

# Data Science for Economic and Social Issues

M. Jahangir Alam  
Texas A&M University

January 30, 2024

First Draft - Textbook

# Preface

This textbook is designed for students and researchers interested in the application of data science in economic and social issues...

First Draft - Textbook

First Draft - Textbook

# Contents

<b>Preface</b>	<b>iii</b>
<b>1 Introduction to Data Science and Economic Analysis</b>	<b>1</b>
1.1 Introduction to Data Science Tools . . . . .	2
1.2 Basics of Economic Policy Analysis . . . . .	2
1.3 Quantitative Methods in Economics . . . . .	2
1.4 Endogeneity and Selection Bias . . . . .	3
1.4.1 Introduction to Endogeneity . . . . .	3
1.4.2 What is Endogeneity? . . . . .	3
1.5 Multicollinearity and Causality Identification . . . . .	4
1.6 Approaches to Address These Problems . . . . .	4
<b>2 Causal Estimation Techniques</b>	<b>5</b>
2.1 Instrumental Variables (IV) . . . . .	6
2.2 Difference-in-Differences (DID) . . . . .	6
2.3 Regression Discontinuity Design (RDD) . . . . .	6
2.4 Propensity Score Matching (PSM) . . . . .	6
2.5 Interrupted Time Series (ITS) . . . . .	6
<b>3 Data Handling and Machine Learning in Economics</b>	<b>7</b>
3.1 Machine Learning Integration in Economics . . . . .	8
3.2 Data Preprocessing and Visualization . . . . .	8
3.3 Introduction to Prophet for Forecasting . . . . .	8
3.4 Introduction to LSTM for Sequence Data Analysis . . . . .	8
3.5 News Sentiment and Stock Price . . . . .	8

First Draft - Textbook

# List of Figures

First Draft - Textbook

First Draft - Textbook



# List of Tables

First Draft - Textbook

First Draft - Textbook

# Chapter 1

## Introduction to Data Science and Economic Analysis

### Contents

---

1.1	Introduction to Data Science Tools . . . . .	2
1.2	Basics of Economic Policy Analysis . . . . .	2
1.3	Quantitative Methods in Economics . . . . .	2
1.4	Endogeneity and Selection Bias . . . . .	3
1.4.1	Introduction to Endogeneity . . . . .	3
1.4.2	What is Endogeneity? . . . . .	3
1.5	Multicollinearity and Causality Identification . . . . .	4
1.6	Approaches to Address These Problems . . . . .	4

---

- 1.1 Introduction to Data Science Tools
- 1.2 Basics of Economic Policy Analysis
- 1.3 Quantitative Methods in Economics

First Draft - Textbook

## 1.4 Endogeneity and Selection Bias

### 1.4.1 Introduction to Endogeneity

Endogeneity is a significant concern in statistical modeling and econometrics. It refers to the scenario where key assumptions of the Classical Linear Regression Model (CLRM) are violated due to the correlation between the explanatory variables and the error term. The CLRM relies on several fundamental assumptions for the validity of the estimates:

- **Linearity:** The relationship between the dependent and independent variables is assumed to be linear.
- **Independence:** Observations are assumed to be independent of each other.
- **Homoscedasticity:** The error term is assumed to have a constant variance, irrespective of the value of the explanatory variables.
- **Normality:** For small sample sizes, it is assumed that the errors are normally distributed, at least approximately, for reliable inference.
- **No Endogeneity:** A critical assumption is that the error term should not be correlated with the independent variables.

Violation of these assumptions, particularly the absence of endogeneity, can lead to significant challenges in identifying causal relationships. Endogeneity can bias the estimates from a regression model, leading to incorrect conclusions about the relationship between the variables.

### 1.4.2 What is Endogeneity?

Endogeneity is a fundamental concept in econometrics that occurs when there is a correlation between an explanatory variable and the error term in a regression model. Consider a standard linear regression model:

$$Y = \beta_0 + \beta_1 X + u$$

In this model,  $Y$  represents the dependent variable,  $X$  is an explanatory variable, and  $u$  is the error term. Endogeneity is present if  $X$  is endogenous, which mathematically means that the covariance between  $X$  and  $u$  is not zero, i.e.,  $\text{Cov}(X, u) \neq 0$ .

This correlation between the explanatory variable and the error term can arise from various sources such as omitted variables, measurement errors, or simultaneous causality. When endogeneity is present, it leads to biased and inconsistent estimators in regression analysis, posing significant challenges to drawing reliable conclusions about causal relationships.

## **1.5 Multicollinearity and Causality Identification**

## **1.6 Approaches to Address These Problems**

First Draft - Textbook

## Chapter 2

# Causal Estimation Techniques

### Contents

2.1	Instrumental Variables (IV) . . . . .	6
2.2	Difference-in-Differences (DID) . . . . .	6
2.3	Regression Discontinuity Design (RDD) . . . . .	6
2.4	Propensity Score Matching (PSM) . . . . .	6
2.5	Interrupted Time Series (ITS) . . . . .	6

- 2.1 Instrumental Variables (IV)
- 2.2 Difference-in-Differences (DID)
- 2.3 Regression Discontinuity Design (RDD)
- 2.4 Propensity Score Matching (PSM)
- 2.5 Interrupted Time Series (ITS)

First Draft - Textbook



## Chapter 3

# Data Handling and Machine Learning in Economics

### Contents

3.1	Machine Learning Integration in Economics . . . . .	8
3.2	Data Preprocessing and Visualization . . . . .	8
3.3	Introduction to Prophet for Forecasting . . . . .	8
3.4	Introduction to LSTM for Sequence Data Analysis . . . . .	8
3.5	News Sentiment and Stock Price . . . . .	8

- 3.1 Machine Learning Integration in Economics
- 3.2 Data Preprocessing and Visualization
- 3.3 Introduction to Prophet for Forecasting
- 3.4 Introduction to LSTM for Sequence Data Analysis
- 3.5 News Sentiment and Stock Price

First Draft - Textbook