

Income Inequality and the Labor Market

Xinyuan Zhu

Jiaxing Senior High School, Jiaxing,
China

Zhuxinyuan080423@outlook.com

Abstract:

Income inequality is a major problem in labor economics, especially in a rapidly developing economy like China. With China's remarkable economic growth in the past few decades, the issue of income distribution has attracted more and more attention from scholars and decision-makers. This article investigates the Simpson paradox, a statistical phenomenon in which the trend of aggregated data conflicts with the subgroup pattern and explores how this paradox distorts the interpretation of the income gap. This paradox is particularly relevant in the context of China. Regional differences and departmental subdivisions create complex data structures that need to be carefully analyzed. By integrating set theory with labor economics, we construct an analytical framework to examine earnings distribution across gender and industry segments within China's labor market. Our methodological approach involves defining the labor market as a universal set composed of multiple mutually exclusive subsets based on demographic and economic characteristics. Through this framework, we systematically analyze how income distribution patterns vary across different population strata and how these variations might be masked in aggregate data. The analysis of large-scale data sets including China's household income items shows that relying entirely on comprehensive indicators such as the Gini coefficient may mask significant inequality within the subgroup. The research process includes collecting data from multiple sources, stratifying samples, and comparative analysis between and within subgroups. These findings highlight the need to implement layered data analysis in economic research and policy formulation, and also provide actionable suggestions to mitigate deviations in inequality measurement.

Keywords: Income inequality; Simpson paradox; Set theory; Labor Market.

1. Introduction

Income inequality is an economic and social challenge facing countries around the world [1]. In China, while the economy is growing rapidly, the income distribution gap is also gradually widening, which has attracted wide attention in the fields of research and policy [2]. Commonly used inequality measurement methods (such as the Gini coefficient) are often based on aggregate data and may ignore differences within different populations or economic groups.

The Simpson paradox was proposed by statistician Simpson in 1951, which means that the trends in grouped data may disappear or even reverse after consolidation [3]. This phenomenon shows that ignoring the data structure may lead to wrong conclusions.

This study adopts the set theory method to regard the labor market as a system composed of different subsets, so as to analyze the income distribution more carefully. By combining the theoretical framework with China's actual data, this article aims to reveal the structural problems behind inequality and provide reference for relevant research and policies. The structure of the article is as follows: the second part reviews the relevant literature, the third part introduces the research methods, the fourth part conducts simulation analysis, the fifth part discusses the meaning of the results, and the sixth part summarizes and puts forward the future research direction.

2. Literature Review

2.1 The Current Situation of Income Inequality in China

The study of income inequality in China reveals significant differences between regions, educational backgrounds and industries. The research of Li et al. and Yue recorded the continuous widening of the urban-rural income gap and the impact of industrial transformation on wage distribution [2,]. These studies usually use overall indicators. These numbers do give some useful pointers, but honestly, they don't really show the whole picture of what's going on within different groups. Check this back in 2022, the average disposable income for city folks hit 49,283 yuan, while for rural residents it was only 20,133 yuan. and the ratio of urban and rural income reached 2.45:1 [5]. Between different industries, the average annual salary of practitioners in the information transmission, software and information technology service industry is 181,006 yuan, while the average annual salary of agriculture, forestry, animal husbandry and fishery is only 54,957 yuan, the difference between the two is more than 3 times.

2.2 Simpson's Paradox in Social Science

Simpson's paradox is widely discussed in the fields of statistics and applications. When studying the employment

rate of college graduates, Wang found that although the employment rate of women in all majors is higher than that of men, the overall data shows that the employment rate of men is higher [6]. Liu Hezhang's research based on the data of the China Health and Nutrition Survey (CHNS) shows that there is a similar phenomenon in the analysis of medical expenditure [7]. From an international perspective, Liu discussed this paradox through many practical cases and emphasized its importance in policy formulation [8]. For example, in the postgraduate admission data of the University of California, Berkeley in 1973, the admission rate of women in individual departments was higher than that of men, but the overall data showed that the admission rate of men was higher.

2.3 Application of Set Theory in Economic Analysis

Set theory is often used in economics to construct models of complex systems and relationships, but its application in inequality research is still relatively limited [9]. In the past, the applications were mainly concentrated in the fields of social selection theory and game theory. This paper expands set theory to labor market analysis and provides a structured method for the decomposition and study of subgroup interaction. By defining the labor market as a complete U and dividing it into mutually exclusive subsets by industry, region, gender and other attributes, we can more accurately analyze the structural characteristics of income distribution. For example, the national labor market is divided into three subsets of eastern, central and western by region. The data shows that the average annual wages of urban non-private units in these three regions in 2022 will be 124,016-yuan, 97,096 yuan and 111,588 yuan respectively, showing obvious regional differences. Recruitment.

3. Methodology

3.1 The Current Situation of Income Inequality in China

The labor market is defined as a finite universal set U , representing all workers. It is partitioned into mutually exclusive subsets based on gender, industry, and region:

gender-based partition:

$$U = M \cup F$$

industry-based partition: $U = I_1 \cup I_2 \cup \dots \cup I_n$

region-based partition: $U = R_1 \cup R_2 \cup \dots \cup R_k$

Income is modeled as a function $U \rightarrow R^+$. The average income of a subset $S \in U$ is:

$$Y_s = 1/S \cdot \sum_{u \in S} u$$

3.2 Application of Set Theory in Economic Analysis

Simpson's paradox presents a counterintuitive statistical

phenomenon: when the trend observed in individual subgroups disappears or even reverses after data aggregation. In labor economics, a typical example shows that when analyzing gender income differences: although men's income may be higher than women's in each specific industry, cross-industry summary data may paradoxically show that women's average income is higher.

This statistical reversal stems from the structural differences in the distribution of labor between industries. Its mathematical basis lies in the weighted calculation method of the subgroup average. Formally speaking, if the overall average income of a gender group is defined as the weighted sum of their average income in each industry, the paradox will arise when the following conditions are met:

Among them, there is a significant difference between the industry distribution for female and the male distribution for male of the female labor force. Specifically, when female workers are concentrated in high-income industries, although their income is still lower than that of their male colleagues in each industry, their overall average income may be higher.

This formal expression emphasizes the importance of examining data at an appropriate subdivision level, because relying only on overall statistics may lead to a fundamentally wrong conclusion about the underlying relationship of the data.

4. Research Result and Discussion

Due to data limitations, we present a hypothetical scenario reflecting patterns identified in CHIP and CLDS surveys:

Assume two sectors:

High-income technology sector (IT)

Low-income service sector (IS)

Income values:

IT for male: 15,000 yuan; female: 14,000 yuan

IS for male: 6,000 yuan; female: 5,500 yuan

Subgroup weights:

90% of men in IT, 10% in IS

10% of women in IT, 90% in IS

Overall averages:

male: $0.9 * 15000 + 0.1 * 6000 = 14100$ yuan

female: $0.1 * 14000 + 0.9 * 5500 = 6350$ yuan

These results show that overall data may exaggerate inequality, but in different cases, the Simpson paradox may also cover up this inequality. These findings are consistent with existing studies on occupational segregation and gender wage gap in China [10].

The set theory method provides a clear way to identify and interpret such patterns, which is conducive to the implementation of more targeted policy interventions. For example, policies aimed at narrowing the gender-income

gap must address both the problem of pay equity within sectors and the problem of distribution imbalances between sectors.

5. Conclusion

This study introduces a set theory framework combined with Simpson's paradox to analyze the income inequality in China's labor market. The main research results include:

If subgroup analysis is not carried out, the overall index may produce misleading conclusions; Structural composition has a significant impact on the observed inequality; Set theory provides a powerful tool for decomposing and analyzing complex data structures.

Limitations include the use of hypothetical data; future research should include micro-level data sets such as CHIP or CLDS. In addition, relevant research can extend the framework to a dynamic environment or incorporate other variables such as education and experience.

References

- [1] Khan, A. R., & Riskin, C. (2021). Inequality and poverty in China in the age of globalization. *Oxford Development Studies*, 49(3), 217-235.
- [2] Li, S., Luo, C., & Wei, Z. (2021). The changing pattern of China's income distribution: New evidence from household survey data. *China Economic Review*, 66, 101585.
- [3] Simpson, E. H. (1951). The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society: Series B*, 13(2), 238-241.
- [4] Yue, X., & Zhang, J. (2022). Regional inequality and labor mobility in China: A structural perspective. *Journal of Chinese Economic and Business Studies*, 20(1), 45-63.
- [5] National Bureau of Statistics of China. (2022). *China Statistical Yearbook 2022*. China Statistics Press.
- [6] Wang, F. (2020). Simpson's Paradox in the context of employment and education: Evidence from China. *Journal of Labor Research*, 41(3), 287-305.
- [7] Liu, C., & Zhang, L. (2021). Health expenditures and insurance coverage: A case of Simpson's Paradox in Chinese health data. *Health Economics Review*, 11(1), 1-12.
- [8] Liu, X. (2019). Understanding Simpson's Paradox: Applications in policy and research. *American Statistician*, 73(1), 82-90.
- [9] Chen, G., & Hamori, S. (2021). Set theory applications in behavioral economics. *Journal of Economic Surveys*, 35(4), 1125-1148.
- [10] Zhang, Y., & Zhao, Z. (2023). Gender wage gaps in China: The role of industry segregation. *Feminist Economics*, 29(1), 112-137.