# CYBER CRIMES IN SAUDI ARABIA

Group Members :

- Raghad Almeataz    - Ahmed Altowairqi    - Mohammed Albesher

- Abdullah Albutih    - Abdulrahman Aljubaylan

Business Understanding

- Cyber crime cases are drastically increasing in Saudi Arabia. At present, Saudi Arabia is one of the top 10 countries in significant cyber attacks.

- In addition, according to specialists in information security,  the financial loss estimated  by Saudi companies as a result of piracy and electronic sabotage ranges from SR 300,000 to more than SR1 million for each case, with bank losses, estimated at more than a billion US dollars.

# Overview

- The 2030 VISION OF SAUDI ARABIA seeks to Increase the Country Income that is not based on Oil that include Information Technology Sector where decreasing the Internet Security Threats to save the lost Money and to use the lost money in other beneficial areas such as investments and Organizations Development and other.

Banks Financial Losses

4.5B

Companies Financial Losses per Case

300K - 1.2M

Number of Victims in 2016

3.6M

All Financial Losses are in SAR Currency - SAUDI RIYALS

# Project Objective

- measure the awareness of Saudi people about the danger of cyber crimes.

- find out the major factors for becoming a potential victim.

- give Solutions for the Main faced issue .

- measure the Risk factors and decrease it .

**The Target Audience**

– Cyber security companies

– Saudi Federation for Cybersecurity, Programming and Drones Institute (SAFCSP)

– Businesses (organizations) that want to measure the awareness of their employees

– Any company that is working remotely, needs to check security applications of their own

–  employees to ensure consistent work

Data Collection & Preparation

# Dataset Overview

- It is a survey conducted using Google forms on Saudis whose age is above than 18. It consists of 64 questions, which indicates the number of columns and 1230 number of responses which indicates the number of rows.

SURVEY

64 Questions

1230 Responses

Rows : 1230    |    Columns : 64

**The applied Approaches :**

- Filling with zeros
- Filling with Values
- Dropping The Null Values
- Filing with Forward Filling
- Filling with Backward Filling

After conducting the different approaches, we continue with the best fit which is the **Dropping The Null values.**
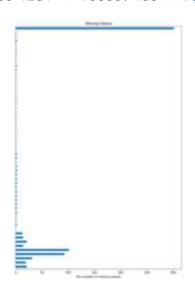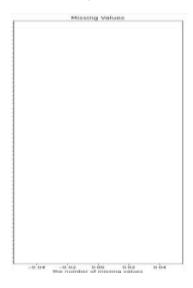
# Data Preprocessing

we dropped the columns that with 60% missing values and more , because they do not have enough data

The Original Dataset                60*1231 =>  73860 / 100 => 739                Dropping Null Values
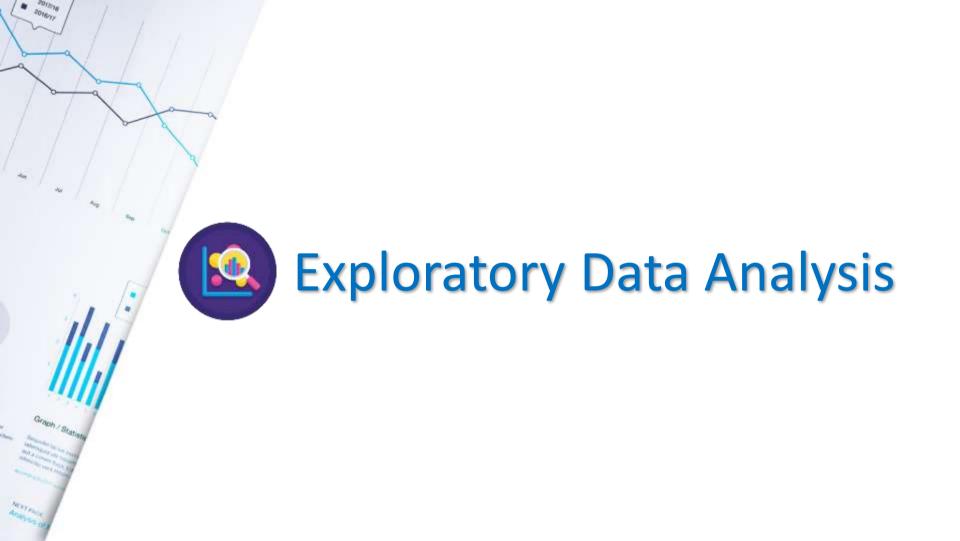
# Data Preprocessing

Trim white spaces, imputing empty values and extracting new columns out of existing ones.

```
 Economic                             1
Information and Libraries             1
Mechanic- Prpducing Department        1
Actuarial Science                     1
 Biology                              1
Name: major, Length: 74, dtype: int64
```

```
Undergraduate (Diploma, BSc)        744
Postgraduate (Master's, PhD)        167
High School                         110
Middle School                         3
                                      1
Name: education, dtype: int64
```

| | updated_about_CC(_Online Sources) | updated_about_CC(_Online Sources)_2 |
|---|---|---|
| 77 | I do not feel that I keep myself updated | No |
| 172 | Internet, website, email bulletins, blogs, etc. | Yes |
| 188 | TV, news, radio, Government websites (e.g. CER… | Yes |
| 193 | Internet service provider ISPs | No |

Exploratory Data Analysis

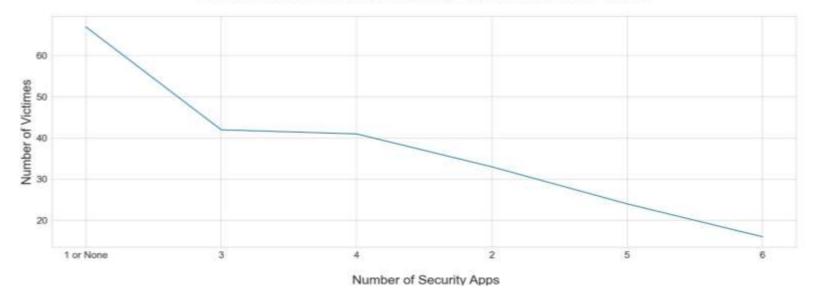❖ **Insight** :

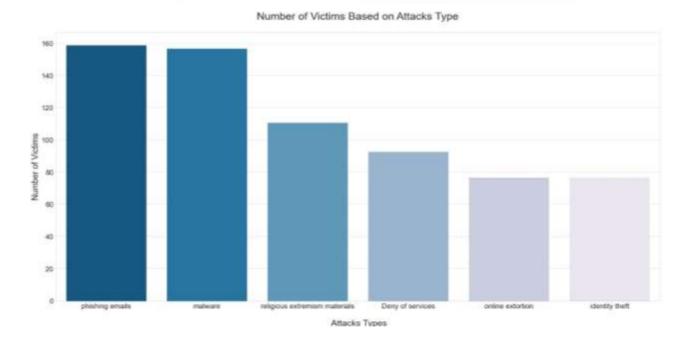• The less security apps the more potential of becoming victim.

The Relationship Between Security Apps & Number of Victims

❖ **Insight** :

• Which type of attacks has the most effect on becoming a victim?

Number of Victims Based on Attacks Type

❖ **Insight** :

- Which purpose has the most effect on becoming a victim?

Number of Victims Based on Purpose of Using the Internet

❖ **Insight** :

- The more good practices you follow the less potential of becoming a victim

**Good Practices:**

1. Careful clicking on links
2. Install software updates
3. Change pass frequently
4. Terms conditions
5. Privacy settings
6. Check legitimacy

The Percentage of Victims Answers about Following The Good Practises

Always or Often

71.3%

14.9%

Seldom or Never

13.8%

Sometimes

❖ **Insight** :

- People always receive threats/attacks

The Percentage of Victims Based on Their Answers
for Receiving Threats/Attacks

❖ **Insight** :

- The more security level the less potential of becoming victim



Number of Victims per Security Level

❖ **Insight** :

- Most people agree with avoiding disclose personal information.

The Percentage of People Answers for 'Avoid Disclosing Personal Info' Question



The Percentage of People who agree with avoid disclosing personal info Based on Their Answers for 'Pass Info' Question

❖ **Insight** :

• What are the top channels that people use to raise their awareness on cyber crimes?

The percentage of channels that people use to raise their awareness about cyber crimes

❖ **Insight** :

• Which party is more responsible for raising the awareness of cyber crimes?

The Responsibility of Each Party in Raising The Awareness

Modeling

| The Dataset Name | Cyber crimes |
|---|---|
| The Target Label | Have you been a victim of a cyber crime ? |
| The Machine Learning Type | Supervised Learning |
| The Modeling Category Type | Classification |
| The Programming Language | Python |
| The Main Library | Scikit-learn |

**The applied Mechanisms :**

- Mapping

- Label encoding

- One Hot encoding

| eservice_usage | device_skill_level | freq_used | connection_type | used_purpose | ... | tele_responsible | user_responsible | education_responsible | government_role | victim |
|---|---|---|---|---|---|---|---|---|---|---|
| Once a day | Intermediate | 1 | Cellular network | 1 | ... | Agree | Agree | Agree | public awareness | No |
| Frequently | Intermediate | 2 | Private Wi-Fi | 5 | ... | Strongly Agree | Strongly Agree | Strongly Agree | global cyber security | No |
| Frequently | Beginner | 2 | Private Wi-Fi | 5 | ... | Strongly Agree | Agree | Strongly Disagree | public awareness | No |

| eservice_usage | device_skill_level | secure_level | check_legitimacy | pass_my_info | clicking_on_banners | privacy_settings | SM_services_protect_info | terms_conditions | change_pass_frequently | |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2 | 1 | 4 | 4 | 2 | 2 | 4 | 4 | 4 | |
| 1 | 2 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | |
| 0 | 1 | 3 | 0 | 1 | 0 | 0 | 3 | 2 | 2 | |

After conducting the different mechanisms and try them all, we continue with
the best fit which is mix between **the Label encoding and one hot encoding.**

**The Challenges :**

- **The label Imbalance**.

- **The Dataset Complexity**.

- **beating the baseline model**.

**The Solution :**

After Conducting many methods to fix the faced issues, we found out that the best solution is to remove around 200 records from the dataset because going with this method would not just resolve the missing values issue but also would resolve the bias in the dataset and makes it more suitable and equivalent. Beating the baseline model in the start was a rough process but after removing the specified records and decreasing the bias in the dataset now the data classes gap minimized which makes more sense.

**The Models Results :**

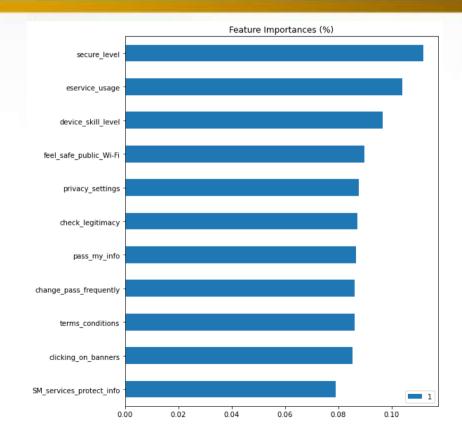| The Baseline | 0.7707317073170732 |
|---|---|
| The Logistic Regression | 0.824390243902439 |
| The Decision Tree | 0.824390243902439 |
| **The Random Forest** | **0.8146341463414634** |
| The SVC | 0.824390243902439 |

**Why did we chose this specific Model ?**

The Random Forest model characterized by Low risk of overfitting and runs efficiently on a large dataset.

We chose the Random Forest model because it's the best model to represent Great accuracy with Strong reliable almost parallel features that all the entered features contributing to the prediction.



Feature Importances (%)

# Tableau Dashboard

The Dashboard

We proposed a cybercrime Dataset that:

- ⚠ Filled with Messy Data

- ⚠ Filled with biased Data

- ⚠ Lack of Specification

After we finished the project , we could :

- ✅ Measure The awareness of people

- ✅ Get The Main factors that Couse cybercrimes

- ✅ Determine if a person is going to be a cybercrime victim or not

**In the near future :**

- Add the project into Different blogs.

- We Share The Results with the Local Community

- We seek to continue to develop the project.

- Publish The work and add more Data.

♡ Acknowledgement

**Mr. Mikio Harman** for his time, effort and dedication to helping us. We learned a lot from him and enjoyed working with him.

**Saudi Digital Academy** for supporting us and provide us with this wonderful Data Science Bootcamp that gave us the Strong Background we need so we can take our Careers to the Next level.

**Coding dojo** for supporting us and provide us with all the Resources we need to Increase our knowledge and work practically with different creative new projects.

# THANK YOU

**Any Questions ?**