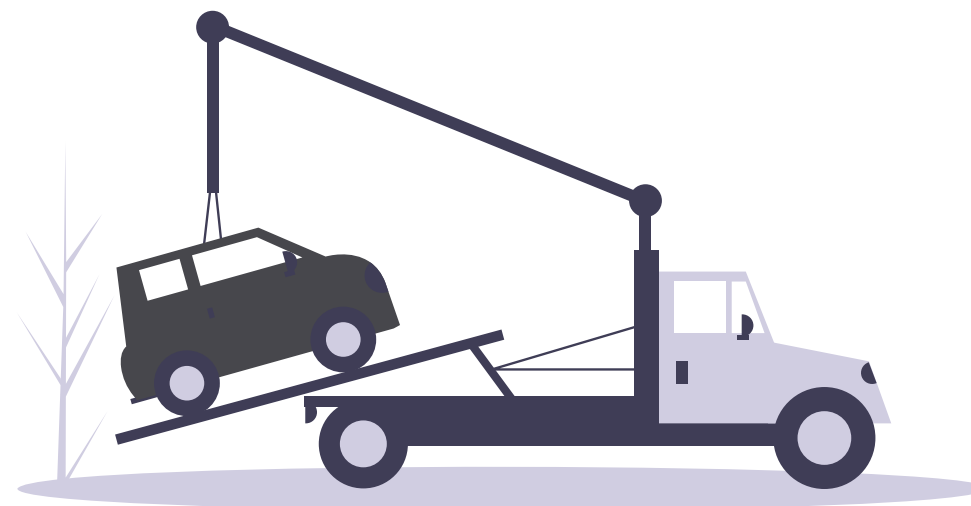# ACCIDENT SEVERITY PREDICTION

MATTEO ORSINI 1795119

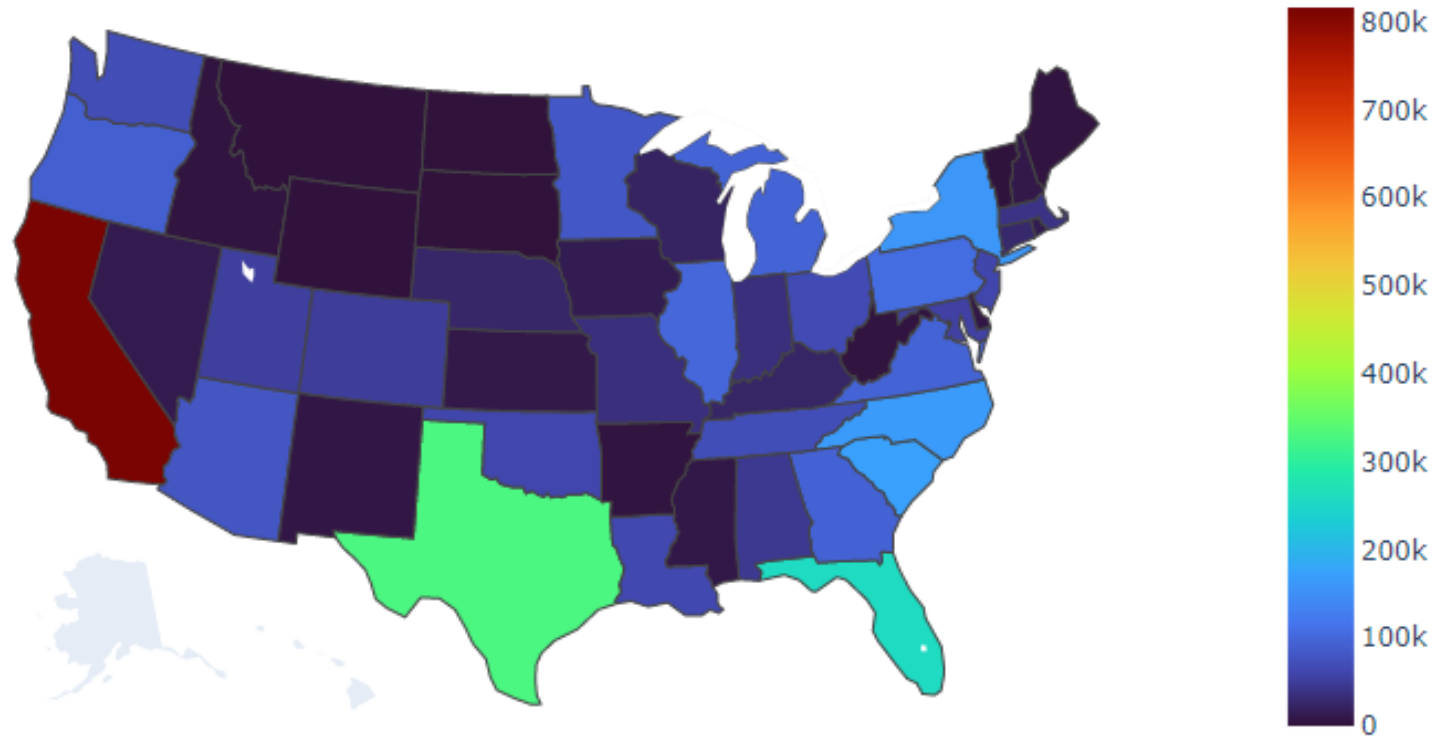FABRIZIO ROSSI 1815023

MINA MAKAR 1804475

# TASK - SEVERITY

- Classify the **impact of an accident** on the traffic (severity)

- The scale of the severity goes from 1 (low impact) to 4 (high impact)

- We used the dataset «US Accidents» from Kaggle

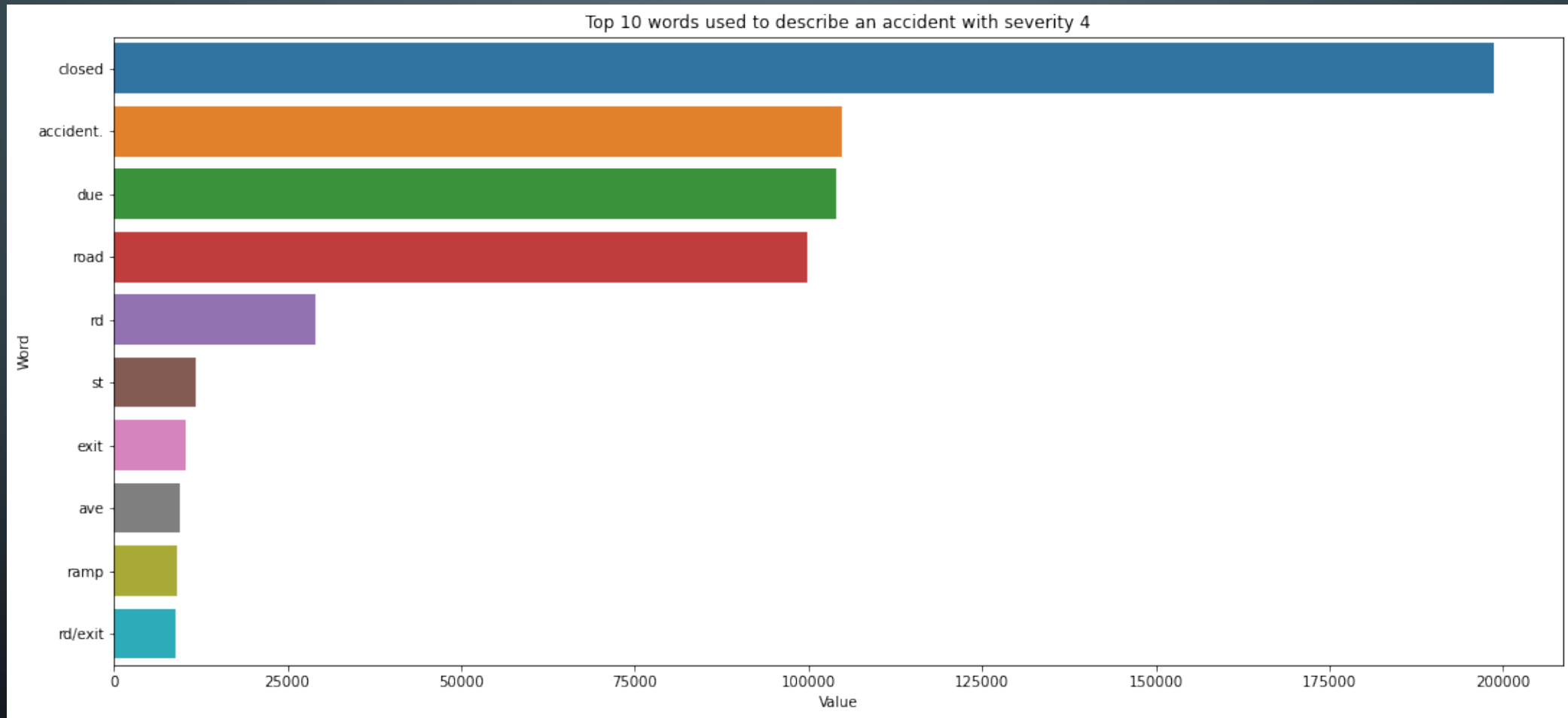- Contains about **3.5 million** accidents collected in the US
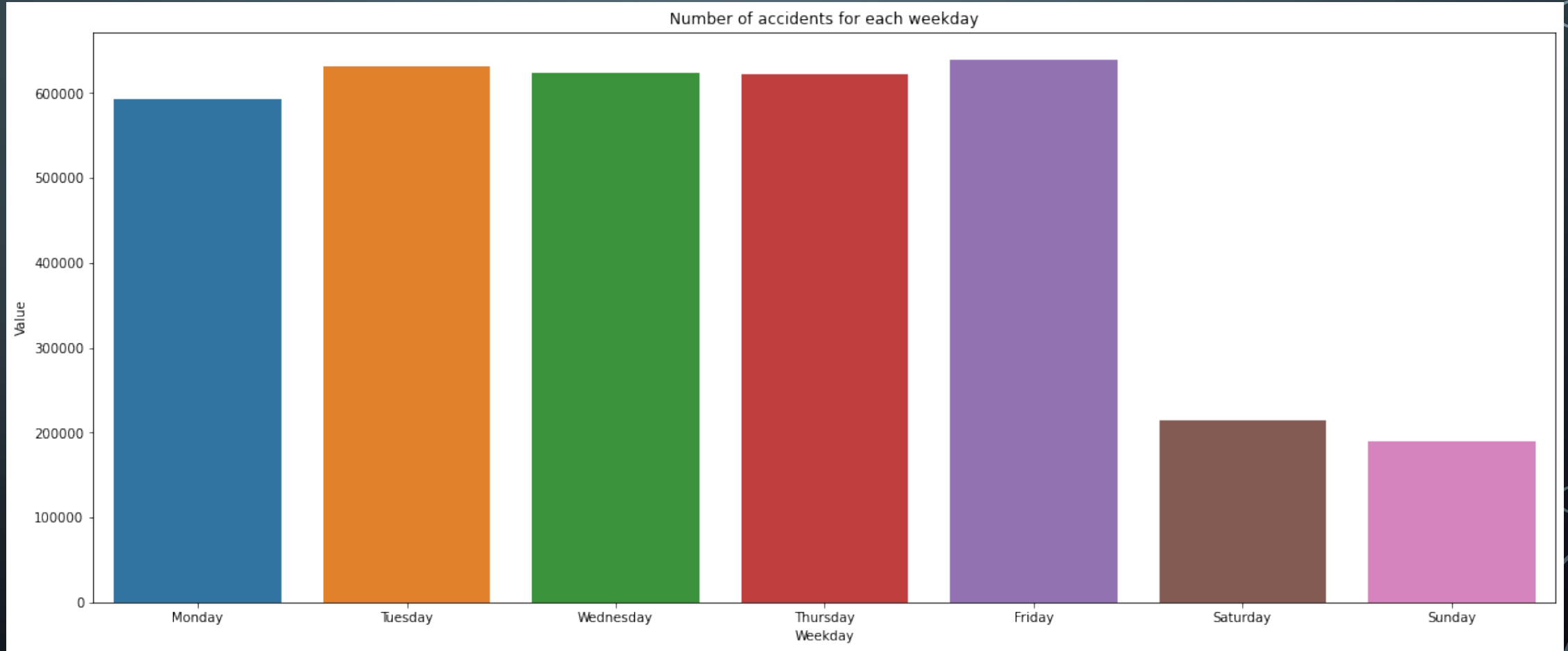
# EXPLORATORY DATA ANALYSIS



Number of US Accidents for each State

# EXPLORATORY DATA ANALYSIS



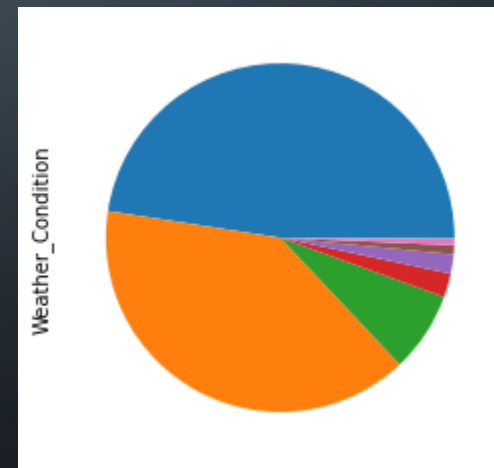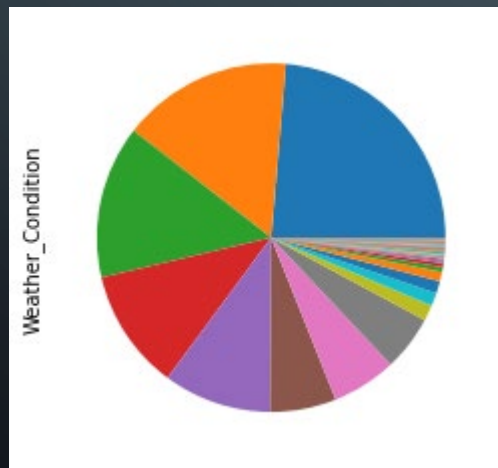Top 10 words used to describe an accident with severity 4
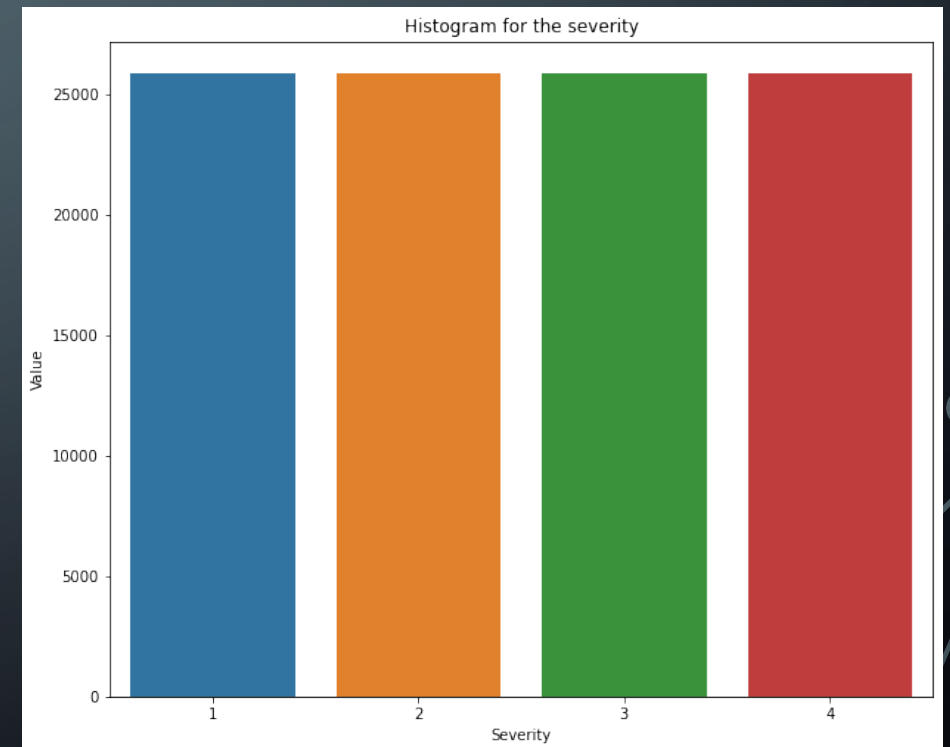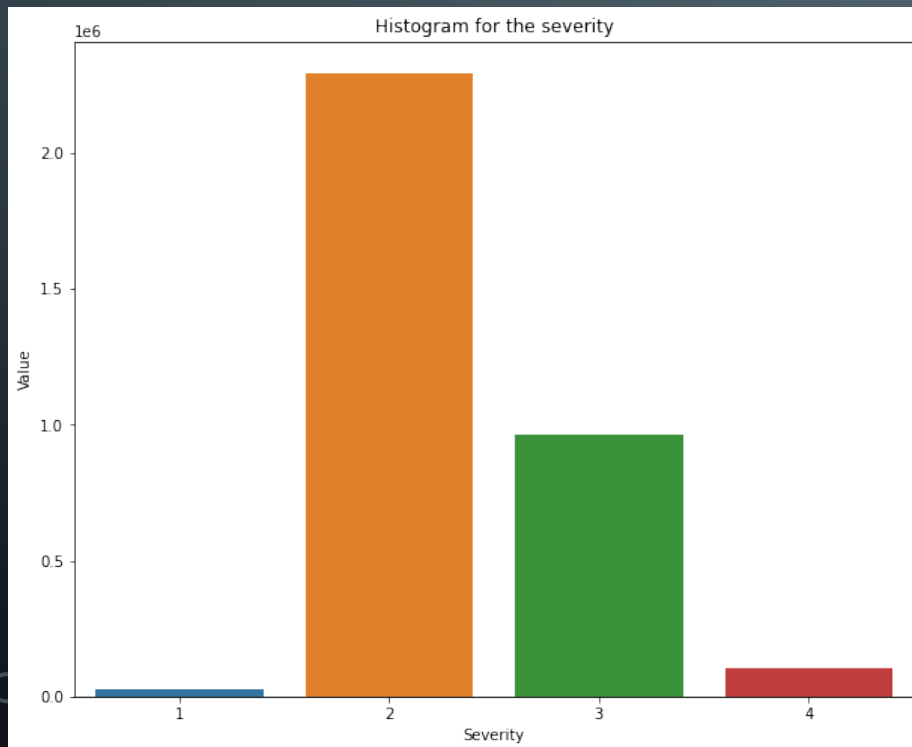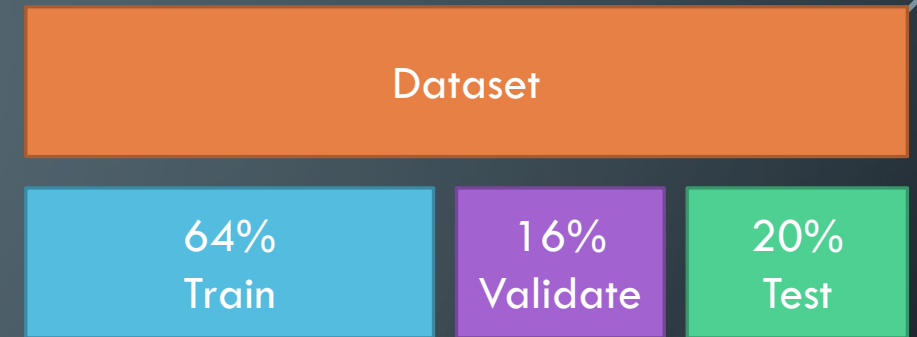
# EXPLORATORY DATA ANALYSIS



Number of accidents for each weekday

# DATA PREPROCESSING

- Reduced number of classes for Weather_Condition (128 → 11) and Wind_Direction (24 → 10)

- Filled missing values with mean for numerical features

- Removed records with missing values for categorical features

- Scaled and encoded features using the one-hot encoding for the categorical features

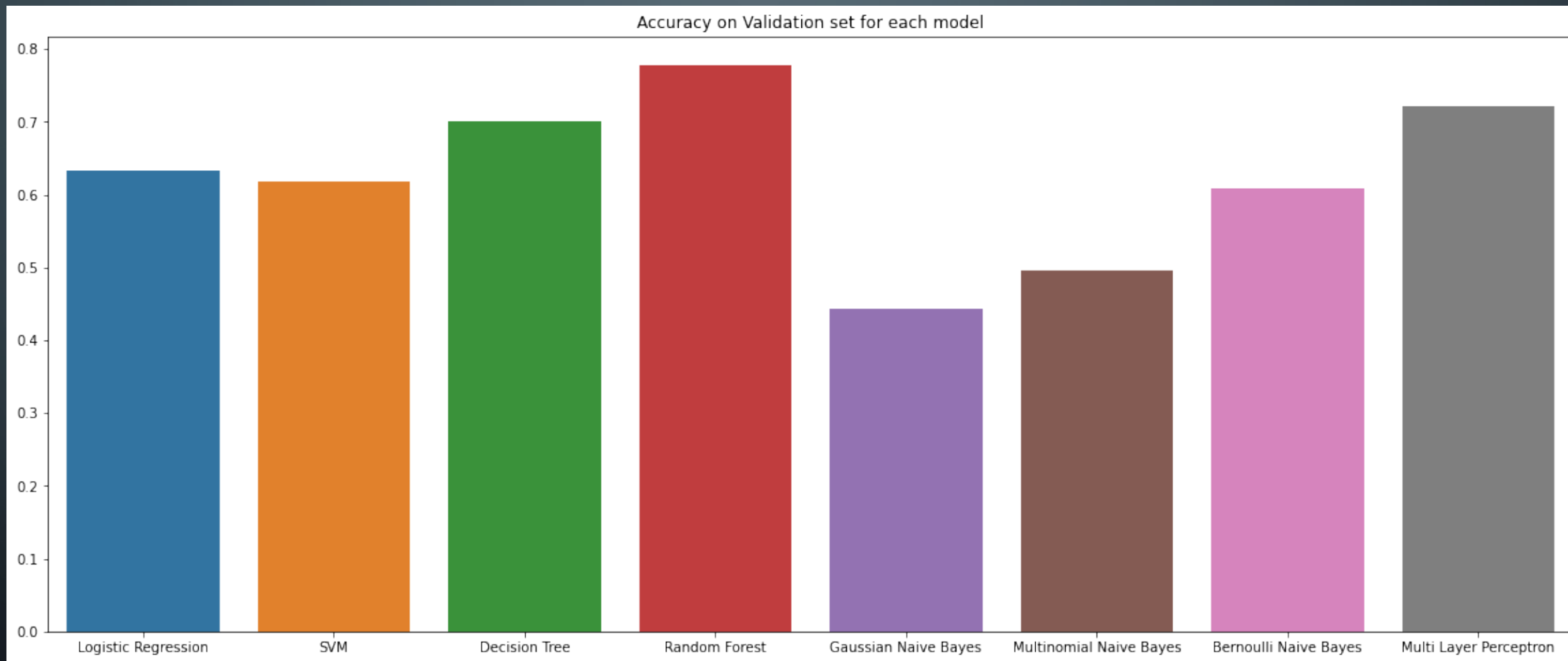- City was encoded using the binary encoder due to the large number of unique values
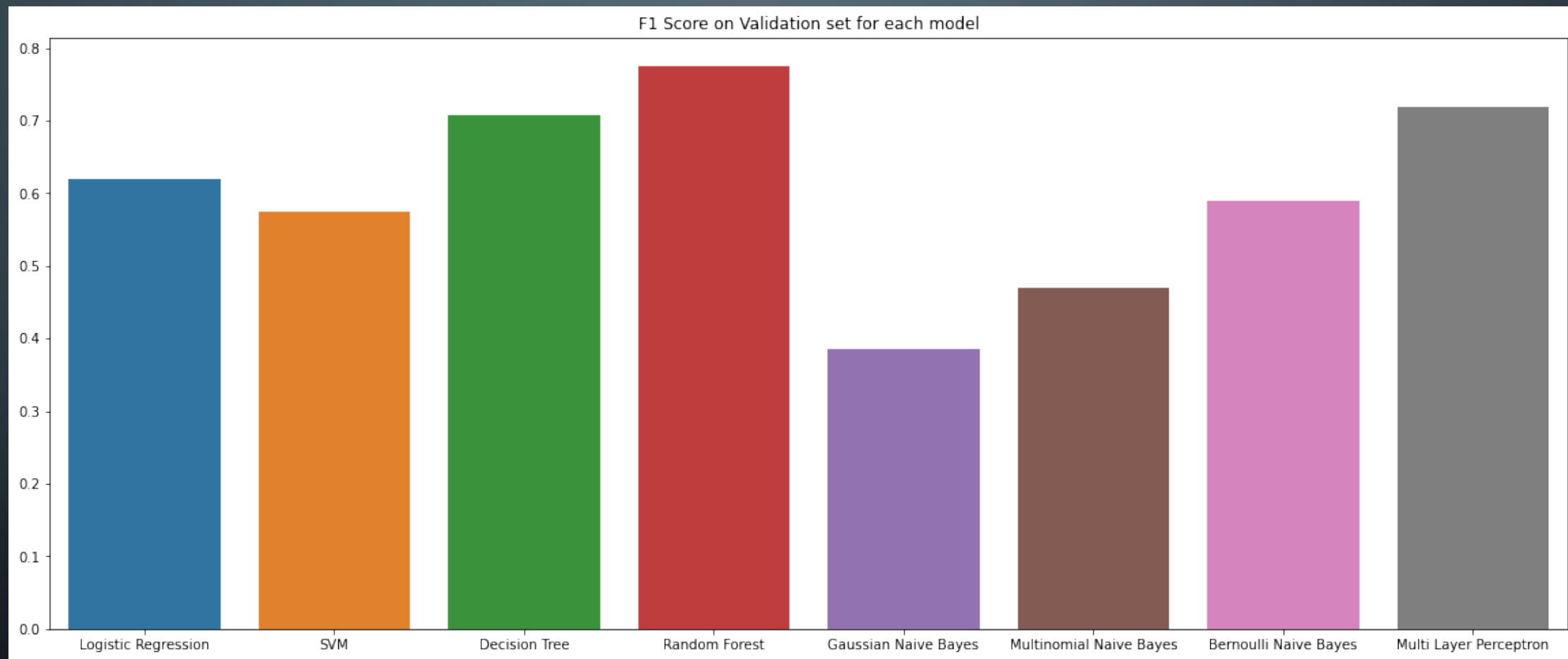
# DATA PREPROCESSING

- We handled the unbalanced dataset problem using the undersampling technique

- We splitted the dataset in **train set, validation set** and **test set**
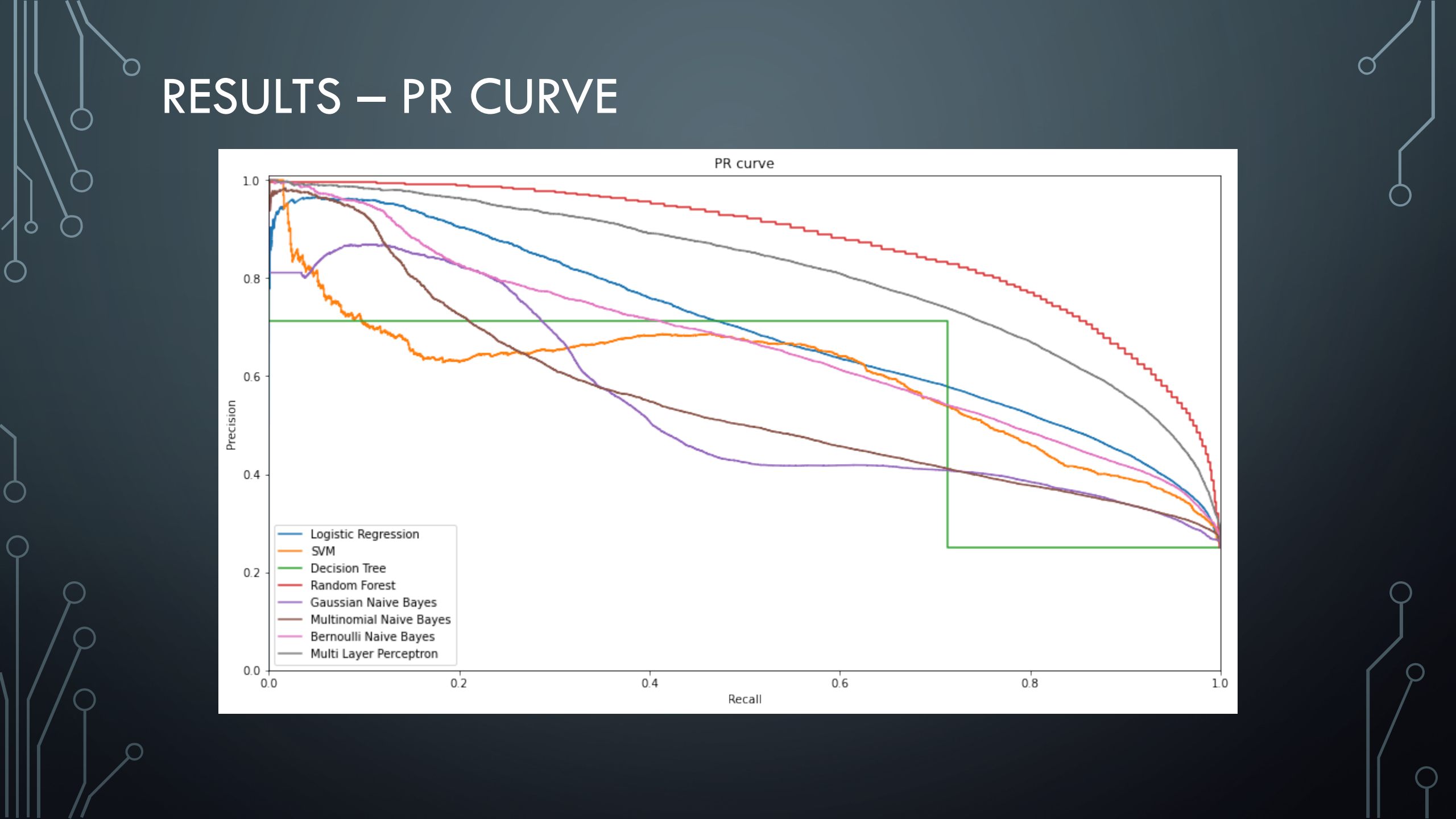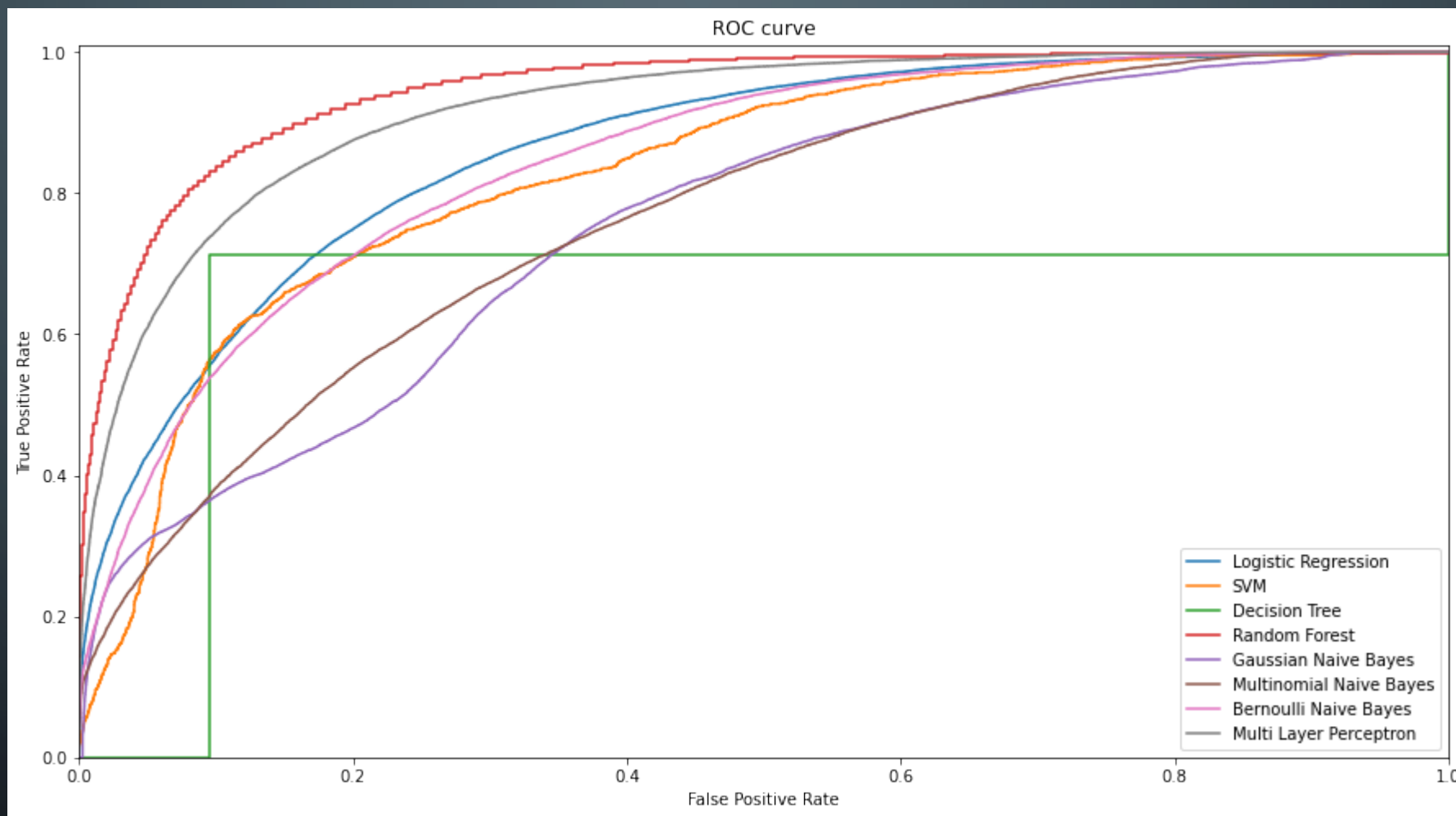
# RESULTS – ACCURACY



Accuracy on Validation set for each model

# RESULTS – F1 SCORE



F1 Score on Validation set for each model

# RESULTS — PR CURVE

# RESULTS – ROC CURVE

# RESULTS – RANDOM FOREST TEST SET



| Class | Precision | Recall | F1-Score |
|---|---|---|---|
| 1 | 0,90 | 0,97 | 0,93 |
| 2 | 0,75 | 0,59 | 0,66 |
| 3 | 0,70 | 0,65 | 0,67 |
| 4 | 0,74 | 0,91 | 0,81 |
| Macro avg | 0,77 | 0,78 | 0,77 |