated by $Y_i = \beta_{i0} + \sum_{s=1}^p \beta_{is} x_{is} + A(\psi_{i0} + \sum_{t=1}^p \psi_{it} x_{it}) + \epsilon$, for p = 10 or 20, and $\epsilon \sim N(0, 0.25)$. The site-specific parameters β_{is} and ψ_{it} are generated under the assumption of varying effects: $\theta_i = (\beta_{i0}, \ldots, \beta_{ip}, \psi_{i0}, \ldots, \psi_{ip}) \sim \text{MVN}(\theta, \Sigma_{\theta})$, where $\theta = (\beta_0, \ldots, \beta_p, \psi_0, \ldots, \psi_p)$ is a vector of common parameters and $\beta_0 = 4$, $\beta_s = 1$ for $s \geq 2$, $\psi_0 = 2.5$, $\psi_1 = -0.5$, $\psi_2 = 2$, $\psi_3 = -1$, and $\psi_t = 0$ for $t \geq 4$. The common optimal ITR is thus $d^{\text{opt}}(x) = I(\psi_0 + \psi_1 x_1 + \psi_2 x_2 + \psi_3 x_2 > 0)$. The between-site variance in the diagonal variance-covariance matrix Σ_{θ} is obtained from the heterogeneity level $I^2 = \frac{\sigma_B^2}{\sigma_B^2 + \sigma_s^2} = 0.1$.

B.3.3 Estimands, methods, and performance metrics

The estimands of interest are the common blip parameters which fully characterize the optimal ITR. All analyses rely on the two-stage IPD meta-analysis, using linear regression as the stage-one model. For all scenarios, we use a Bayesian hierarchical model for the second stage. For the mean parameters in all settings other than the many covariates setting, we use a normal prior with mean 0 and variance 10,000. In the setting of many covariates, we assign the same normal prior to the common main treatment effect parameter but a horseshoe prior to all treatment-covariate interactions, selecting only those whose 95% posterior credible intervals do not include zero. For variance component parameters, three priors with varying levels of informativeness are considered: half-Cauchy priors with location 0 and scale parameters 1, 10, or 100. However, in the many covariates setting, only a half-Cauchy (0,1) prior is used for the variance component parameters.

For all scenarios, 2000 iterations are performed. Measures of performance used to assess the ITR estimation are: (i) the relative bias of estimators of the blip parameters, which represents the difference between the mean of the estimates and the true value, divided by the latter, (ii) the standard deviation of the estimators, (iii) the difference between the value function (dVF) under the true optimal ITR and the value function under the estimated optimal ITR, where the value function with respect to an ITR is the expected outcome if all patients in a population (in our simulation, it is a new cohort of patients of size n = 100,000) were treated according to the ITR, and (iv) the empirical standard deviation of the value function difference when the estimated treatment rule was applied to the same population. For the many covariates setting, these measures are assessed over: (1) a full set of 2000 iterations, and (2) a subset of iterations where the non-zero treatment-covariate interactions are correctly selected. The proportion of selection, calculated as the number of times the covariate is selected divided by the total number of simulation iterations, is also measured. The results are compared with results obtained from a one-stage approach based on the full individual-level data.

B.4 Model details in a sparse data setting in simulation

To illustrate how the proposed two-stage approach can deal with data sparsity, a small simulation is performed and a binary treatment $A \in \{0, 1\}$ is considered for simplicity. A binary covariate X_1 and a categorical covariate X_2 consisting of three levels are assumed: for the *i*-th site,

$$X_{1} \sim \begin{cases} 1 & i = 3s \\ 0 & i = 3s + 1 \\ \text{Bernoulli}(0.5) & i = 3s + 2 \end{cases} \qquad \begin{cases} \text{Multinomial}(1; 0, 0.5, 0.5) & i = 3s \\ \text{Multinomial}(1; 0.5, 0, 0.5) & i = 3s + 1 \\ \text{Multinomial}(1; 1/3, 1/3, 1/3) & i = 3s + 2 \end{cases}$$

where s = 0, ..., 3. The continuous outcome for an individual at site *i* is generated by

$$Y_{i} = \beta_{i0} + \beta_{i1}x_{1} + \beta_{i2}x_{2,2} + \beta_{i3}x_{2,3} + a(\psi_{i0} + \psi_{i1}x_{1} + \psi_{i2}x_{2,2} + \psi_{i3}x_{2,3}) + \epsilon,$$
(B.3)

where the random error ϵ follows a normal distribution with mean 0 and residual variance $\sigma_{\epsilon}^2 = 0.25$. For X_2 , the first category is assumed as the reference level, and two binary indicators $X_{2,2}$ and $X_{2,3}$ are created for the second and third categories of X_2 . As discussed in the main manuscript, both common effect and varying effects settings are explored for the site-specific parameters $\boldsymbol{\theta}_i = (\beta_{i0}, \beta_{i1}, \beta_{i2}, \beta_{i3}, \psi_{i0}, \psi_{i1}, \psi_{i2}, \psi_{i3})$.

Due to the data-generating mechanism, the implied (correctly-specified) linear regression models at the first stage for sites i = 3s, s = 1, 2, 3, are

$$E(Y_i \mid \boldsymbol{x}, a) = \gamma_{i0} + \gamma_{i2} x_{2,2} + a(\xi_{i0} + \xi_{i2} x_{2,2}).$$
(B.4)

Since no patients within sites i = 3s are in the reference category of X_2 , in site-specific analyses, one of $X_{2,2}$ and $X_{2,3}$ will be chosen as the new reference category, and its main effect as well as the interaction effect with the treatment in (B.3) cannot be estimated. Here, without loss of generality, we assume that $X_{2,3}$ is the new reference category. Then, the site-specific main effect estimator of $X_{2,2}$, $\hat{\gamma}_{i2}$, and its interaction effect estimator, $\hat{\xi}_{i2}$, in (B.4) will be biased for β_{i2} and ψ_{i2} respectively, as $\gamma_{i2} = \beta_{i2} - \beta_{i3}$ and $\xi_{i2} = \psi_{i2} - \psi_{i3}$. In addition, since $X_1 = 1$ for all patients within sites i = 3s, the effect of X_1 (i.e., β_{i1}) and its interaction with treatment (i.e., ψ_{i1}) cannot be estimated. The site-specific intercept and main treatment effect estimators in (B.4) (i.e., $\hat{\gamma}_{i0}$ and $\hat{\xi}_{i0}$) are biased for the original parameters β_{i0} and ψ_{i0} in (B.3), as $\gamma_{i0} = \beta_{i0} + \beta_{i1} + \beta_{i3}$ and $\xi_{i0} = \psi_{i0} + \psi_{i1} + \psi_{i3}$. Therefore, to recover the original parameters, the likelihood model for these sites at the second stage should be reparametrized as

$$\hat{\xi}_{i0} \sim N \left(\psi_{i0} + \psi_{i1} + \psi_{i3}, \operatorname{sd}(\hat{\xi}_{i0})^2 \right), \qquad \hat{\xi}_{i2} \sim N \left(\psi_{i2} - \psi_{i3}, \operatorname{sd}(\hat{\xi}_{i2})^2 \right),$$

$$\psi_{i0} \sim N \left(\psi_0, \sigma_{\psi_0}^2 \right), \qquad \psi_{i1} \sim N \left(\psi_1, \sigma_{\psi_1}^2 \right),$$

$$\psi_{i2} \sim N \left(\psi_2, \sigma_{\psi_2}^2 \right), \qquad \psi_{i3} \sim N \left(\psi_3, \sigma_{\psi_3}^2 \right).$$

For sites i = 3s + 1, s = 0, ..., 3, the site-specific linear regression models are

$$E(Y_i \mid \boldsymbol{x}, a) = \gamma_{i0} + \gamma_{i3} x_{2,3} + a(\xi_{i0} + \xi_{i3} x_{2,3}).$$
(B.5)

In these sites, $X_1 = 0$ for all patients, and no patients are in the second category of X_2 . Therefore, $\beta_{i1}, \beta_{i2}, \psi_{i1}, \psi_{i2}$ cannot be estimated. However, the estimators $\hat{\gamma}_{i0}, \hat{\gamma}_{i3}, \hat{\xi}_{i0}$ and $\hat{\xi}_{i3}$ in (B.5) are still be consistent for the parameters $\beta_{i0}, \beta_{i3}, \psi_{i0}$ and ψ_{i3} in (B.3), as there are patients with $X_1 = 0$ and $X_{2,2} = X_{2,3} = 0$. Thus, the likelihood model for these sites at the second stage will be

$$\hat{\xi}_{i0} \sim N\left(\psi_{i0}, \mathrm{sd}(\hat{\xi}_{i0})^2\right), \qquad \hat{\xi}_{i3} \sim N\left(\psi_{i3}, \mathrm{sd}(\hat{\xi}_{i3})^2\right),$$

$$\psi_{i0} \sim N\left(\psi_0, \sigma_{\psi_0}^2\right), \qquad \psi_{i3} \sim N\left(\psi_3, \sigma_{\psi_3}^2\right).$$

For sites i = 3s + 2, s = 0, 1, 2, all levels of all covariates are represented and thus all parameters are estimable. The regression estimators in

$$E(Y_i \mid \boldsymbol{x}, a) = \gamma_{i0} + \gamma_{i1}x_1 + \gamma_{i2}x_{2,2} + \gamma_{i3}x_{2,3} + a(\xi_{i0} + \xi_{i1}x_1 + \xi_{i2}x_{2,2} + \xi_{i3}x_{2,3})$$

will be consistent for the corresponding parameters in (B.3). The likelihood model at the second stage will be

$$\begin{aligned} \hat{\xi}_{i0} &\sim N\left(\psi_{i0}, \operatorname{sd}(\hat{\xi}_{i0})^{2}\right), & \hat{\xi}_{i1} &\sim N\left(\psi_{i1}, \operatorname{sd}(\hat{\xi}_{i1})^{2}\right), \\ \hat{\xi}_{i2} &\sim N\left(\psi_{i2}, \operatorname{sd}(\hat{\xi}_{i2})^{2}\right), & \hat{\xi}_{i3} &\sim N\left(\psi_{i3}, \operatorname{sd}(\hat{\xi}_{i3})^{2}\right), \\ \psi_{i0} &\sim N\left(\psi_{0}, \sigma_{\psi_{0}}^{2}\right), & \psi_{i1} &\sim N\left(\psi_{1}, \sigma_{\psi_{1}}^{2}\right), \\ \psi_{i2} &\sim N\left(\psi_{2}, \sigma_{\psi_{2}}^{2}\right), & \psi_{i3} &\sim N\left(\psi_{3}, \sigma_{\psi_{3}}^{2}\right). \end{aligned}$$

B.5 Additional simulation results

This section presents additional results from the simulations carried out. Figures B.1 and B.2 shows estimates of ψ_0 and the dVF under half-Cauchy (0,10) and half-Cauchy (0,100) priors in the binary treatment setting with small sample size, and Figure B.3 shows simulation results for ψ_1 . Figures B.4 - B.6 present simulation results for the large sample size and binary treatment setting, including the estimation of (ψ_0, ψ_1) and the dVF. Figures B.7 and B.8 present the simulation results in the continuous treatment setting with both small and large sample sizes, including the estimation of (ψ_0, ψ_1, ψ_2) and the dVF. Figures B.9 and

B.10 present simulation results of $(\psi_0, \psi_1, \psi_2, \psi_3)$ in the sparse data setting with the large sample size.

Similar patterns to those in the main manuscript are observed. The results are not sensitive to the different prior choices. The one- and two-stage approaches give similar results. Both provide unbiased estimations of blip function parameters. However, when the heterogeneity is small, the variability of estimators in the two-stage approach is larger than that in the onestage approach, as the one stage approach is able to borrow information across sites when the heterogeneity is small. When the heterogeneity is large, this difference is smaller, and variability in both approaches increase compared to that with small heterogeneity. In both binary and continuous treatment settings, the dVF increases with increasing heterogeneity, suggesting a worse optimal ITR estimation. In the sparse data setting, the dVF is zero or close to zero in all scenarios, regardless of the heterogeneity levels. We only consider binary covariates and binary treatment in the sparse data setting. The indicator function $I(\psi_0 + \psi_1 x_1 + \psi_2 x_{2,2} + \psi_3 x_{2,3} > 0)$ is less sensitive to the errors in the blip function parameter estimation compared with the optimal ITR in the setting of continuous treatment/covariates. Therefore, even if the parameter estimators are more varied with large heterogeneity, the dVF does not change much with different heterogeneity levels, and we (almost) obtain the true optimal ITR in all cases. Also, in all settings, with a larger sample size, we obtain a more precise optimal ITR estimation, as the variability of blip parameter estimation and the dVF are smaller.









Simulation results for the small sample size and the binary treatment setting. Performance of the methods is under different confounding scenarios, heterogeneity levels $(I^2 = 0.1, 0.2, 0.3)$, and prior choices. The triangles represent the assessed over 2000 iterations. The difference in the value function (dVF) between the true and estimated optimal ITR is shown mean of the estimates in each case. Figure B.2:





Simulation results for the small sample size and the binary treatment setting. Performance of the methods is The dashed line of ψ_1 are shown under different confounding scenarios, heterogeneity case. The triangles represent the mean of the estimates in each Estimates (posterior means) levels $(I^2 = 0.1, 0.2, 0.3)$, and prior choices. assessed over 2000 iterations. shows the true value of -0.5. Figure B.3:







Performance of the methods is The dashed line of ψ_1 are shown under different confounding scenarios, heterogeneity case. The triangles represent the mean of the estimates in each Simulation results for the large sample size and the binary treatment setting. Estimates (posterior means) = 0.1, 0.2, 0.3, and prior choices. assessed over 2000 iterations. shows the true value of -0.5. Figure B.5: levels (I^2)



Simulation results for the large sample size and the binary treatment setting. Performance of the methods is 0.1, 0.2, 0.3, and prior choices. The triangles represent the assessed over 2000 iterations. The difference in the value function (dVF) between the true and estimated optimal ITR is shown under different confounding scenarios, heterogeneity levels ($I^2 =$ mean of the estimates in each case. Figure B.6:



Model 🖻 One-stage 🖻 Two-stage

Figure B.7: Simulation results for the small sample size and the continuous treatment setting. Performance of the methods is assessed over 2000 iterations. Estimates (posterior means) of ψ_0 , ψ_1 , ψ_2 , and the difference in the value function (dVF) between the true and estimated optimal ITR are shown under different heterogeneity levels ($I^2 = 0.1, 0.2, 0.3$), and prior choices. The triangles represent the mean of the estimates in each case. The dashed lines show the true values of ψ_0 , ψ_1 , ψ_2 .



Model 🖻 One-stage 🖻 Two-stage

Figure B.8: Simulation results for the large sample size and the continuous treatment setting. Performance of the methods is assessed over 2000 iterations. Estimates (posterior means) of ψ_0, ψ_1, ψ_2 , and the difference in the value function (dVF) between the true and estimated optimal ITR are shown under different heterogeneity levels $(I^2 = 0.1, 0.2, 0.3)$ and prior choices. The triangles represent the mean of the estimates in each case. The dashed lines show the true values of ψ_0 , ψ_1 , ψ_2 .



Model 🖻 One-stage 🖻 Two-stage

Figure B.9: Simulation results for the small sample size and the sparse data setting. Performance of the methods is assessed over 2000 iterations. Estimates (posterior means) of ψ_0 , ψ_1 , ψ_2 , ψ_3 are shown under different heterogeneity levels ($I^2 = 0.1, 0.2, 0.3$) and half-Cauchy (0,10) and half-Cauchy (0,100) priors. The triangles represent the mean of the estimates in each case. The dashed lines show the true values of ψ_0 , ψ_1 , ψ_2 , ψ_3 .



Model 🖻 One-stage 🔳 Two-stage

Figure B.10: Simulation results for the large sample size and the sparse data setting. Performance of the methods is assessed over 2000 iterations. Estimates (posterior means) of ψ_0 , ψ_1 , ψ_2 , ψ_3 are shown under different heterogeneity levels ($I^2 = 0.1, 0.2, 0.3$) and prior choices. The triangles represent the mean of the estimates in each case. The dashed lines show the true values of ψ_0 , ψ_1 , ψ_2 , ψ_3 .

B.6 Analysis of Warfarin data

B.6.1 Visual inspection of the overlap assumption

The overlap assessment was conducted by discretizing the continuous dose into four ordinal dose groups based on the minimum, 25th, 50th, and 75th quantiles, and maximum of the observed doses. Then, following that suggested by Li and Li (2019), a proportional odds logistic model including all potential confounders was used to estimate the generalized propensity score (Imbens, 2000). The distribution of the generalized propensity score is shown in Figure B.11. A moderate lack of overlap was observed, particularly within the dose groups [4.5, 22.8] and [42, 95].



Figure B.11: Distribution of the generalized propensity score by dose group in the Warfarin analysis.

B.6.2 Details of the models in the Warfarin analysis

As discussed in the main paper, the blip functions are quadratic in Warfarin dose. The linear regression model for site i can be explicitly stated as:

$$\begin{split} E(Y_{i} \mid \boldsymbol{x}, a) &= \beta_{i0} + \beta_{i1} \text{Age} + \beta_{i2} \text{Amiodarone} + \beta_{i3} \text{Female} \\ &+ \beta_{i4} \text{VKORC1}(\text{AG}) + \beta_{i5} \text{VKORC1}(\text{AA}) + \beta_{i6} \text{CYP2C9}(12) \\ &+ \beta_{i7} \text{CYP2C9}(\text{other}) + \beta_{i8} \text{Weight} + \beta_{i9} \text{Height} + \beta_{i,10} \text{Non-White} \\ &+ a \times \left\{ \psi_{i0}^{(1)} + \psi_{i1}^{(1)} \text{Age} + \psi_{i2}^{(1)} \text{Amiodarone} + \psi_{i3}^{(1)} \text{Female} + \psi_{i4}^{(1)} \text{VKORC1}(\text{AG}) \\ &+ \psi_{i5}^{(1)} \text{VKORC1}(\text{AA}) + \psi_{i6}^{(1)} \text{CYP2C9}(12) + \psi_{i7}^{(1)} \text{CYP2C9}(\text{other}) \right\} \\ &+ a^{2} \times \left\{ \psi_{i0}^{(2)} + \psi_{i1}^{(2)} \text{Age} + \psi_{i2}^{(2)} \text{Amiodarone} + \psi_{i3}^{(2)} \text{Female} + \psi_{i4}^{(2)} \text{VKORC1}(\text{AG}) \\ &+ \psi_{i5}^{(2)} \text{VKORC1}(\text{AA}) + \psi_{i6}^{(2)} \text{CYP2C9}(12) + \psi_{i7}^{(2)} \text{CYP2C9}(\text{other}) \right\}. \end{split}$$

We assume that $\psi_{it}^{(u)} \sim N(\psi_t^{(u)}, (\sigma_t^{(u)})^2)$, t = 0, ..., 8, u = 1, 2. The parameters of interests are the common blip function parameters $\psi^{(1)} = (\psi_0^{(1)}, ..., \psi_8^{(1)})$ and $\psi^{(2)} = (\psi_0^{(2)}, ..., \psi_8^{(2)})$ which fully characterize the optimal Warfarin dosing. The unknown between-site variability associated with $\psi_{it}^{(u)}$ is denoted by $(\sigma_t^{(u)})^2$.

Tables B.2 and B.3 show site-specific blip function parameter estimates obtained from the stage-one (frequentist) linear regression models and the associated standard deviations for the Warfarin data. Due to data sparsity, some site-specific blip function parameters cannot be estimated in some sites. As demonstrated in the simulation, we need to modify the proposed two-stage model.

For i = 3, the linear regression model in the first stage is

$$E(Y_{i}|\boldsymbol{x}, a) = \gamma_{i0} + \gamma_{i1} \text{Age} + \gamma_{i2} \text{Amiodarone} + \gamma_{i3} \text{Female} + \gamma_{i4} \text{VKORC1}(\text{AG}) + \gamma_{i5} \text{VKORC1}(\text{AA}) + \gamma_{i7} \text{CYP2C9}(\text{other}) + \gamma_{i8} \text{Weight} + \gamma_{i9} \text{Height} + a \times \left\{ \xi_{i0}^{(1)} + \xi_{i1}^{(1)} \text{Age} + \xi_{i2}^{(1)} \text{Amiodarone} + \xi_{i3}^{(1)} \text{Female} + \xi_{i4}^{(1)} \text{VKORC1}(\text{AG}) + \xi_{i5}^{(1)} \text{VKORC1}(\text{AA}) + \xi_{i7}^{(1)} \text{CYP2C9}(\text{other}) \right\} + a^{2} \times \left\{ \xi_{i0}^{(2)} + \xi_{i1}^{(2)} \text{Age} + \xi_{i2}^{(2)} \text{Amiodarone} + \xi_{i3}^{(2)} \text{Female} + \xi_{i4}^{(2)} \text{VKORC1}(\text{AG}) + \xi_{i5}^{(2)} \text{VKORC1}(\text{AA}) + \xi_{i7}^{(2)} \text{CYP2C9}(\text{other}) \right\}.$$
(B.7)

The parameters in equation (B.7) satisfy

$$\gamma_{il} = \beta_{il},$$
 $\xi_{it}^{(1)} = \psi_{it}^{(1)},$ $\xi_{it}^{(2)} = \psi_{it}^{(2)},$
 $\gamma_{i0} = \beta_{i0} + \beta_{i,10}$

for $l \neq 0, 6, 10$ and $t \neq 6$, since all patients in Site 3 are non-White and none carry CYP2C9 genotype 12. The modified likelihood model in the second stage is then

$$\hat{\xi}_{it}^{(u)} \sim N\left(\psi_{it}^{(u)}, \mathrm{sd}(\hat{\xi}_{it}^{(u)})^2\right), \psi_{it}^{(u)} \sim N\left(\psi_t^{(u)}, (\sigma_t^{(u)})^2\right),$$

for $t \neq 6$ and u = 1, 2.

For i = 6, the linear regression model in the first stage is

$$E(Y_{i}|\boldsymbol{x}, a) = \gamma_{i0} + \gamma_{i1} \text{Age} + \gamma_{i2} \text{Amiodarone} + \gamma_{i3} \text{Female} + \gamma_{i4} \text{VKORC1}(\text{AG}) + \gamma_{i7} \text{CYP2C9}(\text{other}) + \gamma_{i8} \text{Weight} + \gamma_{i9} \text{Height} + a \times \left\{ \xi_{i0}^{(1)} + \xi_{i1}^{(1)} \text{Age} + \xi_{i2}^{(1)} \text{Amiodarone} + \xi_{i3}^{(1)} \text{Female} + \xi_{i4}^{(1)} \text{VKORC1}(\text{AG}) + \xi_{i7}^{(1)} \text{CYP2C9}(\text{other}) \right\}$$
(B.8)
$$+ a^{2} \times \left\{ \xi_{i0}^{(2)} + \xi_{i1}^{(2)} \text{Age} + \xi_{i2}^{(2)} \text{Amiodarone} + \xi_{i3}^{(2)} \text{Female} + \xi_{i4}^{(2)} \text{VKORC1}(\text{AG}) + \xi_{i7}^{(2)} \text{CYP2C9}(\text{other}) \right\}.$$

The parameters in equation (B.8) satisfy

$$\begin{aligned} \gamma_{i0} &= \beta_{i0} + \beta_{i5} + \beta_{i,10}, & \xi_{i0}^{(1)} &= \psi_{i0}^{(1)} + \psi_{i5}^{(1)}, & \xi_{i0}^{(2)} &= \psi_{i0}^{(2)} + \psi_{i5}^{(2)}, \\ \gamma_{i4} &= \beta_{i4} - \beta_{i5}, & \xi_{i4}^{(1)} &= \psi_{i4}^{(1)} - \psi_{i5}^{(1)}, & \xi_{i4}^{(2)} &= \psi_{i4}^{(2)} - \psi_{i5}^{(2)}, \\ \gamma_{il} &= \beta_{il}, & \xi_{it}^{(1)} &= \psi_{it}^{(1)}, & \xi_{it}^{(2)} &= \psi_{it}^{(2)}, \end{aligned}$$

for $l \neq 0, 4, 5, 6, 10$ and $t \neq 0, 4, 5, 6$, since all patients in Site 6 are non-White, and none carry VKORC1 genotype GG or CYP2C9 genotype 12. The modified likelihood model is

$$\begin{split} \hat{\xi}_{i0}^{(u)} &\sim N\left(\psi_{i0}^{(u)} + \psi_{i5}^{(u)}, \mathrm{sd}(\hat{\xi}_{i0}^{(u)})^2\right), \\ \hat{\xi}_{i4}^{(u)} &\sim N\left(\psi_{i4}^{(u)} - \psi_{i5}^{(u)}, \mathrm{sd}(\hat{\xi}_{i4}^{(u)})^2\right), \\ \hat{\xi}_{it}^{(u)} &\sim N\left(\psi_{it}^{(u)}, \mathrm{sd}(\hat{\xi}_{it}^{(u)})^2\right), \\ \psi_{il}^{(u)} &\sim N\left(\psi_{l}^{(u)}, (\sigma_{l}^{(u)})^2\right), \end{split}$$

for $t = 1, 2, 3, 7, l \neq 6, u = 1, 2$.

For i = 7, the linear regression model in the first stage is

$$E(Y_{i}|\boldsymbol{x}, a) = \gamma_{i0} + \gamma_{i1} \text{Age} + \gamma_{i3} \text{Female} + \gamma_{i4} \text{VKORC1}(\text{AG}) + \gamma_{i6} \text{CYP2C9}(12) + \beta_{i7} \text{CYP2C9}(\text{other}) + \gamma_{i8} \text{Weight} + \gamma_{i9} \text{Height} + a \times \left\{ \xi_{i0}^{(1)} + \xi_{i1}^{(1)} \text{Age} + \xi_{i3}^{(1)} \text{Female} + \xi_{i4}^{(1)} \text{VKORC1}(\text{AG}) + \xi_{i6}^{(1)} \text{CYP2C9}(12) + \xi_{i7}^{(1)} \text{CYP2C9}(\text{other}) \right\} + a^{2} \times \left\{ \xi_{i0}^{(2)} + \xi_{i1}^{(2)} \text{Age} + \xi_{i3}^{(2)} \text{Female} + \xi_{i4}^{(2)} \text{VKORC1}(\text{AG}) + \xi_{i6}^{(2)} \text{CYP2C9}(12) + \xi_{i7}^{(2)} \text{CYP2C9}(\text{other}) \right\}.$$
(B.9)

The parameters in equation (B.9) satisfy

$$\begin{aligned} \gamma_{i0} &= \beta_{i0} + \beta_{i,10}, & \xi_{it}^{(1)} &= \psi_{it}^{(1)}, & \xi_{it}^{(2)} &= \psi_{it}^{(2)}, \\ \gamma_{il} &= \beta_{il}, \end{aligned}$$

for $l \neq 0, 2, 5, 10$ and $t \neq 2, 5$, since all patients in Site 7 are non-White and none take amiodarone or carry VKORC1 genotype AA. The modified likelihood model is

$$\hat{\xi}_{it}^{(u)} \sim N\left(\psi_{it}^{(u)}, \operatorname{sd}(\hat{\xi}_{it}^{(u)})^2\right),$$
$$\psi_{it}^{(u)} \sim N\left(\psi_t^{(u)}, (\sigma_t^{(u)})^2\right),$$

for $t \neq 2, 5, u = 1, 2$.

For i = 1, 2, 4, 5, 8, 9, 10, 11, 12, 13, no modification is needed, as all levels of all covariates are represented in the sites' data. Therefore, the likelihood model is

$$\hat{\xi}_{it}^{(u)} \sim N\left(\psi_{it}^{(u)}, \mathrm{sd}(\hat{\xi}_{i0}^{(u)})^2\right), \psi_{it}^{(u)} \sim N\left(\psi_t^{(u)}, (\sigma_t^{(u)})^2\right),$$

Table B.2: Site-specific blip function parameter estimates $\hat{\xi}_i^{(1)}$ in the stage-one (frequentist) linear regression models and the associated standard deviations (in parantheses) in the analysis of International Warfarin Pharmacogenetics Consortium data. The results are rescaled by a factor of 1000.

Site	1	2	3	4	5
Intercept	8.51 (11.84)	-24.8(81.93)	27.14(45.78)	10.05(18.13)	5.24(24.95)
Age	-0.63(2.15)	0.95(9.77)	-1.43(6.57)	-1.18(3.34)	-2.35(3.82)
Amiodarone	-40.12(35.86)	16.08(74.61)	-102.99(82.5)	10.28(21.91)	55.34(356.95)
Female	0.39(6.12)	15.67(64.71)	-23.85(17.18)	-1.66(11.45)	-4.77(10.79)
VKORC1 (AG)	$0.85\ (7.98)$	0.62(78.88)	-3.6(38.47)	-4.67(7.67)	$15.81 \ (8.72)$
VKORC1 (AA)	18.7(21.97)	-114.16(399.46)	22.33(37.04)	4.21 (25.11)	8.32(28.86)
CYP2C9 (12)	-11.29(9.87)	$16.21 \ (33.51)$	NA	7.92(13.4)	0.09(16.29)
CYP2C9 (other)	-16.1(17.9)	-243.63(562.22)	17.63(37.39)	11.1 (12.83)	1.41 (18.21)
Site	6	7	8	9	10
Intercept	1.96(41.61)	4.01 (24.56)	-9.02 (32.16)	3.05(14.61)	14.23 (12.31)
Age	6.33(9.76)	-2.46(5.1)	1.16 (4.81)	-4.68 (2.71)	-4.1(1.74)
Amiodarone	-82(155.67)	NÁ	15.18(24.78)	-17.29 (20.42)	-17.27 (19.31)
Female	-19.48 (19.06)	-7.57(16.65)	-4.54 (14.14)	15.03(7.98)	5.7(6.83)
VKORC1 (AG)	-10.33(31.52)	-3.36(22.95)	-3.27(16.99)	13.01(7.18)	3.48(9.28)
VKORC1 (AA)	NA	NA	-28.17(76.13)	-11.57(23)	-7.08(15.56)
CYP2C9 (12)	NA	-61.24(116.28)	31.57(14.58)	13.54(12.28)	-1.79(8.46)
CYP2C9 (other)	-3.04(53.03)	75.87(43.38)	-15.36(29.83)	-12.19(8.7)	4.83(9.23)
Site	11	12	13		
Intercept	9.32 (19.44)	36.21 (29.18)	-10.13 (13.22)		
Age	0.6(3.27)	-9.94 (5.84)	2.18(2.67)		
Amiodarone	-55.53 (58.3)	-77.95 (72.19)	8.87 (21.42)		
Female	5.34(9.88)	-1.68(19.07)	0.17(7.81)		
VKORC1 (AG)	-11.19 (20.17)	33.44 (29.42)	3.54 (9.86)		
VKORC1 (AA)	-20.34 (15.77)	34.67(49.1)	-59.53 (43.27)		
CYP2C9 (12)	-8.15 (14.56)	-95.07 (43.58)	-10.29 (13.55)		
CYP2C9 (other)	40.9 (31.66)	-49.45(62.55)	9.82(21.24)		

for $t = 0, \ldots, 7, u = 1, 2$.

To select the variables that are truly relevant for the treatment decision, a horseshoe prior (Carvalho et al., 2010) is assumed for all treatment-covariate interactions. Specifically, for

Table B.3: Site-specific blip function parameter estimates $\hat{\xi}_i^{(2)}$ in the stage-one (frequentist) linear regression models and the associated standard deviations (in parantheses) in the analysis of International Warfarin Pharmacogenetics Consortium data. The results are rescaled by a factor of 1000.

Site	1	2	3	4	5
Intercept	-0.07 (0.12)	0.23(0.9)	-0.45(0.67)	-0.02 (0.21)	-0.22 (0.31)
Age	0(0.02)	0 (0.16)	0.06(0.12)	-0.01 (0.04)	0.05(0.05)
Amiodarone	0.92(0.7)	-0.22(1.15)	2.75(2.45)	-0.12(0.29)	-1.19(7.7)
Female	$0.03\ (0.07)$	-0.22(0.99)	0.5~(0.36)	0.03(0.14)	$0.02 \ (0.15)$
VKORC1 (AG)	-0.03(0.09)	0.16(1.1)	-0.15(0.48)	$0.01 \ (0.09)$	-0.15(0.1)
VKORC1 (AA)	-0.39(0.41)	3.3(9.72)	-0.82(0.54)	-0.06(0.49)	-0.21 (0.53)
CYP2C9 (12)	0.17(0.13)	-0.26(0.38)	NA(NA)	-0.11(0.19)	-0.1 (0.23)
CYP2C9 (other)	$0.29\ (0.31)$	6.43(14.62)	-0.03(0.89)	-0.2(0.19)	-0.02(0.33)
Site	6	7	8	9	10
Intercept	0.04(0.66)	-0.06 (0.24)	0.07(0.36)	-0.02 (0.18)	-0.22 (0.13)
Age	-0.13 (0.16)	0.03(0.05)	-0.01 (0.05)	0.06(0.03)	0.05(0.02)
Amiodarone	1.18(2.71)	NA (NA)	-0.01 (0.27)	0.2(0.3)	0.24(0.29)
Female	0.32(0.3)	0.07(0.16)	0.05(0.18)	-0.21(0.1)	-0.08(0.07)
VKORC1 (AG)	0.34(0.44)	0.12(0.28)	0.05(0.19)	-0.16(0.08)	0.02(0.1)
VKORC1 (AA)	NA(NA)	NA (NA)	0.86(1.7)	0.19(0.41)	$0.11 \ (0.23)$
CYP2C9(12)	NA (NA)	1.22(1.69)	-0.42(0.17)	-0.17(0.14)	$0.01 \ (0.09)$
CYP2C9 (other)	0.07(1.15)	-0.96(0.55)	0.26(0.53)	0.13(0.09)	0 (0.11)
Site	11	12	13		
Intercept	-0.13 (0.24)	-0.36 (0.35)	0.1(0.14)		
Age	-0.01 (0.04)	0.11(0.07)	-0.03 (0.03)		
Amiodarone	0.89(1.21)	1.05(1.17)	-0.12 (0.31)		
Female	0.01(0.14)	-0.02(0.23)	0.01(0.09)		
VKORC1 (AG)	0.11(0.36)	-0.43(0.34)	-0.05(0.12)		
VKORC1 (AA)	0.27(0.21)	-0.76 (0.84)	1.27(0.96)		
CYP2C9 (12)	0.09(0.21)	1.23(0.6)	0.13(0.17)		
CYP2C9 (other)	-0.89(0.7)	0.9(1.04)	-0.15 (0.34)		
t = 1, ..., 7 and u = 1, 2, we have

$$\psi_t^{(u)} \sim N(0, \tau^2 (\lambda_t^{(u)})^2),$$

 $\lambda_t^{(u)} \sim \text{Half-Cauchy (0,1)},$
 $\tau \sim \text{Half-Cauchy (0,1)},$

where τ and $\lambda_t^{(u)}$ are, respectively, the global and local shrinkage parameters. If the 95% credibility interval of $\psi_t^{(u)}$, t = 1, ..., 7, u = 1, 2, contains zero, the corresponding treatment-covariate interaction will not be selected, suggesting that the associated covariate has no evidence of a tailoring effect on the optimal Warfarin dosing. For $\psi_0^{(u)}$, u = 1, 2, the priors are

$$\psi_0^{(u)} \sim \begin{cases} N(0, 100^2)^+, & u = 1\\ N(0, 100^2)^-, & u = 2 \end{cases}$$

Here, we use truncated priors for $\psi_0^{(1)}$ and $\psi_0^{(2)}$, as a positive dose effect and a negative squared-dose effect on the defined outcome are substantively reasonable and have been found in previous work (Schulz and Moodie, 2021). Regarding the variance component parameters $\sigma_t^{(u)}$, a half-Cauchy (0,1) prior is used. The Bayesian hierarchical model is implemented in RStan (Stan Development Team, 2021, 2020); 2000 posterior samples are drawn from two chains for each parameter.

APPENDIX C

Appendix to Manuscript 3

The supplementary material includes simulation results for all blip parameters other than ψ_{20} and ψ_{21} . Tables C.1 and C.2 show the simulation results of ψ_{h0} and ψ_{h1} for $h = 3, \ldots, 7$, respectively, for the setting where the data are generated under an unstructured betweensite heterogeneity model, while tables C.3 and C.4 show the simulation results of the same set of parameters but for the setting where the data are generated under the assumption of common between-site heterogeneity. We note that not all parameter ψ_{ht} , $h = 3, \ldots, 7$, t = 0, 1 are available for all networks due to the difference in treatment sets.

Table C.1: model. Rel networks, n	Simulation results lative Bias (%, de numbers of sites, an	s for the setting moted by RB) a nd heterogeneity	where the data and standard dev assumptions in	are generated u viations (SD) of the Bayesian hid	nder an unstruc $\hat{\psi}_{h0}, h = 3, \dots$ srarchical model	tured between-s ,7, are reported based on 2000 j	ite heterogeneity 1 across different iterations.
Network	Number of sites	Heterogeneity	$\hat{\psi}_{30}$ RB (SD)	$\hat{\psi}_{40}$ RB (SD)	$\hat{\psi}_{50}$ RB (SD)	$\hat{\psi}_{60}$ RB (SD)	$\hat{\psi}_{70}$ RB(SD)
	1	0	-0.217 (0.422)				
а	c.	Common	-0.037 (0.333)				
	>	Varying	-0.045(0.477)				
	1	0	-0.042(0.277)				
p	ст.	Common	-0.075(0.182)				
)	Varying	-0.047 (0.258)				
	1	0	$-0.031 \ (0.276)$	-0.051(0.304)	-0.124(0.280)		
С	ст.	Common	$0.063\ (0.176)$	$0.051 \ (0.192)$	$0.042\ (0.175)$		
)	Varying	$0.065\ (0.212)$	$0.049\ (0.242)$	$0.043\ (0.215)$		
	1	0	-0.002(0.428)	$0.267\ (0.402)$	$0.050\ (0.406)$	-0.454(0.588)	$-0.269\ (0.592)$
q	67	Common	$0.014\ (0.273)$	$0.154\ (0.259)$	$0.098\ (0.262)$	$0.135\ (0.377)$	-0.074(0.381)
)	Varying	$0.014\ (0.483)$	$0.112\ (0.444)$	$0.089\ (0.455)$	$0.111 \ (0.663)$	-0.045(0.667)
	1	0	-0.050(0.372)	$0.093 \ (0.544)$	$-0.107\ (0.560)$	$0.196\ (0.531)$	
в	~~	Common	$0.050\ (0.227)$	$0.250\ (0.319)$	$-0.051 \ (0.326)$	$0.134\ (0.310)$	
		Varying	$0.057\ (0.262)$	$0.295\ (0.553)$	-0.058 (0.548)	$0.036\ (0.517)$	

Table C.2: model. Rel networks, n	Simulation results lative Bias (%, de numbers of sites, an	s for the setting moted by RB) a nd heterogeneity	where the data and standard dev assumptions in	are generated u viations (SD) of the Bayesian hid	nder an unstruc $\hat{\psi}_{h_1}, h = 3, \dots$ srarchical model	tured between-s ,7, are reported based on 2000 j	ite heterogeneity 1 across different iterations.
Network	Number of sites	Heterogeneity	$\hat{\psi}_{31}$ RB (SD)	$\hat{\psi}_{41}$ RB (SD)	$\hat{\psi}_{51}$ RB (SD)	$\hat{\psi}_{61}$ RB (SD)	$\hat{\psi}_{71}$ RB(SD)
	1	0	$0.317\ (0.203)$				
а	c.	Common	$0.120\ (0.210)$				
	>	Varying	$0.083\ (0.355)$				
	1	0	-0.134(0.136)				
p	ст.	Common	-0.173(0.107)				
)	Varying	-0.280(0.186)				
	1	0	$0.004\ (0.134)$	$0.171 \ (0.125)$	-0.218(0.125)		
C	er.	Common	$0.103\ (0.100)$	$0.318\ (0.104)$	$0.055 \ (0.097)$		
)	Varying	$0.050\ (0.122)$	$0.249 \ (0.128)$	$0.089\ (0.120)$		
	1	0	$0.010\ (0.206)$	$0.228\ (0.158)$	-0.241(0.179)	-1.224(0.274)	-0.585(0.272)
q	ст.	Common	-0.271(0.163)	-0.018(0.146)	-0.136(0.156)	-0.552(0.227)	$-0.541 \ (0.226)$
)	Varying	-0.230(0.351)	-0.004 (0.296)	-0.132(0.323)	-0.600(0.474)	-0.465(0.479)
	1	0	-0.378(0.155)	$-0.207 \ (0.166)$	-0.443(0.191)	-0.253(0.200)	
в	~~	Common	-0.112(0.114)	$0.096\ (0.134)$	-0.127(0.145)	$0.206\ (0.154)$	
		Varying	$-0.104 \ (0.142)$	0.193(0.283)	$-0.167\ (0.309)$	$0.291 \ (0.338)$	

Table C.3: heterogenei different ner	Simulation result ity. Relative Bias tworks, numbers or	is for the setting (%, denoted by f sites, and heter	g where the dat y RB) and stan ogeneity assump	a are generated dard deviations otions in the Bay	. under the assu (SD) of $\hat{\psi}_{h0}$, h esian hierarchica	mption of comr = $3, \ldots, 7$, are l model based or	non between-site e reported across n 2000 iterations.
Network	Number of sites	Heterogeneity	$\hat{\psi}_{30}$ RB (SD)	$\hat{\psi}_{40}$ RB (SD)	$\hat{\psi}_{50}$ RB (SD)	$\hat{\psi}_{60}$ RB (SD)	$\hat{\psi}_{70}$ RB(SD)
	1	0	-0.195(0.413)				
g	ст.	Common	-0.021(0.326)				
	>	Varying	-0.018(0.471)				
	1	0	$-0.039\ (0.271)$				
q	с.	Common	-0.074 (0.177)				
	D	Varying	-0.076(0.233)				
	1	0	-0.022(0.271)	$-0.057\ (0.310)$	-0.126(0.280)		
С	ст.	Common	$0.059\ (0.172)$	$0.055\ (0.195)$	$0.044 \ (0.174)$		
	>	Varying	$0.061\ (0.208)$	$0.055\ (0.248)$	$0.045\ (0.214)$		
	1	0	-0.011(0.418)	$0.268\ (0.413)$	$0.061\ (0.405)$	-0.405(0.581)	-0.376(0.583)
q	67	Common	$0.021 \ (0.267)$	$0.152\ (0.264)$	$0.093 \ (0.262)$	$0.137\ (0.372)$	-0.015(0.375)
)	Varying	$0.030\ (0.477)$	$0.157\ (0.478)$	$0.079\ (0.460)$	$0.081 \ (0.658)$	-0.007 (0.664)
	1	0	$0.094\ (0.387)$	$-0.192\ (0.555)$	$-0.159\ (0.567)$	$0.856\ (0.541)$	
е	ст .	Common	$-0.074 \ (0.247)$	$0.105\ (0.344)$	$-0.009\ (0.347)$	-0.183(0.341)	
		Varying	$-0.048 \ (0.354)$	$0.064\ (0.553)$	$-0.041 \ (0.557)$	-0.075(0.594)	

Table C.4: heterogenei different ne	Simulation result ity. Relative Bias tworks, numbers of	 is for the setting (%, denoted by f sites, and heter 	g where the dat γ RB) and stam ogeneity assump	a are generated dard deviations tions in the Bay	under the assu (SD) of $\hat{\psi}_{h1}$, h esian hierarchica	mption of comr $= 3, \ldots, 7$, are l model based or	non between-site e reported across n 2000 iterations.
Network	Number of sites	Heterogeneity	$\hat{\psi}_{31}$ RB (SD)	$\hat{\psi}_{41}$ RB (SD)	$\hat{\psi}_{51}$ RB (SD)	$\hat{\psi}_{61}$ RB (SD)	$\hat{\psi}_{71}$ RB(SD)
	1	0	$0.279 \ (0.181)$				
g	с.	Common	$0.123\ (0.196)$				
	>	Varying	$0.100\ (0.324)$				
	1	0	-0.108(0.123)				
q	c.	Common	-0.166(0.098)				
	D	Varying	-0.178(0.111)				
	1	0	-0.003(0.123)	$0.135\ (0.140)$	-0.225(0.124)		
С	ст.	Common	$0.106\ (0.095)$	$0.363\ (0.108)$	$0.040 \ (0.096)$		
	>	Varying	$0.047\ (0.115)$	$0.259\ (0.141)$	$0.053\ (0.119)$		
	1	0	-0.016(0.184)	$0.249\ (0.183)$	-0.217(0.179)	-1.122(0.257)	-0.740(0.256)
q	ст.	Common	-0.224(0.152)	-0.049(0.154)	-0.179(0.154)	-0.616(0.218)	-0.453(0.217)
)	Varying	-0.220(0.324)	$0.092\ (0.368)$	-0.193(0.329)	-0.565(0.493)	-0.427 (0.460)
	1	0	$0.252\ (0.193)$	-0.623(0.200)	$-0.744 \ (0.230)$	-1.073(0.252)	
е	er:	Common	$0.060\ (0.145)$	$0.115\ (0.169)$	$0.117\ (0.185)$	$0.335\ (0.193)$	
		Varying	$0.151\ (0.193)$	$0.138\ (0.348)$	$0.116\ (0.366)$	$0.431 \ (0.391)$	

References

- Archer, C., Kessler, D., Lewis, G., Araya, R., Duffy, L., Gilbody, S., Lewis, G., Kendrick, T., Peters, T. J., and Wiles, N. (2024). What predicts response to sertraline for people with depression in primary care? a secondary data analysis of moderators in the PANDA trial. *PLOS ONE*, 19(5):1–12.
- Arjas, E. and Saarela, O. (2010). Optimal dynamic regimes: presenting a case for predictive inference. *The International Journal of Biostatistics*, 6(2).
- Beck, A. T. and Alford, B. A. (2014). Depression: causes and treatment. University of Pennsylvania Press, Philadelphia, 2nd edition.
- Béliveau, A., Goring, S., Platt, R. W., and Gustafson, P. (2017). Network meta-analysis of disconnected networks: how dangerous are random baseline treatment effects? *Research Synthesis Methods*, 8(4):465–474.
- Béliveau, A. and Gustafson, P. (2021). A theoretical investigation of how evidence flows in Bayesian network meta-analysis of disconnected networks. *Bayesian Analysis*, 16(3):803– 823.
- Benrimoh, D., Fratila, R., Israel, S., Perlman, K., Mirchi, N., Desai, S., Rosenfeld, A., Knappe, S., Behrmann, J., Rollins, C., You, R. P., and Aifred Health Team, T. (2018). Aifred health, a deep learning powered clinical decision support system for mental health.

In Escalera, S. and Weimer, M., editors, *The NIPS '17 Competition: Building Intelligent Systems*, pages 251–287, Cham. Springer International Publishing.

- Benrimoh, D., Tanguay-Sela, M., Perlman, K., Israel, S., Mehltretter, J., Armstrong, C., Fratila, R., Parikh, S. V., Karp, J. F., Heller, K., Vahia, I. V., Blumberger, D. M., Karama, S., Vigod, S. N., Myhr, G., Martins, R., Rollins, C., Popescu, C., Lundrigan, E., Snook, E., Wakid, M., Williams, J., Soufi, G., Perez, T., Tunteng, J.-F., Rosenfeld, K., Miresco, M., Turecki, G., Gomez Cardona, L., Linnaranta, O., and Margolese, H. C. (2021). Using a simulation centre to evaluate preliminary acceptability and impact of an artificial intelligence-powered clinical decision support system for depression treatment on the physician-patient interaction. *BJPsych Open*, 7(1):e22.
- Berlin, J. A., Santanna, J., Schmid, C. H., Szczech, L. A., and Feldman, H. I. (2002). Individual patient- versus group-level data meta-regressions for the investigation of treatment effect modifiers: ecological bias rears its ugly head. *Statistics in Medicine*, 21(3):371–387.
- Berry, S. M., Carlin, B. P., Lee, J. J., and Muller, P. (2010). Bayesian adaptive methods for clinical trials. Chapman & Hall/CRC, London.
- Bondell, H. D. and Reich, B. J. (2012). Consistent high-dimensional Bayesian variable selection via penalized credible regions. *Journal of the American Statistical Association*, 107(500):1610–1624.
- Borenstein, M., Hedges, L. V., Higgins, J. P., and Rothstein, H. R. (2021). Introduction to meta-analysis. John Wiley & Sons, Oxford.
- Bowden, J., Tierney, J. F., Simmonds, M., Copas, A. J., and Higgins, J. P. (2011). Individual patient data meta-analysis of time-to-event outcomes: one-stage versus two-stage approaches for estimating the hazard ratio under a random effects model. *Research Syn*thesis Methods, 2(3):150–162.

- Burke, D. L., Ensor, J., and Riley, R. D. (2017). Meta-analysis using individual participant data: one-stage and two-stage approaches, and why they may differ. *Statistics in Medicine*, 36(5):855–875.
- Caldwell, D. M., Dias, S., and Welton, N. J. (2015). Extending treatment networks in health technology assessment: how far should we go? *Value Health*, 18(5):673–681.
- Campbell, M. J. (2000). Cluster randomized trials in general (family) practice research. Statistical Methods in Medical Research, 9(2):81–94.
- Campbell, M. J., Donner, A., and Klar, N. (2007). Developments in cluster randomized trials and statistics in medicine. *Statistics in Medicine*, 26(1):2–19.
- Campbell, M. K., Mollison, J., and Grimshaw, J. M. (2001). Cluster trials in implementation research: estimation of intracluster correlation coefficients and sample size. *Statistics in Medicine*, 20(3):391–399.
- Carvalho, C. M., Polson, N. G., and Scott, J. G. (2010). The horseshoe estimator for sparse signals. *Biometrika*, 97(2):465–480.
- Chakraborty, B. and Moodie, E. E. M. (2013). Statistical methods for dynamic treatment regimes: reinforcement learning, causal inference, and personalized medicine. Springer, New York.
- Chakraborty, B. and Murphy, S. A. (2014). Dynamic treatment regimes. Annual Review of Statistics and Its Application, 1:447–464.
- Chevret, S. (2012). Bayesian adaptive clinical trials: a dream for statisticians only? *Statistics in Medicine*, 31(11-12):1002–1013.
- Chow, S. C. and Chang, M. (2008). Adaptive design methods in clinical trials a review. Orphanet Journal of Rare Diseases, 3(11).

- Chow, S. C., Chang, M., and Pong, A. (2005). Statistical consideration of adaptive methods in clinical development. *Journal of Biopharmaceutical Statistics*, 15(4):575–591.
- Cipriani, A., Furukawa, T. A., Geddes, J. R., Malvini, L., Signoretti, A., McGuire, H., Churchill, R., Nakagawa, A., Barbui, C., and MANGA Study Group (2008). Does randomized evidence support sertraline as first-line antidepressant for adults with acute major depression? a systematic review and meta-analysis. *Journal Clinical Psychiatry*, 69(11):1732–1742.
- Cipriani, A., Higgins, J. P., Geddes, J. R., and Salanti, G. (2013). Conceptual and technical challenges in network meta-analysis. Annals of Internal Medicine, 159(2):130–137.
- Cochran, W. G. (1954). The combination of estimates from different experiments. *Biometrics*, 10(1):101–129.
- Cole, S. R. and Hernán, M. A. (2008). Constructing inverse probability weights for marginal structural models. American Journal of Epidemiology, 168(6):656–664.
- Cooper, N. J., Peters, J., Lai, M. C. W., Juni, P., Wandel, S., Palmer, S., Paulden, M., Conti, S., Welton, N. J., Abrams, K. R., Bujkiewicz, S., Spiegelhalter, D., and Sutton, A. J. (2011). How valuable are multiple treatment comparison methods in evidence-based health-care evaluation? *Value Health*, 14(2):371–380.
- Danieli, C. and Moodie, E. E. M. (2022). Preserving data privacy when using multi-site data to estimate individualized treatment rules. *Statistics in Medicine*, 41(9):1627–1643.
- Demets, D. L. and Ware, J. H. (1980). Group sequential methods for clinical trials with a one-sided hypothesis. *Biometrika*, 67(3):651–660.
- DeMets, D. L. and Ware, J. H. (1982). Asymmetric group sequential boundaries for monitoring clinical trials. *Biometrika*, 69(3):661–663.

- DerSimonian, R. and Kacker, R. (2007). Random-effects model for meta-analysis of clinical trials: an update. *Contemporary Clinical Trials*, 28(2):105–114.
- DerSimonian, R. and Laird, N. (1986). Meta-analysis in clinical trials. Controlled Clinical Trials, 7(3):177–188.
- Dias, S., Welton, N. J., and Ades, A. (2010a). Study designs to detect sponsorship and other biases in systematic reviews. *Journal of Clinical Epidemiology*, 63(6):587–588.
- Dias, S., Welton, N. J., Caldwell, D. M., and Ades, A. E. (2010b). Checking consistency in mixed treatment comparison meta-analysis. *Statistics in Medicine*, 29(7-8):932–944.
- Dmitrienko, A. and Wang, M.-D. (2006). Bayesian predictive approach to interim monitoring in clinical trials. *Statistics in Medicine*, 25(13):2178–2195.
- Donegan, S., Williamson, P., D'Alessandro, U., and Tudur Smith, C. (2013). Assessing key assumptions of network meta-analysis: a review of methods. *Research Synthesis Methods*, 4(4):291–323.
- Donner, A. and Klar, N. (2000). Design and analysis of cluster randomization trials in health research. Wiley, London.
- Donner, A. and Klar, N. (2004). Pitfalls of and controversies in cluster randomization trials. American Journal of Public Health, 94(3):416–422.
- Dwork, C. (2006). Differential privacy. In Bugliesi, M., Preneel, B., Sassone, V., and Wegener,
 I., editors, *International Colloquium on Automata, Languages and Programming*, pages 1–
 12, Berlin, Heidelberg. Springer.
- Dwork, C. (2008). Differential privacy: a survey of results. In Agrawal, M., Du, D., Duan, Z., and Li, A., editors, *Theory and Applications of Models of Computation*, pages 1–19, Berlin, Heidelberg. Springer.

- Eldridge, S. M., Ashby, D., and Kerry, S. (2006). Sample size for cluster randomized trials: effect of coefficient of variation of cluster size and analysis method. *International Journal* of Epidemiology, 35(5):1292–1300.
- Fisher, D., Copas, A., Tierney, J., and Parmar, M. (2011). A critical review of methods for the assessment of patient-level interactions in individual participant data meta-analysis of randomized trials, and guidance for practitioners. *Journal of Clinical Epidemiology*, 64(9):949–967.
- Food and Drug Administration (2010). Guidance for the use of Bayesian statistics in medical device clinical trials. https://www.fda.gov/regulatory-information/search-fda-g uidance-documents/guidance-use-bayesian-statistics-medical-device-clinica l-trials. Accessed: 2024-07-22.
- Freedman, L. S. and Spiegelhalter, D. J. (1989). Comparison of Bayesian with group sequential methods for monitoring clinical trials. *Controlled Clinical Trials*, 10(4):357–367.
- Freedman, L. S., Spiegelhalter, D. J., and Parmar, M. K. B. (1994). The what, why and how of Bayesian clinical trials monitoring. *Statistics in Medicine*, 13(13-14):1371–1383.
- Friedman, L. M., Furberg, C. D., DeMets, D. L., Reboussin, D. M., and Granger, C. B. (2015). Fundamentals of clinical trials. Springer, New York, 5th edition.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models. Bayesian Analysis, 1(3):515–533.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B. (2013). Bayesian data analysis. CRC Press, New York.
- Gold, L. H. (2014). DSM-5 and the assessment of functioning: the World Health Organization Disability Assessment Schedule 2.0 (WHODAS 2.0). The Journal of the American Academy of Psychiatry and the Law, 42(2):173–181.

- Goring, S. M., Gustafson, P., Liu, Y., Saab, S., Cline, S. K., and Platt, R. W. (2016). Disconnected by design: analytic approach in treatment networks having no common comparator. *Research Synthesis Methods*, 7(4):420–432.
- Grayling, M. J., Wason, J. M., and Mander, A. P. (2017). Group sequential designs for stepped-wedge cluster randomised trials. *Clinical Trials*, 14(5):507–517.
- Grayling, M. J., Wason, J. M. S., and Villar, S. S. (2022). Response adaptive intervention allocation in stepped-wedge cluster randomized trials. *Statistics in Medicine*, 41(6):1081– 1099.
- Greenland, S. (1983). Tests for interaction in epidemiologic studies: a review and a study of power. Statistics in Medicine, 2(2):243–251.
- Gsponer, T., Gerber, F., Bornkamp, B., Ohlssen, D., Vandemeulebroecke, M., and Schmidli, H. (2014). A practical guide to Bayesian group sequential designs. *Pharmaceutical Statis*tics, 13(1):71–80.
- Guittet, L., Ravaud, P., and Giraudeau, B. (2006). Planning a cluster randomized trial with unequal cluster sizes: practical issues involving continuous outcomes. BMC Medical Research Methodology, 6(17).
- Hahn, P. R. and Carvalho, C. M. (2015). Decoupling shrinkage and selection in Bayesian linear models: a posterior summary perspective. *Journal of the American Statistical Association*, 110(509):435–448.
- Hahn, P. R., Murray, J. S., and Carvalho, C. M. (2020). Bayesian regression tree models for causal inference: regularization, confounding, and heterogeneous effects (with discussion). *Bayesian Analysis*, 15(3):965–1056.
- Hales, R. J. and Moulton, L. H. (2017). Cluster randomised trials. Chapman & Hall/CRC, London, 2nd edition.

- Hardy, R. J. and Thompson, S. G. (1996). A likelihood approach to meta-analysis with random effects. *Statistics in Medicine*, 15(6):619–629.
- Henssler, J., Bschor, T., and Baethge, C. (2016). Combining antidepressants in acute treatment of depression: A meta-analysis of 38 studies including 4511 patients. *The Canadian Journal Psychiatry*, 61(1):29–43.
- Higgins, J. P., Thomas, J., Chandler, J., Cumpston, M., Li, T., Page, M. J., and Welch, V. A. (2019). Cochrane handbook for systematic reviews of interventions. John Wiley & Sons, Oxford.
- Higgins, J. P. T., Jackson, D., Barrett, J. K., Lu, G., Ades, A. E., and White, I. R. (2012). Consistency and inconsistency in network meta-analysis: concepts and models for multiarm studies. *Research Synthesis Methods*, 3(2):98–110.
- Hougaard, P. (1995). Frailty models for survival data. Lifetime Data Analysis, 1:255–273.
- Ibrahim, J. G., Chen, M.-H., and Sinha, D. (2001). Bayesian survival analysis. Springer, New York, 2nd edition.
- Imbens, G. W. (2000). The role of the propensity score in estimating dose-response functions. Biometrika, 87(3):706–710.
- International Warfarin Pharmacogenetics Consortium (2009). Estimation of the Warfarin dose with clinical and pharmacogenetic data. *New England Journal of Medicine*, 360(8):753–764.
- Kass, R. E. and Raftery, A. E. (1995). Bayes factors. Journal of the American Statistical Association, 90(430):773–795.
- Kennedy, S. H., Andersen, H. F., and Thase, M. E. (2009). Escitalopram in the treatment of major depressive disorder: a meta-analysis. *Current Medical Research and Opinion*, 25(1):161–175.

- Kennedy, S. H., Lam, R. W., McIntyre, R. S., Tourjman, S. V., Bhat, V., Blier, P., Hasnain, M., Jollant, F., Levitt, A. J., MacQueen, G. M., McInerney, S. J., McIntosh, D., Milev, R. V., Müller, D. J., Parikh, S. V., Pearson, N. L., Ravindran, A. V., Uher, R., and the CANMAT Depression Work Group (2016). Canadian Network for Mood and Anxiety Treatments (CANMAT) 2016 clinical guidelines for the management of adults with major depressive disorder: section 3. pharmacological treatments. *The Canadian Journal of Psychiatry*, 61(9):540–560.
- Kerry, S. M. and Martin Bland, J. (2001). Unequal cluster sizes for trials in English and Welsh general practice: implications for sample size calculations. *Statistics in Medicine*, 20(3):377–390.
- Kessler, R. C., van Loo, H. M., Wardenaar, K. J., Bossarte, R. M., Brenner, L. A., Ebert, D. D., de Jonge, P., Nierenberg, A. A., Rosellini, A. J., Sampson, N. A., and et al. (2017). Using patient self-reports to study heterogeneity of treatment effects in major depressive disorder. *Epidemiology and Psychiatric Sciences*, 26(1):22–36.
- Kidwell, K. M., Seewald, N. J., Tran, Q., Kasari, C., and Almirall, D. (2018). Design and analysis considerations for comparing dynamic treatment regimens with binary outcomes from sequential multiple assignment randomized trials. *Journal of Applied Statistics*, 45(9):1628–1651.
- Killip, S., Mahfoud, Z., and Pearce, K. (2004). What is an intracluster correlation coefficient? crucial concepts for primary care researchers. *Annals of Family Medicine*, 2(3):204–208.
- Kirino, E. (2012). Escitalopram for the management of major depressive disorder: a review of its efficacy, safety, and patient acceptability. *Patient Preference and Adherence*, 6:853–861.
- Klar, N. and Donner, A. (2001). Current and future challenges in the design and analysis of cluster randomization trials. *Statistics in Medicine*, 20(24):3729–3740.

- Kontopantelis, E. (2018). A comparison of one-stage vs two-stage individual patient data meta-analysis methods: a simulation study. *Research Synthesis Methods*, 9(3):417–430.
- Koopman, L., van der Heijden, G. J. M. G., Hoes, A. W., Grobbee, D. E., and Rovers, M. M. (2008). Empirical comparison of subgroup effects in conventional and individual patient data meta-analyses. *International Journal of Technology Assessment in Health Care*, 24(3):358–361.
- Kosorok, M. R. and Laber, E. B. (2019). Precision medicine. Annual Review of Statistics and Its Application, 6:263–286.
- Kravitz, R. L., Duan, N., and Braslow, J. (2004). Evidence-based medicine, heterogeneity of treatment effects, and the trouble with averages. *The Milbank Quarterly*, 82(4):661–687.
- Kroenke, K. and Spitzer, R. L. (2002). The PHQ-9: a new depression diagnostic and severity measure. *Psychiatric Annals*, 32(9):509–515.
- Laber, E. B., Lizotte, D. J., Qian, M., Pelham, W. E., and Murphy, S. A. (2014). Dynamic treatment regimes: technical challenges and applications. *Electronic Journal of Statistics*, 8(1):1225–1272.
- Lake, S., Kammann, E., Klar, N., and Betensky, R. (2002). Sample size re-estimation in cluster randomization trials. *Statistics in Medicine*, 21:1337–1350.
- Lan, K. K. G. and DeMets, D. L. (1983). Discrete sequential boundaries for clinical trials. Biometrika, 70(3):659–663.
- Lewandowski, D., Kurowicka, D., and Joe, H. (2009). Generating random correlation matrices based on vines and extended onion method. *Journal of Multivariate Analysis*, 100(9):1989–2001.
- Lewis, R. J. and Berry, D. A. (1994). Group sequential clinical trials: a classical evaluation

of Bayesian decision-theoretic designs. Journal of the American Statistical Association, 89(428):1528–1534.

- Li, F. and Li, F. (2019). Propensity score weighting for causal inference with multiple treatments. *The Annals of Applied Statistics*, 13(4):2389–2415.
- Li, H. and Pati, D. (2017). Variable selection using shrinkage priors. Computational Statistics
 & Data Analysis, 107:107–119.
- Linn, K. A., Laber, E. B., and Stefanski, L. A. (2017). Interactive Q-learning for quantiles. Journal of the American Statistical Association, 112(518):638–649.
- Liu, Y., Wang, Y., Kosorok, M. R., Zhao, Y., and Zeng, D. (2018). Augmented outcomeweighted learning for estimating optimal dynamic treatment regimens. *Statistics in Medicine*, 37(26):3776–3788.
- Logan, B. R., Sparapani, R., McCulloch, R. E., and Laud, P. W. (2019). Decision making and uncertainty quantification for individualized treatments using Bayesian additive regression trees. *Statistical Methods in Medical Research*, 28(4):1079–1093.
- Lu, G. and Ades, A. E. (2006). Assessing evidence inconsistency in mixed treatment comparisons. *Journal of the American Statistical Association*, 101(474):447–459.
- Manatunga, A. K., Hudgens, M. G., and Chen, S. (2001). Sample size estimation in cluster randomized studies with varying cluster size. *Biometrical Journal*, 43(1):75–86.
- McGlothlin, A. E. and Lewis, R. J. (2014). Minimal clinically important difference: defining what really matters to patients. *JAMA*, 312(13):1342–1343.
- Mehltretter, J., Fratila, R., Benrimoh, D., Kapelner, A., Perlman, K., Snook, E., Israel, S., Armstrong, C., Miresco, M., and Turecki, G. (2020a). Differential treatment benefit prediction for treatment selection in depression: a deep learning analysis of STAR*D and CO-MED Data. *Computational Psychiatry*, 4:61–75.

- Mehltretter, J., Rollins, C., Benrimoh, D., Fratila, R., Perlman, K., Israel, S., Miresco, M., Wakid, M., and Turecki, G. (2020b). Analysis of features selected by a deep learning model for differential treatment selection in depression. *Frontiers in Artificial Intelligence*, 2:31.
- Moodie, E. E. M., Coulombe, J., Danieli, C., Renoux, C., and Shortreed, S. M. (2022). Privacy-preserving estimation of an optimal individualized treatment rule: a case study in maximizing time to severe depression-related outcomes. *Lifetime Data Analysis*, 28:1–31.
- Moodie, E. E. M., Dean, N., and Sun, Y. R. (2014). Q-learning: flexible learning about useful utilities. *Statistics in Biosciences*, 6(2):223–243.
- Morris, T. P., Fisher, D. J., Kenward, M. G., and Carpenter, J. R. (2018). Meta-analysis of Gaussian individual patient data: two-stage or not two-stage? *Statistics in Medicine*, 37(9):1419–1438.
- Morris, T. P., White, I. R., and Crowther, M. J. (2019). Using simulation studies to evaluate statistical methods. *Statistics in Medicine*, 38(11):2074–2102.
- Murray, T. A., Yuan, Y., and Thall, P. F. (2018). A Bayesian machine learning approach for optimizing dynamic treatment regimes. *Journal of the American Statistical Association*, 113(523):1255–1267.
- Noda, A., Kraemer, H. C., Taylor, J. L., Schneider, B., Ashford, J. W., and Yesavage, J. A. (2006). Strategies to reduce site differences in multisite studies: a case study of Alzheimer disease progression. *The American Journal of Geriatric Psychiatry*, 14(11):931–938.
- Noma, H., Furukawa, T. A., Maruo, K., Imai, H., Shinohara, K., Tanaka, S., Ikeda, K., Yamawaki, S., and Cipriani, A. (2019). Exploratory analyses of effect modifiers in the antidepressant treatment of major depression: individual-participant data meta-analysis of 2803 participants in seven placebo-controlled randomized trials. *Journal of Affective Disorders*, 250:419–424.

- O'Brien, P. C. and Fleming, T. R. (1979). A multiple testing procedure for clinical trials. Biometrics, 35(3):549–556.
- Paule, R. C. and Mandel, J. (1982). Consensus values and weighting factors. Journal of Research of the National Bureau of Standards (1977), 87(5):377–385.
- Perlman, K., Benrimoh, D., Israel, S., Rollins, C., Brown, E., Tunteng, J.-F., You, R., You, E., Tanguay-Sela, M., Snook, E., Miresco, M., and Berlim, M. T. (2019). A systematic meta-review of predictors of antidepressant treatment outcome in major depressive disorder. *Journal of Affective Disorders*, 243:503–515.
- Piepho, H. P., Williams, E. R., and Madden, L. V. (2012). The use of two-way linear mixed models in multitreatment meta-analysis. *Biometrics*, 68(4):1269–1277.
- Pocock, S. J. (1977). Group sequential methods in the design and analysis of clinical trials. *Biometrika*, 64(2):191–199.
- Popescu, C., Golden, G., Benrimoh, D., Tanguay-Sela, M., Slowey, D., Lundrigan, E., Williams, J., Desormeau, B., Kardani, D., Perez, T., Rollins, C., Israel, S., Perlman, K., Armstrong, C., Baxter, J., Whitmore, K., Fradette, M.-J., Felcarek-Hope, K., Soufi, G., Fratila, R., Mehltretter, J., Looper, K., Steiner, W., Rej, S., Karp, J. F., Heller, K., Parikh, S. V., McGuire-Snieckus, R., Ferrari, M., Margolese, H., and Turecki, G. (2021). Evaluating the clinical feasibility of an artificial intelligence–powered, web-based clinical decision support system for the treatment of depression in adults: longitudinal feasibility study. JMIR Formative Research, 5(10):e31862.
- Puffer, S., Torgerson, D. J., and Watson, J. (2005). Cluster randomized controlled trials. Journal of Evaluation in Clinical Practice, 11(5):479–483.
- Rassen, J., Moran, J., Toh, D., Kowal, M., Johnson, K., Shoabi, A., Hammad, T., Raebel, M., Holmes, J., Haynes, K., et al. (2013). Evaluating strategies for data sharing and analyses in distributed data settings. *Technical report, Mini-Sentinel*.

- Rettie, A. E. and Tai, G. (2006). The pharmocogenomics of Warfarin: closing in on personalized medicine. *Molecular interventions*, 6(4):223–227.
- Riley, R. D. and Fisher, D. J. (2021). Individual participant data meta-analysis: a handbook for healthcare research. John Wiley & Sons, Oxford.
- Riley, R. D., Lambert, P. C., and Abo-Zaid, G. (2010). Meta-analysis of individual participant data: rationale, conduct, and reporting. *British Medical Journal*, 340.
- Robins, J. M. (1997). Causal inference from complex longitudinal data. In Berkane, M., editor, Latent Variable Modeling and Applications to Causality, pages 69–117. Springer, New York.
- Robins, J. M. (2000). Marginal structural models versus structural nested models as tools for causal inference. In Halloran, M. E. and Berry, D., editors, *Statistical Models in Epidemiology, the Environment, and Clinical Trials*, pages 95–133, New York. Springer.
- Robins, J. M. (2004). Optimal structural nested models for optimal sequential decisions. In Lin, D. Y. and Heagerty, P. J., editors, *Proceedings of the Second Seattle Symposium in Biostatistics*, pages 189–326. Springer, New York.
- Rodriguez Duque, D., Stephens, D. A., Moodie, E. E. M., and Klein, M. B. (2022). Semiparametric Bayesian inference for optimal dynamic treatment regimes via dynamic marginal structural models. *Biostatistics*, 24(3):708–727.
- Rosenberger, W. F., Sverdlov, O., and Hu, F. (2012). Adaptive randomization for clinical trials. *Journal of Biopharmaceutical Statistics*, 22(4):719–736.
- Ross, D. A., Dollimore, N., Smith, P. G., Kirkwood, B. R., Arthur, P., Morris, S. S., Addy, H. A., Binka, F., Arthur, P., Gyapong, J. O., and Tomkins, A. M. (1993). Vitamin A supplementation in northern Ghana: effects on clinic attendances, hospital admissions, and child mortality. *The Lancet*, 342(8862):7–12.

- Rubin, D. (1980). Discussion of "Randomization analysis of experimental data in the Fisher randomization test" by D. Basu. *Journal of the American Statistical Association*, 75:591– 593.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688–701.
- Rush, A., Batey, S. R., Donahue, R. M., Ascher, J. A., Carmody, T. J., and Metz, A. (2001). Does pretreatment anxiety predict response to either bupropion SR or sertraline? *Journal of Affective Disorders*, 64(1):81–87.
- Rush, A., Fava, M., Wisniewski, S. R., Lavori, P. W., Trivedi, M. H., Sackeim, H. A., Thase, M. E., Nierenberg, A. A., Quitkin, F. M., Kashner, T., Kupfer, D. J., Rosenbaum, J. F., Alpert, J., Stewart, J. W., McGrath, P. J., Biggs, M. M., Shores-Wilson, K., Lebowitz, B. D., Ritz, L., Niederehe, G., and for the STAR*D Investigators Group (2004). Sequenced treatment alternatives to relieve depression (STAR*D): rationale and design. *Controlled Clinical Trials*, 25(1):119–142.
- Saha-Chaudhuri, P. and Weinberg, C. (2017). Addressing data privacy in matched studies via virtual pooling. BMC Medical Research Methodology, 17(1):1–10.
- Salanti, G. (2012). Indirect and mixed-treatment comparison, network, or multipletreatments meta-analysis: many names, many benefits, many concerns for the next generation evidence synthesis tool. *Research Synthesis Methods*, 3(2):80–97.
- Saville, B. R., Connor, J. T., Ayers, G. D., and Alvarez, J. (2014). The utility of Bayesian predictive probabilities for interim monitoring of clinical trials. *Clinical Trials*, 11(4):485– 493.
- Schmitz, S., Maguire, A., Morris, J., Ruggeri, K., Haller, E., Kuhn, I., Leahy, J., Homer, N., Khan, A., Bowden, J., et al. (2018). The use of single armed observational data to closing

the gap in otherwise disconnected evidence networks: a network meta-analysis in multiple myeloma. *BMC Medical Research Methodology*, 18(66):1–18.

- Schnitzer, M. E., Steele, R. J., Bally, M., and Shrier, I. (2016). A causal inference approach to network meta-analysis. *Journal of Causal Inference*, 4(2):20160014.
- Schulz, J. and Moodie, E. E. M. (2021). Doubly robust estimation of optimal dosing strategies. Journal of the American Statistical Association, 116(533):256–268.
- Shen, J., Golchi, S., Moodie, E. E. M., and Benrimoh, D. (2022). Bayesian group sequential designs for cluster-randomized trials. *Stat*, 11(1):e487.
- Shen, J., Moodie, E. E. M., and Golchi, S. (2024). Sparse two-stage Bayesian meta-analysis for individualized treatments.
- Shi, C., Song, R., Lu, W., and Fu, B. (2018). Maximin projection learning for optimal treatment decision with heterogeneous individualized treatment effects. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(4):681–702.
- Shi, H. and Yin, G. (2019). Control of type I error rates in Bayesian sequential designs. Bayesian Analysis, 14(2):399 – 425.
- Simmonds, M. C. and Higgins, J. P. (2007). Covariate heterogeneity in meta-analysis: criteria for deciding between meta-regression and individual patient data. *Statistics in Medicine*, 26(15):2982–2999.
- Simoneau, G., Moodie, E. E. M., Nijjar, J. S., Platt, R. W., and the Scottish Early Rheumatoid Arthritis Inception Cohort Investigators (2020). Estimating optimal dynamic treatment regimes with survival outcomes. *Journal of the American Statistical Association*, 115(531):1531–1539.
- Smarr, K. L. and Keefer, A. L. (2020). Measures of depression and depressive symptoms. *Arthritis Care & Research*, 72(S10):608–629.

- Smith, C. T. and Williamson, P. R. (2007). A comparison of methods for fixed effects meta-analysis of individual patient data with time to event outcomes. *Clinical Trials*, 4(6):621–630.
- Spicker, D., Moodie, E. E. M., and Shortreed, S. M. (2024). Differentially private outcomeweighted learning for optimal dynamic treatment regime estimation. *Stat*, 13(1):e641.
- Splawa-Neyman, J., Dabrowska, D. M., and Speed, T. P. (1990). On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, 5(4):465–472.
- Stan Development Team (2020). RStan: the R interface to Stan. http://mc-stan.org/. R package version 2.21.2.
- Stan Development Team (2021). Stan modeling language users guide and reference manual. https://mc-stan.org/. version 2.28.
- Stevens, J. W., Fletcher, C., Downey, G., and Sutton, A. (2018). A review of methods for comparing treatments evaluated in studies that form disconnected networks of evidence. *Research Synthesis Methods*, 9(2):148–162.
- Sutton, R. S. and Barto, A. G. (2018). Reinforcement learning: an introduction. MIT press, Cambridge.
- Tanguay-Sela, M., Benrimoh, D., Popescu, C., Perez, T., Rollins, C., Snook, E., Lundrigan, E., Armstrong, C., Perlman, K., Fratila, R., Mehltretter, J., Israel, S., Champagne, M., Williams, J., Simard, J., Parikh, S. V., Karp, J. F., Heller, K., Linnaranta, O., Cardona, L. G., Turecki, G., and Margolese, H. C. (2022). Evaluating the perceived utility of an artificial intelligence-powered clinical decision support system for depression treatment using a simulation center. *Psychiatry Research*, 308:114336.

- Tchetgen Tchetgen, E. J., Robins, J. M., and Rotnitzky, A. (2010). On doubly robust estimation in a semiparametric odds ratio model. *Biometrika*, 97(1):171–180.
- Torgerson, D. J. (2001). Contamination in trials: is cluster randomisation the answer? British Medical Journal, 322(7282):355–357.
- Trivedi, M. H., Kocsis, J. H., Thase, M. E., Morris, D. W., Wisniewski, S. R., Leon, A. C., Gelenberg, A. J., Klein, D. N., Niederehe, G., Schatzberg, A. F., Ninan, P. T., and Keller, M. (2008). REVAMP research evaluating the value of augmenting medication with psychotherapy: rationale and design. *Psychopharmacology Bulletin*, 41(4):5–33.
- Trivedi, M. H., McGrath, P. J., Fava, M., Parsey, R. V., Kurian, B. T., Phillips, M. L., Oquendo, M. A., Bruder, G., Pizzagalli, D., Toups, M., Cooper, C., Adams, P., Weyandt, S., Morris, D. W., Grannemann, B. D., Ogden, R. T., Buckner, R., McInnis, M., Kraemer, H. C., Petkova, E., Carmody, T. J., and Weissman, M. M. (2016). Establishing moderators and biosignatures of antidepressant response in clinical care (EMBARC): rationale and design. Journal of Psychiatric Research, 78:11–23.
- Trivedi, M. H., Rush, A. J., Ibrahim, H. M., Carmody, T. J., Biggs, M. M., Suppes, T., Crismon, M. L., Shores-Wilson, K., Toprac, M. G., Dennehy, E. B., Witte, B., and Kashner, T. M. (2004). The inventory of depressive symptomatology, clinician rating (IDS-C) and self-report (IDS-SR), and the quick inventory of depressive symptomatology, clinician rating (QIDS-C) and self-report (QIDS-SR) in public sector patients with mood disorders: a psychometric evaluation. *Psychological Medicine*, 34(1):73–82.
- Turner, E. L., Li, F., Gallis, J. A., Prague, M., and Murray, D. M. (2017). Review of recent methodological developments in group-randomized trials: part 1—design. American Journal of Public Health, 107(6):907–915.
- van Erp, S., Oberski, D. L., and Mulder, J. (2019). Shrinkage priors for Bayesian penalized regression. *Journal of Mathematical Psychology*, 89:31–50.

- van Schie, S. and Moerbeek, M. (2014). Re-estimating sample size in cluster randomised trials with active recruitment within clusters. *Statistics in Medicine*, 33:3253–3268.
- Ventz, S. and Trippa, L. (2015). Bayesian designs and the control of frequentist characteristics: a practical solution. *Biometrics*, 71(1):218–226.
- Vyas, D. A., Eisenstein, L. G., and Jones, D. S. (2020). Hidden in plain sight reconsidering the use of race correction in clinical algorithms. New England Journal of Medicine, 383(9):874–882.
- Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228– 1242.
- Wallace, M. P. and Moodie, E. E. M. (2015). Doubly-robust dynamic treatment regimen estimation via weighted least squares. *Biometrics*, 71(3):636–644.
- Watkins, C. J. C. H. (1989). Learning from delayed rewards. PhD thesis, King's College, Cambridge, UK.
- Weinberger, M., Oddone, E. Z., Henderson, W. G., Smith, D. M., Huey, J., Giobbie-Hurder, A., and Feussner, J. R. (2001). Multisite randomized controlled trials in health services research: scientific challenges and operational issues. *Medical Care*, 39:627–634.
- White, I. R., Barrett, J. K., Jackson, D., and Higgins, J. P. T. (2012). Consistency and inconsistency in network meta-analysis: model estimation using multivariate meta-regression. *Research Synthesis Methods*, 3(2):111–125.
- Yin, G., Lam, C. K., and Shi, H. (2017). Bayesian randomized clinical trials: from fixed to adaptive design. *Contemporary Clinical Trials*, 59:77–86.
- Zhang, B., Tsiatis, A. A., Davidian, M., Zhang, M., and Laber, E. (2012). Estimating optimal treatment regimes from a classification perspective. *Stat*, 1(1):103–114.

- Zhang, C., Chen, J., Fu, H., He, X., Zhao, Y.-Q., and Liu, Y. (2020). Multicategory outcome weighted margin-based learning for estimating individualized treatment rules. *Statistica Sinica*, 30:1857–1879.
- Zhang, C., Mayo, M. S., Wick, J. A., and Gajewski, B. J. (2021). Designing and analyzing clinical trials for personalized medicine via Bayesian models. *Pharmaceutical Statistics*, 20(3):573–596.
- Zhao, Y., Zeng, D., Rush, A. J., and Kosorok, M. R. (2012). Estimating individualized treatment rules using outcome weighted learning. *Journal of the American Statistical* Association, 107(499):1106–1118.
- Zhou, T. and Ji, Y. (2023). On Bayesian sequential clinical trial designs. The New England Journal of Statistics in Data Science, 2(1):136–151.
- Zhou, X., Liu, S., Kim, E. S., Herbst, R. S., and Lee, J. J. (2008). Bayesian adaptive design for targeted therapy development in lung cancer — a step toward personalized medicine. *Clinical Trials*, 5(3):181–193.
- Zhou, X., Mayer-Hamblett, N., Khan, U., and Kosorok, M. R. (2017). Residual weighted learning for estimating individualized treatment rules. *Journal of the American Statistical* Association, 112(517):169–187.
- Zhu, H. and Yu, Q. (2017). A Bayesian sequential design using alpha spending function to control type I error. Statistical Methods in Medical Research, 26(5):2184–2196.
- Zhu, L., Yu, Q., and Mercante, D. E. (2019). A Bayesian sequential design for clinical trials with time-to-event outcomes. *Statistics in Biopharmaceutical Research*, 11(4):387–397.
- Zhu, R., Zhao, Y.-Q., Chen, G., Ma, S., and Zhao, H. (2017). Greedy outcome weighted tree learning of optimal personalized treatment rules. *Biometrics*, 73(2):391–400.

Zou, G. Y., Donner, A., and Klar, N. (2005). Group sequential methods for cluster randomization trials with binary outcomes. *Clinical Trials*, 2(6):479–487.