

# Exploratory Data Analysis

Aayush Shrestha   Anmol Jha

Kathmandu University  
Department of Computational Mathematics

MATH 252 Progress Presentation

# Outline

1. Introduction
2. Objective
3. Progressed Work
4. Work to Complete

# Exploratory Data Analysis

Exploratory Data Analysis (EDA) is a method of analyzing data using statistical summaries and graphical representations.

EDA is a critical step in any Data Analysis or Data Science project as it provides a better understanding of data set variables and their relationships, and what data can reveal beyond formal modeling or hypothesis testing tasks. It aids in determining how to best manipulate data sources to obtain the required answers, making it easier for data scientists to discover patterns, detect anomalies, determine test hypotheses, and validate assumptions.

# Dashboard

A dashboard is basically a GUI for Data Visualization. A dashboard is a good choice if you need to summarize and present a lot of information on a single window.

A dashboard is one of the major tools for presenting data as it gives at-a-glance views of key performance indicators (KPI) relevant to a particular objective.

A dashboard helps Data Analysts and Data Scientists perform many data-related tasks and also provides a visual aid for other stakeholders to understand data and make accurate data-based decisions.

## Goals for the Project

- Familiarity with data cleaning, feature extraction, and data conversion.
- Describing the data using statistical methods such as mean, median, standard deviation, correlation, and so on.
- Implementation of different plots in R.
- Implementation of a dashboard.
- Understanding and explaining the trends and outliers in the data based on the domain.

## Initial View

## SAMPLE DATA

|                                |   |   |                |   |                                  |
|--------------------------------|---|---|----------------|---|----------------------------------|
| <b>Patient Id</b>              | PID0x81d5   | <b>Test 1</b>                                     | 0              | <b>Heart Rate</b><br>(rates/min)                        | Normal                           |
| <b>Patient Age</b>             | 7   | <b>Test 2</b>                                     | 0              | <b>Respiratory Rate</b><br>(breaths/min)                | Tachypnea                        |
| <b>Genes in mother's side</b>  | Yes   | <b>Test 3</b>                                     | 0              | <b>History of anomalies in previous pregnancies</b>     | Yes                              |
| <b>Inherited from father</b>   | Yes   | <b>Test 4</b>                                     | 1              | <b>No. of previous abortion</b>                         | 2                                |
| <b>Maternal gene</b>           | Yes   | <b>Test 5</b>                                     | 0              | <b>Birth defects</b>                                    | Singular                         |
| <b>Paternal gene</b>           | Yes   | <b>Parental consent</b>                           | Yes            | <b>White Blood cell count (thousand per microliter)</b> | 7.785072684                      |
| <b>Blood cell count (mcl.)</b> | 4.743937401   | <b>Follow-up</b>                                  | High           | <b>Blood test result</b>                                | slightly abnormal                |
| <b>Patient First Name</b>      | Irene   | <b>Gender</b>                                     | Female         | <b>Assisted conception IVF/ART</b>                      | Yes                              |
| <b>Family Name</b>             | Trainer   | <b>Birth asphyxia</b>                             | No record      | <b>Symptom 1</b>  | 1                                |
| <b>Father's name</b>           | Isaul   | <b>Autopsy shows birth defect (if applicable)</b> | No             | <b>Symptom 2</b>  | 1                                |
| <b>Mother's age</b>            | 31  | <b>Place of birth</b>                             | Institute      | <b>Symptom 3</b>  | 1                                |
| <b>Father's age</b>            | 61  | <b>Folic acid details (peri-conceptual)</b>       | No             | <b>Symptom 4</b>  | 0                                |
| <b>Institute</b>               | New England Medical Center  | <b>H/O serious maternal illness</b>               | No             | <b>Symptom 5</b>  | 1                                |
| <b>Location of Institute</b>   | 818 HARRISON AV<br>SOUTH END, MA 02118<br>(42.335925371008436,<br>-71.07378404259959) | <b>H/O radiation exposure (x-ray)</b>             | Not applicable | <b>Genetic Disorder</b>                                 | Single-gene inheritance diseases |
| <b>Status</b>                  | Deceased  | <b>H/O substance abuse</b>                        | Yes            | <b>Disorder Subclass</b>                                | Cystic fibrosis                  |

Figure: Sample raw data.

# Initial Pre-processing

First of all, the data was changed from an xlsx file to a csv file, and the column names were modified to better suit programming. Also, all the rows with any NA values were removed. After completing all that, the dataset now contains 45 columns and 6706 observations.

## EDA Work

Exploring the columns and their types, the columns with data related to the patient but not the disorder were removed. Then, a statistical summary of the numeric columns was created, and the columns having no variance were also removed. We are left with 32 columns.

Furthermore, based on graphical methods such as stacked bar graphs, box and whisker plots, and domain knowledge, a few more columns were disregarded based on their distribution and implication.



# Dashboard

## Remaining Work

- Develop a better dashboard with more and better visual aids and information.
- Further reduce the dimensionality of the parameters based on more robust statistical and domain knowledge.
- Finally, create a properly processed dataset to be used for future classification projects.

# THANK YOU!

Please share if you have any queries.  
We would like to show the code now.