

# Exploratory Data Analysis On Genetic Data

Aayush Shrestha   Anmol Jha

Kathmandu University  
Department of Computational Mathematics

MATH 252 Final Presentation  
Supervisor: Mr. Kiran Shrestha  
Special Thanks to Mr. Simon Shrestha

# Outline

1. Introduction
2. Methodology
3. Completed Work
4. Conclusion

# Exploratory Data Analysis

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

EDA is a critical process that 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 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.

# Project

Our project consist of the task to perform EDA on a Kaggle dataset so as to find the significant attributes that contribute to the classification of the target genetic disorder.

Also we were to create a dashboard including all those Key performance indicators.

## Initial View

## SAMPLE DATA

<b>Patient Id</b>	PID0x81d5	<b>Test 1</b>	0	<b>Heart Rate</b> (rates/min)	Normal
<b>Patient Age</b>	7	<b>Test 2</b>	0	<b>Respiratory Rate</b> (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 Name</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 SOUTH END, MA 02118 (42.335925371008436, -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.

# Non Graphical

Initially we used the summary statistics, count, minimum, maximum, average, standard deviation etc to describe the numeric columns.

Also correlation between every pair of numeric attributes were computed to find highly correlated attributes.

# Graphical

We used R programming packages to create visualization related to the distribution of an attribute on its own and also to demonstrate the relation between attributes. Here we mainly used ggplot package to create stacked bar graph and stem and leaf plots.

Also we used a heat-map plot to better visualize the correlation matrix created.



## Domain Based

Based on the further continuation of the project as a classification of genetic disorder based on genetic inheritance the attributes were considered for their significance to the project and promptly discarded.

# Dashboard

We used `https://visual.is` as a proper dashboard creation tool. This let us interactively choose different plots and data that we want to display in our dashboard.

## Initial Pre-processing

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. The dataset now contains 45 columns and 6706 observations.

```
#importing package  
library(readxl)  
library(tidyverse)
```

```
#read the excel file  
df <- read_excel("train.xlsx")
```

```
#rename the columns eliminating obstacle character  
names(df) <- gsub(" ", "_", names(df))  
names(df) <- gsub("\\(", "_", names(df))  
names(df) <- gsub("\\)", "_", names(df))  
names(df) <- gsub("\\/", "_", names(df))  
names(df) <- gsub("'", "_", names(df))
```

```
#rows with any na value are dropped  
df <- df %>% drop_na()
```

```
#save as csv  
write.csv(df, "train.csv", row.names = FALSE)
```

Figure: Preprocessing Code.

## Non Graphical EDA

Summary statistics showed that the columns of Test 1, Test 2,, Test 3, Test 4, Test 5 have no variance in them hence they are removed.

	Test_1	Test_2	Test_3	Test_4	Test_5
Min.	0	0	0	1	0
1st Qu.	0	0	0	1	0
Median	0	0	0	1	0
Mean	0	0	0	1	0
3rd Qu.	0	0	0	1	0
Max.	0	0	0	1	0

Figure: Summary Statistics

# Non Graphical EDA

Also a correlation heatmap showed that there are no highly correlated attribute

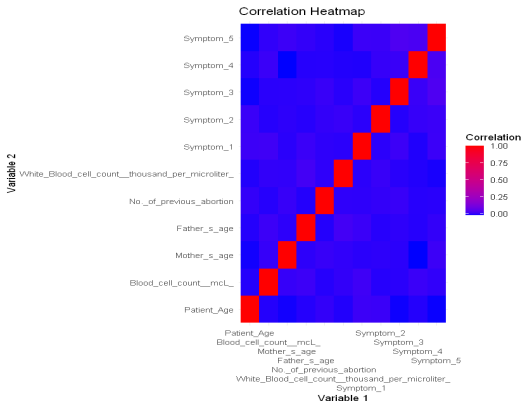


Figure: Correlation Heatmap

## Graphical EDA

We find that the Genetic disorder attribute is dependent on the Disorder Subclass attribute. Hence The Genetic Disorder attribute is removed.

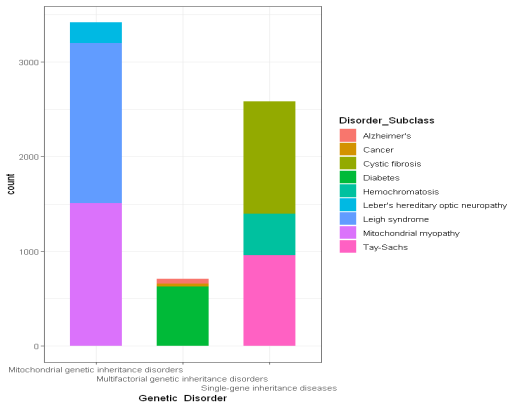


Figure: Dependency between Target Attributes

## Graphical EDA

For categorical attributes, a stacked bar graph is created with Disorder subclass attribute on x-axis. By this we found that most of such attributes divided the dataset almost symmetrically along the Disorder Subclass and hence is removed.

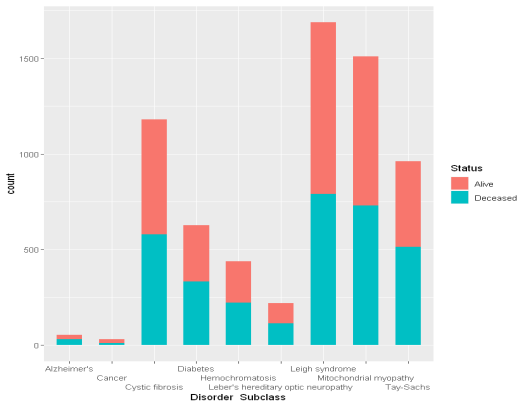


Figure: Disorder Subclass and Status

## Graphical EDA

For continuous attributes, a box and whisker plot is created with Disorder subclass attribute on x-axis. Here the data were not as symmetrical.

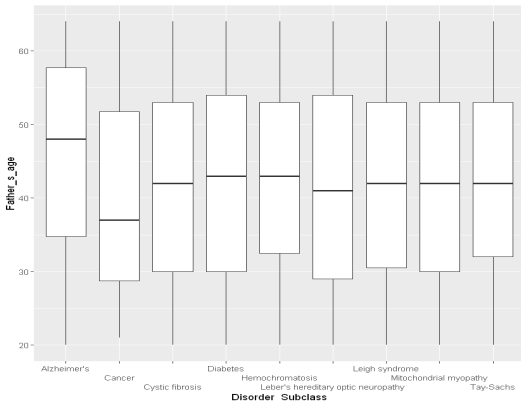


Figure: Disorder Subclass and Patient's Father's age



## Graphical EDA + Domain Information

Given the domain of study to be classification of disorder due to genetic inheritance, the following attribute were considered significant and have shown reliable variance.

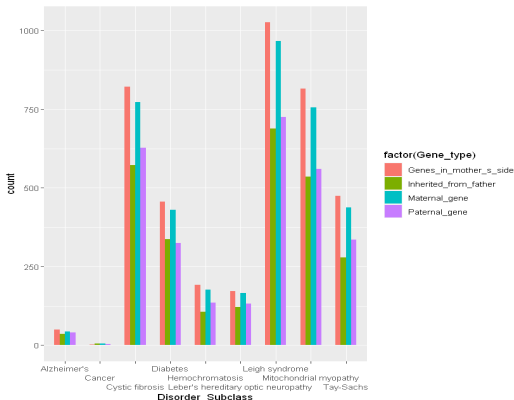


Figure: Disorder Subclass and Genetic Presence (YES)

# Graphical EDA + Domain Information

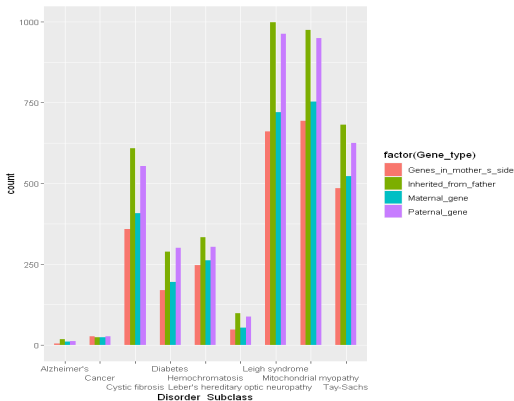


Figure: Disorder Subclass and Genetic Presence (NO)

# Graphical EDA + Domain Information

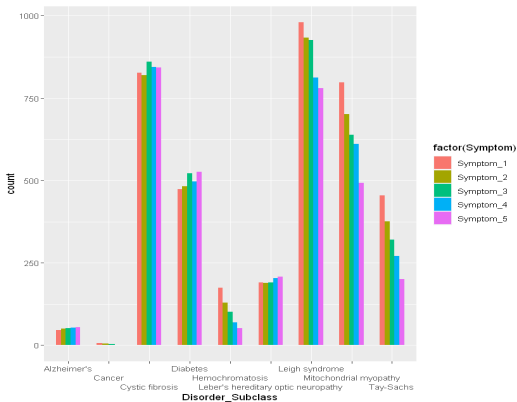


Figure: Disorder Subclass and Symptom Detected (True)

# Graphical EDA + Domain Information

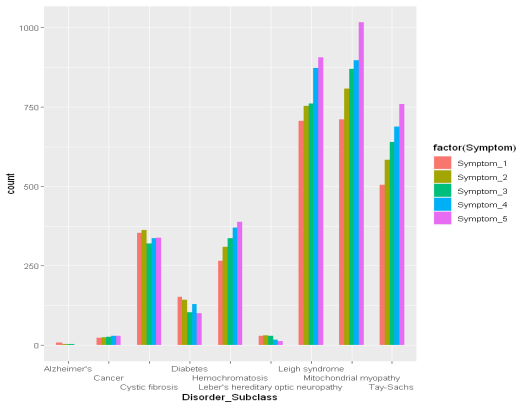


Figure: Disorder Subclass and Symptom Detected (False)

## Domain Based

Attributes like Patients Name, birth address, History of alcohol which doesn't have much significance to the domain of genetic inheritance were also removed.

All this left us with 16 columns with one of them being the target column.

# Dashboard

We have a basic dashboard framework ready that gives us basic visualization of the data.

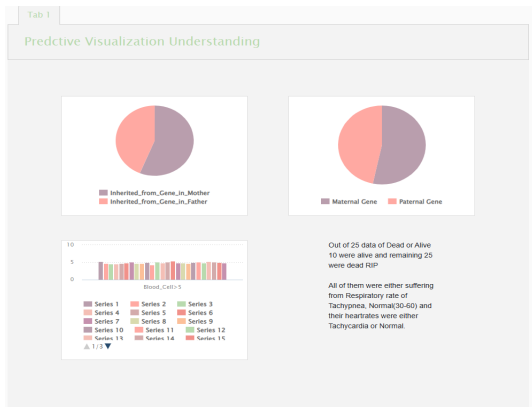


Figure: Dashboard Sample

## Conclusion

In this EDA project, we applied various statistical analysis, graphical visualization, and domain knowledge to reduced the dataset from 45 to 16 columns.

Patient Age	Genes in mother's side
Inherited from father	Maternal gene
Paternal gene	Blood cell count mcL
Mother's age	Father's age
Number of previous abortion	White Blood cell count
Symptom 1	Symptom 2
Symptom 3	Symptom 4
Symptom 5	Disorder Subclass

## Discussion and Limitations

The dashboard created is entirely dependent on the capability of the tool used.

Given that the data is from a Kaggle dataset, there is a possibility of it being synthetic, as such there are certain trends and outliers that isn't explainable from a real world perspective.

The methodology discussed here is still useful and backed by proper citable sources. Hence there are still proper real world application of this project and or its workflow.



## Demonstration

Now we would like to  
demonstrate the Project

# THANK YOU!

Please share if you have any queries.

We would like to show the code now.