# Three Essays on Occupational Tasks, City Size Wage Gap, and Multidimensional Mismatch

Qi Xu

Department of Economics

McGill University, Montreal

December, 2020

*A thesis submitted to McGill University in partial fulfillment of the requirements of the degree of Ph.D. in Economics*

# Abstract

This dissertation is composed of three chapters that provide a comprehensive analysis of within-occupation task differences and worker-occupation mismatch in the U.S. The first chapter employs a unique U.S. online job posting data set and presents two interesting empirical facts: first, tasks of narrowly defined occupations change spatially, and occupations in large cities are more complex than in small cities; second, large cities pay 8.5% higher than small cities and around 2/3 of the wage premium is within occupation. The estimation results show that the within occupation task difference is the main source of the city size wage premium. Chapter 2 investigates the multidimensional skill mismatch of the male workers among two surveys of NLSY79 and NLSY97. I document the empirical facts that the aggregate mismatch rate is higher in NLSY97 than NLSY79 in three dimensions of verbal, math, and social; The mismatch rate varies across locations in both surveys. The last chapter studies the cyclicality of multidimensional skill mismatch between workers and their occupations in NLSY79 and NLSY97. I provide strong evidence that the mismatch rate is procyclical in NLSY97.

# Abrégé

Cette thèse est composée de trois chapitres qui fournissent une analyse complète de la différence entre les tâches au sein de la profession et de l'inadéquation entre les travailleurs et les professions aux États-Unis. Le premier chapitre utilise un ensemble de données unique sur les offres d'emploi en ligne aux États-Unis et présente deux faits empiriques intéressants: premièrement, les tâches de les professions restreintes changent dans l'espace et les professions dans les grandes villes sont plus complexes que dans les petites villes; Deuxièmement, les grandes villes paient 8,5% de plus que les petites villes et environ 2/3 de l'avantage salarial est au sein de la profession. Les résultats de l'estimation montrent que la différence entre les tâches au sein de la profession est la principale source de l'avantage salarial lié à la taille de la ville. Le chapitre 2 étudie l'inadéquation multidimensionnelle des compétences des travailleurs masculins dans deux enquêtes de NLSY79 et NLSY97. Je documente les faits empiriques selon lesquels le taux d'inadéquation global est plus élevé dans NLSY97 que NLSY79 dans trois dimensions verbale, mathématique et sociale; Le taux d'inadéquation varie d'un endroit à l'autre dans les deux enquêtes. Le dernier chapitre étudie la cyclicité de l'inadéquation multidimensionnelle des compétences entre les travailleurs et leurs professions dans NLSY79 et NLSY97. Je fournis des preuves solides que le taux d'inadéquation est procyclique dans NLSY97.

# Contribution of Authors

All three papers in this thesis are original and distinct contributions to the fields of Labour Economics. The first chapter is joint work with Ben Bradley from Burning Glass Technologies. He provided the U.S. online job posting data set. I analyzed the data set, built the model, performed the empirical analysis, and wrote the paper. Chapter 2 and chapter 3 are solely authored by me.

# Acknowledgements

I am incredibly grateful to my supervisor Professor Rui Castro for his support, mentorship, and guidance on my research. I would like to thank my co-supervisor Professor Theodore Papageorgiou and my committee member Professor Fabian Lange. I am deeply grateful for their insightful visions and valuable advice.

I thank Professor Francesco Amodio, Professor Francisco Alvarez-Cuadrado, Professor Saraswata Chaudhuri, and Professor Laura Lasio for their assistance in job market. I also would like to thank Professor Markus Poschke for his extensive advice in the macro advising group.

I would like to thank my friends Xintong Wang and Siyu Ma. I really appropriate their friendship and support.

Finally, I would like to thank my loving, caring, and supportive husband Jie Yang. I cannot overemphasize the support from him. His love, encouragement, and sacrifice helped me get through the difficult times.

# Introduction

A wealth of quantitative evidence documents that wages are higher in large cities than in small cities. Understanding such city size wage gap is one of the most important questions in labor and urban economics. In the first chapter of my dissertation, **Occupational Task Difference and the City Size Wage Premium**, I investigate the reasons behind city size wage premium. This chapter examines the relationship between occupational task inputs and the city size wage gap. Our analysis focuses on how the nature and complexity of an occupation vary across local labor markets of different sizes and how this change affects the city size wage gap. Our evidence uncovers that occupational tasks vary spatially, and occupations in large cities are more complex than in small cities. We also show that the occupational task difference plays an important role in explaining the city size wage gap.

Large cities pay 8.5% higher than small cities. Around 2/3 of the wage premium is within occupation and only 1/3 is explained by different occupations in large vs small cities. Our goal is to understand within occupation wage differentials. Using data on US online job postings (unique data from Burning Glass Technology), we investigate the structure of within occupation task requirements in large vs small cities. The main advantage of this data set is that it comprises ample job posting information, especially the job requirement keywords at various location levels, which enables us to outline the task requirement within occupation across cities. We are,

to the best of our knowledge, the first to build a spatial data set with comprehensive information on occupational task requirements within occupations across cities.

To process the task keywords in the enormous data set, we adopt the natural language processing method. We measure and classify the task keywords in job postings and construct a new data set detailing task compositions for 6-digit occupations. In this data set, we aggregate the task requirements at the job posting level into five categories at occupational levels. We find that task distributions do differ across cities. On average and for the same occupation, more tasks are required in large cities versus small cities. For three out of four 6-digit occupations, large cities demand more tasks than small cities. The shares of nonroutine analytic tasks, cognitive tasks, and interactive tasks are higher in large cities; small cities emphasize more manual tasks. Within detailed 6-digit occupations, aggregate tasks as well as task compositions are different across cities. We hypothesize that these occupational task differences are the main source of the within occupation city size wage gap

We develop a supply-demand framework and derive an explicit causal link between occupational wage gaps and task distributions, total factor productivities, and task biased productivities. Our model is estimated, employing both Cobb-Douglas and CES production functions. We use the data of five tasks and the average wage of each occupation in large and small cities for the years 2007 and 2010 through 2016 and estimate the parameters of the production function. Based on the estimation result, we decompose the city size wage gap and evaluate the ability of task inputs to interpret the wage gaps. We find that occupational task inputs account for half of occupational city size wage gaps.

This paper helps to understand the regional disparities, which have risen since the 1980s and have attracted greater interest. Our research suggests that jobs in large cities are significantly more complex in terms of task requirements and this task difference accounts for a large part of the wage differentials.

Labor market frictions cause inefficient assignments of workers to jobs. A growing body of literature focuses on measuring this 'mismatch' in the labor market. Chapter two, **A Comparison of Multidimensional Skill Mismatch between NLSY79 and NLSY97**, compares a measure of 3-dimensional skill mismatch between the two cohorts in NLSY79 and NLSY97. Following Guvenen et al. (2020), I infer the abilities of workers from NLSY79/NLSY97 and the occupational skill requirements from O*NET. I then measure the extent of mismatch between worker abilities and skill requirements in each cohort. I find that the mismatch rate in NLSY97 is higher than in NLSY79. In the new cohort, both the positive and negative mismatch rates are higher than the old cohort. I also find that the mismatch varies across locations in both cohorts. The mismatch rate is higher in rural areas than in urban areas. It's also higher in non-MSA (less densely-populated) areas than in MSA areas.

Chapter three, **The Cyclicality of Multidimensional Skill Mismatch: Evidence from NLSY79 and NLSY97**, studies the cyclicality of the multidimensional skill mismatch rate. I investigate the effect of labor market conditions on the average match quality of the worker-occupation pairs. I find that the average mismatch between workers and their occupations is procyclical in NLSY97, meaning that in economic downturns the mismatch rate is lower. This result suggests an important role for the so-called "cleansing effect" of recessions, i.e. that low-quality matches are destroyed and only high-quality matches survive in recessions. I also provide evidence that the procyclicality of mismatch is more pronounced for highly educated workers in NLSY97.

# Table of Contents

# List of Figures

# List of Tables

# Chapter 1

# Occupational Task Difference and City Size Wage Premium

Qi XU[1]
Ben BRADLEY[2]

## Abstract

This paper focuses on how the nature and complexity of an occupation vary across local labor markets of different sizes and how this change affects the city size wage gap. Using a unique U.S. online job posting data, we investigate the task requirements within detailed six-digit occupations in large and small cities. Our evidence shows that (i) more tasks are required in large vs small cities, and (ii) the shares of nonroutine analytic tasks, cognitive tasks, and interactive tasks are higher in large cities; small cities emphasize more manual tasks. Within the narrowly defined occupations, aggregate tasks as well as task compositions are therefore different across cities. Then we examine the relationship between the occupational city size wage

---

[1]Department of Economics, McGill University
[2]Burning Glass Technologies

gaps, task input differences, total factor productivity, and task biased productivities in a supply-demand framework. Based on the estimation and decomposition results, we show that occupational task differences and task biased productivities account for nearly half of the within occupation city size wage difference.

## 1.1   Introduction

A wealth of quantitative evidence documents that wages are higher in large cities than in small cities (Roback (1982), Glaeser and Mare (2001), Baum-Snow and Pavan (2011), and Roca and Puga (2017)). Our paper revisits this issue and confirms that large cities pay 8.51% higher wages than small cities.[3] We also find that 65% of the city size wage premium comes from within occupation wage differences. Nine out of ten occupations pay higher wages in large cities than small cities. Therefore, within occupation city size wage gaps play a big part in aggregate wage differences across cities.

This paper formally explores the importance of the mechanisms that may have generated wage gaps across cities of different sizes. In particular, we inspect the impact of variations in the nature of occupations across cities on the city size wage gap. We use task inputs to measure the nature of each occupation. To analyze the structure of occupational task requirements in local labor markets of different sizes, we utilize data on US online job postings from Burning Glass Technology. The main advantage of this data set is that it contains ample job posting information, especially the job requirement keywords at various locations, which enables us to examine the task requirement within occupation across cities. We are, to the best of our

---

[3]This essay is from a productivity perspective and focuses on nomial wages. In the future research about welfare, we will focus on the real wages.

2

knowledge, the first to build a spatial data set with comprehensive information on occupational task requirements within detailed occupations across cities.

To process the task keywords in the enormous data set, we adopt the natural language processing method from linguistics. We measure and classify the task keywords in job postings and construct a new data-set of task compositions for 6-digit SOC ( Standard Occupational Classification) occupations. We aggregate the task requirements at the job posting level into five categories at 6-digit occupational levels. We first find that, on average, more tasks are required in large cities versus small cities. For three out of four 6-digit occupations, large cities demand more tasks than small cities. Second, the shares of nonroutine analytic tasks, cognitive tasks, and interactive tasks are higher in large cities; small cities emphasize more manual tasks.

Having documented the patterns of occupational wage gaps and task compositions across cities, we hypothesize that these occupational task differences are the main source of the city size wage gap. In order to examine how much the city size wage gap can be explained by occupational task differences, we offer a supply-demand framework and derive an explicit causal link between occupational wage gaps, task distributions, TFPs, and task biased productivities. We think of the wage difference across local labor markets as being caused by the interaction of task inputs and factor demand schedules in the context of the task-based production process. This model employs either Cobb-Douglas or the CES production functions that incorporate five tasks as inputs of production, and allow for both total factor productivity and task-biased productivities.

In the estimation part, we use the data of the five tasks and the average wage of each occupation in large and small cities for the years 2007 and 2010 through 2016 and estimate the parameters of the production function. By decomposition, we assess the abilities of total factor productivity, task inputs, and task biased productivi-

3

ties to explain the city size wage gap. Our evidence indicates that occupational task differences across cities are the main source of the city size wage gap. In particular, we simulate the model by shutting down the channel of total factor productivity. In the framework with the Cobb-Douglas production function, only the occupational task inputs difference could explain 52.36% of the city size wage gap.

This paper contributes to several strands of literature. The first examines the existence of city size wage gaps. In Baum-Snow and Pavan (2011), the nominal wage premium for medium-sized cities over smaller areas is 7% if individual fixed effects are controlled for. Papageorgiou (2013) also confirms this fact that workers in more highly populated areas are paid significantly higher wages. In countries outside of the U.S., there is also evidence about wage premiums across cities. Roca and Puga (2017) use rich administrative data for Spain and show a positive relationship between earnings and cities. Workers in Madrid (the most populous city in Spain) earn 46% more than workers in Santiago de Compostela (the median-sized city). Combes et al. (2008) document that wages in Paris are on average 15% higher than in large French cities such as Lyon or Marseille and 35% higher than in mid-sized French cities. This paper verifies the existence of the city size wage gap in the U.S. Our evidence shows that large cities pay 8.5% higher wages than small cities. Differently from the literature, we look at the within occupation wage premium. Our results show that about 2/3 of the overall 8.5% wage premium between large and small cities is within occupation. This motivates us to seek the within-occupation reasons that could explain the city size wage gap.

The literature is also related to the literature on job/occupation tasks. Autor et al. (2003) start the discussion of task content for individual occupations. They document that from 1977 to 1991, there is less input of routine tasks and more input of nonroutine cognitive tasks within occupation and education groups. However, they exploit information from the 1977 and 1991 edition of *Dictionary of Occupational*

4

*Titles* (DOT), which is constant across cities and could only capture the occupational task variations at the two time points.

Spitz-Oener (2006) takes the exploration one step further. She shows how skill requirement changes within occupation, using a survey data set from Germany. She finds that occupations are more complex today than in 1979 as more tasks are required for individual occupations. This paper analyzes the within occupation skill variation over time. Motivated by this paper, our paper scrutinizes the occupational task levels and compositions in different locations. This paper follows the task classifications in her paper and measures aggregate and constructs a spatial data set detailing occupational task compositions across locations.

The third strand literature is about task/skill distributions across cities. Eeckhout et al. (2014) discover that the average skills of workers are constant across city size but the distribution of skills varies across cities. In their paper, only one-dimensional instead of multi-dimensional skills are measured. They do not consider the within-occupation task composition disparity between large and small cities. The evidence presented in our paper improves considerably on the evidence offered by Eeckhout et al. (2014) as we examine the task compositions variations within very detailed occupations across cities.

Finally, this paper contributes to the literature using online job posting data to analyze the labor market. More and more scholars like Papageorgiou (2013), Deming and Kahn (2018), and Hershbein and Kahn (2018) employ Burning Glass online data to dig the copious and novel job information.

This paper is organized into six sections. In the next section, we discuss the data set and measurement of variables. In section 1.3, we present the patterns of city size wage premium and occupational task distributions. In Section 1.4, we lay out the theoretical framework. In sections 1.5 and 1.6, we discuss the estimation approach and the empirical results. In section 1.7, we conclude.

## 1.2 Data Set and Measurement of Variables

### 1.2.1 Data Set

The primary data set of this paper is U.S. online job posting data provided by Burning Glass Technology (hereafter BGT). BGT collects and processes information from nearly 40,000 online websites including job boards and company postings, which includes the near-universe of all online vacancy postings. They parsed and developed the job descriptions into a series of systematic and user-friendly tables. The sample used in this paper is for the years 2007 and 2010 through 2016. The data for years 2008 and 2009 is current unavailable from Burning Glass Technology.

BGT collects information on job title, job date, industry and occupation identifiers (6-digit Standard Occupational Classification (SOC) codes), geography variables, education requirement, salary and job type, and most importantly, keywords about task requirements for each vacancy. Every job advertisement comes with several keywords as task requirement, such as problem solving, detail-oriented, supervisory, etc. These keywords describe the tasks workers should conduct for a particular job and occupation. In 2016, more than $23,000,000$ jobs were posted and there were more than 12,000 different keywords for all the postings. On average, each posting demands 9 keywords to describe the activities that workers should perform for this job.

The second data set used in this paper is from Occupational Employment Statistics (OES) of the Bureau of Labor Statistics (BLS). OES produces employment and wage estimates annually for the 6-digit SOC occupations in different locations. We obtain the average wages for each occupation by location in the years of 2007, 2010-2016.

Both data-sets provide the geographical information of the jobs and occupations. In this paper, we use Core Based Statistics Area (CBSA), which refers collectively to

6

metropolitan and micropolitan statistical areas, as the geographical indicator. Each metropolitan statistical area contains at least one urbanized area of 50,000 or more inhabitants. Each micropolitan statistical area contains at least one urban cluster of at least 10,000 but less than 50,000 population. CBSAs are defined in terms of population as well as density. The metropolitan areas include relatively more concentrated areas than micropolitan areas. We refer to a metropolitan statistical area as a large city and a micropolitan statistical area as a small city. As shown in Figure 1.1, red bars represent small cities and blue bars are large cities. The overlap refers to the metropolitan areas with a lower population but a larger density, or the micropolitan areas with a higher population but a smaller density.



**Figure 1.1:** Sizes of Large and Small Cities

Table 1.1 presents the summary statistics about online job postings in 2016. There are 366 metropolitan statistic areas and 576 micropolitan statistic areas. The number of job postings in metro areas is also 15 times more than in micro areas. The average

population of metropolitan areas ("large cities") is around 700,000. The average
population of micropolitan areas ("small cities") is 53,000. As shown in Table 1.2,
"communicate" is the most popular keyword in all postings, as it appears in more
than 28% of all ads. The other top 10 keywords are "writing", "customer service",
"organisation", "sales", "M.S. excel", "problem solving", "planning", "team work",
"scheduling". There is a slight difference between the top 10 keywords in large and
small cities. 'Physical Demand' appears in 12.08% of all ads in small cities but is not
the top 10 keywords in large cities. In general, these 10 keywords are more about
"interacting" and "analyzing" and not about "manual work", which leads us to
investigate the task structure of job postings: what kind of tasks do they emphasize?
Does task structure vary within occupations? Do large and small cities highlight
different tasks? These questions will be answered in the next section.

**Table 1.1:** Job Posts Statistics in 2016

|  | Metro Stat Area | Micro Stat Area |
| --- | --- | --- |
| Number of Areas | 366 | 576 |
| Job Posts | 22,283,752 | 1,549,445 |
| Job Posts per capita | 0.086 | 0.053 |
| Avg Popupation | 705,000 | 53,000 |

Compared to the Dictionary of Occupational Titles (DOT), which defines the task
contents of occupations and the definitions are invariable across locations, the most
important advantage of BGT data is that it enables us to track the task requirements
within occupation and at the job vacancy level. The geographical variables provide
us a chance to see the variations of task contents within occupation and across lo-
cations. This property makes Burning Glass data a perfect source to analyze the
task contents, the complexity, and even the nature of an occupation in different local
labor markets.

**Table 1.2:** Top 10 Task Keywords in 2016

| Top 10 Keywords | Freqency | Top 10 Keywords in Large Cities | Frequency | Top 10 Keywords in Small Cities | Frequency |
|---|---|---|---|---|---|
| Communicate | 28.84% | Communicate | 29.88% | Communicate | 21.17% |
| Writing | 17.34% | Writing | 18.00% | Customer Srvc | 16.23% |
| Customer Srvc | 15.31% | Customer Srvc | 15.46% | Writing | 12.31% |
| Organisation | 12.39% | Organisation | 12.98% | Physical Demand | 12.08% |
| Sales | 12.02% | M.S. Excel | 12.84% | Organisation | 11.06% |
| M.S. Excel | 11.85% | Sales | 12.04% | Sales | 10.93% |
| Problem Solve | 11.03% | Problem Solve | 11.61% | Scheduling | 9.99% |
| Planning | 10.51% | Planning | 11.13% | Computer Skills | 8.91% |
| Team Work | 10.49% | Teamwork | 10.83% | Supervisory Skills | 8.76% |
| Scheduling | 10.07% | Scheduling | 10.25% | Retail Setting | 8.34% |

## 1.2.2 Measurement of Occupational Tasks

The main advantage of BGT data is that it includes keywords from the text descriptions in advertisements as well as geographical variables. For example, in a posting of 2016, the keywords for one job of plumbers are: 'Repair', 'Hand Tools', 'Inspection', and 'Test Equipment'. Another posting for plumbers demands 'Piping Repair', 'Plumbing Maintenance', 'Customer Service', and 'Bilingual' to fulfill the job duties. These keywords contain the necessary information for a job, such as, skills, work styles, languages, etc., which indicate what kind of tasks workers should perform on the job.

In this section, we use the BGT data to construct a new data set for the task compositions of six-digit Standard Occupational Classification (SOC) occupations. 6-digit occupations are not at a general level but belong to a very detailed classification. These 6-digit occupations refer to specialized jobs. For example, drivers under

this classification system include heavy truck drivers, bus drivers, taxi drivers, light truck drivers, etc. The occupation plumber is also categorized as plumbers and helpers of plumbers. Therefore, the task keywords in the posting of each 6-digit occupation could reflect its task requirement to a considerable extent. Another feature of this data set is that it captures variations of occupational tasks across time (Table 1.17 in the Appendix shows the task compositions in large and small cities for 2007 and 2016) and locations. We will discuss the within occupational task difference across locations for 2016 in next the section and compare it with the empirical patterns of 2007 in the Appendix.

**Job Task Measurement**

Before quantifying occupational tasks, we first measure the task contents of each job posting based on the task classification in the literature. Autor et al. (2003) proposes a basic category to classify tasks into five groups: nonroutine analytic task, nonroutine interactive task, routine cognitive task, routine manual task, and nonroutine manual task. The abbreviations NA, NI, RC, RM, NM are used for simplicity. According to Spitz-Oener (2006), each of the five task groups is assigned with a sequence of keywords that describe the activities of that task (see Table 1.3).

In the BGT data, every job is characterized by a series of keywords. To evaluate the task composition of each posting based on the classification of Spitz-Oener (2006), we need to map the keywords of each job posting in BGT data onto the five categories.

10

**Table 1.3:** Task Classification

| Five Categories | Basic Words |
|---|---|
| Non-routine Analytic | ′researching′, ′analyzing′, ′evaluating′, ′planning′, ′designing′, ′sketching′, ′research′ ′devising rule′, ′interpreting rule′ |
| Non-routine Interactive | ′negotiating′, ′lobbying′, ′coordinating′, ′organizing′, ′teaching′, ′selling′, ′buying′ ′advertising′, ′entertaining′, ′presenting′, ′managing′, ′advising′, |
| Routine Cognitive | ′calculating′, ′bookkeeping′, ′correcting′, ′measuring′, ′calculate′, ′corrections′, ′measurement′ |
| Routine Manual | ′operating′, ′controlling machines′, ′equipping machines′ |
| Non-routine Manual | ′repairing′, ′renovating houses/apartments/machines/vehicles′ ′restoring art/monuments′, ′serving′, ′accommodating′ |

Source: Spitz-Oener (2006)

We employ Natural Language Processing method from computer science and artificial intelligence to measure the task keywords of each job.

Step 1: Measure the semantic distance (semantic similarity) between keywords in BGT data and the basic keywords in Spitz-Oener (2006). The Wu-Palmer Similarity is introduced in this step. It returns a score denoting to what extent the senses of two words are similar, based on the depth of the two senses in WordNet [4] taxonomy, and that of their Least Common Subsumer: $Sim_{WP} = 2 * \frac{N}{N_1 + N_2}$. The similarity score $0 < Sim_{WP} <= 1$. The score is 1 if the two words are the same. The larger the score is, the more similar the two words are.



where least common subsumer (LCS) is the most specific common ancestor deepest in the taxonomy of two words ($C_1$ and $C_2$). Semantically, it represents the commonality of $C_1$ and $C_2$.R (in the above figure) is the ontology root. In Wu-Palmer Similarity, $N_1$ and $N_2$ are the numbers of arcs between the words $C_1$, $C_2$ and the ontology root R. N is the number of arcs between the LCS and the ontology root R. (Wu and Palmer, 1994).

For example, $C_1$ is 'automobile', $C_2$ is 'boat', LCS is 'vehicle', and R is 'object'. The LCS of "boat" and "automobile" is "vehicle", which is the most recent ancestor of $C_1$ and $C_2$ in Wordnet taxonomy. Using the Natural Language Toolkit (NLTK)

---

[4]WordNet is a large lexical database of English developed by Princeton University. In WordNet, words are grouped into sets of cognitive synonyms (synsets).

of Python, we find that the Wu-Palmer similarity score between 'automobile' and 'boat' is 0.696.

Select the keyword "communicate" in one posting. We calculate the similarities between this keyword and all the basic words in the five categories. For example, for the words "communicate" and "researching" (basic word of Non-routine Analytic task in Table 1.3), we find that their Wu-Palmer similarity score is 0.13. For the words "communicate" and "organizing" (basic words of Non-routine Interactive task in Table 1.3), the similarity score is 0.25. These two similarity scores indicate that "communicate" is more semantically related to "organize" than to "research".

Step 2: Having got the similarity scores between each keyword in all the postings and the basic words in the five categories, we need to classify the keywords into the basic groups. Now we choose the keyword "instruct" as an example. From step 1, we obtain the Wu-Palmer similarities between "instruct" and every basic word in NA, NI, RC, RM, and NM. Then, we define the similarity score between "instruct" and the five groups as the highest score within each group. As a result, the largest semantic similarity between "instruct" and NA group is 0.25; the largest semantic similarity between "instruct" and NI group is 0.73; the largest semantic similarity between "instruct" and RC group is 0.2; the largest semantic similarity between "instruct" and RM group is 0.18; the largest semantic similarity between "instruct" and NM group is 0.18. Thus, "instruct" enters the group of the Nonroutine Interactive task. In the end, the keywords of each job posting get into one of the five task categories.

Step 3: Based on the similarity scores of the keywords, we count the numbers of five tasks for each job posting. $N_{Xp}$ ($X = NA,\ NI,\ RC,\ RM,\ NM$) denotes the number of $X$ task that a job $p$ demands.

$$N_{Xp} = number\ of\ keywords\ for\ job\ p\ that\ are\ in\ task\ group\ X$$

For example, if the keywords of one posting are: "communication", "edit", and "planning", then $N_{NA} = 2$, $N_{NI} = 1$, $N_{RC} = 0$, $N_{RM} = 0$, $N_{Nm} = 0$. This job needs two nonroutine analytic tasks, one nonroutine interactive task, and doesn't need other tasks. In the end, each posting is endowed with a measurement. This five-dimensional vector measures the amounts of five tasks that each job requires.

**Occupational Task Measurement**

Next, we construct the measurement of task composition for each occupation. In total, there are 840 6-digit occupations. Each job posting is mapped to an occupation. We define the number of $X$ ($X = NA, NI, RC, RM, NM$) tasks for each occupation $j$ as the average number of tasks per vacancy under that occupation:

$$N_{Xj} = \frac{\sum N_{Xp}}{number\ of\ vacancies\ for\ each\ occupation\ j}$$

## 1.3 City Size Wage Premium and Occupational Task Distributions

In this section, we measure task inputs for the detailed occupations, in both large and small cities. We also document a series of empirical facts of the occupational task distributions. Additionally, we investigate the city size wage gap and zoom in to the within occupation wage difference.

### 1.3.1 Occupational City Size Wage Premium

The existence of city size wage gap has been well documented by Roback (1982), Glaeser and Mare (2001), and Baum-Snow and Pavan (2011). We revisit this issue from the perspective of occupational wages and obtain a similar outcome. In 2016, the average hourly wage in large cities is 8.51% higher than in small cities.

The average wage gap may result from differences along two margins. The first margin is the different occupational structure of employment in large and small cities. Large cities may have more highly-paid occupations, such as doctors, engineers, and lawyers. The second margin is the wage difference between cities within occupations. For the same occupations, like a plumber, large cities pay higher than small cities. The city size wage gap is shifted through the two channels. Intuitively, within occupation effect means the average wage gap between large and small cities if we keep the occupational employment distribution the same across cities; between occupation effect is the average wage gap when we allow the occupational employment distribution to differ across cities but the occupational wage to be constant. We use the Marshall-Edgeworth-Type decomposition to quantify the levels of the two margins. $W_L$ and $W_S$ are weighted average wages in large and small cities. $o$ denotes occupation. $\omega_{oL}$ and $\omega_{oS}$ are the weights of occupation $o$'s employ-

15

ment in large and small cities. The wage gap is decomposed into two parts: within occupation effect reflects the within occupation wage difference between cities; between occupation effect indicates the wage gap attributable to the difference in the occupational distribution of employment.

$$W_L - W_S = \underbrace{\sum (W_{oL} - W_{oS}) \frac{\omega_{oL} + \omega_{oS}}{2}}_{Within\ Occ\ Effect} + \underbrace{\sum (\omega_{oL} - \omega_{oS}) \frac{W_{oL} + W_{oS}}{2}}_{Btw\ Occ\ Effect}$$

Table 1.4 shows the wage premium and the decomposition in 2016. The average hourly wage in a large city is 21.66 and in a small city is 19.01. The weighted wage premium (Weight is the employment of each occupation) is 8.51%.

**Table 1.4:** Decomposition of City Size Wage Premium in 2016

|  | 2016 | | |
|---|---|---|---|
|  | Large City | Small City | Premium |
| Overall | 21.66 | 19.01 | 8.51% |
| Btw Occ | –– | –– | 34.97% |
| Within Occ | –– | –– | 65.03% |

The result of this table evinces the prominent city size wage premium. At the same time, within occupation effect plays a major role in pushing up the premium. Around 65% of the city size wage gap results from the within occupation wage difference.[5] If we assume the occupational employment distribution is the same across cities, i.e., there are the same numbers of doctors, lawyers, and engineers in large and small cities, wage differences within doctors (or lawyers, or engineers)

---

[5]I also conduct the decomposition analysis by finer gradations and the results are consistent with the original decomposition analysis: extra large cities ($population >= 1,000,000$), large cities ($300,000 =< population < 1,000,000$), medium cities ($50,000 < population < 300,000$). The wage gap between extra large and large cities is 7.57%, among which 74% is within occupation. The wage gap between large and medium cities is 3.42%, among which 72% is within occupation.

account for 65% of the average wage gap. Within occupation wage difference is the first important motivation for this paper. Next, we look deeper to investigate the wage gap within each detailed occupation.

The information in Figure 1.2 describes the fact that there's a wage gap between cities within the 6-digit detailed occupations. Each red bar represents the occupation's wage gap between large and small cities. Almost nine out of ten occupations pay higher wages in large cities than in small cities.

Table 1.5 lists the top 10 occupations with the highest positive wage gaps and bottom 10 occupations with the highest negative wage gaps. We omit the occupations with few employments in large or small cities, and the occupations listed in the table are all with more than 80 of employment in both large and small cities. For the top 10 occupations, wages are higher in large cities; for the bottom 10 occupations, wages are higher in small cities. What is the reason behind the within occupation wage difference across cities? When we compare the occupations in the upper panel and the lower panel, like, historians VS mining machine operators, and real estate brokers VS plasterers and stucco masons, the intuitive impression we have from the titles is that the way to fulfill job duties is different between top 10 and bottom 10 occupations. According to the SOC definition, the top 10 occupations require tasks like research, analyze, interpret, plan, create, sell, arrange, etc. These are related to analytic and interactive tasks. For the bottom 10 occupations, the tasks in the SOC definition are operate, load, drill, shaper, perform manual labor, etc, which seem to be manual tasks. In the next section, we will probe deeper, using the abundant online job posting data to inspect the task contents of the detailed occupations in both large and small cities.

17

**Figure 1.2:** City Size Wage Premium for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the wage gap in 2016.

**Table 1.5:** Occupations with Highest Positive and Negative Wage Gaps in 2016

| Top 10 |
|---|
| Historians; Choreographers; Administrative Law Judges, Adjudicators; Designers; Advertising and Promotions Managers; Camera Operators; Real Estate Brokers; Radio and Television Announcers; Fire inspectors and investigators; Actors |

| Bottom 10 |
|---|
| Mining Machine Operators; Woodworkers; Manufactured building and mobile home installers Nuclear engineers; Tank car, truck, and ship loaders; Explosives workers, ordnance handling experts, and blasters; Continuous mining machine operators; Plasterers and stucco masons; Forest and conservation workers; Mine cutting and channeling machine operators |

*Notes: Occupations listed above are with significant numbers of employment*

## 1.3.2 Occupational Task Compositions

In this paper, we utilize the near-universal U.S. online job posting data to construct a new spatial data set with the measurements for occupational task compositions. We find out that the occupational task distributions vary across locations.

**Table 1.6:** Tasks of Plumbers, Pipefitters, and Steamfitters in Large and Small cities

| 47-2152: Plumbers, Pipefitters, and Steamfitters | |
|---|---|
| Number of Tasks: 104 in large cities VS 28 in small cities | |
| Type of Tasks | |
| Common Tasks | 'Repair' 'Hand Tools' 'Inspection' 'Test Equipment' 'Schematic Diagrams' 'Plumbing' 'Soldering' 'Piping Repair' 'Plumbing Maintenance' 'Plumbing Repairs' 'Valve Installation' 'Problem Solving' 'Project Management' 'Communication Skills' 'Welding' 'Piping Systems' 'PipeFitting' 'Pipe Installation', etc. |
| Only in Large Cities | 'Bilingual' 'Teaching' 'Music' 'Microsoft Excel' 'Microsoft Outlook' 'Microsoft Word' 'Machine Tools' 'Welding Equipment' 'Microsoft 'Project' 'Leadership' 'Customer Contact', etc. |

We use an example to start our exploration. The occupation is "Plumbers, Pipefitters, and Steamfitters" (hereafter plumbers). The task requirements of the online job postings are summarized in Table 1.6. The plumbers conduct more tasks in large cities. In large cities, the number of unique task keywords in all postings of "Plumbers" is 104, while in small cities the number is only 28. On top of that, the task types of the same occupation, plumbers, vary across cities. The common

**Table 1.7:** Tasks in Large and Small Cities of 2016

|  | Large City | Small City | Task Gap |
|---|---|---|---|
| Overall | 8.47 | 6.93 | 1.54 |
| Nonroutine Analytic | 1.37 (**16.17%**) | 1.07 (15.44%) | 0.30 |
| Nonroutine Interactive | 1.10 (**12.99%**) | 0.83 (11.98%) | 0.27 |
| Routine Cognitive | 1.19 (**14.05%**) | 0.83 (12.12%) | 0.35 |
| Nonroutine Manual | 1.63 (19.24%) | 1.36 (**19.62%**) | 0.27 |
| Routine Manual | 3.17 (37.43%) | 2.83 (**40.84%**) | 0.34 |

tasks, like "repair", "plumbing", "installation", appear in both large and small cities. These tasks are traditional tasks for a plumber. In addition to the common tasks, there are special task keywords such as "bilingual", and "Microsoft" that only exist in large cities. The task content difference of "plumber" points out that a plumber S in a small city tends to conduct traditional activities on the job. Another plumber L, who works in a large city, needs to tackle a situation using distinct abilities from conventional ones. Imagine that Plumber S provides pipe plumbing services for the households in a community. The main tasks are "plumbing" and "communication." Meanwhile, Plumber L works in a big complex in Montreal. He should be bilingual, coordinate with other plumbers, and use some software on a computer to record the cleaning progress. This example reflects that the task requirements of an occupation with the same title could be different in different cities. Thus, occupations in large cities tend to be more complicated.

The average number of tasks per vacancy within occupation is 8.47 in large cities and 6.93 in small cities. Table 1.7 shows the aggregate task levels and task compositions in large and small cities. The average number of tasks a job vacancy entails is 1.54 more tasks in large cities. Furthermore, we also change the cutoff of large vs small cities and prove that our result is robust.

If we look in detail at Table 1.7, we find two pieces of information. First, occupational task requirements vary across large and small cities. The numbers of the

aggregate task and the five compositions per posting are all larger in large cities than in small cities. The last column shows the positive gaps. Second, the compositions, ie. the structure of five tasks differ across locations. The share of nonmanual (NA, NI, and RC) tasks in large cities is 44%, compared to 39% in small cities. Manual tasks (NM and RM) dominate in small cities with a proportion of 61%.

Table 1.8 shows the decomposition of task gaps into between occupation effect and within occupation effect, using the Marshall-Edgeworth-Type decomposition method. It presents the two effects for the overall task as well as the five compositions. Similar to the decomposition of wage gaps between large and small cities, the two channels that affect the overall task gaps are: task differences within occupation and occupational distribution distinction across cities. 49% of the overall task gap is attributable to the within occupation effect. The within occupation effect is more important for the nonroutine analytic and the routine manual tasks.

**Table 1.8:** Shift-Share Analysis of Task Gaps Between Large and Small Cities in 2016

|  | Total Gap | Btw-Occ | Within-Occ |
|---|---|---|---|
| Overall | 1.54 | 50.96% | 49.04% |
| Nonroutine Analytic | 0.31 | 55.33% | 44.67% |
| Nonroutine Interactive | 0.27 | 60.43% | 39.57% |
| Routine Cognitive | 0.35 | 59.50% | 40.50% |
| Nonroutine Manual | 0.27 | 67.97% | 32.03% |
| Routine Manual | 0.34 | 53.89% | 46.11% |

Overall, within occupation aggregate and compositional task differences are the substantial driving forces behind the city size task gap. Even though the occupational identifiers in our data set are at the elaborate 6-digit classification level, they still cannot fully and correctly demonstrate the essential characteristics of the occu-

pations. The significant within occupation effect suggests that the nominal identical occupations in different locations have various task compositions. In fact, jobs with exactly the same titles might involve distinctive tasks if they are posted in different cities. In reality, the nature of an occupation has evolved along with the development of cities.

Nearly 50% of the city size task gap is caused by within occupation differences. Figure 1.3 exhibits some intriguing features about within occupation task distributions. This figure is sorted by values of occupational task differences. Three out of four occupations have positive task gaps between large and small cities. When we look deeper at the five task differences as shown in Figure 1.11-Figure 1.15 in the Appendix, we find that for the majority of the occupations, the five tasks are more required in large cities than small cities.

Table 1.9 displays the top 10 occupations with the highest task gaps and the bottom 10 occupations with the lowest task gaps. Similar to Table 1.5, we exclude the occupations with less than 80 vacancies in both large and small cities. Again we could compare the occupations in the upper panel and lower panel. In the top 10 occupations, there are experts in science, workers conducting crushing machines, managers, practitioners, coordinators, and supervisors. In the bottom 10 occupations, we have workers conducting multiple machines, maintenance workers, anthropologists and archeologists; electronics engineers (except computer), technicians and repairers. Even though the distinction between the occupations in the two panels is not as obvious as the distinction in Table 1.5, we still get the impression that the tasks conducted in the top 10 occupations are more about analytic, cognitive, and interactive activities, and the tasks performed in the bottom 10 occupations are more related to manual tasks.

**Figure 1.3:** City Size Task Difference for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the occupational task gaps.

**Table 1.9:** Occupations with Highest Positive and Negative Task Gaps in 2016

| Top 10 |
| --- |
| Biochemists and Biophysicists; Crushing, Grinding, and Polishing Machine Setters, Operators, and Tenders; Advertising and Promotions Managers; Purchasing Managers; Health Diagnosing and Treating Practitioners; Instructional Coordinators; Medical Scientists (except epidemiologists); Pesticide Handlers, Sprayers, and Applicators (vegetation); First-Line Supervisors of Farming, Fishing, and Forestry Work; Natural Sciences Managers; |

| Bottom 10 |
| --- |
| Zoologists and Wildlife Biologists; Multiple Machine Tool Setters, Operators, and Tenders; Soil and Plant Scientists; Maintenance Workers (machinery); Anthropologists and Archeologists; Electronics Engineers (except computer); Occupational Health and Safety Technicians; Telecommunications Equipment Installers and Repairers (except line installers); Demonstrators and Product Promoters; Home Appliance Repairers |

*Notes: Occupations listed above are with significant numbers of vacancies*

Having documented the patterns of occupational city size wage gap and task difference, we wonder whether these two patterns are related somehow. Next, we show two simple scatters to roughly describe the effect of tasks on the wage gaps. In Figure 1.4, we find that the occupational city size wage gap is positively related to total task difference and nonmanual task difference. The raw correlation coefficient for the left figure is 0.008. The raw correlation coefficient for the right figure is 0.016. Both are significant at the 95% confidence interval. From these figures, we see a simple relationship that occupations with larger task differences and larger nonmanual task differences across cities tend to pay higher wages in large cities.



**Figure 1.4:** Scatters of Wage Gap VS Task Difference and Wage Gap VS Nonmanual Task Difference

Notes: The associates in both cases are significant at the 95% confidence interval.

## 1.4 Theoretical Framework

The preceding sections have discussed the patterns of occupational wages and tasks in both large and small cities. Our evidence has unveiled the compelling facts that the occupational city size wage gap exists and has gone hand in hand with a considerable gap in task requirements. On top of that, the wage gap is positively related to total task differences and nonmanual task differences. For the occupations with larger task differences and higher nonmanual task differences, their wage gaps are more likely to be larger.

In order to examine how much the city size wage gap can be explained by occupational task difference, we offer a supply-demand framework and propose the mechanisms that could generate variations in occupational wage gaps across local labor markets of different sizes. In particular, we derive an explicit causal link between occupational wage gaps and task distributions, TFPs, and task biased productivities.

We begin our exploration of the occupational city size wage gap with a simple supply-demand framework. We think of the wage difference across local labor markets as being caused by the interaction of task inputs and its associated factor demand schedules in the context of the task-based production process. This is a partial equilibrium analysis as we do not identify the determinant of task supplies. We take the occupational task distributions as given.

The framework involves a task-based production of $K$ task inputs. Output by occupation $j$ in city $c \in \{small, \, large\}$ is denoted by the task-based production function:

$$Y_{jt}^c = F(\mathbf{X_{jt}^c}, \mathbf{Z_{jt}^c})^\theta \tag{1.1}$$

where

$\mathbf{X_{jt}^c}$ : vector of task inputs $(5 \times 1)$

27

$\mathbf{Z^c_{jt}}$ : vector of demand shifters ($5 \times 1$)

$\theta$: degree of homogeneity.

Individual workers are characterized by $(\mathbf{x^c_{ijt}}; h^c_{ijt})$. $\mathbf{x^c_{ijt}}$ is a $5 \times 1$ vector and denotes the endowments of worker $i$ in occupation $j$ at time $t$ in city $c$. $h^c_{ijt}$ is the working hour of worker $i$ in occupation $j$.

We assume that the market is competitive, where each occupation minimizes cost over all factors (tasks). Accordingly, the earnings of each individual are the sum of task returns:

$$\omega^c_{ijt} = \rho'_{\mathbf{jt}} \mathbf{x^c_{ijt}} h^c_{ijt} \tag{1.2}$$

where $\rho^c_{\mathbf{jt}} = D(\mathbf{X^c_{jt}}, \mathbf{Z^c_{jt}})$, denoting the task demand ($5 \times 1$).

In aggregate, all workers' endowment equals to total task input:

$$\sum \mathbf{x^c_{ijt}} h^c_{ijt} = \mathbf{X^c_{jt}} \tag{1.3}$$

$$\sum h^c_{ijt} = H^c_{jt} \tag{1.4}$$

(1.4) is the sum of all workers' working hours for a specific occupation.

The wage of an occupation is the average earning of workers in that occupation:

$$W^c_{jt} = \frac{\sum \omega^c_{ijt}}{\sum h^c_{ijt}} \tag{1.5}$$

By combining (1.2), (1.3), (1.4), and (1.5), we link the occupational wage with task inputs, demand shifters, and the total working hours for that occupation, and obtain equation (1.6), which will be estimated in the next section.

$$W^c_{jt} = \frac{\theta (Y^c_{jt})^{\frac{1}{\theta}}}{H^c_{jt}} \tag{1.6}$$

28

We introduce $\theta$ to denote the degree of homogeneity of the production function, which allows the possibility of decreasing returns to scale, as we observe a congestion effect $H_j^c$ in the wage equation 1.6. With this framework, our purpose is to give an explanation for the city size wage gap based upon task inputs gaps and productivity differences. We assume that the production function is identical across all occupations, cities, and over time, in terms of the degree of homogeneity and elasticity parameters. By estimating Equation 1.6, we infer the parameters of the production function. Based on the estimation result, the next step is to evaluate how much the wage gap could be explained by task differences across cities. Furthermore, we will evaluate the ability of the five task differences to interpret the city size wage gap.

We start with a Cobb-Douglas production function:

$$Y_{jt}^c = Z_{jt}^c[(NA_{jt}^c)^{\alpha_1}(RM_{jt}^c)^{\alpha_2}(NI_{jt}^c)^{\alpha_3}(NM_{jt}^c)^{\alpha_4}(RC_{jt}^c)^{\alpha_5}]^\theta \tag{1.7}$$

where $Z_{jt}^c$ is the total factor productivity of occupation $j$ in city $c$ at time $t$. $\alpha_1 - \alpha_5$ are output elasticities and $\sum_i^5 \alpha_i = 1$. The parameters $\alpha_i$ are constant over time and across cities. $NA_{jt}^c$, $RM_{jt}^c$, $NI_{jt}^c$, $NM_{jt}^c$ and $RC_{jt}^c$ are the inputs of nonroutine analytic task, routine manual task, nonroutine interactive task, nonroutine manual task and routine cognitive task. This specification imposes a unitary cross-input substitution elasticity and excludes a role of task-biased productivities and allows us to focus on the examination of how far one can go toward explaining the city size wage gap simply using TFP and task inputs.

In addition, we also introduce a generalized CES production function 1.8. The elasticities of substitution between inputs are $\frac{1}{1-\sigma}$. If $\sigma = 0$, the function degenerates to a Cobb-Douglas function. In addition to total factor productivity $Z_{jt}^c$ and the five task inputs, CES production function allows for task-biased productivities. The important feature of CES is that it characterizes the task-biased productivities $\lambda_1$, $\lambda_2$,

$\lambda_3, \lambda_4, \lambda_5$. Each of these productivities favors a specific task. The effect of certain task gaps on wage gaps could be magnified by the corresponding task-biased productivities. The use of the CES production function emphasizes the role of productivities towards different tasks, which can be recovered from estimations. This specification enables us to check how the productivities towards different tasks impact on the wage gaps.

$$Y_{jt}^c = Z_{jt}^c[\lambda_1(NA_{jt}^c)^\sigma + \lambda_2(RM_{jt}^c)^\sigma + \lambda_3(NI_{jt}^c)^\sigma + \lambda_4(NM_{jt}^c)^\sigma + \lambda_5(RC_{jt}^c)^\sigma]^{\frac{\theta}{\sigma}} \quad (1.8)$$

## 1.5  Estimation

This section outlines how we estimate the parameters of the model presented in the previous section and then discuss how the model captures the mechanisms behind the occupational city size wage gaps.

The general function we will estimate is Equation 1.6. First, we employ Equation 1.7 to estimate the model under the assumption of Cobb-Douglas production function. By plugging Equation 1.7 into the wage function 1.6, we derive that:

$$\begin{aligned}
lnW_{jt}^c &= ln\theta + lnZ_{jt}^c + \alpha_1\theta ln(\frac{NA_{jt}^c}{H_{jt}^c}) + \alpha_2\theta ln(\frac{RM_{jt}^c}{H_{jt}^c}) + \alpha_3\theta ln(\frac{NI_{jt}^c}{H_{jt}^c}) \\
&+ \alpha_4\theta ln(\frac{NM_{jt}^c}{H_{jt}^c}) + \alpha_5\theta ln(\frac{RC_{jt}^c}{H_{jt}^c}) + (\theta - 1)lnH_{jt}^c \\
&= ln\theta + lnZ_{jt}^c + \alpha_1\theta ln(na_{jt}^c) + \alpha_2\theta ln(rm_{jt}^c) + \alpha_3\theta ln(ni_{jt}^c) \\
&+ \alpha_4\theta ln(nm_{jt}^c) + \alpha_5\theta ln(rc_{jt}^c) + (\theta - 1)lnH_{jt}^c
\end{aligned} \quad (1.9)$$

where $na_{jt}^c$, $rm_{jt}^c$, $ni_{jt}^c$, $nm_{jt}^c$, and $rc_{jt}^c$ denote the per hour task inputs on occupation $j$, at time $t$ in city $c$. We shall eventually estimate the following equation:

$$lnW_{jt}^c = \beta_{jt}^c + \beta_1 ln(na_{jt}^c) + \beta_2 ln(rm_{jt}^c) + \beta_3 ln(ni_{jt}^c) + \beta_4 ln(nm_{jt}^c) + \beta_5 ln(rc_{jt}^c) + \beta_6 lnH_{jt}^c + \epsilon_{jt}^c$$

$$(1.10)$$

In Equation 1.10, $lnW_{jt}^c$ is the logarithm of occupation $j's$ wage in time $t$ and city $c$. The task inputs are all in the logarithms. $\beta_{jt}^c$ includes occupation fixed effect, time fixed effect, and city fixed effect, capturing the role of productivity. $lnH_{jt}^c$ is the working hours for each occupation. It captures the congestion effect, indicating how crowded an occupation is. We can estimate Equation 1.10 by ordinary least squares using panel data for the years 2007, and 2010-2016.

In the production process, the choices of task inputs are correlated with residual, which causes the endogeneity problem. When estimating the production function, endogeneity can occur both at the occupational level and at the local economy level, which could potentially bias the estimated coefficients obtained from ordinary least squares. We try two ways to tackle this issue. First, we employ the fixed effects method (specifications (1) and (2) in Table 1.10 and Table 1.11). In specification (1), we employ the occupational fixed effect. In specification (2), we allow the productivities of the large and small cities to be different. Second, following the spirit of Arellano and Bond (1991), we employ lagged dependent variables to include information of past periods (specification (3) in Table 1.10 and Table 1.11).

Table 1.10 displays the regression results for the Cobb-Douglas production function. The term $\beta_{jt}^c$ includes occupation effect, year effect, and city fixed effect. Columns (1) and (2) exhibit the estimates of fixed effect estimations, without and with city fixed effect respectively. For the estimation of column (1), the productivity is occupation-specific but not city-specific. From column (1) to column (2), we allow large cities to be more productive than small cities, which increases $R^2$ by 6.6%. To

31

make use of additional dynamic information, in the regression of Column (3), we follow Arellano and Bond (1991) and use the value of the dependent variable in the previous period as a predictor for the current value of the dependent variable. In this scenario, Cobb-Douglas production function could explain 83% of the occupational wages. The estimates of $\theta$, degree of homogeneity, are all significantly smaller than one in the three experiments. This shows that the production function employed in this experiment should exhibit decreasing returns to scale. The increase of task inputs leads to a less proportional output. The output elasticities of the five tasks $\alpha_1 - \alpha_5$ are all significantly positive. The elasticity of routine manual task $\alpha_2$ is around 50%, dominating the elasticities of the other four tasks. This is consistent with the empirical facts about task inputs. The elasticity of Nonroutine Analytic task is 0.09, the smallest of the five. However, the nonroutine analytic task input is not the smallest one. This could suggest that task counts of these five task types are in some sense in different unites.

**Table 1.10:** Estimation Results for Cobb-Douglas Production Function

|  | (1) | (2) | (3)* |
|---|---|---|---|
| $\alpha_1$ (Nonroutine Analytic) | 0.083 | 0.128 | 0.090 |
|  | (0.00) | (0.00) | (0.00) |
| $\alpha_2$ (Routine Manual) | 0.434 | 0.440 | 0.517 |
|  | (0.00) | (0.00) | (0.00) |
| $\alpha_3$ (Nonroutine Interactive) | 0.133 | 0.096 | 0.103 |
|  | (0.00) | (0.00) | (0.00) |
| $\alpha_4$ (Nonroutine Manual) | 0.271 | 0.249 | 0.125 |
|  | (0.00) | (0.00) | (0.00) |
| $\alpha_5$ (Routine Cognitive) | 0.080 | 0.087 | 0.115 |
|  | (0.00) | (0.00) | (0.00) |
| $\theta$ | 0.935 | 0.904 | 0.909 |
|  | (0.00) | (0.00) | (0.00) |
| Year Fixed Effect | Yes | Yes | Yes |
| Occ Fixed Effect | Yes | Yes | Yes |
| City Fixed Effect |  | Yes | Yes |
| $R^2$ | 71.13% | 77.73% | 82.75% |
| Observations | 7069 | 7069 | 7069 |

Note: p-values in brackets; * Arellano and Bond (1991)

**Table 1.11:** Estimation Results for CES Production Function

| | (1) | (2) | (3)* |
|---|---|---|---|
| $\lambda_1$ (Nonroutine Analytic) | 0.266 | 0.267 | 0.217 |
| | (0.00) | (0.00) | (0.00) |
| $\lambda_2$ (Routine Manual) | -0.002 | -0.030 | -0.046 |
| | (0.96) | (0.42) | (0.39) |
| $\lambda_3$ (Nonroutine Interactive) | 0.180 | 0.175 | 0.162 |
| | (0.00) | (0.00) | (0.00) |
| $\lambda_4$ (Nonroutine Manual) | 0.322 | 0.318 | 0.351 |
| | (0.00) | (0.00) | (0.00) |
| $\lambda_5$ (Routine Cognitive) | 0.235 | 0.270 | 0.316 |
| | (0.00) | (0.00) | (0.00) |
| $\theta$ | 0.918 | 0.862 | 0.866 |
| | (0.00) | (0.00) | (0.00) |
| $\sigma$ | -0.0009 | -0.0123 | -0.0181 |
| | (0.00) | (0.00) | (0.00) |
| Year Fixed Effect | Yes | Yes | Yes |
| Occ Fixed Effect | Yes | Yes | Yes |
| City Fixed Effect | | Yes | Yes |
| $R^2$ | 79.68% | 83.77% | 86.07% |
| Observations | 7069 | 7069 | 7069 |

Note: p-values in brackets; * Arellano and Bond (1991)

In order to include the effect of task-biased productivities on the wage, we then resort to a CES production function 1.8. By combining 1.8 and 1.6, we derive Equation 1.11. Analogous to Equation 1.9, $na_{jt}^c$, $rm_{jt}^c$, $ni_{jt}^c$, $nm_{jt}^c$, and $rc_{jt}^c$ denote the per hour task inputs on occupation $j$, at time $t$ in city $c$.

$$W_{jt}^c = \frac{\theta Z_{jt}^c [\lambda_1 (na_{jt}^c)^\sigma + \lambda_2 (rm_{jt}^c)^\sigma + \lambda_3 (ni_{jt}^c)^\sigma + \lambda_4 (nm_{jt}^c)^\sigma + \lambda_5 (rc_{jt}^c)^\sigma]^{\frac{\theta}{\sigma}}}{(H_{jt}^c)^{1-\theta}} \qquad (1.11)$$

Table 1.11 reports the estimation results for the CES production function. Similarly to the estimation of Cobb-Douglas function, the first two columns present the estimates of fixed effect estimations and the estimation of the last column deals with the endogeneity issue following Arellano and Bond (1991). $\lambda_1 - \lambda_5$ measure the task-biased productivities. In the three specifications, the estimates of $\lambda_2$ (productivities that favor routine manual tasks) are insignificant. The estimates of $\lambda_1$, $\lambda_3$, $\lambda_4$ and $\lambda_5$ are all significantly positive. Besides, we estimate $\theta$ to be significantly smaller than 1, which is consistent with the estimate of Cobb-Douglas production function. In column (3), the estimate of $\sigma$ is -0.0181, indicating that the elasticity substitution between inputs is $\frac{1}{1-\sigma} = 0.982$.

## 1.6  Results

Figures 1.5 and 1.6 show the graph of log hourly wages from the actual data vs those predicted by our model in large and small cities respectively. In general, these two figures show that the model fits occupational wage data in both large and small cities quite well.



**Figure 1.5:** Actual and logWage Predicted by Cobb-Douglas Framework by City

Note: The line in the figure is the 45 degree line.

**Figure 1.6:** Actual and logWage Predicted by CES Framework by City

Note: The line in the figure is the 45 degree line.

Based on the estimation results, we next scrutinize the extent to which occupational task difference can explain the city size wage gap. In the framework with Cobb-Douglas production function, the wage gap is decomposed into three factors: TFP differences, task inputs differences, and congestion differences. By shutting down the channels of TFP and congestion difference, we derive the contribution of

task input difference to the wage difference.

$$\Delta lnW_j = \underbrace{(lnZ_j^L - lnZ_j^S)}_{TFP} + \underbrace{\sum_{i=1}^{5} \theta\alpha_i(lnT_{ij}^L - lnT_{ij}^S)}_{Task\ Difference} + \underbrace{(\theta-1)(lnH_j^L - lnH_j^S)}_{Congestion} \quad (1.12)$$

In the framework of CES production function, as task-biased productivities are introduced, the city size wage gap is decomposed into four parts: TFP difference, task inputs difference, congestion difference, and task-biased productivity effect. In this experiment, task biased productivities $\lambda_1$ to $\lambda_5$ work together with task intensity difference. The difference of task intensities will be magnified by the task-biased productivities as the last term shown in Equation 1.13.

$$\Delta lnW_j = \underbrace{(lnZ_j^L - lnZ_j^S)}_{TFP} + \underbrace{\sum_{i=1}^{5} \theta\lambda_i(lnT_{ij}^L - lnT_{ij}^S)}_{Task\ Difference} + \underbrace{(\theta-1)(lnH_j^L - lnH_j^S)}_{Congestion}$$
$$+ \underbrace{\frac{\theta\sigma}{2}\sum_{i=1}^{5}\sum_{k=1}^{5}\lambda_i\lambda_k[(ln\frac{T_{ij}^L}{T_{kj}^L})^2 - (ln\frac{T_{ij}^S}{T_{kj}^S})^2]}_{Task\ Intensity\ with\ TBP} \quad (1.13)$$

Figure 1.7 displays the occupational actual wage difference (blue line), fitted wage difference (red line) and the wage difference contributed by task inputs difference (grey line) if we assume that TFPs and congestions are constant across cities. Figure 1.8 shows the actual wage difference (blue line), fitted wage difference (red line), contributions by task inputs (grey line), and task-biased productivities (orange line) if we assume that TFPs and congestions are constant across cities.

**Figure 1.7:** City Size Wage Gap Decomposition in Cobb-Douglas Framework

**Figure 1.8:** City Size Wage Gap Decomposition in CES Framework

The following Table 1.12 displays the predicted wage gaps in the framework of Cobb-Douglas production function and CES production function. Our predictions about the overall wage difference in Cobb-Douglas and CES are 8.88% and 8.93% respectively, which are close to the actual wage difference 8.5%. We also find that task inputs difference accounts for more than half of the predicted city size wage gap without allowing for task-biased productivities. In addition, if productivities favor specific tasks, the task difference could explain 27.60% of the predicted wage gap and around 7.51% is explained by those task-biased productivities.

**Table 1.12:** Prediction Results

|  | Data | Predicted | |
|---|---|---|---|
|  |  | C-D | CES |
| $\Delta lnWage$ | 8.5% | 8.88% | 8.93% |
| Shut down TFP and Congestion | | | |
| Occ Task Difference |  | 52.36% | 27.60% |
| Task-Biased Productivities |  | — | 7.51% |

## 1.7 Conclusion

This paper focuses on how the nature and complexity of an occupation vary in local labor markets of different sizes, and how this change impacts the city size wage gap. We employ a near universe of US online job posting data set that contains more than 20 million job ads. Using this data set, we investigate the structure of within occupation requirements in large vs small cities. We find that task distributions differ across cities. Within occupation, 1.54 more tasks are required in large cities. Large cities require relatively more analytic, interactive, and cognitive tasks, while small cities emphasize more manual tasks.

Then we lay out a framework to empirically investigate the extent to which occupational task inputs difference can explain the city size wage gap. We first estimate the production function and then decompose the city size wage gap into several factors. Our results indicate that task inputs difference could interpret half of the city size wage in the case of a Cobb-Douglas production function. In the case of a CES production function, task inputs difference, as well as task-biased productivities account for more than 35% of the wage gap.

This paper helps understand U.S. regional disparities, which have risen since the 1980s and have attracted greater interest. Our research suggests that jobs in large cities are significantly more complex in terms of task requirements and this task difference accounts for a large part of the wage differentials.

## 1.8 Appendix: Details of Burning Glass Job Posting Data Set

**Table 1.13:** The Numbers of Job Posts and Tasks in 2007 and 2016

|      | Job Posts  | Task Keywords |
|------|-----------|---------------|
| 2007 | 9,575,975  | 11,017        |
| 2016 | 23,883,197 | 12,214        |

**Table 1.14:** Top 10 Keywords and Frequencies in 2007

| Top 10 Keywords  | Freqency |
|------------------|----------|
| Communicate      | 16.11%   |
| Sales            | 10.45%   |
| Writing          | 9.58%    |
| Customer Service | 7.89%    |
| M.S. Excel       | 7.59%    |
| Research         | 6.03%    |
| Organization     | 5.99%    |
| Planning         | 5.96%    |
| Project Manage   | 5.94%    |
| Computer Skills  | 5.34%    |

**Figure 1.9:** City Size Wage Premium for 6-digit Occupations in 2007

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the wage gap in 2007.

**Figure 1.10:** City Size Task Difference for 6-digit Occupations in 2007

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the occupational task gaps.

**Figure 1.11:** Gap of Nonroutine Analytic Task for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the NA task gap in 2016.

46

**Figure 1.12:** Gap of Nonroutine Interactive Task for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the NI task gap in 2016.

**Figure 1.13:** Gap of Nonroutine Manual Task for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the NM task gap in 2016.

**Figure 1.14:** Gap of Routine Cognitive Task for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the RC task gap in 2016.

**Figure 1.15:** Gap of Routine Manual Task for 6-digit Occupations in 2016

Notes: The horizontal axis is by 6-digit occupations. They are sorted by the RM task gap in 2016.

**Figure 1.16:** City Size Task Difference for 3-digit Occupations in 2016



**Figure 1.17:** Task Compositions in Large and Small Cities for 2007 and 2016

**Figure 1.18:** US Metropolitan and Micropolitan CBSAs

# Chapter 2

# A Comparison of Multidimensional Skill Mismatch between NLSY79 and NLSY97

Qi XU

## Abstract

This paper studies the multidimensional skill mismatch of the male workers a-mong two surveys: National Longitudinal Survey of Youth 79 and 97. Using the survey data and O*NET data set, I construct and compare the worker-occupation match qualities of the two cohorts based on the measure proposed by Guvenen et al. (2020). I document the following facts about the match quality differences: on av-erage, the mismatch rate is higher in NLSY97 than NLSY79; workers are both more over-qualified and more under-qualified in NLSY97 than NLSY79; mismatch rate is higher in NLSY97 than NLSY79 in all dimensions (verbal, math, and social skills); mismatch rate varies across locations in both cohorts. The average mismatch rate is

higher in rural areas than in urban areas. It's also higher in non-MSA areas than in MSA areas.

## 2.1 Introduction

Due to labor market frictions, workers cannot be allocated to perfectly matched occupations. The term skill mismatch is used to describe this situation where the worker's abilities exceed or cannot meet the skills required by the occupations. Guvenen et al. (2020) propose an empirical measure of skill mismatch. It is the distance between the set of skills required by an occupation and the set of abilities possessed by a worker of that occupation. A worker is perfectly matched to his occupation only if his abilities are ideally aligned with the skills required by occupations. Otherwise, he's mismatched to his occupation, either over-qualified or under-qualified.

With the development of the social progress, technical progress, and economic conditions, the distribution of characteristics, the abilities on the worker side, and the skills required by occupations have changed over time, which will affect the match quality of worker-occupation pairs. Given that skill mismatch is directly linked to labor market outcomes, like wage and occupational attainment, it is essential to study the change of skill mismatch over time. For example, how is the skill match quality of American youth compared to the previous generation? Is the new generation better matched to occupations than the previous generation or not? Have people become more over-qualified or more under-qualified? The answers to these questions can be obtained by looking at the differences between the two generations. In this paper, I compare the skill match qualities of the two generations from two cohorts of the National Longitudinal Survey of Youth, 1979 (NLSY79) and the National Longitudinal Survey of Youth, 1997 (NLSY97). There are several aspects of differences between the two cohorts (NLSY79 and NLSY97). Firstly, workers are

younger in NLSY97. The NLSY79 spans from 1979 to 2010 and respondents are aged from 16/23 years old to 47/54 years old. The new cohort (NLSY97) spans from 1997 to 2017 and the workers are aged from 14/18 years old to 34/38 years old. Since mismatch tends to occur early in the career, being younger (cohort NLSY97) should result in a higher mismatch rate. Secondly, the abilities of respondents might differ between the two cohorts since the education attainment and family background are different over time. Lastly, the labor market conditions, which affect the worker-occupation match qualities, are different during the periods of the two surveys. The average unemployment rate during the period of NLSY79 is 6.39%, higher than that of NLSY97 (5.64%). For the workers who entered the labor market at the beginning of the 1980s, the average unemployment rate of this period (1980-1983) is about 8%. For the workers who entered the market around 1997, the average unemployment rate between 1997 to 2000 is only 4.3%.

This paper is mainly descriptive. The main objectives are to describe the characteristics of workers and occupations, measure the multidimensional skill mismatch rate, and compare the mismatch rate of workers between the cohorts of NLSY79 and NLSY97 following the methods of Guvenen et al. (2020). To construct the measures, I combine the data on the worker side and occupation side from NLSY79/97. NLSY79/NLSY97 contains detailed information on occupation, employment history, and wages of each respondent. The respondents of the surveys also participated in the tests of Armed Services Vocational Aptitude Battery (ASVAB) at the beginning of the surveys, which measures the abilities of individuals. The occupation side data is the Occupational Information Network (O*NET) data set from the U.S. Department of Labor. It characterizes the skills required by each occupation. After merging the data sets on two sides, I build two panel data sets that describe the attributes of individuals and their occupations over the period of NLSY79 and NLSY97. The

measure of mismatch is based on the combined data sets along three skill/ability dimensions: math, verbal, and social skills/abilities.

The main findings regarding the changes in mismatch between the cohorts of NLSY79 and NLSY97 are as follows:

(1) On average, the aggregate mismatch rate is higher in NLSY97 than NLSY79. In particular, the aggregate mismatch rate is higher in NLSY97 than NLSY79 by year, by age, or by region of residence.

(2) The aggregate positive and negative mismatch rates are both larger in NL-SY97, meaning that over-qualified workers are more over-qualified and under-qualified workers are more under-qualified in NLSY97 than NLSY79.

(3) The aggregate mismatch rate is higher in NLSY97 along the three skill dimensions of verbal, math, and social skills.

(4) The mismatch rate varies across locations in both cohorts. Guvenen et al. (2020) don't study the variation of mismatch across locations. I find that the mismatch rate is lower in urban areas relative to rural areas. In addition, the mismatch rate is higher in non-MSA areas than in MSA areas in both cohorts of NLSY79 and NLSY97.

This paper relates to the literature on measuring skill mismatch. There are a number of different approaches to measure skill mismatch and no commonly accepted measurement to date. Many papers measure mismatch along one dimension such as years of education and field of study. For example, Duncan and Hoffman (1981) utilize the data from the Panel Study of Income Dynamics (1976) and measure mismatch as the difference between the actual education attainment and the self-reported educational requirement on their jobs. Kiker et al. (1997) denote the workers with educational attainments greater than the modal educational level for

their specific occupation as over-qualified and those whose educational attainments are below the mode in his occupation as under-qualified. Bauer (2002) uses a large German panel data set and measures mismatch as the discrepancy between the educational attainment of an individual and the mean value of education within his occupation. Frei and Sousa-Poza (2012) measure mismatch based on the subjective survey question of Switzerland. One downside of the previous measures of mismatch based on educational attainment is that they cannot capture the variations of mismatch among workers of the same education level and same occupational requirement. Because workers of the same educational attainment still vary significantly in skills. For instance, college students with different fields of study (e.g. computer science and economics) on the same jobs could end up with a different level of mismatch since they are endowed with totally dissimilar skills. Accordingly, Liu et al. (2016) use the field of study in college to proxy the skills of workers. Mismatch occurs if a worker is not matched to an industry that values his skill (field of study).

The availability of data sets on multiple abilities of workers allows more papers to measure the multidimensional mismatch rate. Fredriksson et al. (2018) use data of workers' four cognitive talents and four noncognitive talents from Statistics Sweden and the Swedish War Archives. They design and measure the mismatch as differences between workers' talents and the average talents of the tenured workers on the same occupation. Lise and Postel-Vinay (2020) estimate a model of on-the-job search with multidimensional skills, which infers mismatch through the lens of a structural model by matching data moments. Guvenen et al. (2020), use the same data sets (NLSY79 and O*NET) as Lise and Postel-Vinay (2020), and propose an empirical measure of multidimensional skill mismatch, which is directly measurable with micro data.

My paper follows the empirical measure of Guvenen et al. (2020). I replicate their results about match quality in NLSY79, apply the measure to the new cohort NLSY97, and compare the mismatch rate of workers between the two generations. I find that the aggregate mismatch rate is higher in NLSY97 than NLSY79 even if the age is controlled for. Additionally, the aggregate mismatch rate differs across locations in both cohorts.

The rest of the paper is organized as follows. In Section 2.2, I introduce the data sources and measures. In Section 2.3 I present the empirical facts of mismatch rate in NLSY79 and NLSY97. Section 2.4 concludes.

## 2.2    Data Sources and Measures

### 2.2.1    Data sources

The main sources of data for this paper are the NLSY79, NLSY97, and O*NET Data. The NLSY79 and NLSY97 track two nationally representative samples of individuals in two periods: 1979-2016 and 1997-now. Since I follow the measurement of Guvenen et al. (2020), I will use the same sample selection criteria as them and restrict the analysis to males of the two cohorts. In the following sections, I employ the term "skill" to characterize the requirements of occupations and the term "ability" to describe the features of workers.

The National Longitudinal Surveys of Youth are a set of surveys designed to gather information at multiple points in time on the labor market activities of each worker. The NLSY79 Cohort is a longitudinal project that follows the lives of a sample of 9,964 American youth born between 1957 to 1964. This cohort spans from 1979 to 2016. The respondents were interviewed annually from 1979 to 1994 and biennially thereafter. The NLSY97 follows a sample of 8,984 individuals born be-

tween 1980 to 1984. The respondents were 12 to 17 years old when first interviewed in 1997. The sample members were interviewed annually from 1997 to 2011 and biennially thereafter.

The NLSY79 and NLSY97 comprise detailed information on employment, wages, and occupations of each respondent. They also include survey questions about the data on workers' personalities and attitudes, which I use to measure the "social" abilities of workers. For example, respondents were asked about *How much do you feel that difficult or cooperative describes you as a person?*. In addition, all respondents took the Armed Services Vocational Aptitude Battery (ASVAB) test at the start of each survey. The ASVAB test is a timed multi-aptitude test, which I use to measure workers' "math" and "verbal" abilities. The version of the ASVAB taken by NLSY79 respondents includes ten component tests [1] and the version of the ASVAB taken by NLSY97 respondents includes ten component tests and two speeded subsets to measure the workers' vocational aptitudes.[2] To be consistent with Guvenen et al. (2020), I select four components: Word Knowledge, Paragraph Comprehension, Arithmetic Reasoning, and Mathematics Knowledge of ASVAB test for the two cohorts to measure workers' verbal and math abilities. I use the scores of Word Knowledge and Paragraph Comprehension to measure the verbal abilities and scores of Arithmetic Reasoning and Mathematics Knowledge to calculate the math abilities.

The Occupational Information Network (O*NET) is the primary source of occupational information in the U.S., which contains a rich set of occupation-specific descriptors that describe the characteristics of the work, including but not limited

---

[1] ASVAB tests taken by NLSY79 respondents: Arithmetic Reasoning, Mathematics Knowledge, Paragraph Comprehension, Word Knowledge, General Science, Numerical Operations, Coding Speed, Automotive and Shop Information, Mechanical Comprehension, and Electronics Information.

[2] ASVAB tests taken by NLSY97 respondents: Arithmetic Reasoning, Assembling Objects, Auto Information, Coding Speed, Electronics Information, General Science, Mathematics Knowledge, Mechanical Comprehension, Numerical Operations, Paragraph Comprehension, Shop Information, Word Knowledge

to skill requirements and knowledge to conduct the tasks of each occupation. The O*NET provides the scores of the level of importance for 277 descriptors [3]. Again following Guvenen et al. (2020), I choose 26 descriptors that are consistent with ASVAB tests and 6 descriptors about social skills to calculate the verbal, math, and social skills of each occupation.

I follow the approach of Guvenen et al. (2020) for the sample selection criteria. In addition to the criteria listed in their paper, I limit the samples to the individuals with valid geographic information. The number of the remaining individuals and observations after applying the sample selection criteria are 1,992 individuals and 33,364 observations in NLSY79 and 2,422 individuals and 15,410 observations in NLSY97.

## 2.2.2   Mapping between skills and abilities

Table 2.1 describes the information of the three-dimensional skills and abilities of NLSY79 and NLSY97. In this paper, I use skills to denote the inputs required by occupations. Abilities are the characteristics of the worker side. The three dimensions are 'Verbal', 'Math', and 'Social'. Measurements of workers' abilities are from NLSY79 and NLSY97. Measurements of skills required by occupations are from O*NET. To pair the skills and abilities, I mainly follow the method of Guvenen et al. (2020) except that I use alternative measurements for social ability on the worker side since the Rotter/Rosenberg Scales, which Guvenen et al. (2020) use to measure social skill, are not included in NLSY97.

---

[3]I use the O*NET database version 4.0, the same as Guvenen et al. (2020)

**Verbal and Math**

**Skills required by occupations**

I use the same version of the O*NET database (Version 4.0) for the two cohorts. The 26 descriptors for verbal and math skills and 6 descriptors for social skills are listed in Table 2.1. For the 26 descriptors of verbal and math skills, the Defense Manpower Data Center (DMDC) assigned 26 relatedness scores to the ASVAB components (Word Knowledge, Paragraph Comprehension, Arithmetic Reasoning, Math Knowledge). The first step is to convert the 26 descriptors to four scores that are comparable to the four ASVAB components according to the relatedness scores. Each of the four scores is the weighted average of the 26 O*NET descriptors. By now, I obtain four scores, two comparable to Word Knowledge and Paragraph Comprehension, the other two comparable to Arithmetic Reasoning, Math Knowledge.

**Table 2.1:** Data Sources of Skills of Occupations and Abilities of Workers

| | Worker's Ability in NLSY79 | Worker's Ability in NLSY97 | Occupational Skill in O*NET |
|---|---|---|---|
| Verbal | ASVAB: Word Knowledge/Paragraph Comp. | ASVAB: Word Knowledge/Paragraph Comp. | 26 Verbal and Math Skills |
| Math | ASVAB: Arithmetic Reasoning/Math Knowledge | ASVAB: Arithmetic Reasoning/Math Knowledge | |
| Social | Rotter/Rosenberg Scale | Personality Scale: Open/Extraverted | 6 Social Skills |

The second step is to reduce the 4 scores which are comparable to ASVAB components to 2 categories (verbal and math) by principal component analysis. The score of verbal skill is the first principle component of the two scores related to Word Knowledge and Paragraph Comprehension, and the score of math skill is that of the two scores related to Math Knowledge and Arithmetic Reasoning.

**Table 2.2:** List of Skills in O*NET

| Verbal and Math Skills | |
|---|---|
| 1. Oral Comprehension | 14. Operation and Control |
| 2. Written Comprehension | 15. Equipment Maintenance |
| 3. Deductive Reasonning | 16. Troubleshooting |
| 4. Inductive Reasoning | 17. Repairing |
| 5. Information Ordering | 18. Computers and Electronics |
| 6. Mathematical Reasoning | 19. Engineering and Technology |
| 7. Number Facility | 20. Building and Construction |
| 8. Reading Comprehension | 21. Mechanical |
| 9. Mathematics Skill | 22. Mathematics Knowledge |
| 10. Science | 23. Physics |
| 11. Technology Design | 24. chemistry |
| 12. Equipment Selection | 25. Biology |
| 13. Installation | 26. English Language |
| Social Skills | |
| 1. Social Perceptiveness | 4. Negotiation |
| 2. Coordination | 5. Instructing |
| 3. Persuasion | 6. Service Orientation |

Source: Guvenen et al. (2020)

**Abilities of workers**

At beginning of each survey, most respondents participated in the ASVAB tests. The measurements of verbal abilities and math abilities of workers are the ASVAB test scores. The score of verbal ability is the first principle component of Word Knowledge and Paragraph Comprehension, and the score of math ability is that of Math Knowledge and Arithmetic Reasoning.

**Social**

**Skills required by occupations**

Socials skill on occupation side are from 6 descriptors of O*NET database listed in Table 2.1. These descriptors are Social Perceptiveness, Coordination, Persuasion, Negotiation, Instructing, Service Orientation. Again, I use the first principal component method to convert the 6 descriptors to one single dimension.

**Abilities of workers**

According to Guvenen et al. (2020), for NLSY79, I use the Rotter Locus of Control Scale and Rosenberg Self-Esteem Scale to measure the (non-cognitive) social ability of respondents. The Rotter Locus of Control Scale measures the extent to which individuals believe they have control over their lives through self-motivation or self-determination (internal control) as opposed to the extent that the environment (that is, chance, fate, luck) controls their lives (external control). The Rosenberg self-esteem describes a degree of approval or disapproval toward oneself (Rosenberg (1965)). I take the first component of the two measures.

For NLSY97, the Rotter/Rosenberg scales are not available. The comparable measurement of social ability in NLSY97 comes from two survey questions about the personality traits of respondents. In the questions, the respondents rated how well the following paired traits applied to them: extraverted or enthusiastic; open or complex. These two survey questions are comparable to Rotter/Rosenberg to the

extent that they both measure the relationship between the respondents themselves and the outside world. I also take the first component of the two scales to measure the social abilities of workers. This measure of social ability in NLSY97 correlates with the two other abilities as well as in NLSY79 because the correlations between the social ability and verbal/math abilities are quite close in NLSY79 and NLSY97. (shown in left panel of Table 2.3)

**Table 2.3:** Correlations among Ability and Skill Scores in NLSY79 and NLSY97

| Worker's Ability | | (a) Worker Ability | | | (b) Occupational Skill Requirement | | |
|---|---|---|---|---|---|---|---|
| | | Verbal | Math | Social | Verbal | Math | Social |
| NLSY79 | Verbal | 1.00 | | | 0.37 | 0.34 | 0.35 |
| | Math | 0.78 | 1.00 | | 0.44 | 0.40 | 0.35 |
| | Social | 0.30 | 0.27 | 1.00 | 0.13 | 0.11 | 0.16 |
| NLSY97 | Verbal | 1.00 | | | 0.17 | 0.15 | 0.19 |
| | Math | 0.84 | 1.00 | | 0.17 | 0.16 | 0.17 |
| | Social | 0.26 | 0.29 | 1.00 | 0.03 | 0.03 | 0.03 |

Table 2.3 reports the correlations among ability and skill scores in NLSY79 (upper panel) and NLSY97 (lower panel). The left panel (a) displays the correlations between workers' verbal, math, and social ability scores in the two cohorts. The right panel (b) shows the correlations between workers' ability scores and the corresponding skill requirements of their occupations. People may wonder if the scale of the measures could affect the correlations. In Table 2.4, I also display the correlations between ability and skill percentile ranks in both cohorts. The correlations are lower in NLSY97 than NLSY79 in terms of both scores and ranks. The correlations indicate that:

1. The correlations among the three abilities of workers are similar in two cohorts. In NLSY79, the correlation between verbal ability and math ability is 0.78; the correlation between verbal ability and social ability is 0.30; the correlation between math ability and social ability is 0.27. In NLSY97, the three correlations are 0.84, 0.26, and 0.29. The correlations of NLSY79 and NLSY97 are close, meaning that the variables to measure the abilities of the two cohorts, especially the social ability, are comparable.

2. The correlations between worker's abilities and occupational skills (right panel) are lower in three dimensions of verbal, math, and social in NLSY97 than NLSY79, suggesting that the mismatch rate is higher in NLSY97.

3. In NLSY79, workers with higher math abilities sort into occupations with higher skill requirements in all dimensions. In NLSY97, however, workers with higher math abilities lost this advantage.

**Table 2.4:** Correlations between Ability and Skill Percentile Ranks in NLSY79 and NLSY97

|  |  | Occupational Skill Requirement | | |
|---|---|---|---|---|
| Worker's Ability | | Verbal | Math | Social |
| | Verbal | 0.37 | 0.33 | 0.35 |
| NLSY79 | Math | 0.43 | 0.40 | 0.35 |
| | Social | 0.15 | 0.14 | 0.18 |
| | Verbal | 0.20 | 0.17 | 0.23 |
| NLSY97 | Math | 0.20 | 0.18 | 0.22 |
| | Social | 0.06 | 0.05 | 0.07 |

Figure 2.1 displays the distributions of verbal, math, and social abilities in NL-SY79 (left panel) and NLSY97 (right panel). For verbal ability, the distributions of NLSY79 and NLSY97 are both left-skewed, with more observations higher than the rest. From NLSY79 to NLSY97, the distribution of verbal ability becomes less skewed and more individuals tend to be with average verbal abilities. For math ability, the distribution of NLSY79 is close to uniform and the distribution of NL-SY97 is almost normal. For the workers in NLSY97, their math abilities are more concentrated. For social ability, the distribution of NLSY79 is about normal. Workers' abilities are concentrated to the average level. The distribution of NLSY97 is left-skewed, signifying that more workers are with higher abilities than the rest.

Figure 2.2 displays the distributions of verbal, math, and social skills in NLSY79 (left panel) and NLSY97 (right panel). The skill requirement of each occupation is invariant between the two cohorts, especially for the distributions of verbal skill and math skill. The distribution of social skill changes a little bit. In NLSY97, the distribution is more right-skewed.

**Figure 2.1:** Histograms of Worker's Abilities in NLSY79 and NLSY97

**Figure 2.2:** Histograms of Occupational Skills in NLSY79 and NLSY97

### 2.2.3 Empirical Measure of Mismatch

To measure the extent to which a worker's abilities are matched to the skills required by his occupation, I use the mismatch measurement proposed by Guvenen et al. (2020). This measurement is defined as the difference between a worker's abilities and the skills of his occupation. If the mismatch rate is higher, the worker and his occupation are less well matched. If the mismatch rate is zero, the abilities of workers are perfectly consistent with the skills required by their occupations.

**Mismatch**

Each occupation $o$ is characterized as a set of three skills (verbal, math, and social), $j = 1, 2, 3$, denoted as $S_{oj}$. Every worker $i$ is characterized as a vector of three abilities (verbal, math, and social), $j = 1, 2, 3$, denoted as $A_{ij}$. I use $R(S_{oj})$ and $R(A_{ij})$ to denote the corresponding percentile ranks of the occupational skill requirements and worker abilities. The mismatch between a worker $i$ and his occupation $o$ is the weighted sum of the absolute value of the difference in each of the three dimensions between a worker's abilities and skill requirements. Same as Guvenen et al. (2020), the weights $\omega_j$ are the factor loadings from the first principal component of the set of absolute values of differences $\{|R(A_{ij}) - R(S_{oj})|\}_{j=1}^{3}$.

$$m_{io} \equiv \sum_{j=1}^{3}\{\omega_j \times |R(A_{ij}) - R(S_{oj})|\} \tag{2.1}$$

**Positive and negative mismatch**

The mismatch rate $m_{io}$ is always positive if the worker is not perfectly matched to his occupation. Then I introduce two measurements from Guvenen et al. (2020): positive mismatch and negative mismatch. Positive mismatch measures the part where some of a worker's abilities exceed the skill requirement of his occupation. Negative mismatch measures the part where some of a worker's abilities don't meet

the occupational skill requirement. Equations 2.1, 2.2, and 2.3 imply that $m_{i,o} = m_{i,o}^+ + (-m_{i,o}^-)$.

$$m_{io}^+ \equiv \sum_{j=1}^{3} \omega_j \max[R(A_{ij}) - R(S_{oj}), 0] \tag{2.2}$$

$$m_{io}^- \equiv \sum_{j=1}^{3} \omega_j \min[R(A_{ij}) - R(S_{oj}), 0] \tag{2.3}$$

## 2.3 Empirical Facts of Mismatch

In this section, I present the empirical findings of how the worker-occupation skill mismatch rate and its components change across the cohorts of NLSY79 and NL-SY97. In general, I show evidence that the mismatch rate is higher in the cohort of NLSY97. Additionally, I provide evidence about mismatch rate by year, by age, and by location. Table 2.5 displays the descriptive statistics for mismatch in NLSY79 and NLSY97. The mismatch rate in NLSY97 is systematically higher than NLSY79 by different groups (education, race, and industry). It is also noticeable that the mismatch of $\geq$4-year college graduates is lower than <4-Year College graduates in both cohorts. This is because the highly educated workers are more skilled and more specialized. They have a relatively clearer target of occupation when entering the labor market. Therefore they tend to work in the occupations that are more suitable for them.

I speculate that one reason that the mismatch rate is higher in NLSY97 is overeducation. Overeducation is more and more an issue over time. The workers in NLSY97 are more overeducated and therefore the mismatch rate is higher. Another reason why there is more mismatch is the higher occupational switching rate in NL-SY79 than NLSY97. Workers in NLSY79 switch their occupations more frequently will end up with a relatively lower mismatch rate.

### 2.3.1 Aggregate Mismatch and Components of Mismatch

Figure 2.3 shows the aggregate mismatch rates, verbal/math/social mismatch rates, and positive/negative mismatch rates in the cohorts of NLSY79 and NLSY97. The aggregate mismatch rate and its components all go up from the old cohort to the new cohort. The aggregate mismatch rate is 0.26 in NLSY79 and 0.3 in NLSY97 respectively. The mismatch rate of verbal skill is 0.26 in NLSY79 and 0.29 in NLSY97.

**Table 2.5:** Descriptive Statistics for Mismatch in NLSY79 and NLSY97

| | Mismatch | |
| Group Name | NLSY79 | NLSY97 |
|---|---|---|
| *All Observations* | 0.258 | 0.292 |
| | | |
| *By Education* | | |
| ≥4-Year College | 0.231 | 0.249 |
| <4-Year College | 0.268 | 0.297 |
| | | |
| *By Race* | | |
| Hispanic | 0.261 | 0.286 |
| Black | 0.247 | 0.278 |
| Non-Black, Non-Hispanic | 0.259 | 0.295 |
| | | |
| *By Industry* | | |
| Agriculture, Forestry, Fishing, and Hunting | 0.263 | 0.381 |
| Mining | 0.274 | 0.235 |
| Construction | 0.274 | 0.291 |
| Manufacturing | 0.258 | 0.288 |
| Transportation, Communications, and other Utilities | 0.242 | 0.267 |
| Wholesale and Retail Trade | 0.262 | 0.292 |
| Financial, Insurance, and Real Estate | 0.230 | 0.257 |
| Business and Repair Service | 0.267 | 0.357 |
| Personal Service | 0.274 | 0.307 |
| Entertainment and Recreation Services | 0.317 | 0.344 |
| Professional and Related Services | 0.241 | 0.280 |
| Public Administration | 0.247 | 0.291 |

Note: The data source is NLSY79, NLSY97, and O*NET. The measurement of mismatch is proposed by Guvenen et al. (2020). I use the crosswalk of Guvenen et al. (2020) to convert the different industry classifications in NLSY79 and NLSY97 into the Census 1970 One-Digit Industry Code. Guvenen et al. (2020) normalize the final mismatch measure (to have a mean of the minimum mismatch rate and a standard deviation of one) in their paper. The measure of mismatch rate is not normalized in this paper in order to be compared across two cohorts.

Note: The data source is NLSY79, NLSY97, and O*NET. The sample selection criteria are from Guvenen et al. (2020).

**Figure 2.3:** Components of Mismatch for NLSY79 and NLSY97

The mismatch rate of math skill is 0.25 in NLSY79 and 0.29 in NLSY97. The mismatch rate of social skill is 0.29 in NLSY79 and 0.31 in NLSY97. If I decompose the mismatch rate into two parts, positive mismatch and negative mismatch, we can see that workers are more over-qualified and also more under-qualified in NLSY97 than NLSY79. The positive mismatch rate is 0.13 in NLSY79 and 0.15 in NLSY97. The negative mismatch rate is -0.13 in NLSY79 and -0.14 in NLSY97.

## 2.3.2 Mismatch by Year

Figure 2.4 displays the trends of annual average mismatch rate by year for the cohorts of NLSY79 and NLSY97. First of all, in general, the aggregate mismatch rate goes down over time in the two cohorts, from 0.281 to 0.248 in NLSY79 and from 0.320 to 0.262 in NLSY97. The mechanism is that over time workers gradually switch to occupations which are better matched to their abilities. This leads to a decreasing mismatch rate over the year. Second, the mismatch rate of NLSY97 is completely and thoroughly above that of NLSY79, which indicates a worse matching quality in the new cohort. Here the 'year' partially stands for 'age' of individuals so next, I will show the figure of mismatch by age group to isolate the year and cohort effects.
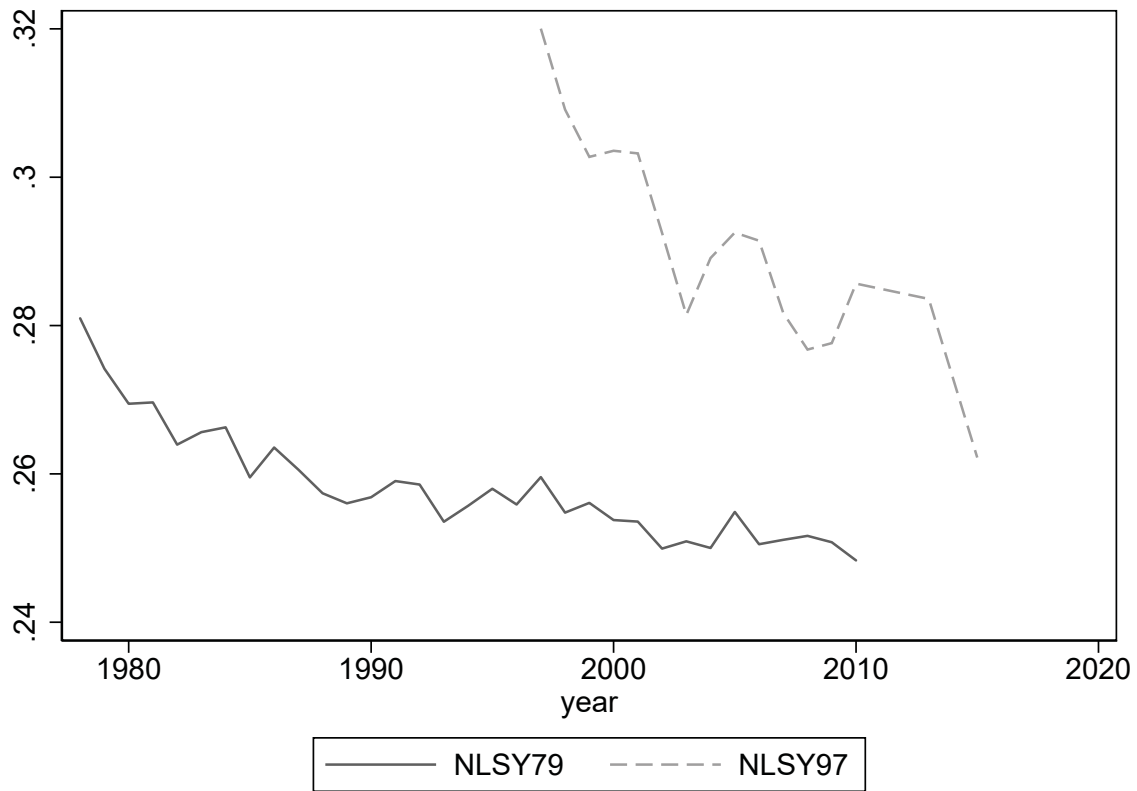
Note: The data source is NLSY79, NLSY97, and O*NET. The sample selection criteria are from Guvenen et al. (2020).

**Figure 2.4:** Mismatch by Year for NLSY79 and NLSY97

### 2.3.3 Mismatch by Age and by Labor Market Experience

The upper panel of Figure 2.5 shows the mismatch rate by age group for the two cohorts of NLSY79 and NLSY97. The upper panel of Figure 2.5 suggests that, even though I control the age for each cohort, there's still a mismatch gap between the two trends (which is roughly 0.04/0.3). For most of the age groups, the workers in the old cohort are better matched to their occupations than the workers in the new cohort. The youngest workers in both samples (after selection) are at age 16. The number of workers aged 16 in NLSY79, however, is under 50. Therefore I drop

this group and present a figure of mismatch by age group with the starting age of 17. The lower panel of Figure 2.5 presents the mismatch rate by the group of labor



**Figure 2.5:** Mismatch by Age an Experience in NLSY79 and NLSY97

market experience for the two cohorts of NLSY79 and NLSY97. With the same labor

market experience, workers in NLSY79 are better matched to their occupations than workers in NLSY97.

## 2.3.4   Mismatch Rate by Locations

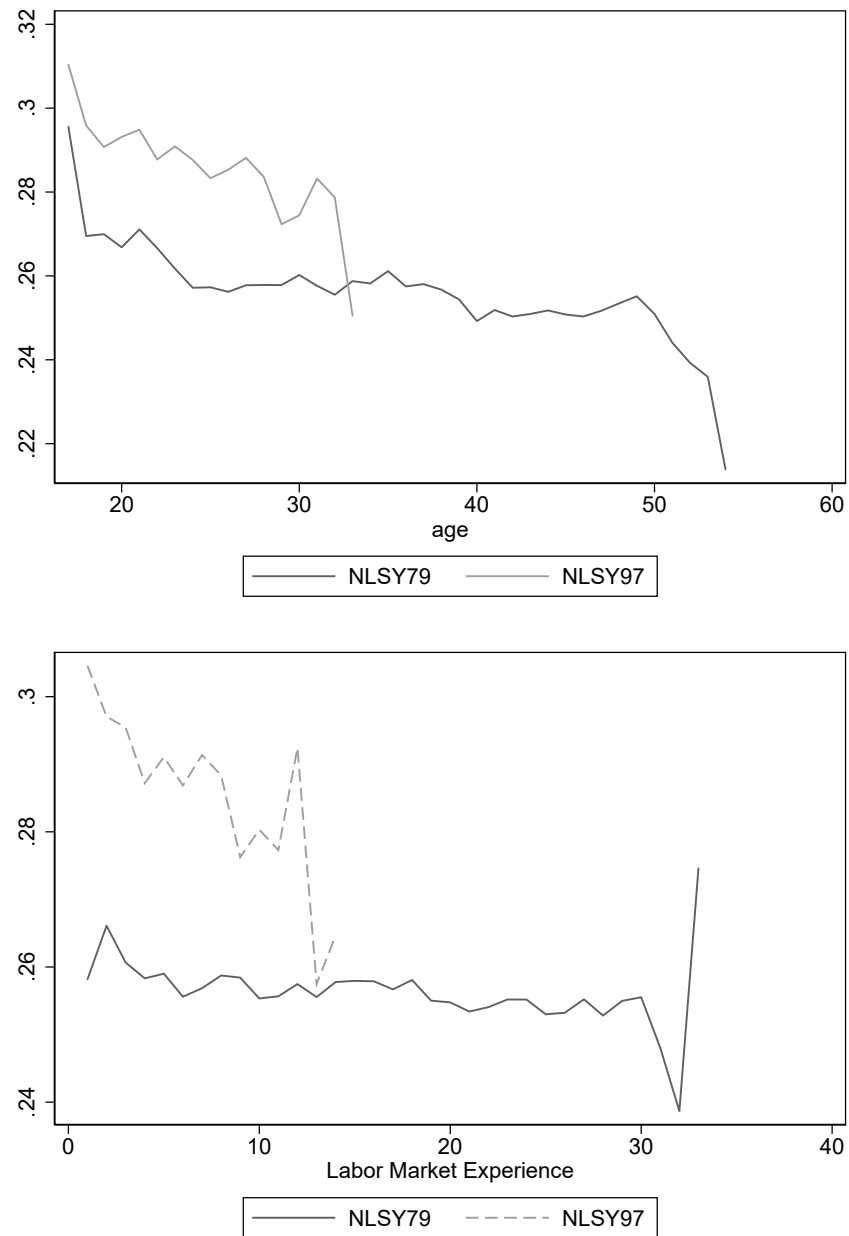

Note: The data source is NLSY79, NLSY97, and O*NET. The sample selection criteria are from Guvenen et al. (2020).

**Figure 2.6:** Mismatch Rate in Rural and Urban Areas for NLSY79 and NLSY97

Figure 2.6 and Figure 2.7 show the mismatch rate by locations. From these two figures, I see the skill mismatch rate varies across different regions in both cohorts. In particular, the mismatch rate is slightly higher in rural areas than in urban areas for the two cohorts. In NLSY79, the mismatch rate in rural and urban areas is 0.259 and 0.257 respectively. In NLSY97, the mismatch rate in rural and urban areas is 0.297 and 0.291 respectively.

Figure 2.7 indicates that the mismatch rate in MSA (Metropolitan Statistical Area) is lower than not in MSA. In NLSY79, the mismatch rate is higher in central city of

MSA than not in central city but in MSA. In NLSY97, the mismatch rate is lower in central city of MSA than other areas. Papageorgiou (2020) provides a view to understanding this mismatch variation between MSAs and non-MSAs. He shows the evidence that the number of occupations is higher in large cities (MSAs) than small cities (MSAs). Cities with double the size have approximately 70 more occupations. Thus workers in larger cities have more occupational options and are able to form better occupational matches. My result is consistent with his evidence as I show that the mismatch rate is higher in non-MSAs (small cities). The reason could be that workers with certain sets of abilities cannot find their right occupation matches in non-MSAs since the perfect occupations do not exist in non-MSAs. While in MSAs, which are large cities with more kinds of occupations available, it's relatively easier for a worker to find an occupation that fits his abilities.
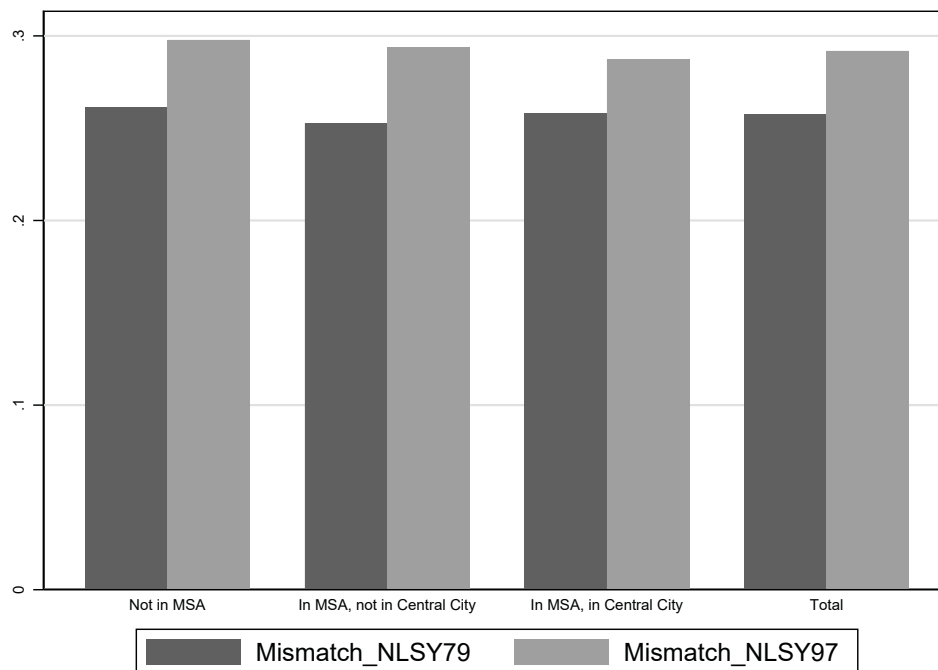


Note: The data source is NLSY79, NLSY97, and O*NET. The sample selection criteria are from Guvenen et al. (2020).

**Figure 2.7:** Mismatch Rate in MSA for NLSY79 and NLSY97

## 2.4   Conclusion and Discussion

In this article, I measure and compare the multidimensional skill mismatch between the two cohorts of NLSY79 and NLSY97. Following the measure proposed by Guvenen et al. (2020), I describe the abilities of workers in NLSY79/NLSY97 and the skills of occupations in O*NET, combine NLSY79/NLSY97 and O*NET and construct the worker-occupation data set, measure the multidimensional skill mismatch for each individual, and compare the aggregate mismatch and mismatch rate by demographic groups between NLSY79 and NLSY97.

I exploit NLSY97 to find the variables that are comparable to variables in Guvenen et al. (2020) to describe the abilities of workers. The abilities and skills are along three dimensions: verbal, math, and social. The correlations among abilities in NLSY79 and NLSY97 are quite close, laying the foundation for comparing the measure of mismatch between the two cohorts. The correlations between the abilities of workers and the skills of occupations provide a big picture of the match qualities in NLSY79 and NLSY97. I see that the correlations between abilities and skills are significantly higher in NLSY79 than NLSY97, indicating that workers in NLSY79 are sorted into the occupations that are more suitable for their abilities.

Next, I measure the aggregate mismatch rate, positive mismatch rate, and negative mismatch rate in the two cohorts. Interestingly, I find the aggregate mismatch rate is higher in NLSY97 than NLSY79. If I take a closer look at the mismatch by demographic groups, I find that the mismatch rate in NLSY97 is higher by year, by age, and even by location.

Then, I also provide the evidence that workers in NLSY97 are more over-qualified and more under-qualified than workers in NLSY79 since the positive mismatch rate and the absolute number of negative mismatch rate are both higher in NLSY97.

Lastly, I find a new fact about the mismatch variations across locations. In NLSY79 and NLSY97, the mismatch rate is higher in rural areas than in urban areas and also higher in non-MSA areas than in MSA areas. This interesting evidence could be a good motivation to explore the behavior and mechanism of job matching across cities.

As a descriptive paper, this chapter suggests some directions for future research. For example, based on the solid results about the mismatch rate during the periods of NLSY79 and NLSY97, future studies could explore the macrodynamics of mismatch and investigate the effect of labor market conditions on mismatch rate. Future research can also focus on how mismatch affects wage and does this effect changes between the two cohorts. In addition, mismatch varying across locations is a good motivation for studies about job search and match across cities.

# Chapter 3

# The Cyclicality of Multidimensional Skill Mismatch: Evidence from NLSY79 and NLSY97

Qi XU

## Abstract

This paper examines the cyclicality of multidimensional skill mismatch between workers and their occupations in NLSY79 and NLSY97. Using the data sets of NLSY79, NLSY97, and O*NET, I estimate the effect of labor market conditions on mismatch. Overall, I find that the mismatch rate is procyclical in the cohort of NLSY97. I also provide strong evidence that the effect of unemployment on mismatch is adjusted by educational level. In NLSY97, the mismatch of higher educated workers is more procyclical.

## 3.1 Introduction

How does the match quality between workers and their occupations vary when the economy is in the downturn? Will the worker-occupation match qualities get worse in recessions? Theoretical work on the cyclical behavior of mismatch does not reach a consensus. One view is that mismatch is countercyclical due to the sullying effect of recessions (Barlevy (2002)). The sullying effect means that in the slack period the workers reallocate to better matches more slowly. Thus the mismatch rate is relatively higher during recessions than expansions. An alternative view is that mismatch is procyclical as a result of the cleansing effect (Caballero and Hammour (1994)). The mechanism of cleansing effect suggests that recessions accelerate the destruction of less efficient matches and high-quality matches remain. In economic downturns, the mismatch rate gets higher since the bad worker-job pairs will disappear.

The goal of this paper is to empirically analyze the cyclicality of the mismatch rate based on micro data (NLSY79, NLSY97, and O*NET). I combine the worker side data sets NLSY79/NLSY97 with the occupation side data set O*NET. Using the combined data set, I calculate the measure of multidimensional skill mismatch rate, which is proposed by Guvenen et al. (2020). According to Guvenen et al. (2020), the multidimensional mismatch rate is the distance between the set of skills required by occupations and the set of workers' abilities. The combined data set features the characteristics of the worker-occupation pairs and covers multiple business cycles during the periods of NLSY79 and NLSY97. I then estimate the effect of labor market conditions on the mismatch rate in the two cohorts.

I address three main questions in this paper. First, what is the effect of cyclical labor market conditions on the worker-occupation mismatch rate? The improvement from previous research (Duncan and Hoffman (1981), Kiker et al. (1997), Bauer (2002), and Liu et al. (2016)) is that I use the direct measure of multidimensional skil-

l mismatch rate instead of using wage/educational level/field of study to measure the mismatch rate. Second, how does the effect of recessions vary by different educational levels? Third, do the answers to the two questions change from the cohort NLSY79 to NLSY97?

To answer these questions, the first step is to measure the worker-occupation mismatch rate for each individual in NLSY79 and in NLSY97. The detailed methods and results are displayed in Chapter 2. In aggregate, the mismatch rate level is higher in NLSY97 than in NLSY79. Workers in NLSY97 are more over-qualified and also more under-qualified. I then build and estimate an empirical specification and investigate the channel through which the cyclical labor market conditions affect the mismatch rate. I find that the effect of the unemployment rate, which is the proxy of labor market conditions, is significantly negative. If the unemployment rate increases, the mismatch rate will decrease, meaning that the mismatch rate is procyclical.

If I decompose the mismatch rate into two parts: positive mismatch and negative mismatch [1] (see definition and results from Chapter 2), I obtain the following results: in NLSY79, positive mismatch is countercyclical and negative mismatch is procyclical. In recessions, it's not easy for over-qualified workers to find a better match. Therefore the positive mismatch rate increases, meaning that the positive mismatch is countercyclical. The under-qualified workers tend to lose their jobs in recessions as they are less skilled. Thus the negative mismatch rate decreases and is procyclical. This result suggests that the sullying effect plays a role in the part of over-qualification while the cleansing effect plays a role in the under-qualification. Overall, the cleansing effect outweighs the sullying effect so that the mismatch rate is procyclical in NLSY79. In NLSY97, positive mismatch is countercyclical while

---

[1]Positive mismatch measures the part where some of a worker's abilities exceed the skill requirement of his occupation. Negative mismatch measures the part where some of a worker's abilities don't meet the occupational skill requirement.

negative mismatch is procyclical. In aggregate, the mismatch rate in NLSY97 is procyclical due to the cleansing effect.

Regarding the second question about the educational level, in addition to the negative effect of educational level on mismatch rate, I find that the effect of unemployment on mismatch is adjusted by the educational level. The procyclicality or countercyclicality of mismatch differs among workers with different educational levels. In NLSY79, the average marginal effect of unemployment rate on positive mismatch rate is positive and increases as the educational level goes up, meaning that being highly educated will make the positive mismatch rate more countercyclical in NLSY79. In NLSY97, the average marginal effect of unemployment rate on mismatch rate is negative and decreases as educational level goes up, i.e. higher educational level increases the procyclicality of mismatch. The cleansing effect gets stronger for workers with higher levels of education.

When comparing the cyclicality of the mismatch rate between NLSY79 and NLSY97, I present the evidence that mismatch is procyclical in both cohorts (insignificant in NLSY79). It differs when the mismatch is decomposed into positive and negative mismatch. Positive mismatch is countercyclical and negative mismatch is procyclical. Additionally, the way that educational level adjusts the effect of unemployment is different. Education lowers the procyclicality of mismatch in NLSY79 and increases the procyclicality of mismatch in NLSY97.

This paper is related to literature about the effect of recessions on mismatch. The two effects of recessions are *cleansing effect* and *sullying effect*. The view of the cleansing effect dates back to Schumpeter (1939). He studies the effect of business cycles on the allocation of resources and proposes that recessions will yield a more efficient allocation of resources by removing the bad investments. Caballero and Hammour (1994) present the evidence about the cleansing effect of recessions. In recessions, the outdated units are most likely to be scrapped. Lise and Robin (2017)

show evidence that during the periods of economic depression, low-type workers are fired, especially those matched with high-type firms; low-type firms hire less; medium-/high-type firms hire relatively more medium-/high-type workers.

On the other hand, there are many papers discussing the sullying effect of recessions. Barlevy (2002) develops a match model with on-the-job search. He points out that in an economic downturn, fewer vacancies are created so workers have a more difficult time moving into jobs that they are best suited for. Kahn (2010) finds that workers who graduate in a worse economy and cannot find good matches are in low-level occupations. Barnichon and Zylberberg (2019) provide the evidence that the underemployed (workers who are overqualified for their jobs) rate is counter-cyclical. In their model, high-skill workers are underemployed in order to avoid competition from other high-skill workers and find a job more easily. In recession, the high-skill workers move down the job ladder to reduce the aggregate shock. Thus the underemployment rate increases.

In this paper, I use the direct measure of multidimensional skill mismatch from micro data sets to investigate which of the sullying and the cleansing effects of recessions on mismatch dominate in the cohorts of NLSY79 and NLSY97. The rest of this paper is organized as follows: In Section 3.2 I introduce the data source and the basic descriptive statistics. In Section 3.3 I discuss the empirical specification and the empirical results. Section 3.4 concludes.

## 3.2 Data Source

In order to estimate the causal effect of labor market conditions on the mismatch rate, I combine the National Longitudinal Survey of Youth 79 and 97 (NLSY79, NLSY97) and the Occupational Information Network (O*NET). NLSY79/NLSY97 includes abundant information on employment history and scores of abilities of each respondent. O*NET is the source of occupational information, especially the scores of the required skills. Based on these data sets, I measure the multidimensional skill mismatch rate for each worker in the cohorts of NLSY79 and NLSY97. The mismatch rate is along three dimensions: verbal, social, and math. The details about how to combine the data sets, how to select the samples, how to construct the measure of skill mismatch, and the statistics of mismatch rate are displayed in Chapter 2. The unemployment rate is an appropriate proxy for labor market conditions in previous literature (Kahn (2010); Oreopoulos et al. (2012); Altonji et al. (2016)). The data of the unemployment rate is published by the Bureau of Labor Statistics (BLS).

In table 3.1, I report the descriptive statistics for the samples of NLSY79 (1979-2010) and NLSY97 (1997-2015). The sample selection criterion is the same as Guvenen et al. (2020) and described in Chapter 2. After selection, the total number of individuals is 1,992 in NLSY79 and 2,422 in NLSY97. The total number of observations is 33,364 in NLSY79 and 15,410 in NLSY97. The average unemployment rate during the period of NLSY79 is 6.39%, higher than that of NLSY97. The average age at the time of interviews is 31.64 for NLSY79 and 22.09 for NLSY97, indicating that the cohort NLSY97 is a younger sample. The average highest grade completed at age 22 and age 33 both increase from NLSY79 to NLSY97. The share of African-American and Hispanic also rises from NLSY79 to NLSY97. The average mismatch rate, positive mismatch, and (absolute value of) negative mismatch are all increase from NLSY79 to NLSY97.

**Table 3.1:** Descriptive Statistics of the Samples NLSY79 and NLSY97

| Statistics | NLSY79 | NLSY97 |
|---|---|---|
| Total Number of Observations | 33,364 | 15,410 |
| Total Number of Individuals | 1,992 | 2,422 |
| | | |
| Average Unemployment Rate (original) | 6.39% | 5.64% |
| Average Age at Time of Interview | 31.64 | 22.09 |
| | | |
| Average Highest Grade at Age 22 | 12.4 | 12.7 |
| Average Highest Grade at Age 33 | 13.8 | 14.1 |
| | | |
| Percentage of African-American | 11.4% | 15.32 % |
| Percentage of Hispanic | 6.87 % | 11.64 % |
| | | |
| Average Mimsatch | 0.258 | 0.292 |
| Average Positive Mismatch | 0.128 | 0.145 |
| Average Negative Mismatch | -0.129 | -0.141 |

Note: The data source is NLSY79, NLSY97, and O*NET. The measurement of mismatch and sample criteria are from Guvenen et al. (2020). Guvenen et al. (2020) rescale the mismatch measure in their paper. The measure is not rescaled in this paper in order to be compared across two cohorts.
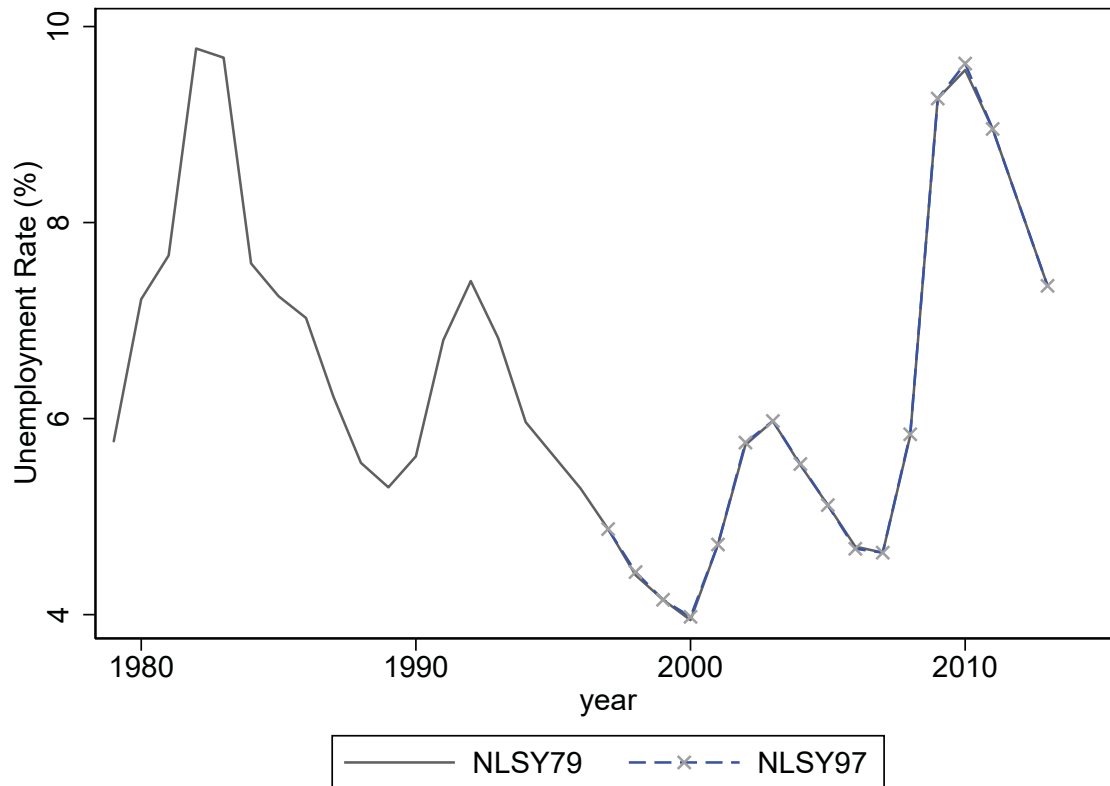
## 3.3 Empirical Analysis

In this section, I study the causal effect of labor market conditions on skill mismatch rate and examine which of the sullying and cleansing effects of recessions on mismatch dominate in NLSY79 and in NLSY97. First I show the labor market conditions (unemployment rate as proxy) during the periods of the two cohorts NLSY79 and NLSY97. Then I present an empirical specification about the relationship between mismatch rate and unemployment. In the last part, I show the empirical results.

### 3.3.1 Labor Market Conditions

The unemployment rate is an appropriate proxy for labor market tightness, which is a key macroeconomic factor that affects each individual's occupational choice. Figure 3.1 plots the annual national unemployment rate from 1979 to 2013, which covers the periods of NLSY79 (1979-2010) and NLSY97 (1997-2015). The average unemployment rate during the period of NLSY79 is 6.39%, higher than that of NLSY97 (5.64%). For the workers who entered the labor market at the beginning of the 1980s, the unemployment rate increased from around 6% to almost 10% and the average unemployment rate of this period is about 8%. For the workers who entered the labor market around 1997, the unemployment rate decreases from around 5% in 1997 to 4% in 2000. The average unemployment rate during this period is only 4.3%, much lower than the average rate at the beginning of the 1980s. The four recessions are: 1980-1983, 1990-1992, 2000-2003, 2008-2010.

Figure 3.2 and Figure 3.3 display the evolution of aggregate mismatch rate and the detrended national unemployment rate for the cohorts of NLSY79 and NLSY97. The time span is 32 years (1979-2010) for NLSY79 and 19 years (1997-2015) for NLSY97. There are two points that can be learned from the figures. First of all, the

general trend of aggregate mismatch rate is that it decreases over time in both cohorts. It decreases from 0.274 to 0.248 in NLSY79 and from 0.320 to 0.262 in NLSY97.



**Figure 3.1:** Unemployment Rate in NLSY79 and NLSY97

Second, the mismatch rate is procyclical in NLSY79 and NLSY97. In economic downturns, the mismatch rate decreases at a faster speed than other periods as the unemployment rate increases. (see 1980-1983, 1990-1992, 2008-2010 in Figure 3.2 and 2000-2003, 2008-2010 in Figure 3.3) The mechanism behind this phenomenon suggests the cleansing effect of recessions. In recessions, low-quality matches are destroyed while only high-quality matches remain, which leads to a lower mismatch rate.

**Figure 3.2:** Aggregate Mismatch Rate and Unemployment in NLSY79



**Figure 3.3:** Aggregate Mismatch Rate and Unemployment in NLSY97

## 3.3.2  Empirical Specification and Results

To examine the effect of labor market conditions on match quality and how the effect varies across different educational levels, I run the following regression:

$$m_{i,t} = \beta_0 + \beta_1 Age_{i,t} + \beta_2 Unempl_t + \beta_3 EDU_{i,t}$$
$$+ \beta_4 EDU_{i,t} * Unempl_t + \beta_5 X_{i,t} + \epsilon_{i,t}$$

(3.1)

In equation 3.1, $m_{i,t}$ is the 3-dimension (verbal, math, social) mismatch rate of individual i with his occupation at time t. Following Chapter 2 and Guvenen et al. (2020), the mismatch rate $m_{i,t}$ is defined as the weighted sum of the difference in each of the three dimensions between worker's abilities and occupational skill requirements. $Age$ is the age of individual i at time t. $Unempl_t$ is the annual national unemployment rate (detrended) at time t. $EDU_{i,t}$ is the highest grade completed of individual i at time t. $X_{i,t}$ is a set of control variables: race and one-digit industry. $\epsilon_{i,t}$ is the error term that captures the unobserved determinants of mismatch.

In this specification, I include age, unemployment rate, educational level, their interaction term, as well as other control variables. The mismatch rate tends to decrease as age goes up since workers switch their occupations and gradually find the occupation which is a better match for them. As the educational level increases, the mismatch rate is expected to go down since the highly educated workers are more skilled and specialized to find a better match than the lowly educated workers. I include the interaction term between unemployment and educational level to measure how the education adjusts the effect of unemployment on mismatch.

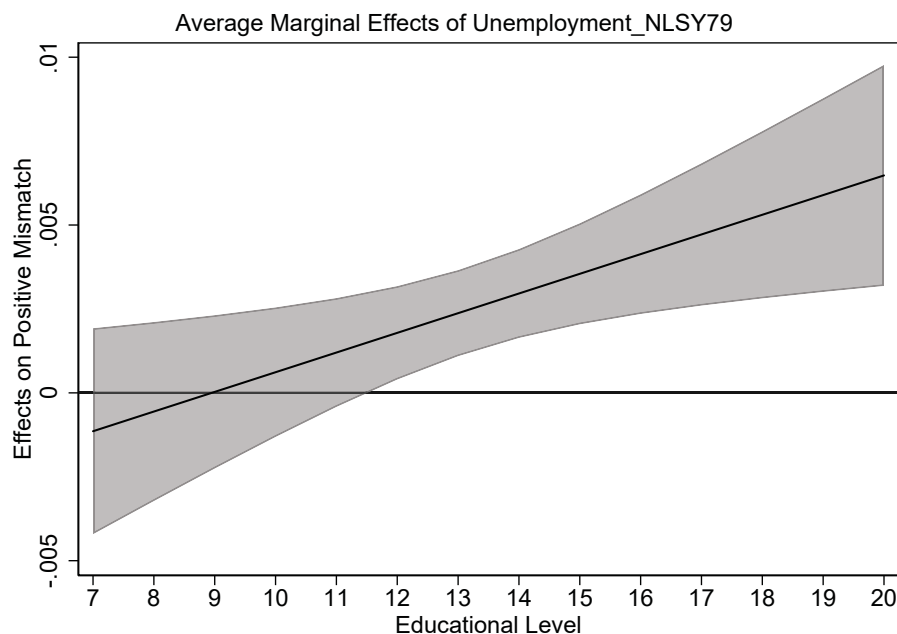**Table 3.2:** The Effect of Unemployment on Mismatch in NLSY79

| NLSY79 | (1) Mismatch | (2) Mismatch | (3) Positive Mismatch | (4) Positive Mismatch | (5) Negative Mismatch | (6) Negative Mismatch |
|---|---|---|---|---|---|---|
| Age | -0.000162 | -0.000226* | -0.00235*** | -0.00238*** | 0.00219*** | 0.00215*** |
| | (-1.34) | (-1.85) | (-19.84) | (-19.91) | (18.75) | (18.31) |
| Unemployment | 0.000509 | -0.0174*** | 0.00263*** | -0.00525* | -0.00213*** | -0.0121*** |
| | (0.77) | (-4.97) | (4.02) | (-1.70) | (-3.35) | (-3.71) |
| EDU | -0.00273*** | -0.00256*** | 0.00782*** | 0.00789*** | -0.0105*** | -0.0104*** |
| | (-6.73) | (-6.29) | (20.56) | (20.68) | (-28.52) | (-28.19) |
| EDU×Unemployment | | 0.00133*** | | 0.000586** | | 0.000745*** |
| | | (5.31) | | (2.56) | | (3.31) |
| Constant | 0.301*** | 0.301*** | 0.108*** | 0.108*** | 0.193*** | 0.194*** |
| | (35.96) | (36.02) | (13.91) | (13.93) | (25.60) | (25.62) |
| Observations | 33364 | 33364 | 33364 | 33364 | 33364 | 33364 |

$t$ statistics in parentheses

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: The data source is NLSY79, NLSY97, and O*NET. The measurement of mismatch is proposed by Guvenen et al. (2020). I use the crosswalk of Guvenen et al. (2020) to convert the different industry classifications in NLSY79 and NLSY97 into the Census 1970 One-Digit Industry Code. Guvenen et al. (2020) rescale the mismatch measure in their paper. The measure is not rescaled in this paper in order to be compared across two cohorts. The control variables are race, one-digit industry, and year effects.

93

Table 3.2 presents the OLS estimates of equation 3.1 using the data of NLSY79. Columns (1), (3), and (5) display the estimation results without the interaction term between educational level and unemployment rate. In column (1), the effect of unemployment on the aggregate mismatch is insignificant. In columns (3) and (5), I decompose the mismatch into positive and negative mismatch and examine the effects of unemployment on them. Column (3) shows a positive relationship between positive mismatch and unemployment. The positive mismatch increases with the unemployment rate, i.e. positive mismatch is countercyclical. The sullying effect plays a role and the mechanism suggests that workers cannot switch to a better job match in economic slumps. Column (5) shows a negative relationship between negative mismatch and unemployment. The negative mismatch rate is procyclical. This can be explained by the cleansing effect: low-quality matches are destroyed in recessions and only high-quality matches remain.



**Figure 3.4:** Average Marginal Effect of Unemployment on Positive Mismatch in NLSY79

Columns (2), (4), and (6) present the estimation results with the interaction term between educational level and unemployment rate. In column (2), the coefficient of the interaction term is positive and significant, indicating that the increase in the educational level will lower the negative effect of unemployment rate on mismatch.

In column (4), the coefficient of the interaction term is also positive and significant. The increasing educational level will improve the positive effect of the unemployment on mismatch. Figure 3.4 shows the average marginal effects of unemployment rate on positive mismatch conditional on the educational level (grade) with 95% CIs. The average marginal effects of unemployment on positive mismatch is positive and increase with the educational levels. Higher educational level will increase the countercyclicality of positive mismatch.

In column (6) of table 3.2, the coefficient of the interaction term is significantly positive. Negative mismatch is procyclical and higher educational level can offset part of the procyclicality of mismatch. Figure 3.5 shows the average marginal effects of unemployment rate on negative mismatch conditional on the educational level (grade) with 95% CIs. From figure 3.5, I find that under grade 16, the average marginal effect of unemployment rate on mismatch is negative and the negative effect tends to zero as grade increases to grade 16.

From Table 3.2, I find that the mismatch rate decreases with age and education level. Positive mismatch increases as the educational level goes up while negative mismatch decreases as the educational level goes up.
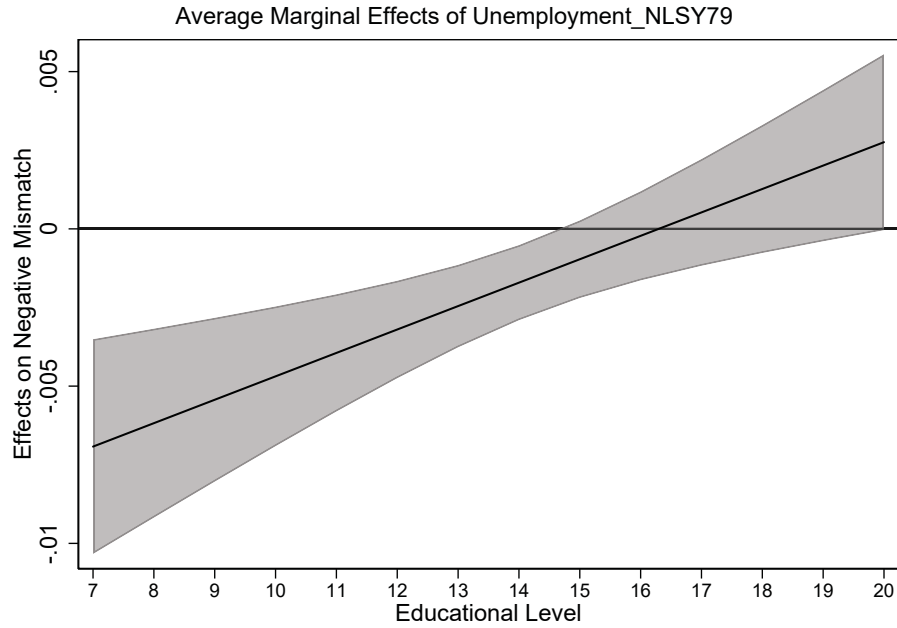
**Table 3.3:** The Effect of Unemployment on Mismatch in NLSY97

| NLSY97 | (1) Mismatch | (2) Mismatch | (3) Positive Mismatch | (4) Positive Mismatch | (5) Negative Mismatch | (6) Positive Mismatch |
|---|---|---|---|---|---|---|
| Age | -0.000269 | -0.000767 | -0.0109*** | -0.0116*** | 0.00982*** | 0.00995*** |
| | (-0.44) | (-1.23) | (-19.72) | (-20.43) | (16.16) | (15.95) |
| Unemployment | -0.00240* | 0.0145*** | 0.00214* | 0.0234*** | -0.00371*** | -0.00830* |
| | (-1.74) | (3.15) | (1.77) | (5.90) | (-2.67) | (-1.69) |
| EDU | -0.00185*** | -0.00123* | 0.00878*** | 0.00956*** | -0.0107*** | -0.0109*** |
| | (-2.59) | (-1.67) | (13.19) | (13.59) | (-14.34) | (-14.36) |
| EDU ×Unemployment | | -0.00130*** | | -0.00163*** | | 0.000352 |
| | | (-3.89) | | (-5.41) | | (1.01) |
| Constant | 0.294*** | 0.299*** | 0.238*** | 0.244*** | 0.0700*** | 0.0687*** |
| | (20.12) | (20.38) | (17.08) | (17.61) | (5.17) | (5.03) |
| Observations | 14791 | 14791 | 14791 | 14791 | 15410 | 15410 |

$t$ statistics in parentheses

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: The data source is NLSY79, NLSY97, and O*NET. The measurement of mismatch is proposed by Guvenen et al. (2020). I use the crosswalk of Guvenen et al. (2020) to convert the different industry classifications in NLSY79 and NLSY97 into the Census 1970 One-Digit Industry Code. Guvenen et al. (2020) rescale the mismatch measure in their paper. The measure is not rescaled in this paper in order to be compared across two cohorts. The control variables are race, one-digit industry, and year effects.

**Figure 3.5:** Average Marginal Effect of Unemployment on Negative Mismatch in NLSY79

Table 3.3 presents the OLS estimates of equation 3.1 using the data of NLSY97. Columns (1), (3), and (5) display the estimation results without the interaction term between educational level and unemployment rate. Column (1) shows a negative relationship between mismatch and unemployment. The mismatch will decline if unemployment increases, i.e. mismatch is procyclical. Similar to the estimation results of NLSY79, the cleansing effect of recessions plays a role. In recessions, bad matches are scrapped and only high-quality matches remain. Column (3) shows a positive relationship between positive mismatch and unemployment. Therefore positive mismatch is countercyclical. The estimate of the unemployment rate in Column (5) is negative. Thus negative mismatch is also procyclical in NLSY97. Overall, the mismatch rate of NLSY97 is procyclical.

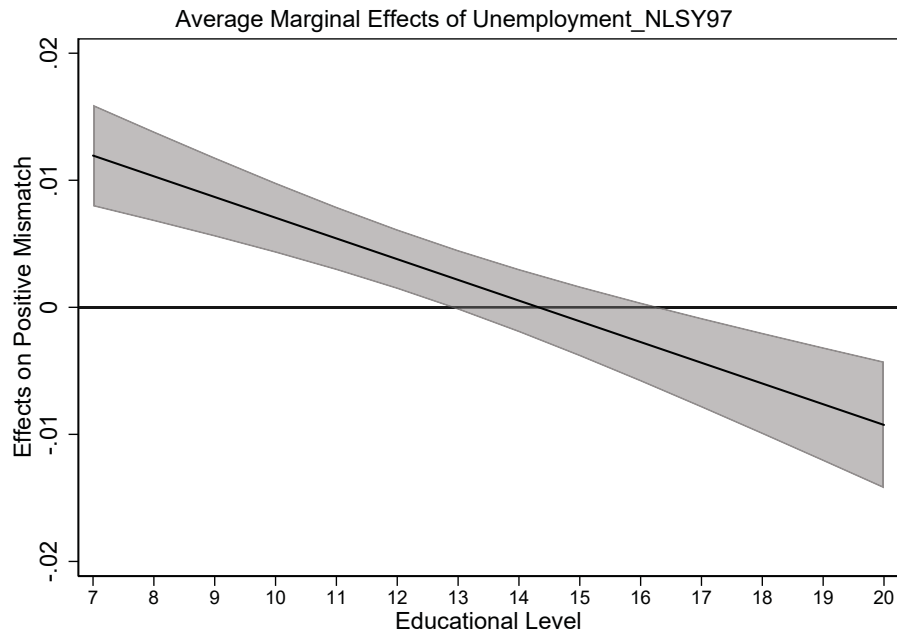**Average Marginal Effects of Unemployment_NLSY97**

**Figure 3.6:** Average Marginal Effect of Unemployment on Mismatch in NLSY97

Columns (2), (4), and (6) present the estimation results with the interaction term between educational level and unemployment rate. In column (2), the coefficient of the interaction term is negative and significant, indicating that the increase in the educational level will improve the negative effect of unemployment rate on mismatch. Figure 3.6 shows the average marginal effects of the unemployment rate on mismatch conditional on the educational level (grade) with 95% CIs. The average marginal effect of unemployment on mismatch is negative and the absolute value gets larger as the educational level goes up. Higher educational level will increase the procyclicality of mismatch.

In column (4), the coefficient of the interaction term is also negative and significant. The increasing educational level will decrease the effect of unemployment on positive mismatch. Figure 3.7 shows the average marginal effects of the unemployment rate on positive mismatch conditional on the educational level (grade) with 95% CIs. The average marginal effects of unemployment on positive mismatch are

firstly positive and then negative as the educational level goes up. Higher educational level will decrease the countercyclicality of positive mismatch.



**Figure 3.7:** Average Marginal Effect of Unemployment on Positive Mismatch in NLSY97

In NLSY79, positive mismatch is countercyclical and its countercyclicity increases with educational level. In NLSY97, positive mismatch is also countercyclical but its countercyclicity decreases with educational level. My understanding about this difference is that in both cohorts sullying effect of recessions outweighs the cleansing effect. The speed of workers to reallocate to better matches is slow. In NLSY97, however, for the higher educated workers, the cleansing effect gets more important than in NLSY79. Therefore the countercyclicality of positive mismatch decreases.

**Table 3.4:** Cyclicality of Mismatch in NLSY79 and NLSY97

| | NLSY79 | | |
| --- | --- | --- | --- |
| | Mismatch | Positive Mismatch | Negative Mismatch |
| Cyclicality | Insignificant | Counter- | Pro- |
| EDU | —- | Counter-↑ as EDU ↑ | Pro-↓ as EDU ↑ |
| | NLSY97 | | |
| | Mismatch | Positive Mismatch | Negative Mismatch |
| Cyclicality | Pro- | Counter- | Pro- |
| EDU | Pro-↑ as EDU ↑ | Counter-↓ as EDU ↑ | —- |

Note: In NLSY79, the effect of unemployment on mismatch is insignificant. In NLSY97, the effect of the interaction term between education and unemployment on negative mismatch is insignificant.

In column (6) of table 3.3, the coefficient of the interaction term is statistically insignificant. Negative mismatch is procyclical and educational level doesn't affect its procyclicality. In aggregate, the mismatch of NLSY79 is procyclical.

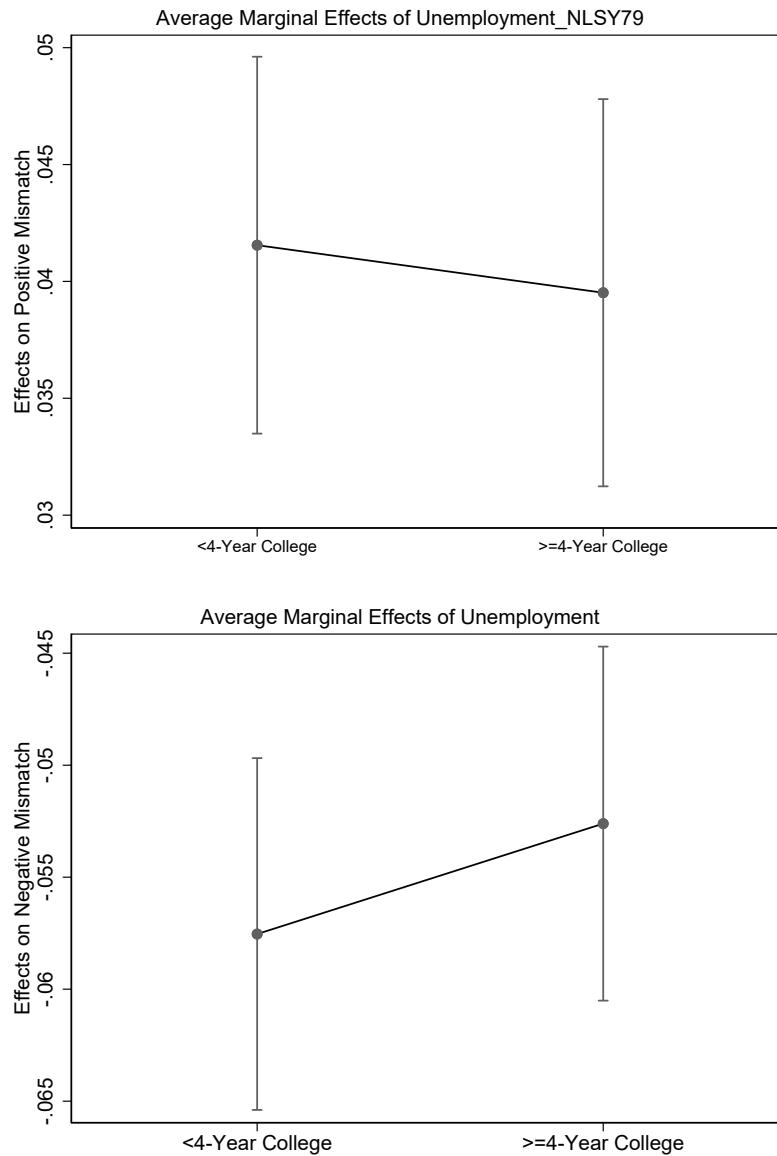Table 3.4 summarizes the cyclicality of mismatch, positive mismatch, and negative mismatch in NLSY79 and NLSY97. In NLSY79, the effect of unemployment on mismatch is insignificant. In NLSY97, the effect of the interaction term between education and unemployment on negative mismatch is insignificant. Mismatch is procyclical in NLSY97. Higher education level increases the procyclicality of mismatch in NLSY97. In NLSY79, higher education level increases the countercyclical of positive mismatch and lowers the procyclicality of negative mismatch.

Table 3.5 and Table 3.6 display the estimation results for NLSY79 and NLSY97 where I use the dummy variable '≥4-Year College or not' as the indicator of educational level. The estimation results are consistent with Table 3.2 and Table 3.3.

Figure 3.8 shows the average marginal effects of the unemployment rate on positive mismatch and negative mismatch in NLSY79. Being college graduates will decrease the procyclicality of negative mismatch but increase the countercyclicality of positive mismatch.



**Figure 3.8:** Average Marginal Effect of Unemployment in NLSY79 (4-Year College as Education dummy)

**Table 3.5:** The Effect of Unemployment on Mismatch in NLSY79 (4-Year College as Education dummy)

| NLSY79 | (1) Mismatch | (2) Mismatch | (3) Positive Mismatch | (4) Positive Mismatch | (5) Negative Mismatch | (6) Negative Mismatch |
|---|---|---|---|---|---|---|
| Age | -0.0000525 | -0.0000780 | -0.00201*** | -0.00198*** | 0.00196*** | 0.00190*** |
| | (-0.44) | (-0.65) | (-17.11) | (-16.74) | (16.96) | (16.34) |
| Unemployment | 0.000394 | -0.000506 | 0.00250*** | 0.00364*** | -0.00211*** | -0.00415*** |
| | (0.60) | (-0.64) | (3.80) | (4.72) | (-3.30) | (-5.28) |
| ≥ 4-Year College | -0.0349*** | -0.0344*** | 0.0106*** | 0.00999*** | -0.0455*** | -0.0444*** |
| | (-17.10) | (-16.78) | (4.90) | (4.57) | (-24.64) | (-24.05) |
| ≥ 4-Year College×Unemployment | | 0.00335** | | -0.00424*** | | 0.00759*** |
| | | (2.52) | | (-3.04) | | (6.55) |
| Constant | 0.267*** | 0.268*** | 0.188*** | 0.186*** | 0.0796*** | 0.0815*** |
| | (36.26) | (36.31) | (26.09) | (25.90) | (12.07) | (12.32) |
| Observations | 33364 | 33364 | 33364 | 33364 | 33364 | 33364 |

$t$ statistics in parentheses

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: The data source is NLSY79, NLSY97, and O*NET. The measurement of mismatch is proposed by Guvenen et al. (2020). I use the crosswalk of Guvenen et al. (2020) to convert the different industry classifications in NLSY79 and NLSY97 into the Census 1970 One-Digit Industry Code. Guvenen et al. (2020) rescale the mismatch measure in their paper. The measure is not rescaled in this paper in order to be compared across two cohorts. The control variables are race, one-digit industry, and year effects.

**Table 3.6:** The Effect of Unemployment on Mismatch in NLSY97 (4-Year College as Education dummy)

| | (1) Mismatch | (2) Mismatch | (3) Positive Mismatch | (4) Positive Mismatch | (5) Negative Mismatch | (6) Negative Mismatch |
|---|---|---|---|---|---|---|
| Age | 0.000619 | 0.000550 | -0.00884*** | -0.00905*** | 0.00854*** | 0.00870*** |
| | (1.04) | (0.90) | (-16.08) | (-16.13) | (14.52) | (14.39) |
| Unemployment | -0.00281** | -0.00244 | 0.00106 | 0.00219* | -0.00301** | -0.00384** |
| | (-2.05) | (-1.58) | (0.87) | (1.65) | (-2.16) | (-2.41) |
| $\geq$ 4-Year College | -0.0473*** | -0.0459*** | 0.0148*** | 0.0190*** | -0.0600*** | -0.0630*** |
| | (-8.67) | (-7.65) | (2.66) | (3.02) | (-13.17) | (-13.33) |
| $\geq$ 4-Year College×Unemployment | | -0.00154 | | -0.00471* | | 0.00344 |
| | | (-0.60) | | (-1.86) | | (1.48) |
| Constant | 0.256*** | 0.258*** | 0.290*** | 0.295*** | -0.0191 | -0.0225 |
| | (17.62) | (17.41) | (21.32) | (21.36) | (-1.36) | (-1.57) |
| Observations | 14791 | 14791 | 14791 | 14791 | 15410 | 15410 |

$t$ statistics in parentheses

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: The data source is NLSY79, NLSY97, and O*NET. The measurement of mismatch is proposed by Guvenen et al. (2020). I use the crosswalk of Guvenen et al. (2020) to convert the different industry classifications in NLSY79 and NLSY97 into the Census 1970 One-Digit Industry Code. Guvenen et al. (2020) rescale the mismatch measure in their paper. The measure is not rescaled in this paper in order to be compared across two cohorts. The control variables are race, one-digit industry, and year effects.

## 3.4 Conclusion

In this paper, by combining the data sets of NLSY79/NLSY97 and O*NET, I measure the 3-dimensional skill mismatch between workers and their occupations and examine the cyclicality of this mismatch rate in the cohorts of NLSY79 and NLSY97. I also examine whether the effect of recessions varies by different educational levels. Then I compare the cyclicality of mismatch in the two cohorts.

Using the combined data set, I find a strong procyclical pattern of the 3-dimensional skill mismatch in the cohort of NLSY97. In NLSY79 and NLSY97, positive mismatch is countercyclical and negative mismatch is procyclical.

I then provide evidence that the effect of unemployment on mismatch is adjusted by the educational level. In NLSY79, the procyclicality of negative mismatch is weakened by higher educational level. The countercyclicality of positive mismatch increases as the educational level goes up. In NLSY97, the procyclicality of mismatch is strengthened by higher educational level and the countercyclicality of positive mismatch decreases with educational level. I also use the dummy variable '$\geq$4-Year College or not' as the indicator of educational level and I find the results about the procyclical pattern in both cohorts are robust.

# Bibliography

Joseph G Altonji, Lisa B Kahn, and Jamin D Speer. Cashier or consultant? entry labor market conditions, field of study, and career success. *Journal of Labor Economics*, 34(S1):S361–S401, 2016.

Manuel Arellano and Stephen Bond. Some tests of specification for panel data: Monte carlo evidence and an application to employment equations. *The review of economic studies*, 58(2):277–297, 1991.

David H Autor, Frank Levy, and Richard J Murnane. The skill content of recent technological change: An empirical exploration. *The Quarterly journal of economics*, 118(4):1279–1333, 2003.

Gadi Barlevy. The sullying effect of recessions. *The Review of Economic Studies*, 69(1): 65–96, 2002.

Regis Barnichon and Yanos Zylberberg. Underemployment and the trickle-down of unemployment. *American Economic Journal: Macroeconomics*, 11(2):40–78, 2019.

Thomas K Bauer. Educational mismatch and wages: a panel analysis. *Economics of education review*, 21(3):221–229, 2002.

Nathaniel Baum-Snow and Ronni Pavan. Understanding the city size wage gap. *The Review of economic studies*, 79(1):88–127, 2011.

Ricardo J Caballero and Mohamad L Hammour. The cleansing effect of recessions. *The American Economic Review*, pages 1350–1368, 1994.

Pierre-Philippe Combes, Gilles Duranton, and Laurent Gobillon. Spatial wage disparities: Sorting matters! *Journal of urban economics*, 63(2):723–742, 2008.

David Deming and Lisa B Kahn. Skill requirements across firms and labor markets: Evidence from job postings for professionals. *Journal of Labor Economics*, 36(S1): S337–S369, 2018.

Greg J Duncan and Saul D Hoffman. The incidence and wage effects of overeducation. *Economics of education review*, 1(1):75–86, 1981.

Jan Eeckhout, Roberto Pinheiro, and Kurt Schmidheiny. Spatial sorting. *Journal of Political Economy*, 122(3):554–620, 2014.

Peter Fredriksson, Lena Hensvik, and Oskar Nordström Skans. Mismatch of talent: Evidence on match quality, entry wages, and job mobility. *American Economic Review*, 108(11):3303–38, 2018.

Christa Frei and Alfonso Sousa-Poza. Overqualification: permanent or transitory? *Applied Economics*, 44(14):1837–1847, 2012.

Edward L Glaeser and David C Mare. Cities and skills. *Journal of labor economics*, 19 (2):316–342, 2001.

Fatih Guvenen, Burhan Kuruscu, Satoshi Tanaka, and David Wiczer. Multidimensional skill mismatch. *American Economic Journal: Macroeconomics*, 12(1):210–44, 2020.

Brad Hershbein and Lisa B Kahn. Do recessions accelerate routine-biased technological change? evidence from vacancy postings. *American Economic Review*, 108 (7):1737–72, 2018.

Lisa B Kahn. The long-term labor market consequences of graduating from college in a bad economy. *Labour economics*, 17(2):303–316, 2010.

Billy Frazier Kiker, Maria C Santos, and M Mendes De Oliveira. Overeducation and undereducation: evidence for portugal. *Economics of Education Review*, 16(2): 111–125, 1997.

Jeremy Lise and Fabien Postel-Vinay. Multidemensional skills, sorting, and human capital accumulation. *American Economic Review*, 110(8):2328–2376, 2020.

Jeremy Lise and Jean-Marc Robin. The macrodynamics of sorting between workers and firms. *American Economic Review*, 107(4):1104–35, 2017.

Kai Liu, Kjell G Salvanes, and Erik Ø Sørensen. Good skills in bad times: Cyclical skill mismatch and the long-term effects of graduating in a recession. *European Economic Review*, 84:3–17, 2016.

Philip Oreopoulos, Till Von Wachter, and Andrew Heisz. The short-and long-term career effects of graduating in a recession. *American Economic Journal: Applied Economics*, 4(1):1–29, 2012.

Theodore Papageorgiou. Worker sorting and agglomeration economies. *Working paper*, 2013.

Theodore Papageorgiou. Occupational matching and cities. Technical report, 2020.

Jennifer Roback. Wages, rents, and the quality of life. *Journal of political Economy*, 90 (6):1257–1278, 1982.

Jorge De La Roca and Diego Puga. Learning by working in big cities. *The Review of Economic Studies*, 84(1):106–142, 2017.

Morris Rosenberg. Rosenberg self-esteem scale (rse). *Acceptance and commitment therapy. Measures package*, 61(52):18, 1965.

Joseph Alois Schumpeter. *Business cycles*, volume 1. McGraw-Hill New York, 1939.

Alexandra Spitz-Oener. Technical change, job tasks, and rising educational demands: Looking outside the wage structure. *Journal of labor economics*, 24(2):235–270, 2006.

Zhibiao Wu and Martha Palmer. Verbs semantics and lexical selection. In *Proceedings of the 32nd annual meeting on Association for Computational Linguistics*, pages 133–138. Association for Computational Linguistics, 1994.