

# Internet Downtime Prediction Analysis using Machine Learning

DATA 245 - Machine Learning Technologies

## Group 6

Poojan Gagrani

Bhavik Patel

Kashish Thakur

Yuti Khamker



# Introduction

- The Internet is a substantial part of our lives and every industry relies on the internet and its services to carry out daily tasks.
- Data Source for the purpose of the project, is owned by Mozilla, which is a telemetry dataset.
- The raw data is cleaned and preprocessed followed by the EDA to explore the correlation outliers in the dataset.
- Four different ML models are combined in an ensemble model to predict the internet outages and shutdowns around the world.



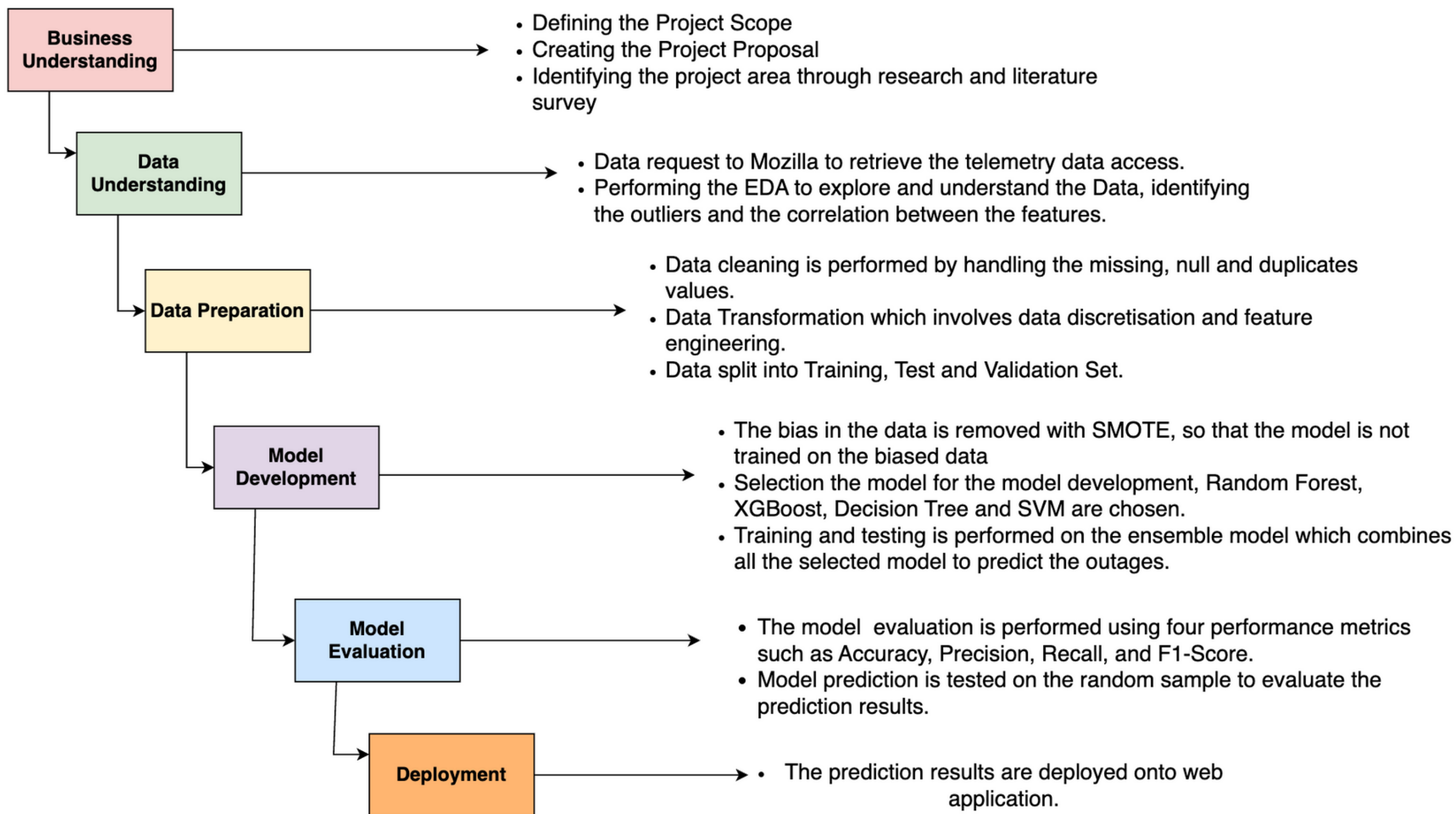
# Motivation

The development and deployment of advanced machine learning models, can achieve enhanced accuracy in predicting outage durations. This not only helps to reduce downtime but also lowers costs for businesses and Internet service providers (ISPs).

Data-driven decision-making can benefit stakeholders by enabling preemptive solutions to network problems and enhanced service-level agreements (SLAs).

An improved user experience will result from shorter and more predictable disruptions.

# Methodology



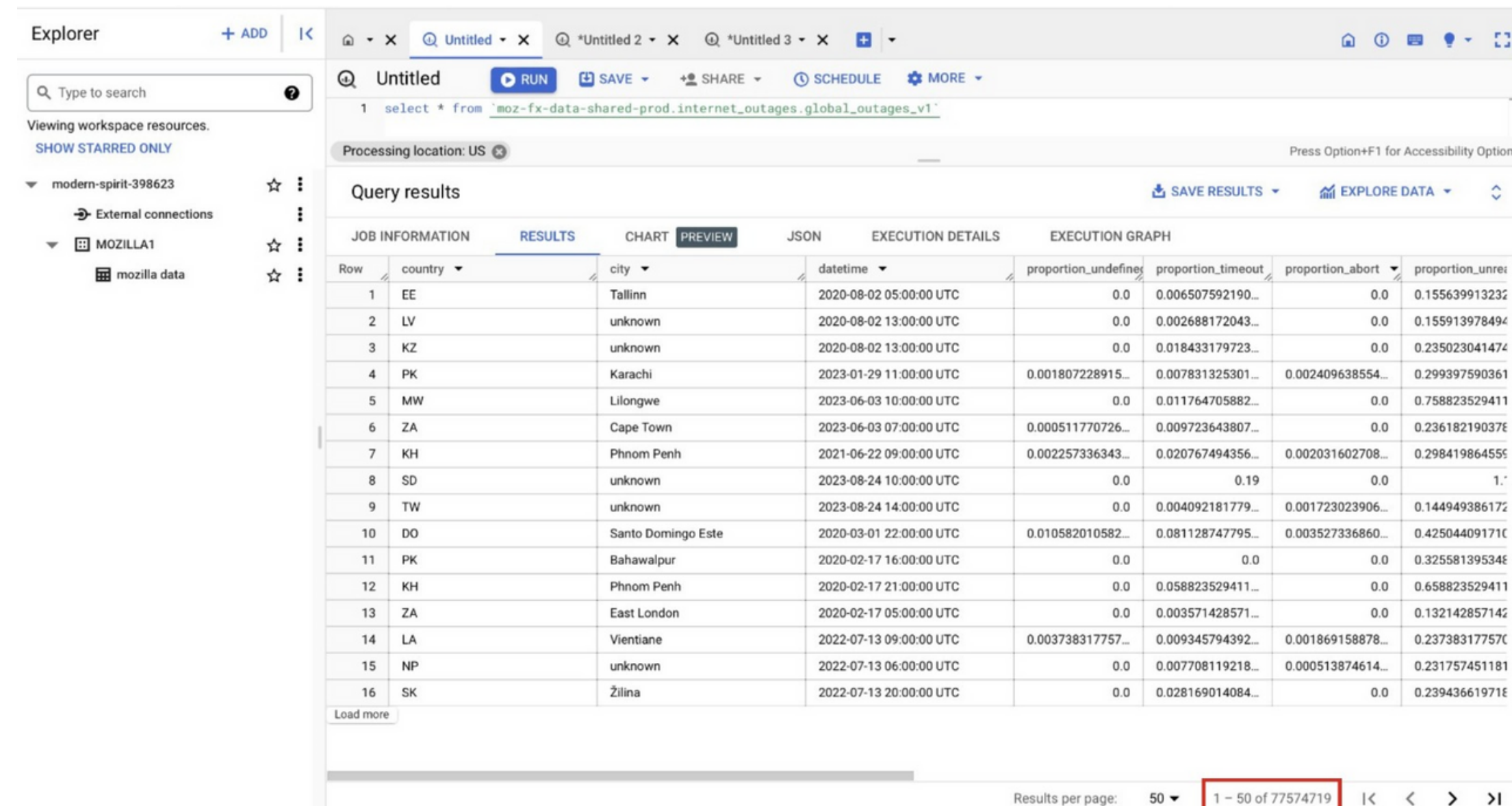
# Literature Survey

Name	About the Dataset	Model Used	Results
Yadwad et al. (2022) Fault Prediction for Network Devices Using Service Outage Prediction Model.	IDM Dataset as been used which contains information about Customer trouble tickets.	Bayesian networks Model and Hidden Markov model	The results imply that Hidden Markov Models are particularly strong as using this probabilistic method improves prediction accuracy when compared to other prediction techniques.
Xu, Yu et al. (2021) Intelligent Outage Probability Prediction for Mobile IoT Networks Based on an IGWO-Elman Neural Network.	IOT sensor Data is collected through Monte-Carlo simulation	IGWO-Elman algorithm	The prediction accuracy of IGWO-Elman is better than algorithms.
Basikolo et al. (2023) Towards zero downtime: Using machine learning to predict network failure in 5G and beyond.	The dataset used is from commerical network testbed.	Random Forest Regressor, Support Vector Regression (SVR), Bayesian ridge and their proposed SVR.	The SVR model, proposed in this study, can swiftly predict the occurrence of a network failure event within the next ten minutes with an f1-score exceeding 0.9 in just ten seconds.
Chen et al. (2019) Outage Prediction and Diagnosis for Cloud Service Systems	The outage dataset is collected from a Microsoft cloud system	Bayesian network and Gradient Boosting Tree based classification model	Using the AirAlertFull the Precision is 71.11, Recall is 100.00, and F1-Score is 83.17



# Data Source

- The source of data is the telemetry dataset owned by Mozilla.
- It consists of the aggregated metrics that correlates to internet outages for different countries in the world.
- The data is available in the Big Query, which is used as the repository where data is fetched and accessed using Google Colab.

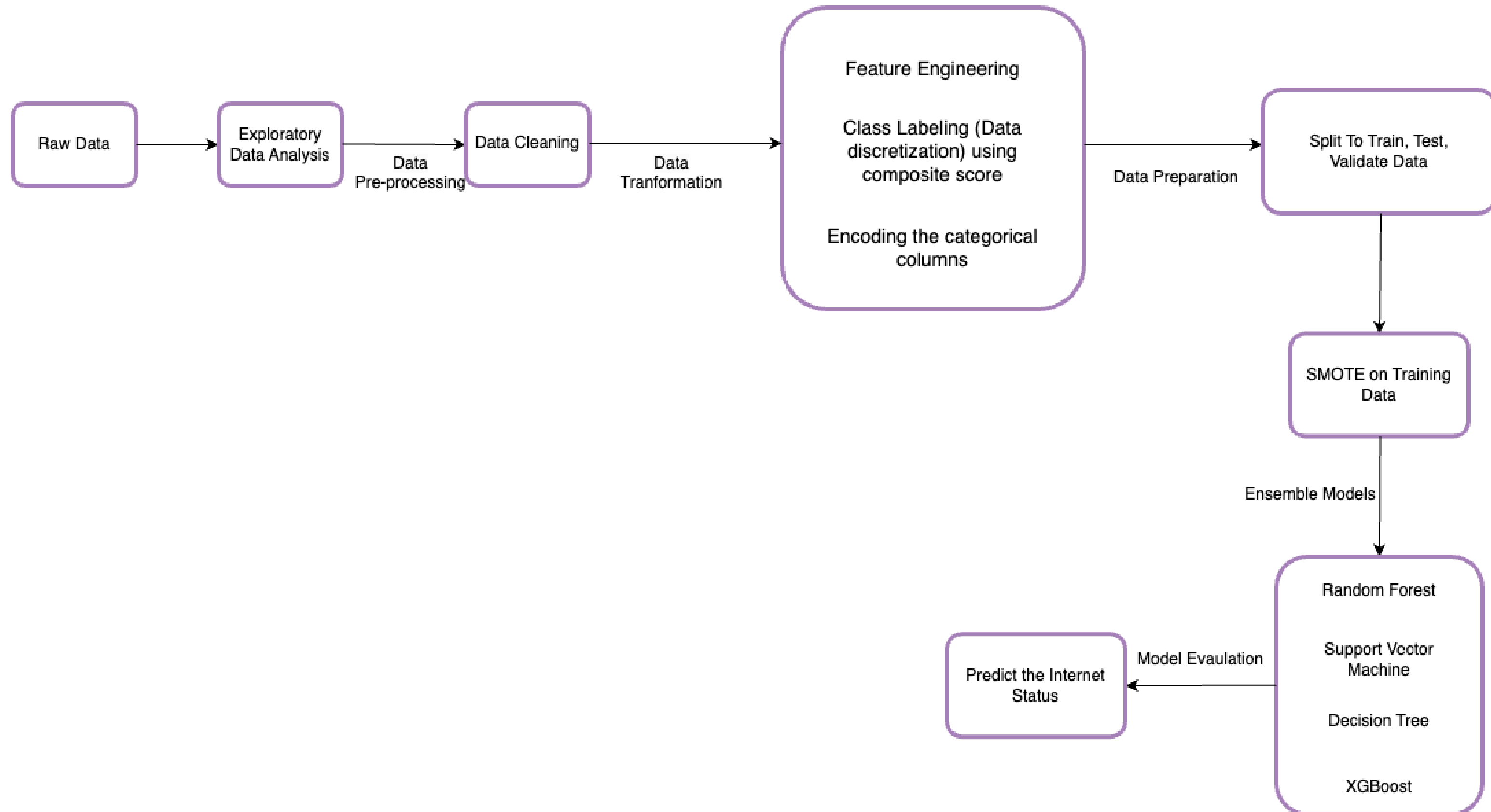


The screenshot displays the Google BigQuery interface. On the left, the 'Explorer' panel shows a workspace with a folder named 'modern-spirit-398623' containing an 'External connections' link and a 'MOZILLA1' folder with a 'mozilla data' table. The main panel shows a query titled 'Untitled' with the following SQL: `select * from 'moz-fx-data-shared-prod.internet_outages.global_outages_v1'`. The query results are displayed in a table with 8 columns: Row, country, city, datetime, proportion\_undefined, proportion\_timeout, proportion\_abort, and proportion\_unre. The table contains 16 rows of data, showing outages for various countries and cities. The bottom of the interface shows 'Results per page: 50' and '1 - 50 of 77574719'.

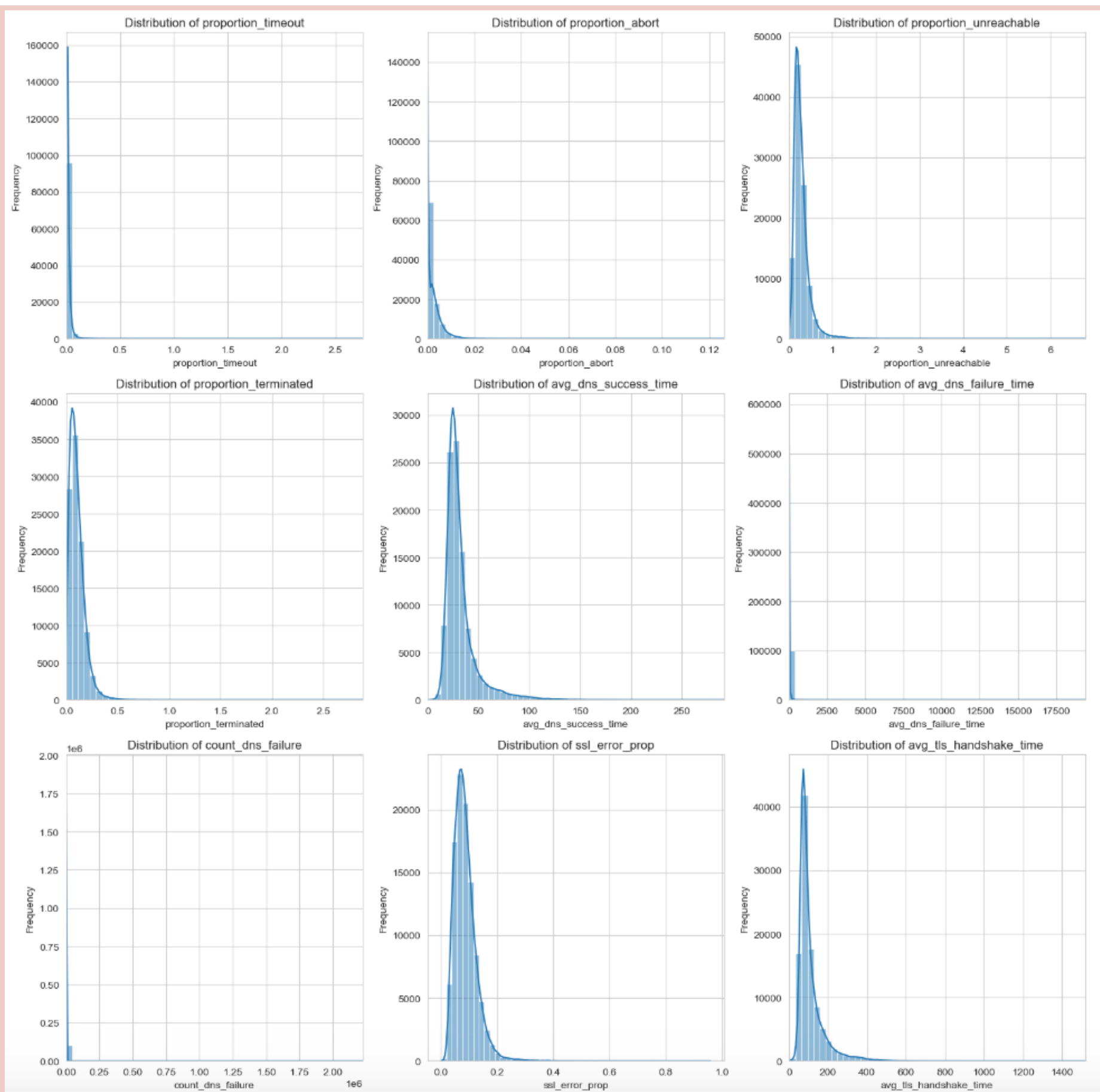
Row	country	city	datetime	proportion_undefined	proportion_timeout	proportion_abort	proportion_unre
1	EE	Tallinn	2020-08-02 05:00:00 UTC	0.0	0.006507592190...	0.0	0.155639913232
2	LV	unknown	2020-08-02 13:00:00 UTC	0.0	0.002688172043...	0.0	0.155913978494
3	KZ	unknown	2020-08-02 13:00:00 UTC	0.0	0.018433179723...	0.0	0.235023041474
4	PK	Karachi	2023-01-29 11:00:00 UTC	0.001807228915...	0.007831325301...	0.002409638554...	0.299397590361
5	MW	Lilongwe	2023-06-03 10:00:00 UTC	0.0	0.011764705882...	0.0	0.758823529411
6	ZA	Cape Town	2023-06-03 07:00:00 UTC	0.000511770726...	0.009723643807...	0.0	0.236182190378
7	KH	Phnom Penh	2021-06-22 09:00:00 UTC	0.002257336343...	0.020767494356...	0.002031602708...	0.298419864556
8	SD	unknown	2023-08-24 10:00:00 UTC	0.0	0.19	0.0	1.0
9	TW	unknown	2023-08-24 14:00:00 UTC	0.0	0.004092181779...	0.001723023906...	0.144949386172
10	DO	Santo Domingo Este	2020-03-01 22:00:00 UTC	0.010582010582...	0.081128747795...	0.003527336860...	0.425044091710
11	PK	Bahawalpur	2020-02-17 16:00:00 UTC	0.0	0.0	0.0	0.325581395348
12	KH	Phnom Penh	2020-02-17 21:00:00 UTC	0.0	0.058823529411...	0.0	0.658823529411
13	ZA	East London	2020-02-17 05:00:00 UTC	0.0	0.003571428571...	0.0	0.132142857142
14	LA	Vientiane	2022-07-13 09:00:00 UTC	0.003738317757...	0.009345794392...	0.001869158878...	0.237383177570
15	NP	unknown	2022-07-13 06:00:00 UTC	0.0	0.007708119218...	0.000513874614...	0.231757451181
16	SK	Žilina	2022-07-13 20:00:00 UTC	0.0	0.028169014084...	0.0	0.239436619718

**Data Source URL:**[https://wiki.mozilla.org/Mozilla\\_Network\\_Outages\\_Data\\_Project](https://wiki.mozilla.org/Mozilla_Network_Outages_Data_Project)

# Project Data Flow



# Exploratory Data Analysis

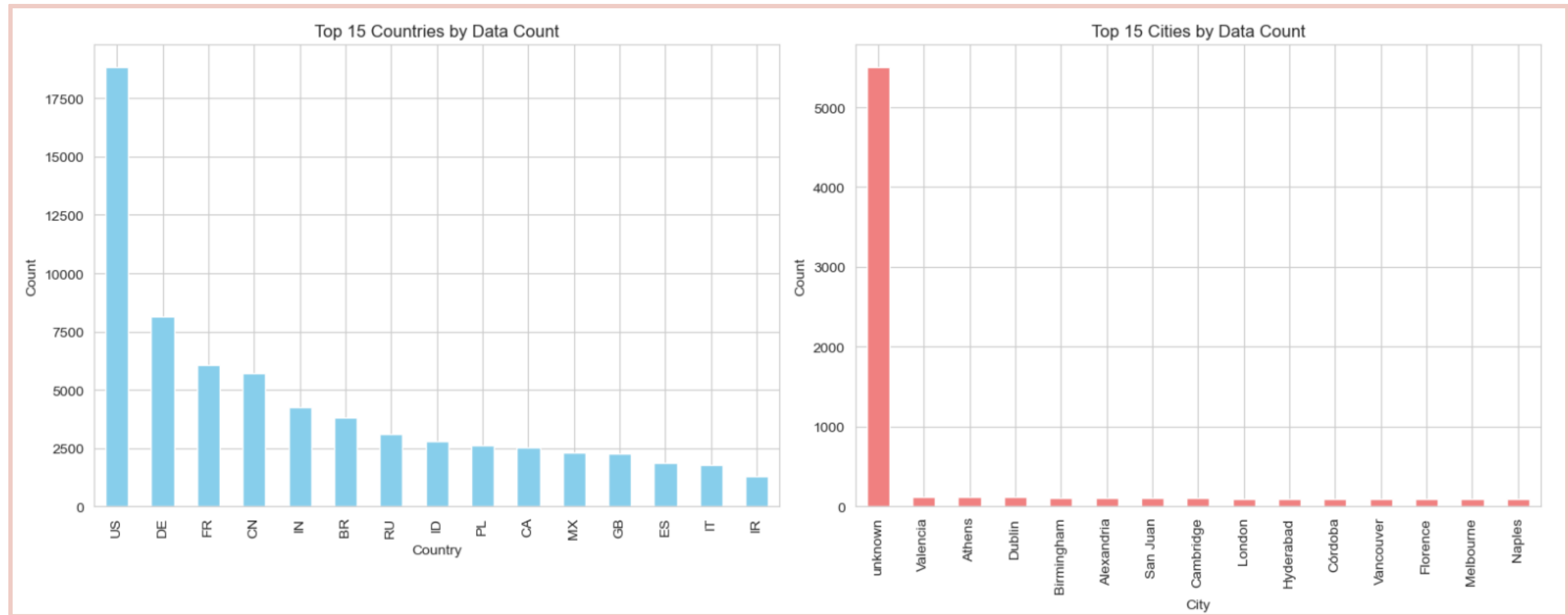


## Identifying the Target Feature

- The histogram depicts the distribution of continuous features.
- Most values range between 0 and 1, with possible outliers exceeding 1.



# Exploratory Data Analysis



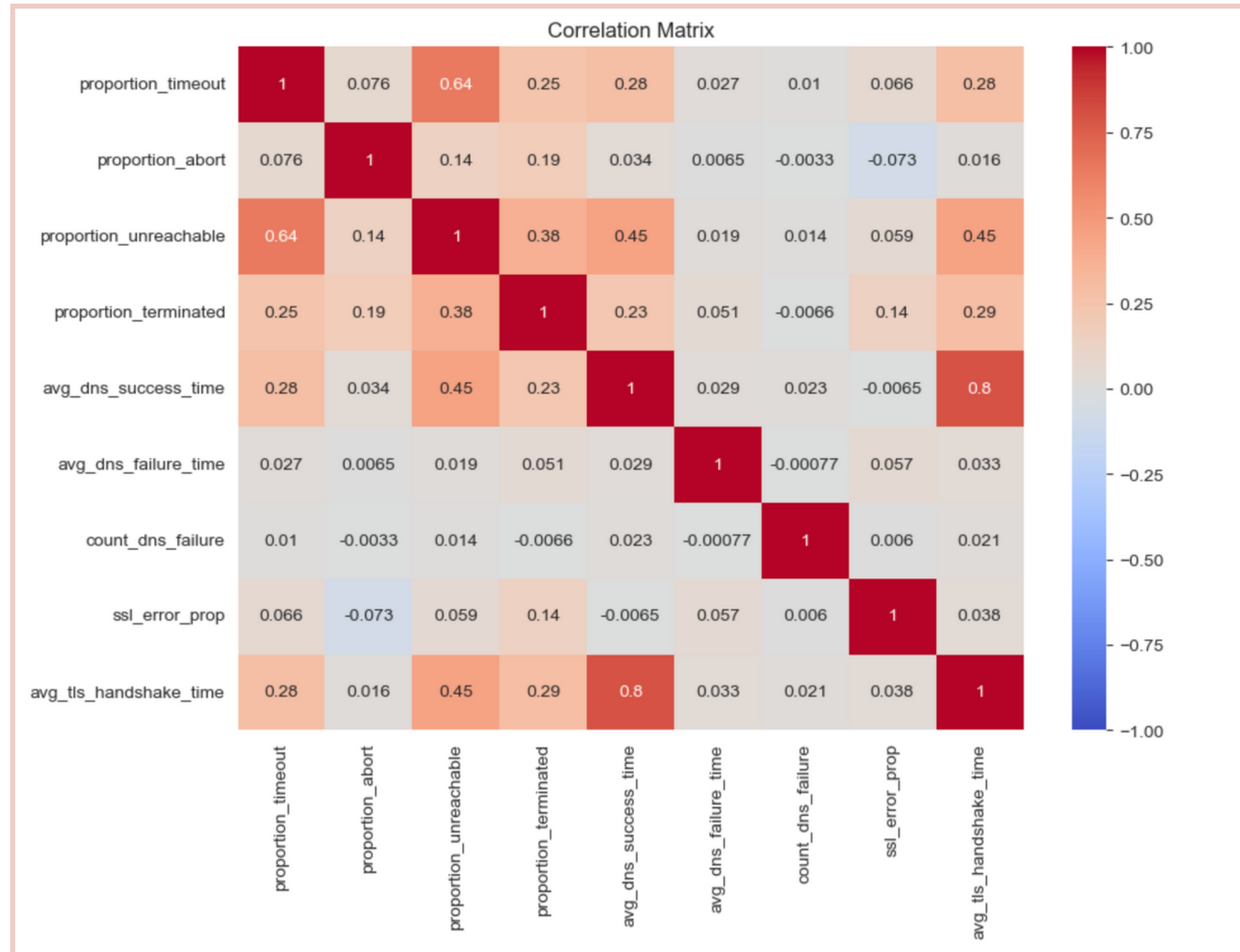
**Identifying the most common country and city in the Dataset**

To analyse the categorical features in the dataset, the plot displays the count of datapoints of top 15 countries and cities

# Exploratory Data Analysis

## Finding the correlation between the numerical features

- The heat-map above is used to understand the correlation between continuous features in the dataset.
- We can see that most of the features have positive correlation.



# Data Cleaning

```
df.isnull().sum()
country          246
city              0
datetime         0
proportion_undefined  0
proportion_timeout  0
proportion_abort  0
proportion_unreachable  0
proportion_terminated  0
proportion_channel_open  0
avg_dns_success_time  0
missing_dns_success  1
avg_dns_failure_time 10
missing_dns_failure  1
count_dns_failure   295
ssl_error_prop      1
avg_tls_handshake_time 696
dtype: int64
```

Identified Null Values

```
[ ] df.isnull().sum()
country          0
city              0
datetime         0
proportion_undefined  0
proportion_timeout  0
proportion_abort  0
proportion_unreachable  0
proportion_terminated  0
proportion_channel_open  0
avg_dns_success_time  0
missing_dns_success  0
avg_dns_failure_time  0
missing_dns_failure  0
count_dns_failure   0
ssl_error_prop      0
avg_tls_handshake_time 0
dtype: int64
```

Removed Null Values

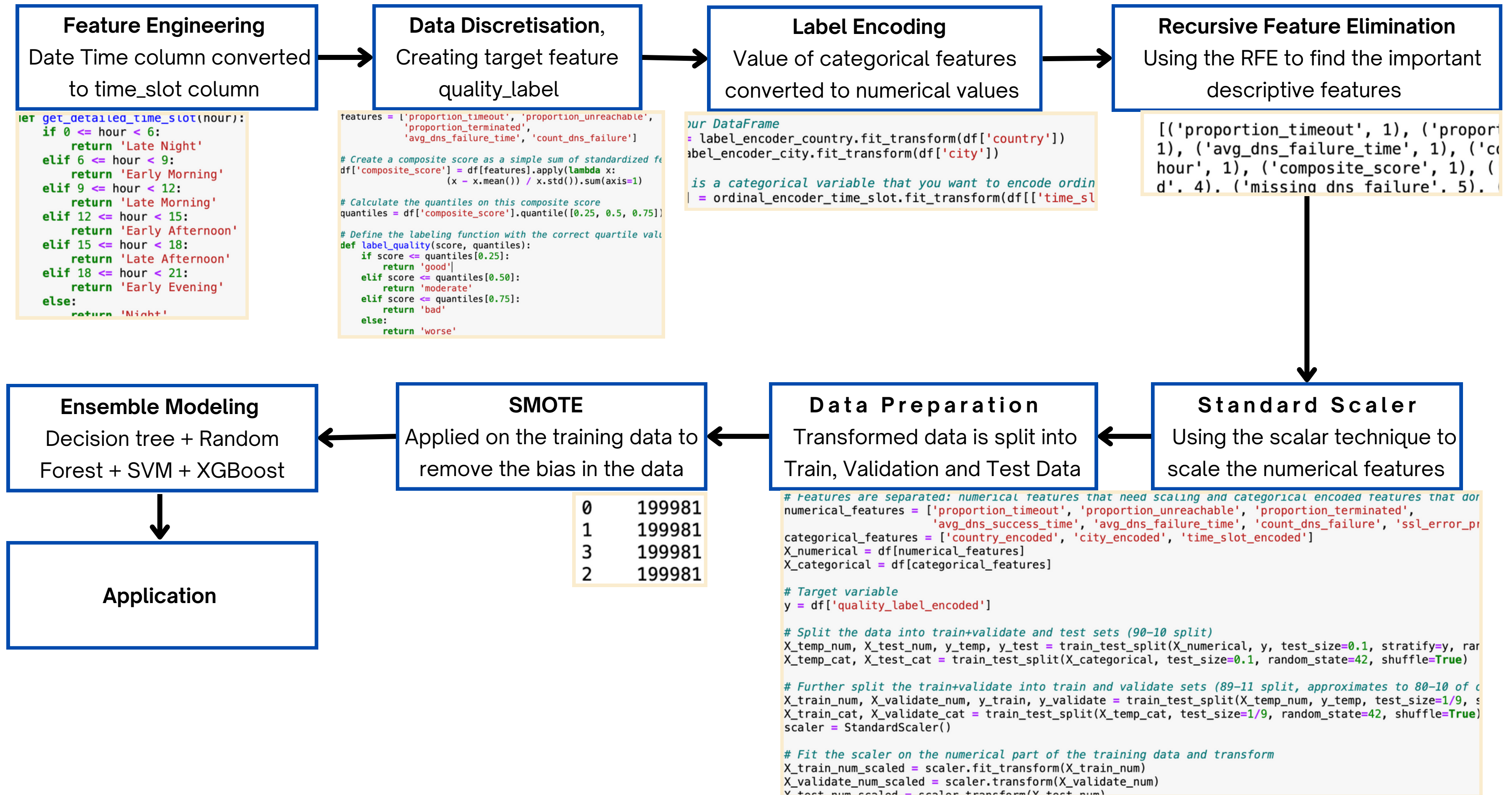
	country	city	datetime	proportion_undefined	proportion_timeout	proportion_abort
42	MM	unknown	2023-03-15 05:00:00+00:00	0.000000	0.009988	0.002497
43	MU	unknown	2023-03-15 10:00:00+00:00	0.000000	0.022642	0.007547
51	NL	unknown	2023-03-15 01:00:00+00:00	0.000000	0.008701	0.000322
54	NZ	unknown	2023-03-15 06:00:00+00:00	0.000000	0.011407	0.007605
98	IR	unknown	2023-02-12 03:00:00+00:00	0.000220	0.049934	0.002782
...	...	...	...	...	...	...
11905423	BE	unknown	2023-07-12 21:00:00+00:00	0.000000	0.008821	0.002786
11905447	CY	unknown	2023-07-12 06:00:00+00:00	0.000000	0.000000	0.000000
11905474	CD	unknown	2023-09-14 21:00:00+00:00	0.000000	0.052632	0.000000
11905497	BG	unknown	2023-06-06 06:00:00+00:00	0.000562	0.005807	0.000375
11905511	CH	unknown	2023-06-06 23:00:00+00:00	0.000000	0.010029	0.001910

649844 rows x 16 columns

Handling unkown values

The identified unknown values are not dropped as it consist of legitimate data values

# Data Transformation and Model building



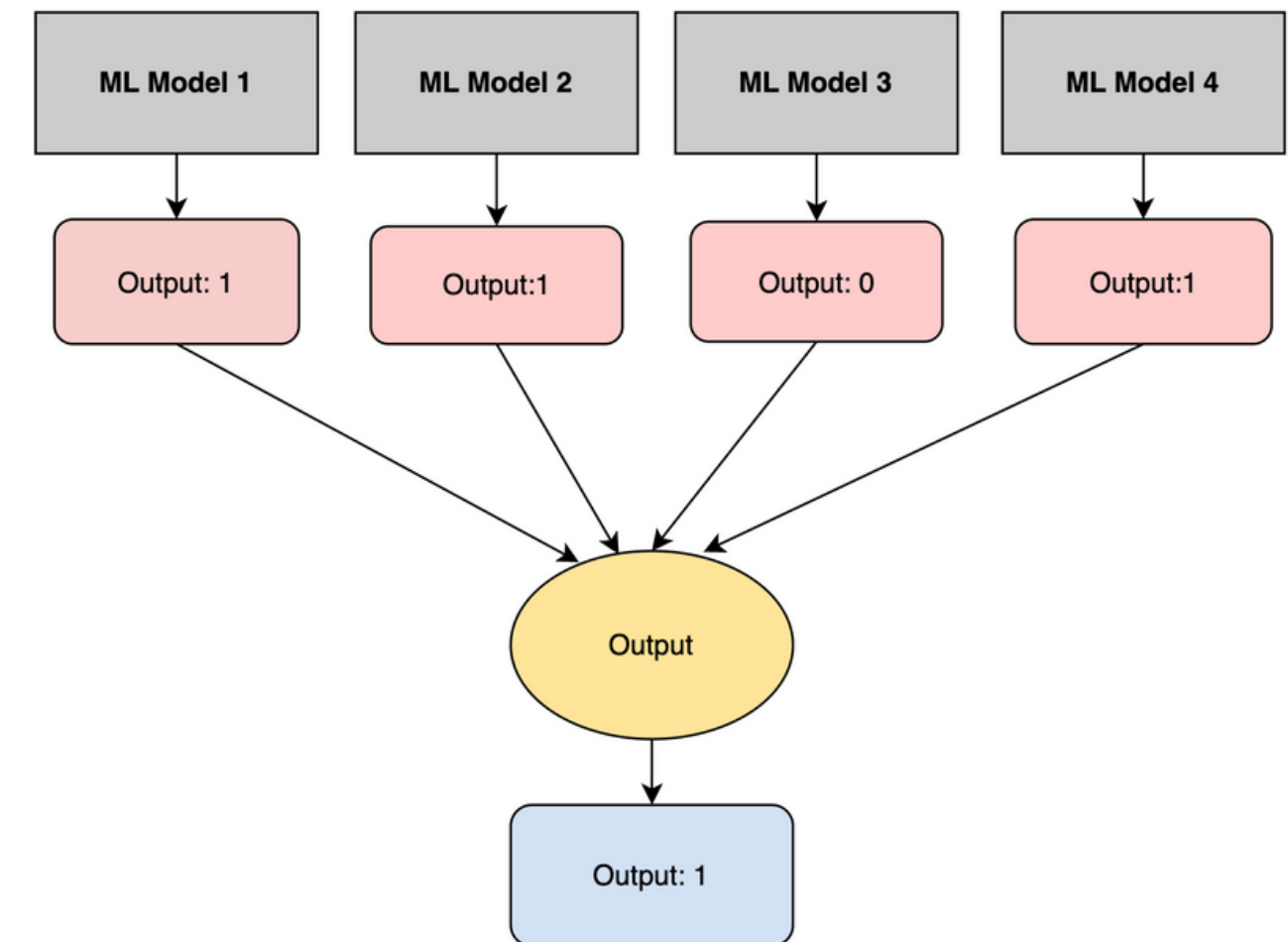
# Model Comparision

Models	Accuracy	Precision	F1 Score	Training Tlme
Random Forest	96%	98%	98%	30 minutes
Support Vector Machine	71.8%	72.3%	72.3%	45 minutes
Decision Tree	93.42%	96%	96%	15 minutes
XGBoost	97%	99%	99%	25 minutes
Ensemble Model	96.28%	96%	96%	20 minutes



# Ensemble Model

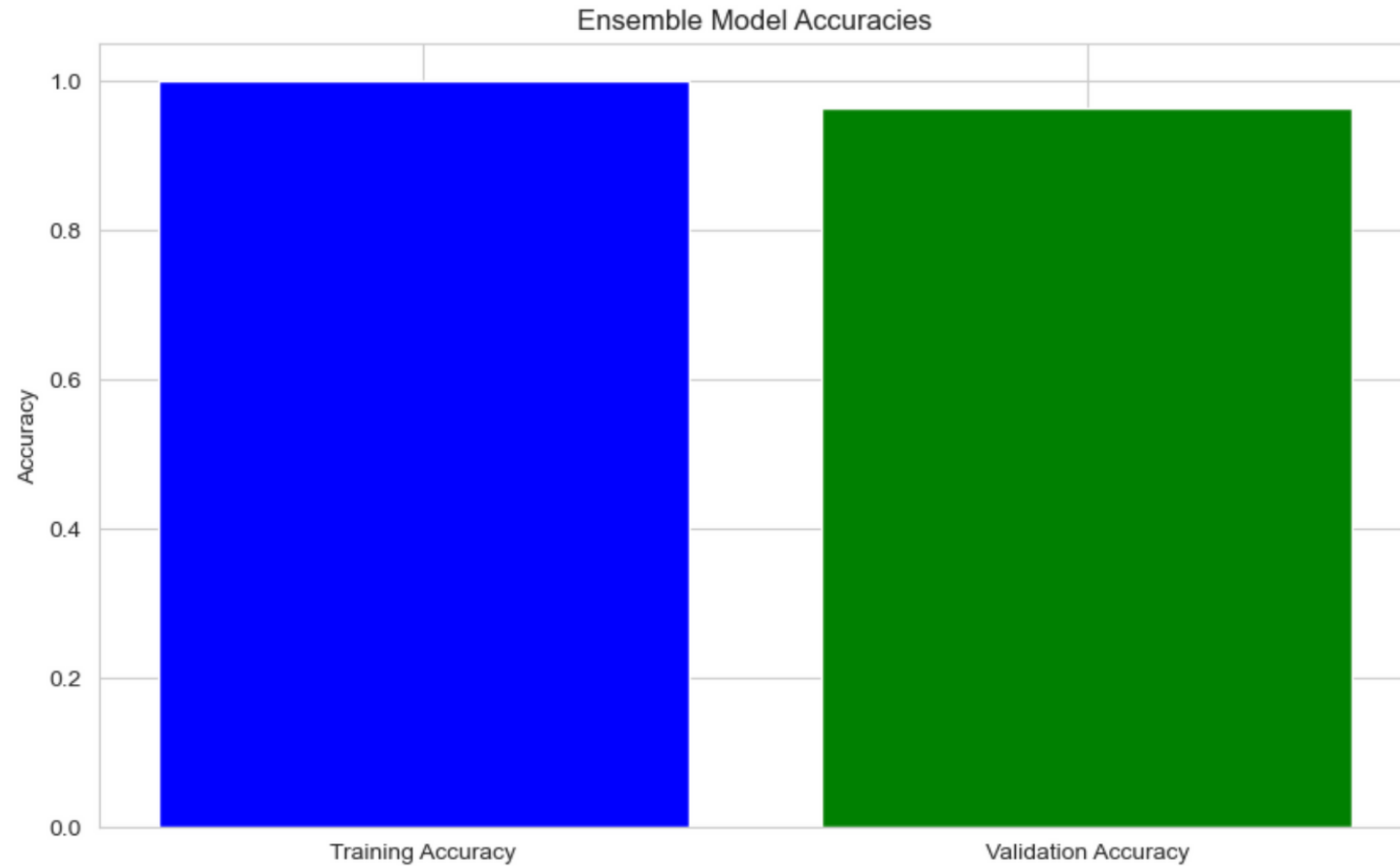
- The ensemble model is created by combining the four machine learning algorithms which are Support Vector Machine, Random Forest, Decision Tree, and XGBoost Algorithm.
- The prediction from all the combined algorithms is integrated to generate the final outcome.
- The ensemble model uses a Hard voting technique that performs the prediction based on the majority voting.



**Depiction of Hard Voting in Ensemble Model**

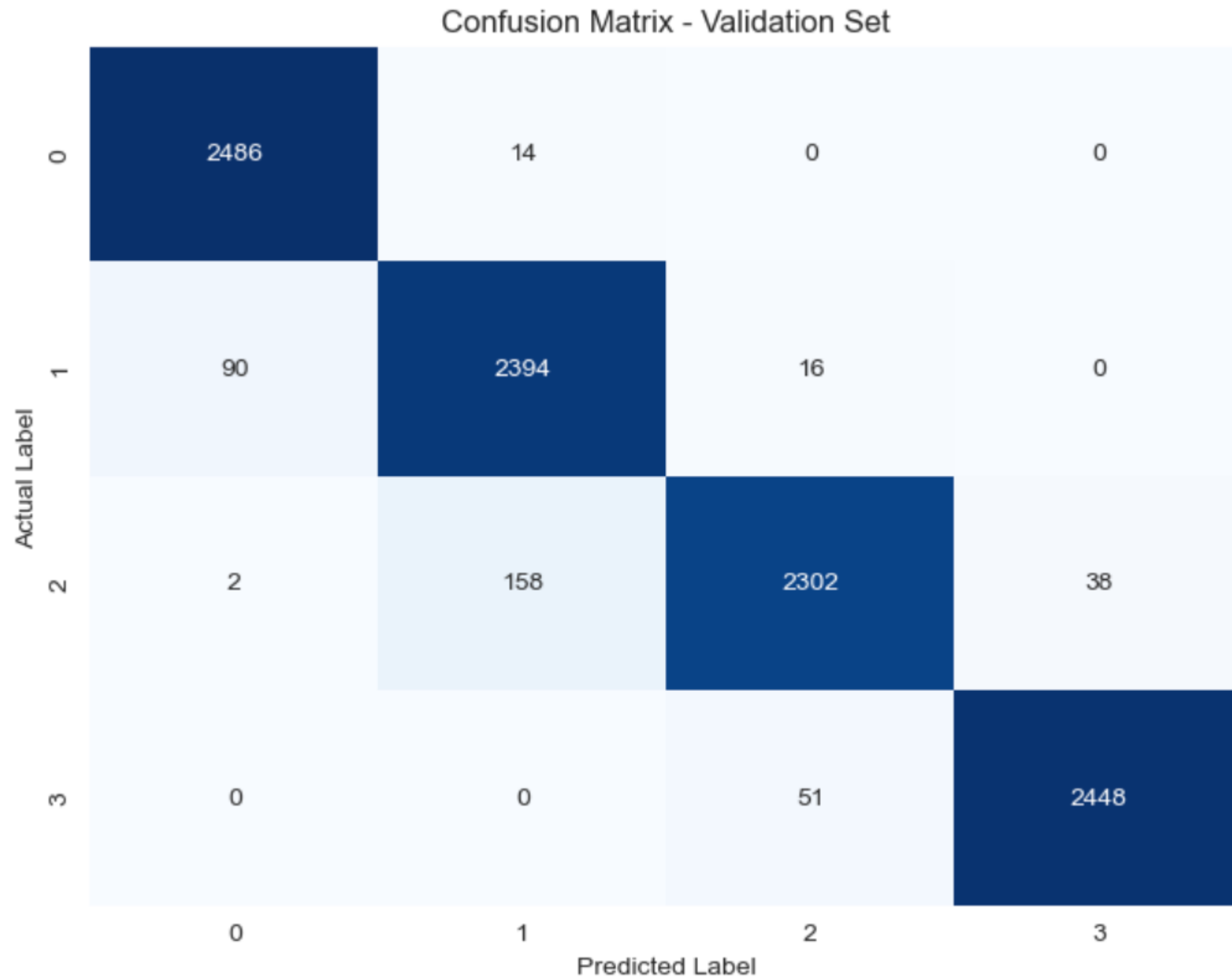


# Model Performance (cont.d)



Plotting Training and Validation accuracies

# Model Performance



The confusion matrix shows the majority of the labels that are predicted were actually from that labeled class

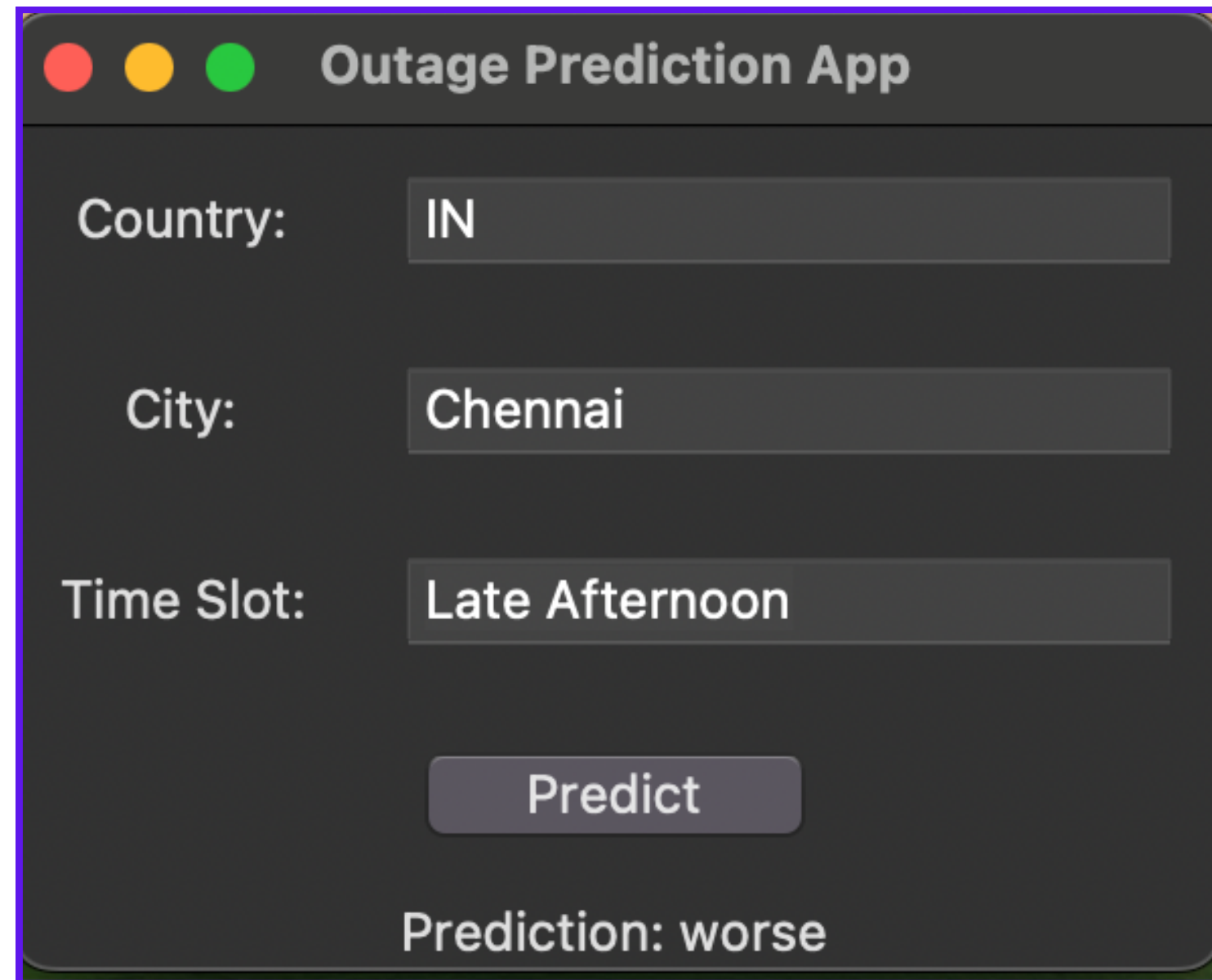
# Conclusion

By precisely forecasting internet outages using spatial and temporal data, this project greatly improves network management and holds out the prospect of better service reliability and customer satisfaction in the telecom industry. By addressing present network management issues, the project lays foundations for upcoming improvements in service quality and predictive maintenance.

# Future Scope

- Incorporating real-time data feeds to update the predictions dynamically.
- Conducting a cost-benefit analysis to understand the financial impact of outages.
- Develop models to detect anomalies in network traffic that could indicate potential issues leading to outages.

# Deployment



The image shows a screenshot of a web application titled "Outage Prediction App". The interface is dark-themed with a purple border. It features three input fields: "Country:" with the value "IN", "City:" with the value "Chennai", and "Time Slot:" with the value "Late Afternoon". Below these fields is a "Predict" button. At the bottom of the form, the text "Prediction: worse" is displayed.

Field	Value
Country:	IN
City:	Chennai
Time Slot:	Late Afternoon
Predict	
Prediction:	worse

Application which accepts Country, City and Time slot as parameter and predicts the outage

Thank you!



*Q & A*