# EAST WEST UNIVERSITY

## CSE303: Statistics for Data Science
## [Spring 2023]

# Project Report

**Course Code** : CSE303
**Course Title** : Statistics for Data Science
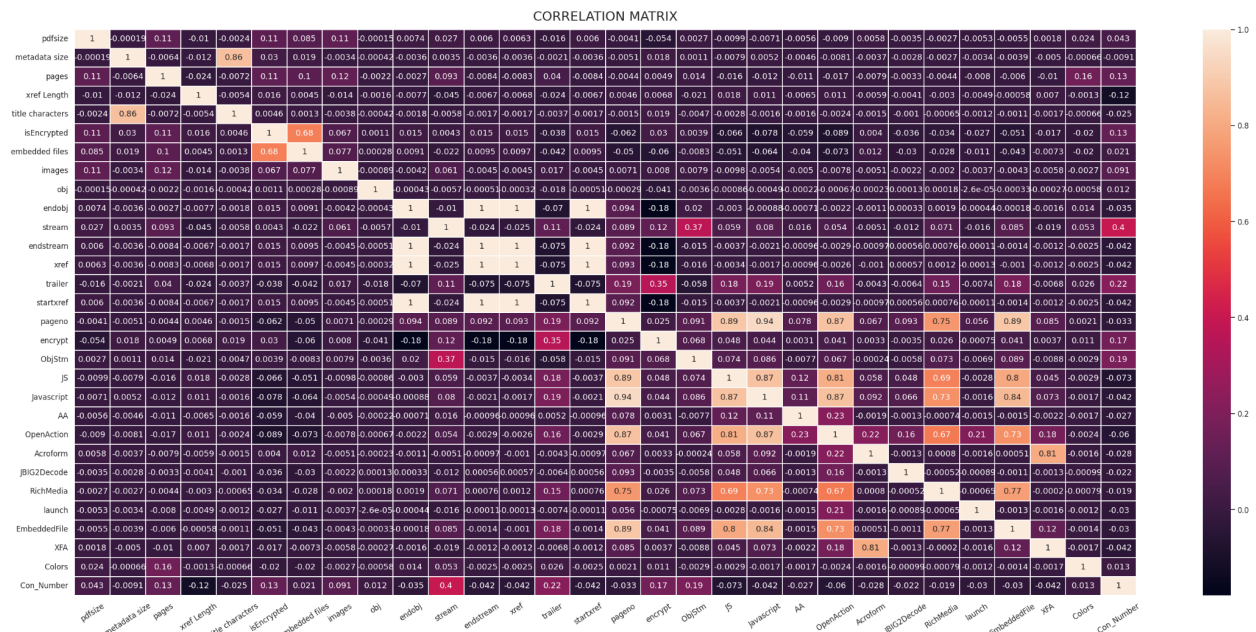**Section** : 02
**Group Number** : 09

## Submitted by:

| Student ID | Student Name | Contribution Percentage |
|---|---|---|
| 2020-1-60-184 | Saiful Islam | 90% |
| 2020-1-60-186 | Sajib Khan | 90% |
| 2020-1-60-273 | Eva Islam | 85% |
| 2020-1-60-211 | Zarin Tasnim Nuzhat | 85% |

# 1. Introduction

The Canadian Institute for Cybersecurity's PDF Malwear Dataset, which was acquired, contains details on PDFs that were attacked maliciously in 2022.The dataset can be used to investigate the variables that affected the prevalence of malicious acts among PDF files and to create prediction models that can calculate the likelihood of malicious acts or malicious assaults based on PDF properties.The 10023 rows and 33 columns of the dataset are provided in CSV format. The columns include pdf size, metadata size, pages, xref Length, title characters, isEncrypted, embedded files, images, obj, endobj, stream, endstream, xref, trailer,  startxref, pageno, encrypt, objstm, js, javascript, AA, OpenAction, Acroform, jBig2Decode, RichMedia, launch, EmbeddedFile, XFA, colors, con_number and whether or not the PDF is under attack by any malware.To manage missing values and outliers, some cleaning of the dataset was necessary.The dataset only contains data on a portion of the PDF files that were attacked by Malware, therefore any analysis or modeling should account for the possibility of selection bias. Additionally, there weren't many missing values or duplicate values in the dataset.

# 2. Exploratory Data Analysis



We have create a correlation data table where,
Positive correlation range: <0.4
Weak correlation range: >0.0999

No correlation range: <=0.0999
Negative correlation range: >0

| Columns | Positive correlation | Weak positive correlation | No correlation | Negetive correlation |
|---|---|---|---|---|
| pdfsize | 1 | 3 | 11 | 15 |
| metadata size | 2 | 0 | 6 | 22 |
| pages | 1 | 6 | 4 | 19 |
| xref Length | 1 | 0 | 8 | 21 |
| title characters | 2 | 0 | 3 | 25 |
| isEncrypted | 2 | 3 | 13 | 12 |
| embedded files | 2 | 1 | 12 | 15 |
| images | 1 | 2 | 8 | 19 |
| obj | 1 | 0 | 5 | 24 |
| endobj | 4 | 0 | 8 | 18 |
| stream | 2 | 4 | 13 | 12 |
| endstream | 4 | 0 | 6 | 20 |
| xref | 4 | 0 | 6 | 20 |
| trailer | 1 | 9 | 4 | 16 |
| startxref | 4 | 0 | 6 | 20 |
| pageno | 6 | 1 | 15 | 8 |
| encrypt | 1 | 14 | 6 | 9 |

| objstm | 1 | 2 | 13 | 14 |
|---|---|---|---|---|
| js | 6 | 2 | 7 | 15 |
| javascript | 6 | 1 | 8 | 15 |
| AA | 1 | 3 | 4 | 22 |
| OpenAction | 6 | 6 | 4 | 14 |
| Acroform | 2 | 1 | 9 | 18 |
| jBig2Decode | 1 | 1 | 8 | 20 |
| RichMedia | 6 | 1 | 9 | 14 |
| launch | 1 | 1 | 1 | 27 |
| EmbeddedFile | 6 | 2 | 4 | 18 |
| XFA | 2 | 2 | 6 | 20 |
| colors | 1 | 1 | 7 | 21 |
| con_number | 1 | 6 | 5 | 18 |

## 3. Data Preprocessing

- **Data collection:** We have collected mour dataset for an online source. Our online source is Canadian Institute for Cybersecurity. From this site we have worked on the CIC-Evasive-PDFMal2022 dataset.
- **Data cleaning:** To process our dataset we have to find out the duplicate and missing data. In our dataset we have a total number of 10026 columns and 33 rows.In the dataset we did not find any duplicate data and we did have only 3 missing data; we have dropped the missing values from the data set. Then we have 10023 columns and 33 rows.
- **Data reduction:** To reduce the noise from diffirent features we replaced the appropriate data with respect to the features.In total we had 17 columns where we found noise.
- **Data storage:** We store the processed data in a spreadsheet and for our coding we have the dataset as a csv file that is easily accessible for analysis.
- **Data analysis:** We have  used machine learning algorithms to identify patterns and visualization tools to identify the most common types of malware.To predict the malwear we made a model using logistic regression and linear regression.

# 4. Machine Learning Models

**Logistic Regression:** A statistical model called logistic regression is used to forecast a binary result (0 or 1) from one or more predictor variables.The link between the predictor variables and the outcome variable is modeled using a logistic function.Modeling the link between a dependent variable and one or more independent variables is done using this form of regression analysis.The logistic regression model converts a linear combination of the input data into a probability score between 0 and 1 using a logistic function, often known as a sigmoid function. The logistic function's formula is as follows:

$g(z) = 1 / (1 + e^{\wedge}(-z))$

The objective is to minimize the difference between the true binary outcomes and the expected probabilities, which is measured by the function.

**Support Vector Machine(SVM):** Support vector machine, or SVM for short, is a well-liked machine learning technique that can be applied to both classification and regression applications. The objective of SVM in classification jobs is to identify a hyperplane that can most effectively classify the data points.The margin—the separation between the hyperplane and the nearest data points from each class—is maximized by selecting the hyperplane. Support vectors are utilized to define the hyperplane; these are the points that are closest to the hyperplane.By utilizing various kernel functions, SVM can handle both linearly and non-linearly separable data. Better class separation is possible because of the kernel function's mapping of the input features into a higher-dimensional space where the data is linearly separable.

In our project we have used logistic regression and linear regression machine learning models.

- **Logistic regression:** With the help of a confusing matrix and performance measure we predicted the dataset and found out accuracy, f-score,precision and recall.
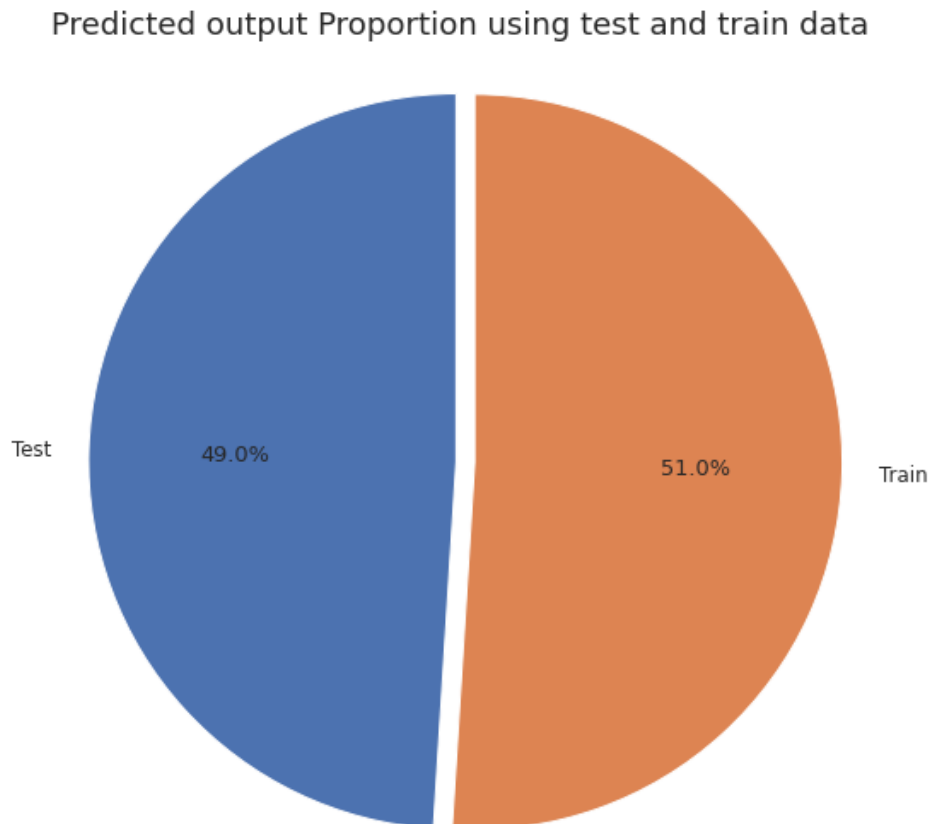From this we can tell how perfectly the model will result in future data.

- **Linear regression:** With the help of regression we can predict the data with MAE,MSE and R2 score.

In the algorithm we have used grid search to get the best parameter from our dataset.
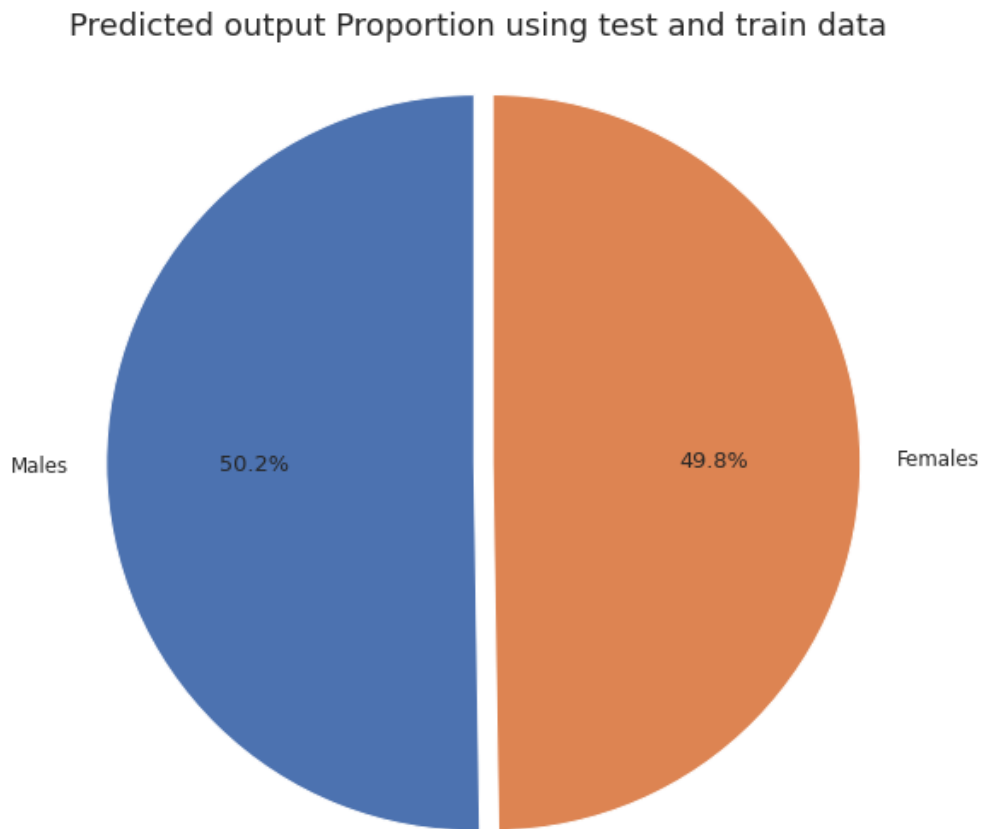
# 5. Performance Evaluation and Discussion

We calculated accuracy, precision, F1-score, and recall to assess the effectiveness of our model using logistic regression.90% is the score and 90% is the accuracy of the logistic regression model using default parameters. As a result, the model is too well fitted.

Predicted output Proportion using test and train data



In the following pie chart we have found the difference between predicted output of test and train data.

To cope up with the overfit problem we have used 5 fold cross validation.

### Predicted output Proportion using test and train data



We used five fold cross validation to address this issue. After cross-validation, the mean score is 91%, while the minimum score is 90%. The outcome didn't significantly improve.90% of the model's predictions for these test and training data are favorable.