

CSC 5741 (2020/21) Data Mining and Warehousing Lecture 5: Exploratory Data Analysis

Lighton Phiri
Department of Library & Information Science
University of Zambia
<http://lis.unza.zm/~lightonphiri>

Lecture Series Outline

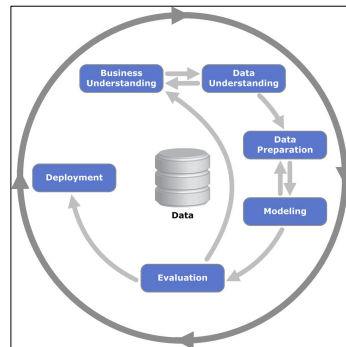
- Introduction
- Exploratory Data Analysis
- Descriptive Statistics
- Graphical Techniques
- Jupyter Notebook Walkthrough

May 30, 2021

CSC 5741 (2020/21) L05 - 2

Introduction (1/3)

- The Cross-industry standard process for data mining (CRISP-DM) is a model commonly used to highlight approaches in data mining
 - CRISP-DM segments a data mining project into six phases with no strict order of execution
 - Surveys conducted suggest CRISP-DM is the most widely used methodology

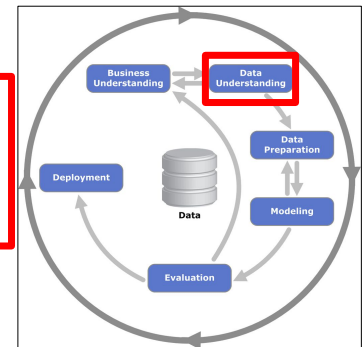


May 30, 2021

CSC 5741 (2020/21) L05 - 3

Introduction (2/3)

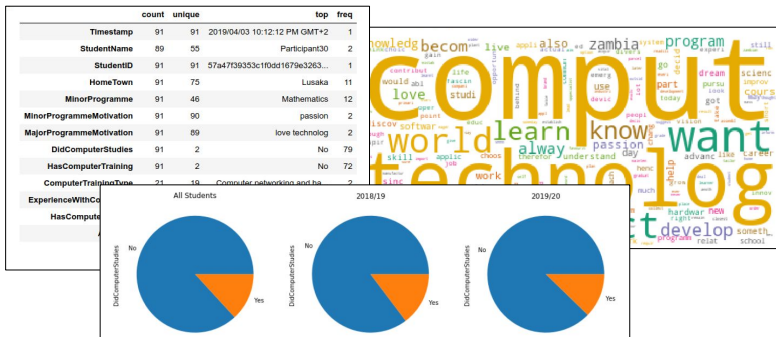
- Identify data sources
- Extract/collect required data
- Described and explore the data collected to gain some sense of what insights to derive
- Ascertain quality of data collected



May 30, 2021

CSC 5741 (2020/21) L05 - 4

Introduction (3/3)

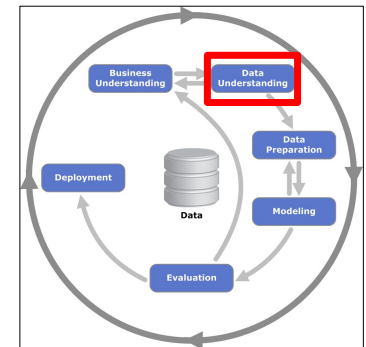


May 30, 2021

CSC 5741 (2020/21) L05 - 5

EDA Provides Insights From Datasets

- The purpose of this EDA is to find insights from datasets and/or data sources
 - Instrumental for setting the stage for data cleaning and transformation—output subsequently fed to machine learning algorithms.
 - Standard practice: **Data Understanding -> Data Preparation**

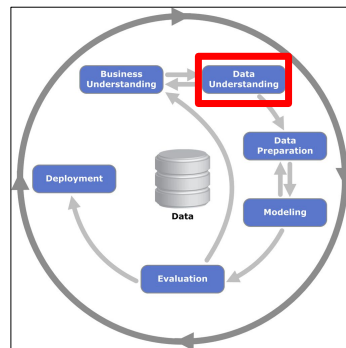


May 30, 2021

CSC 5741 (2020/21) L05 - 6

Goals and Objectives of EDA Process

- Various techniques are employed during EDA in order to achieve the following broad objectives:
 - Gain comprehensive insight of datasets
 - Identify important data characteristics
 - Identify outliers and anomalies
 - Determine correlations of various data characteristics

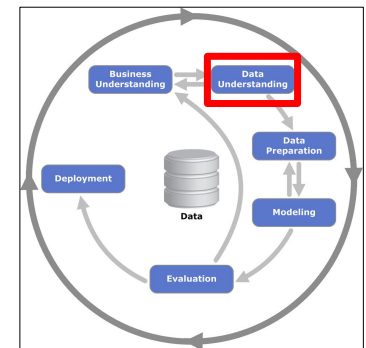


May 30, 2021

CSC 5741 (2020/21) L05 - 7

Outcomes of EDA Process

- Outcome of EDA
 - Important data attributes
 - Determine attribute characteristics—type of attribute, distribution and statistics (min, mode, median, mean)
 - Understand relationships between the different variables
 - DoB vs Age
 - List of outliers



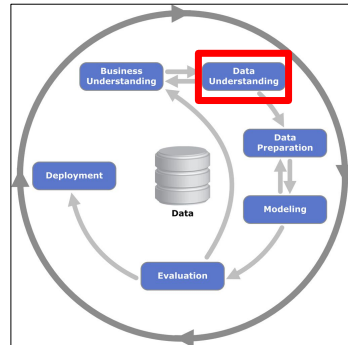
May 30, 2021

CSC 5741 (2020/21) L05 - 8

Exploratory Data Analysis (4/6)

- **Leading questions asked during EDA process**

- What are the different types of data attributes (categorical, continuous, ordinal)?
- How is the data distributed (normal vs non-normal)?
- Is there a correlation between data attributes and outcome?
- What are the most important data attributes?



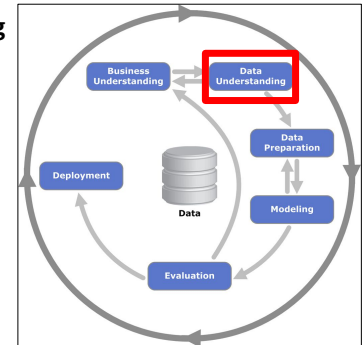
May 30, 2021

CSC 5741 (2020/21) L05 - 9

Exploratory Data Analysis (5/6)

- **Leading questions asked during EDA process**

- What must be done to data attributes with missing values?
- Do datasets have outliers?



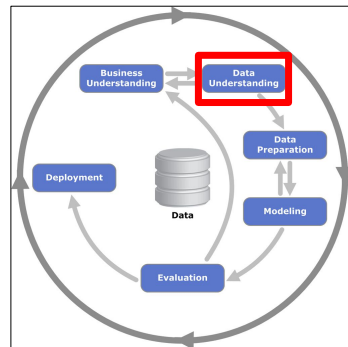
May 30, 2021

CSC 5741 (2020/21) L05 - 10

Exploratory Data Analysis (6/6)

- **A graphical approach to EDA is generally effective, although summary tables could also be used**

- Bar plots for categorical variables and aggregate data
- Line plots for continuous variables
- Histograms for continuous variables



May 30, 2021

CSC 5741 (2020/21) L05 - 11

Q & A Session

- **Comments, concerns and complaints?**

May 30, 2021

CSC 5741 (2020/21) L05 - 12

Lecture Series Outline

- Exploratory Data Analysis
- Descriptive Statistics
- Graphical Techniques
- Jupyter Notebook Walkthrough

Lecture Series Outline

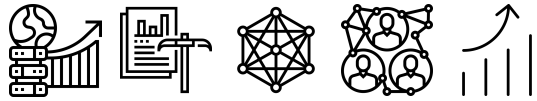
- Part I: Academic Talk
- Part III: Exploratory Data Analysis

Bibliography

- [1] Witten, I. H., Frank, E., Hall, M. A., Pal, C. J. (2017) Data Mining: Practical Machine Learning Tools and Techniques. Chapter 2
<https://www.cs.waikato.ac.nz/ml/weka/book.html>
- [2] NIST/SEMATECH e-Handbook of Statistical Methods. Exploratory Data Analysis. Chapter 1
<https://www.itl.nist.gov/div898/handbook/index.htm>
- [3] Seltman H. J. Experimental Design and Analysis. Exploratory Data Analysis. Chapter 4
<https://www.stat.cmu.edu/~hseltman/309/Book/chapter4.pdf>



✉ csc5741@unza.zm
🔗 <http://bit.ly/39HTdTK>
▶ <http://bit.ly/2kK2ZkA>



CSC 5741 (2020/21)
Data Mining and Warehousing
Lecture 5: Exploratory Data Analysis

Lighton Phiri
Department of Library & Information Science
University of Zambia

<http://lis.unza.zm/~lightonphiri>