



Big Data Course Project

SHAKIBA SADAT
MIRBAGHERI

JULY 2023

Introduction

The purpose of this project is to analyze Bitcoin Heist Ransomware Address Dataset from UCI using Spark library.

- Setting up the VM
- Creating spark session
- Loading the dataset
- Preprocessing the dataset
- Creating RFM features
- Creating ML algorithms
- Evaluation of the models

Setting Up The VM

1. Installing Oracle and Vagrant
2. Creating a vagrant file and setting up the Ubuntu/jammy64
3. Creating and starting the VM
4. Connecting to the VM with vagrant user
5. Setting up the ports
6. Creating the link to access Jupyter notebook

```
PS C:\Users\Lenovo\desktop\test> vagrant init ubuntu/jammy64
```

```
PS C:\Users\Lenovo\desktop\test> vagrant up --provision
```

```
PS C:\Users\Lenovo\desktop\test> vagrant ssh -- -L 8888:localhost:8888
```

Creating spark session

Entry point to PySpark for loading and creating the dataset

```
[1] !pip install pyspark
!pip install findspark Python

... Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting pyspark
  Downloading pyspark-3.3.1.tar.gz (281.4 MB)
    281.4/281.4 MB 4.2 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Collecting py4j==0.10.9.5
  Downloading py4j-0.10.9.5-py2.py3-none-any.whl (199 kB)
    199.7/199.7 KB 15.8 MB/s eta 0:00:00
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.3.1-py2.py3-none-any.whl size=281845512 sha256=a346fedf564ab862fadbd80c96a9ad0191128291ee3d22cb8bea93cec0dfef74
  Stored in directory: /root/.cache/pip/wheels/43/dc/11/ec201cd671da62fa9c5cc77078235e40722170ceba231d7598
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9.5 pyspark-3.3.1
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting findspark
  Downloading findspark-2.0.1-py2.py3-none-any.whl (4.4 kB)
Installing collected packages: findspark
Successfully installed findspark-2.0.1
```

```
[ ] os.environ["JAVA_HOME"] = "C:\Program Files\Java\jdk-19"
os.environ["SPARK_HOME"] = 'C:\\Users\\Nima_S_H\\Desktop\\BDAVM\\spark'
findspark.init() Python
```

A large, irregular pink brushstroke shape with a textured, hand-painted appearance, serving as a background for the section header.

Data Understanding and Preprocessing

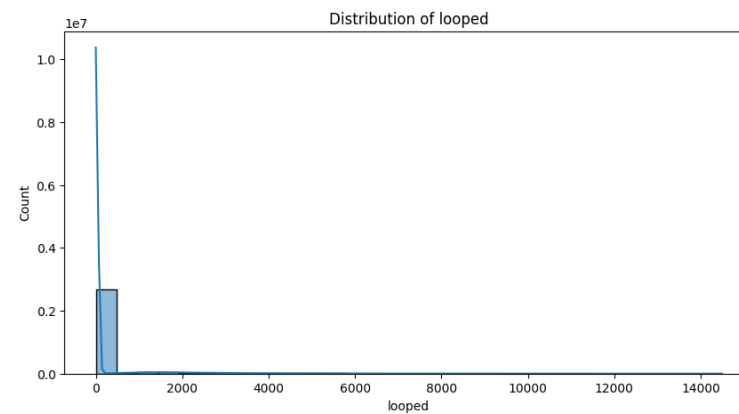
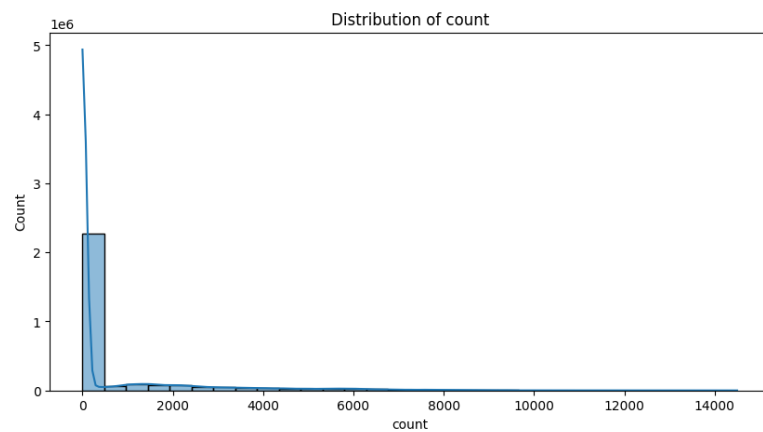
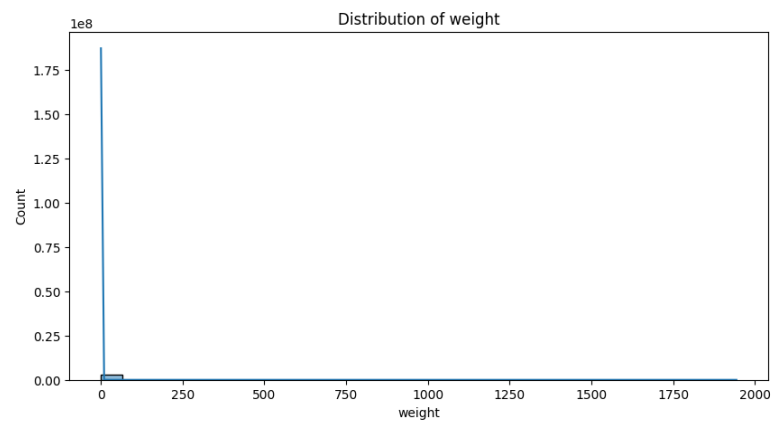
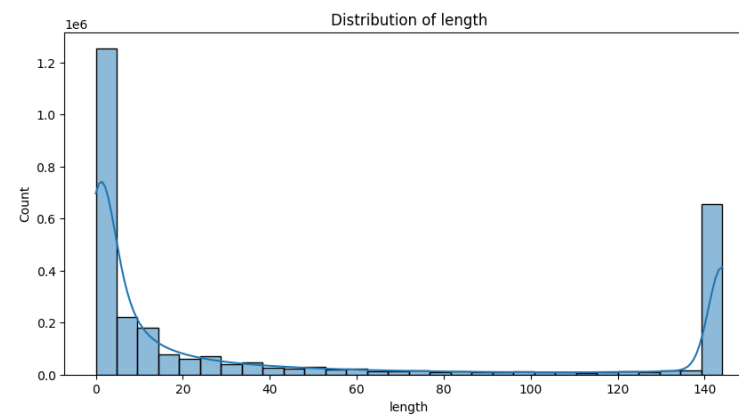
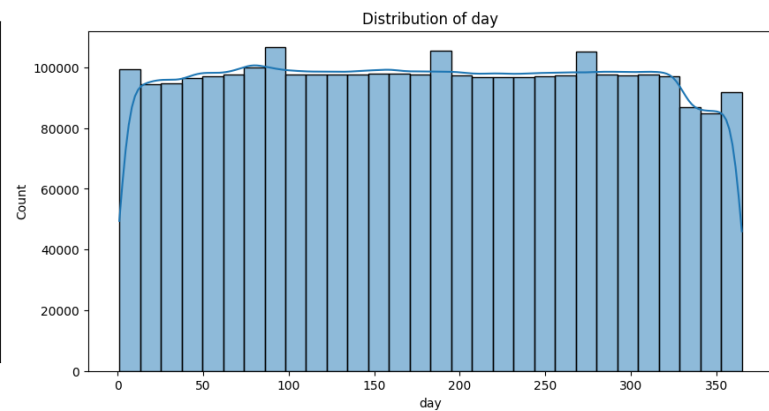
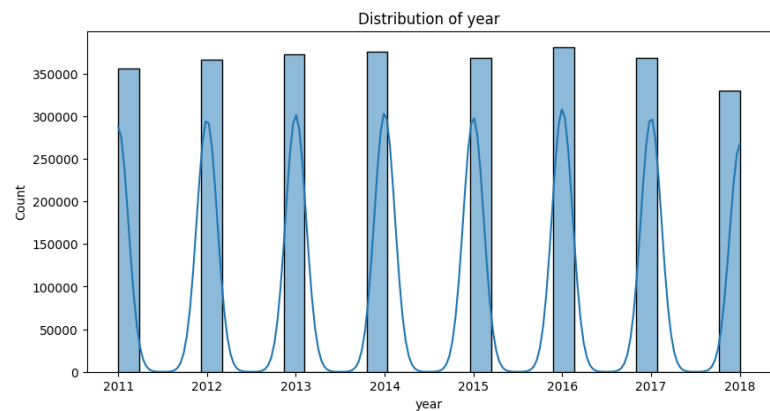
The dataset is quite substantial, containing 2,916,697 rows and 10 columns. Upon conducting a cleanliness check for any null or 'NaN' values, I found dataset to be in excellent shape with no missing values, which streamlined the preprocessing stage.

A large, irregular pink brushstroke shape on the left side of the slide, serving as a background for the title.

Attributes

- **Address:** A string representing a Bitcoin address.
- **Year:** An integer indicating the year of the observation.
- **Day:** An integer representing the day of the year (1 to 365).
- **Length:** An integer value indicating a measurement (specific context not provided).
- **Weight:** A floating-point value (float) associated with the data point (specific context not provided).
- **Count:** An integer representing a quantity (specific context not provided).
- **Looped:** An integer representing a looped value (specific context not provided).
- **Neighbors:** An integer indicating the number of linked addresses.
- **Income:** An integer representing the amount of Satoshi received (smallest unit of Bitcoin).
- **Label:** A categorical string indicating the ransomware family name or "white" if not known to be ransomware.

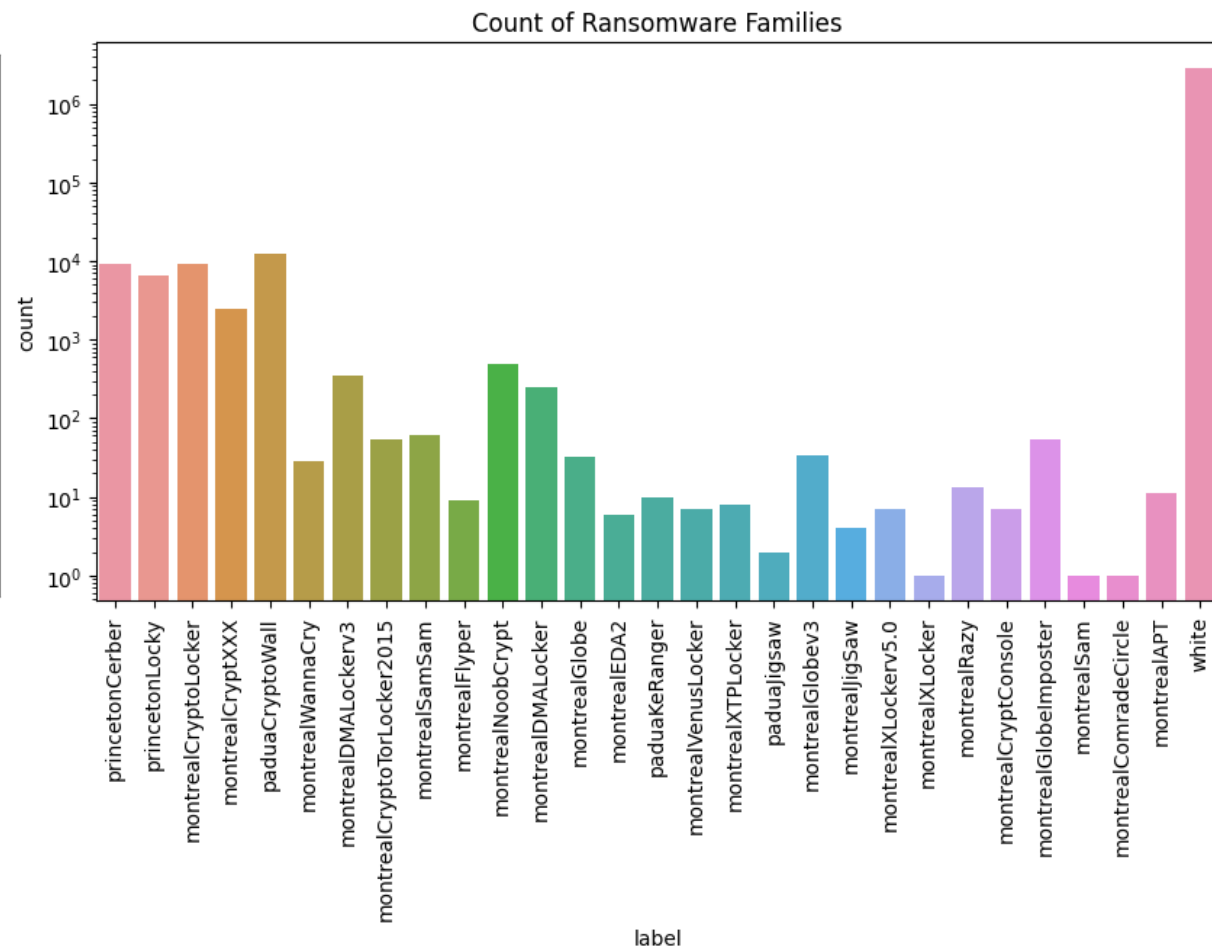
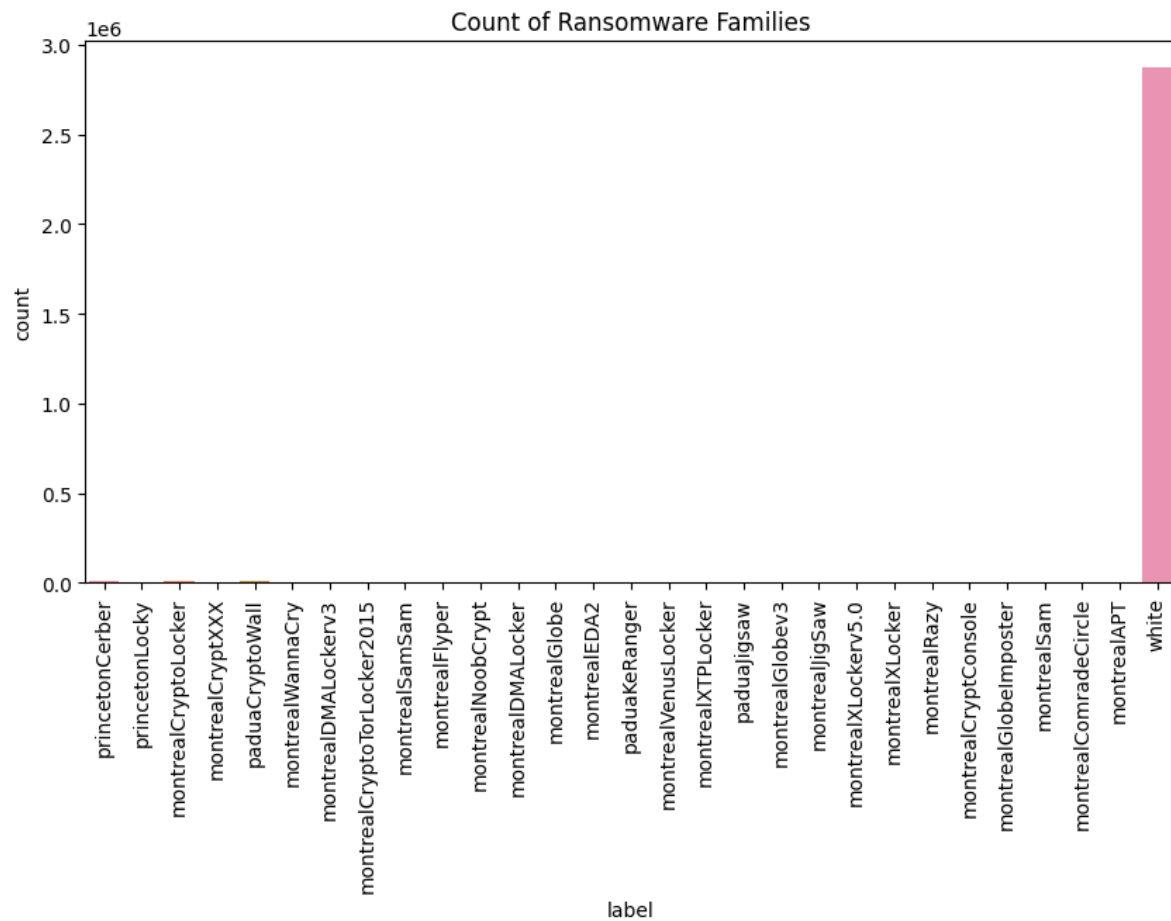
Chart



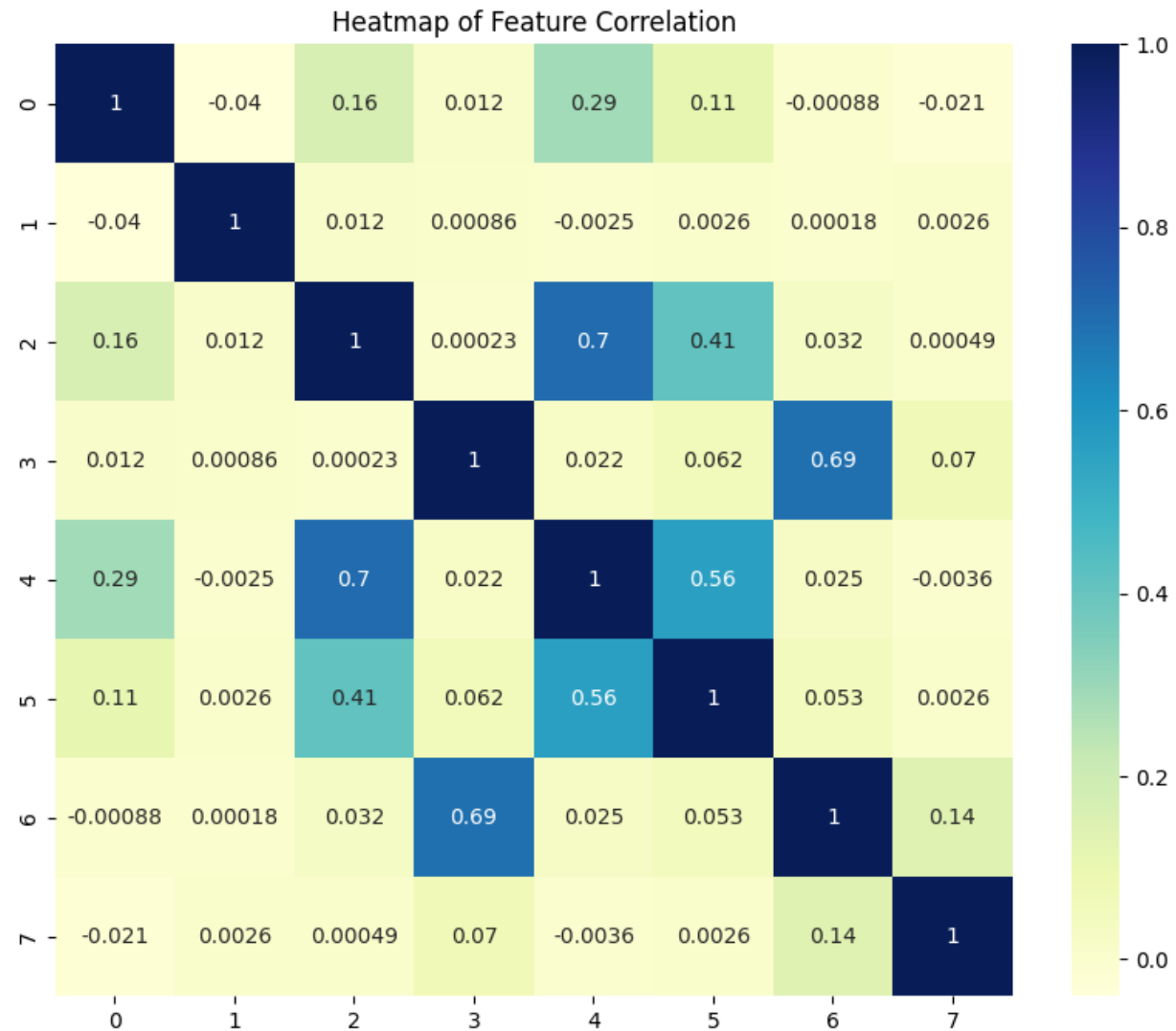
top 10 labels

Label	count
white	2875284
paduaCryptoWall	12390
montrealCryptoLocker	9315
princetonCerber	9223
princetonLocky	6625
montrealCryptXXX	2419
montrealNoobCrypt	483
montrealDMALockerv3	354
montrealDMALocker	251
montrealSamSam	62

Chart



Display Data



A large, irregular pink brushstroke shape that serves as a background for the title text. It has a soft, painterly texture with varying shades of pink.

Transition to Unsupervised Learning

K-means

Feature Scaling

- In this phase, I prepare my dataset for unsupervised machine learning, specifically, K-means clustering. For optimal performance of my clustering algorithm, I need to standardize our features to ensure they're on the same scale.

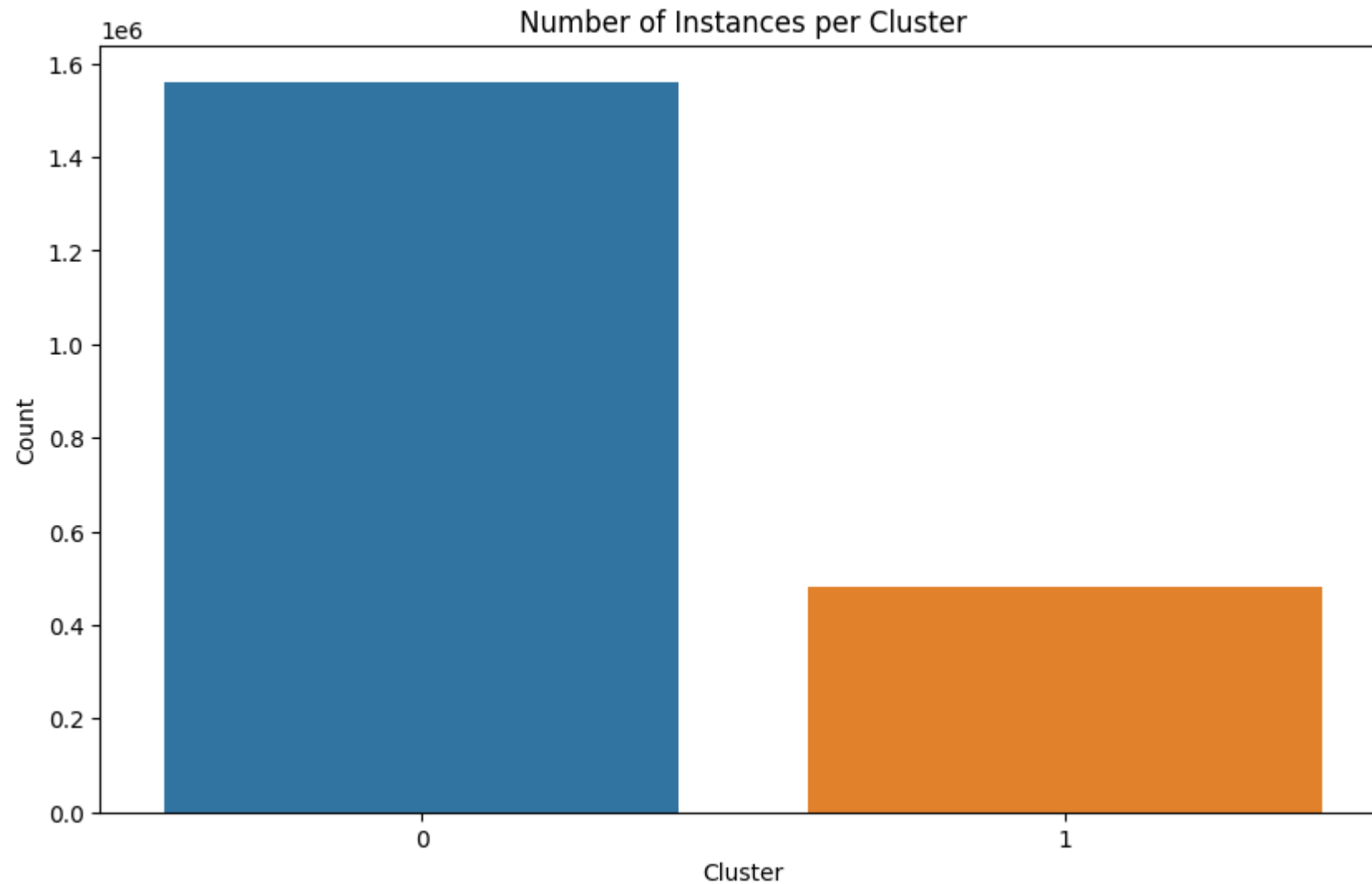
Assembling for Clustering

- To get my data ready for machine learning, I use the **VectorAssembler**, a transformer that combines a given list of columns into a single vector column. In this case, it's transforming my features into one single vector column, which I have labeled as 'features'.

Majority Label Identification

- Since K-means is an unsupervised learning algorithm, I need a way to evaluate its performance. Here, I group our predictions by the 'prediction' and 'label' columns and count the number of instances.
- Then join this DataFrame with my original DataFrame that contains my K-means predictions and labels. This allows me to compare the original label to the majority label for each cluster.

K-means

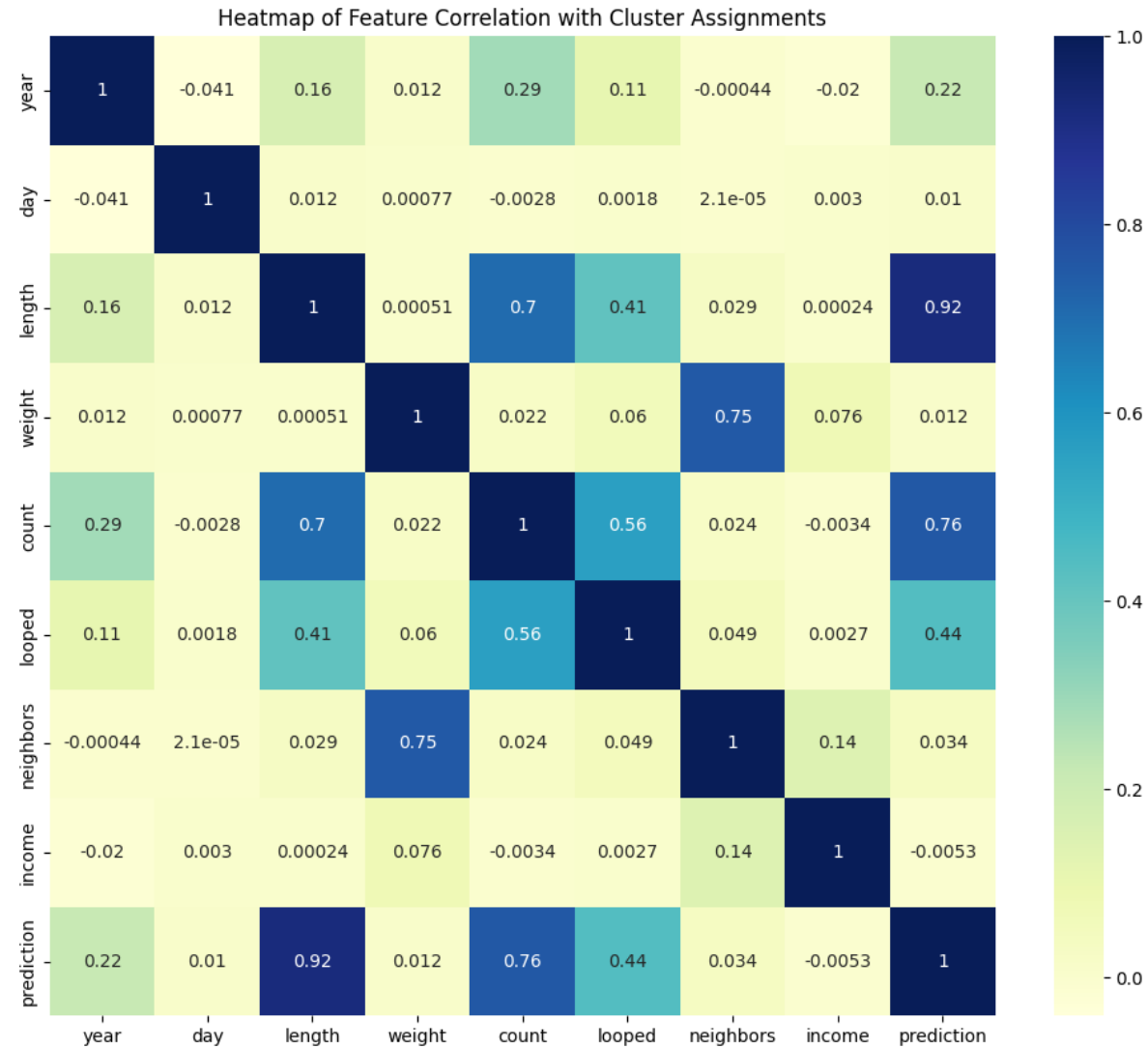


Accuracy Evaluation

84.80% for training

85% for test

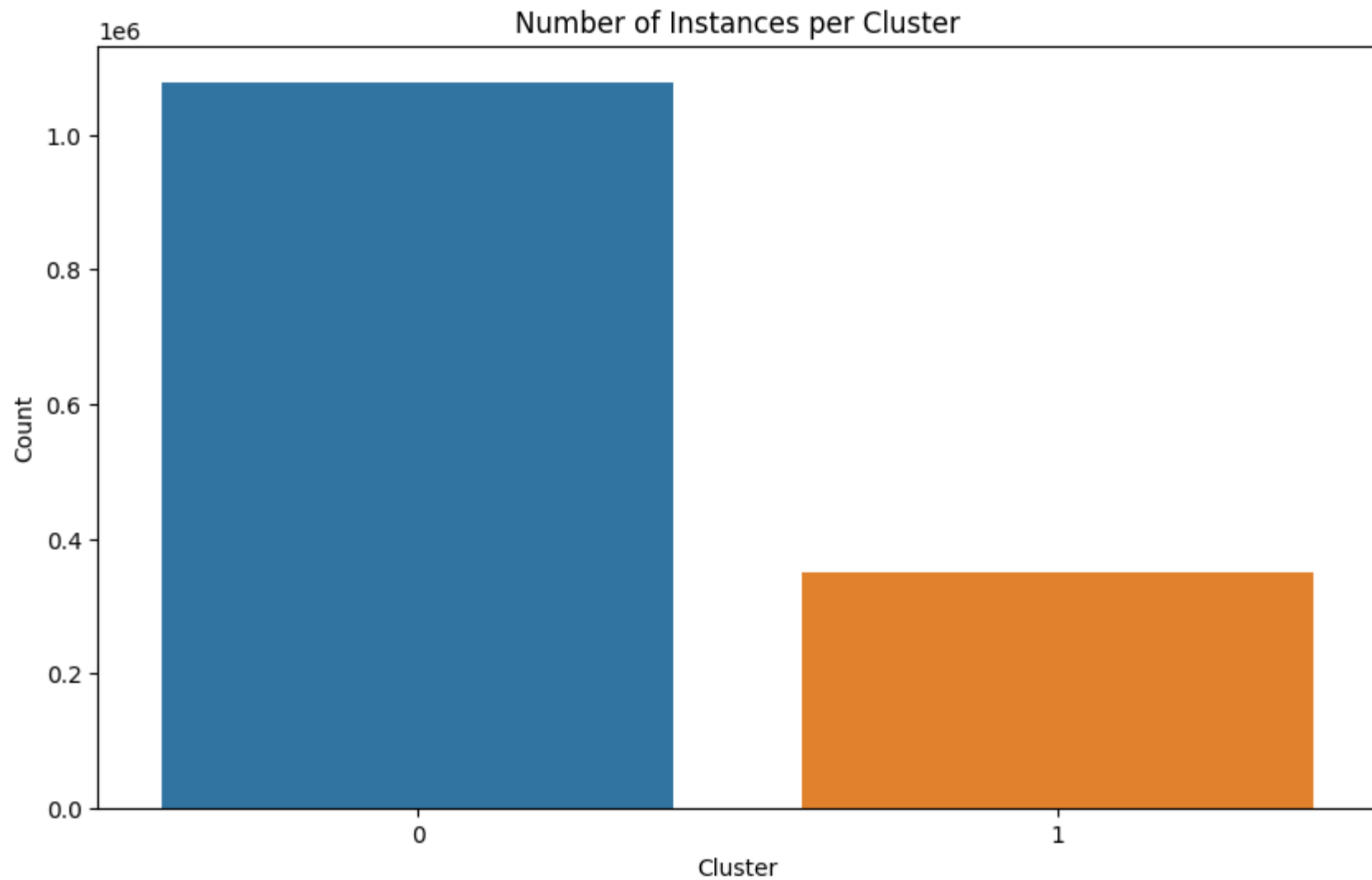
Analyzing Feature Correlation with Cluster Assignments



GMM

84.91% for training

15% for test



A large, irregular pink brushstroke shape on the left side of the slide, serving as a background for the 'Concolusion' text.

Concolusion

The K-means clustering method gave balanced accuracy results between the training and test sets, which indicates a good generalization of the model. On the other hand, the GMM, while it performed slightly better on the training set, had a significantly lower accuracy on the test set, suggesting a potential overfitting issue.

This comparison gives us valuable insights into the strengths and weaknesses of each method, and how they perform with this specific dataset. For future work, it might be beneficial to investigate the reasons behind the low test accuracy of the GMM and adjust the model's parameters or the dataset accordingly.

Thank you