## Libraries Used:
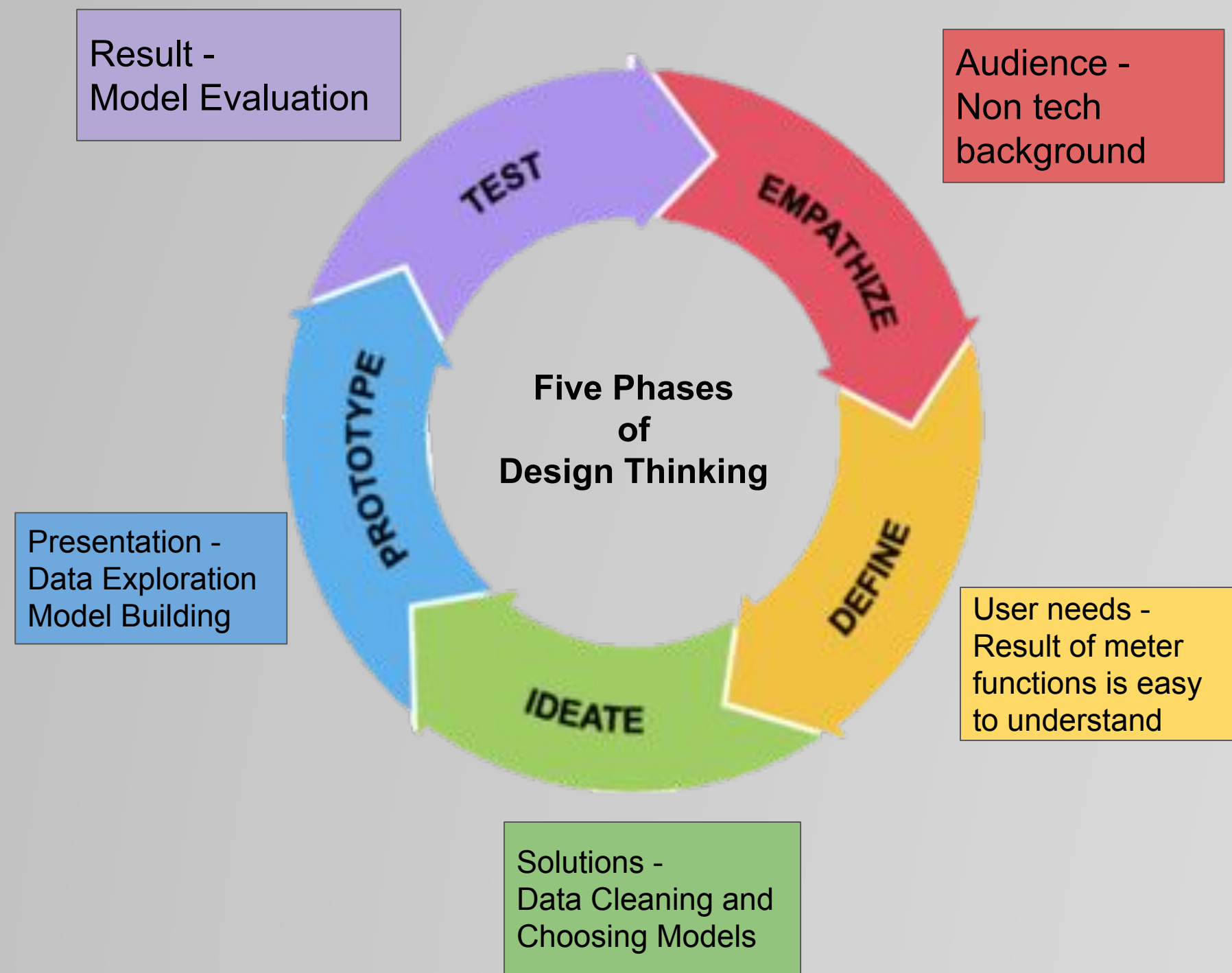- Pandas
- Numpy
- Matplotlib
- Seaborn
- Sklearn
- Datetime
- dmba

**Analyze challenge with Design Thinking**



Result - Model Evaluation

Audience - Non tech background

Presentation - Data Exploration Model Building

Five Phases of Design Thinking

User needs - Result of meter functions is easy to understand

Solutions - Data Cleaning and Choosing Models

## Dimensionality Reduction Techniques Used:
- Principal Component Analysis
- Linear Discriminant Analysis

### PCA



Analysis provided us with 16 components that explained 95% of the variance.
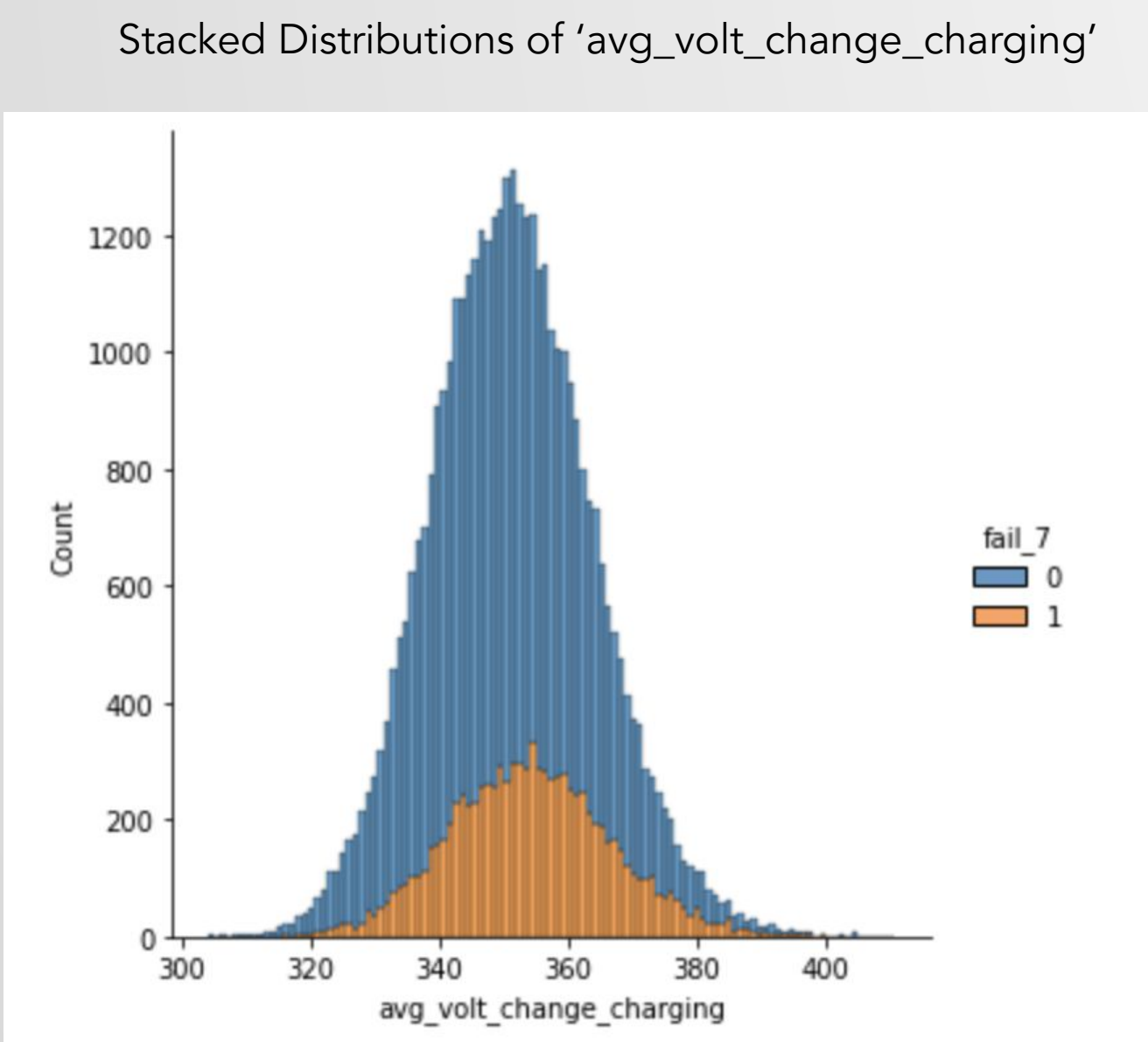
### LDA

```
[[5149    0]
 [1618    0]]
Accuracy 0.7608984779074922
```

- Ran LDA using one linear discriminant
- Gives us ~76% accuracy
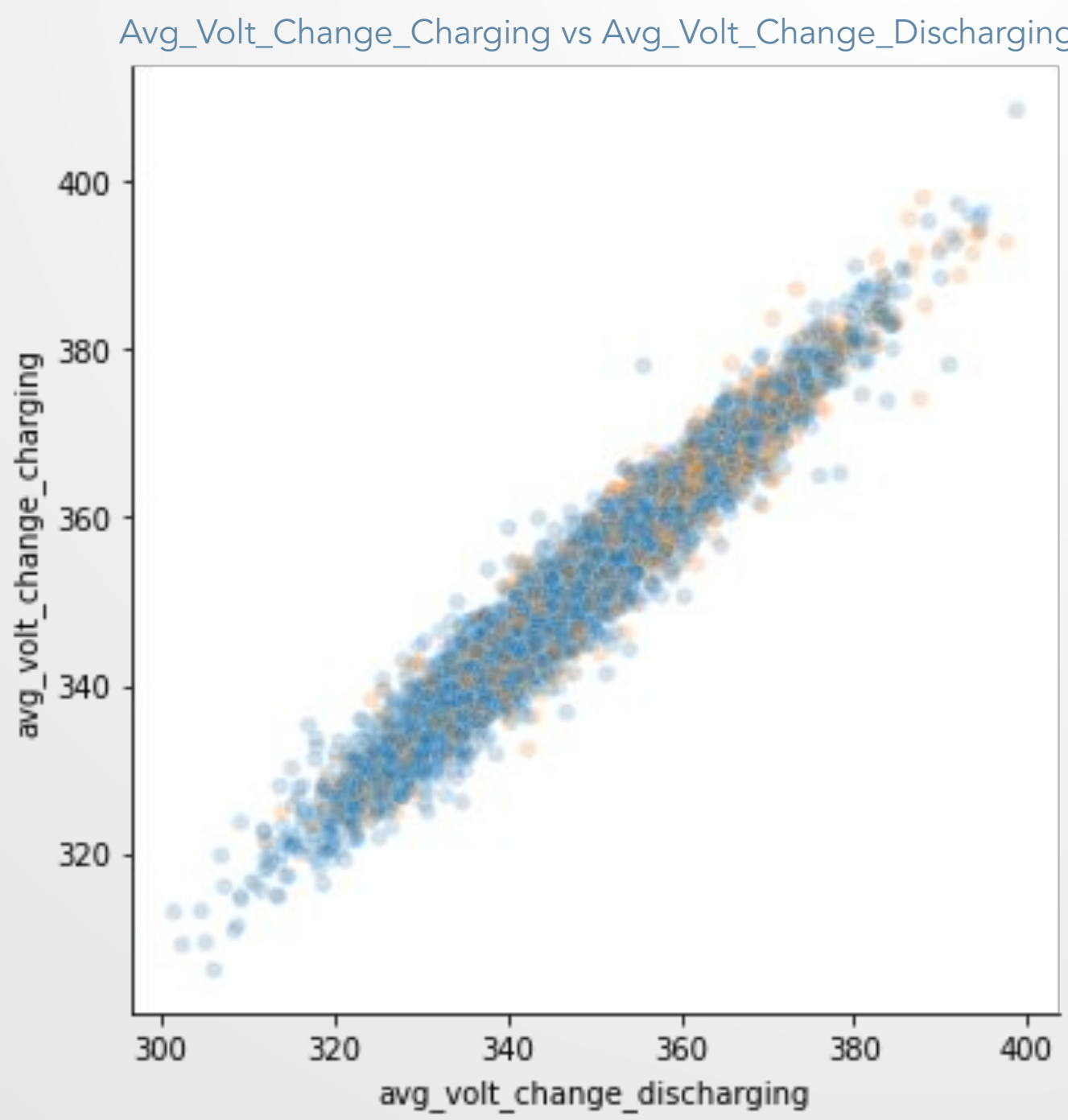- Chose to go with PCA for our dimensionality reduction

# Data Exploration


Correlation Matrix


Max_Voltage_Day Distributions among Working and Failed Meters

By looking at the correlation matrix we noticed 'max_voltage_day' also had a relatively high correlation with the 'fail_7' column and wanted to take a look at the distributions of 'max_voltage_day' day values across meters that did and didn't fail to see if there were any significant differences in the distributions. From the image above we noticed that the meters that failed have comparatively higher counts in the larger max_voltage_day values.


Stacked Distributions of 'avg_volt_change_charging'


Avg_Volt_Change_Charging vs Avg_Volt_Change_Discharging

Here we produced a stacked distribution of avg voltage charging broken down by working and failed meters where both distributions follow normal distribution centered around 350 volts.

For this plot we noticed that both the 'avg_volt_change_charging' and 'avg_volt_change_discharging' had relatively high correlation values to the fail_7 column in the correlation metric so we decided to create a scatter plot. Unfortunately the data was a little too dense to extract any meaningful insight so instead we plotted a sample of the data and noticed while relatively uniform it does seem that higher values in both axis correlate to a larger proportion of failed meters.
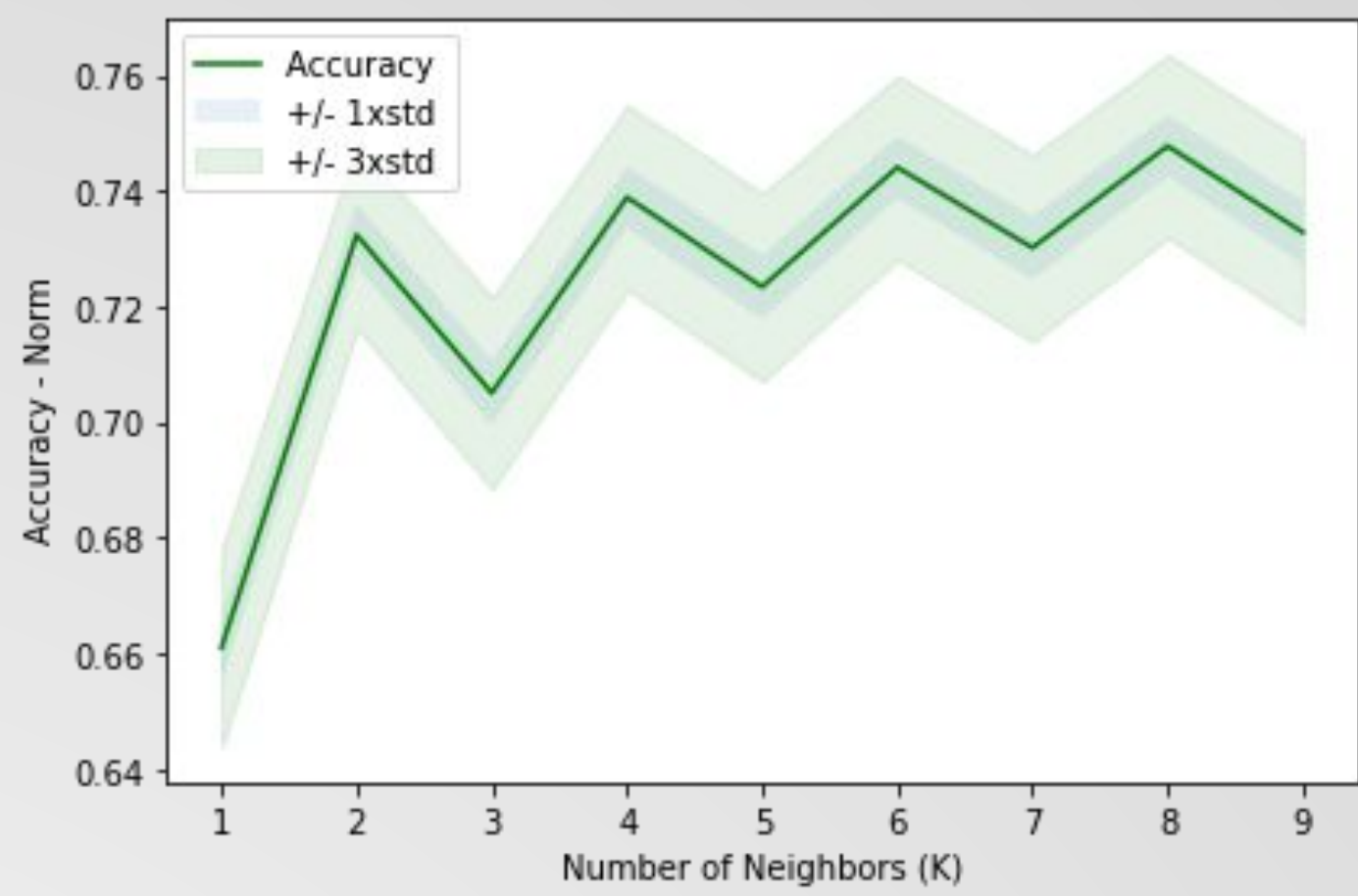
## Summary

As a classification problem, we first looked at the accuracy scores where all our models returned values in the range from 0.73 to 0.77 except Logistic Regression model with normalized dataset, whose accuracy is only 0.38. The Logistic Regression fit with the PCA's has highest accuracy 0.7616.

Since none of the models were able to produce a significantly better accuracy score than the others, we then looked at the F1 scores to which we found our Logistic Regression model fit with our normalized dataset returning the largest value 0.42, but since it returned the lowest accuracy score we are not going to consider it as our choice. The Naive Bayes model fit with PCA's retuned strong precision and recall score compared to the other models. Finally, the Naive Bayes fit with PCA's returned the best F1 score of the remaining models still being considered. Therefore, the Naive Bayes fit with PCA's is our best fit model.
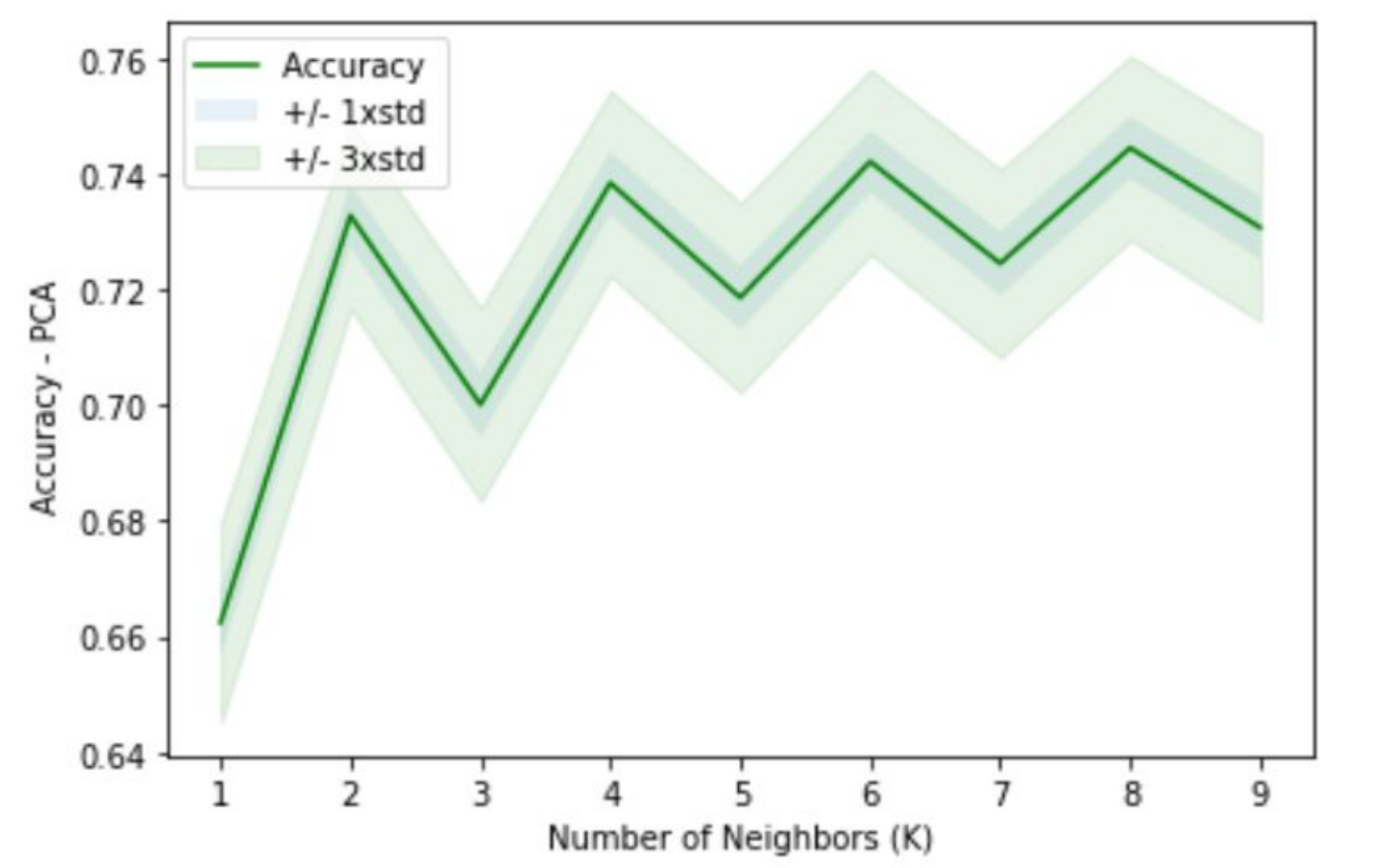
# Model Building

- K-Nearest Neighbors
- Decision Tree
- Naive Bayes
- Logistic Regression

## K-Nearest Neighbor Model



With normalized dataset, the best K is 8 with accuracy around 74.8%



With PCA dataset, the best K is 8 with accuracy around 74.4%

## Logistic Regression



## Naive Bayes


Normalized Dataset

```
Confusion Matrix (Accuracy 0.6198)

        Prediction
Actual    0      1
  0    12105   8337
  1     1952   4671

Accuracy: 0.7608984779074922
Precision score: 0.3590867158671587
Recall score: 0.7052695153253813
F1 score: 0.4758799857368448
```
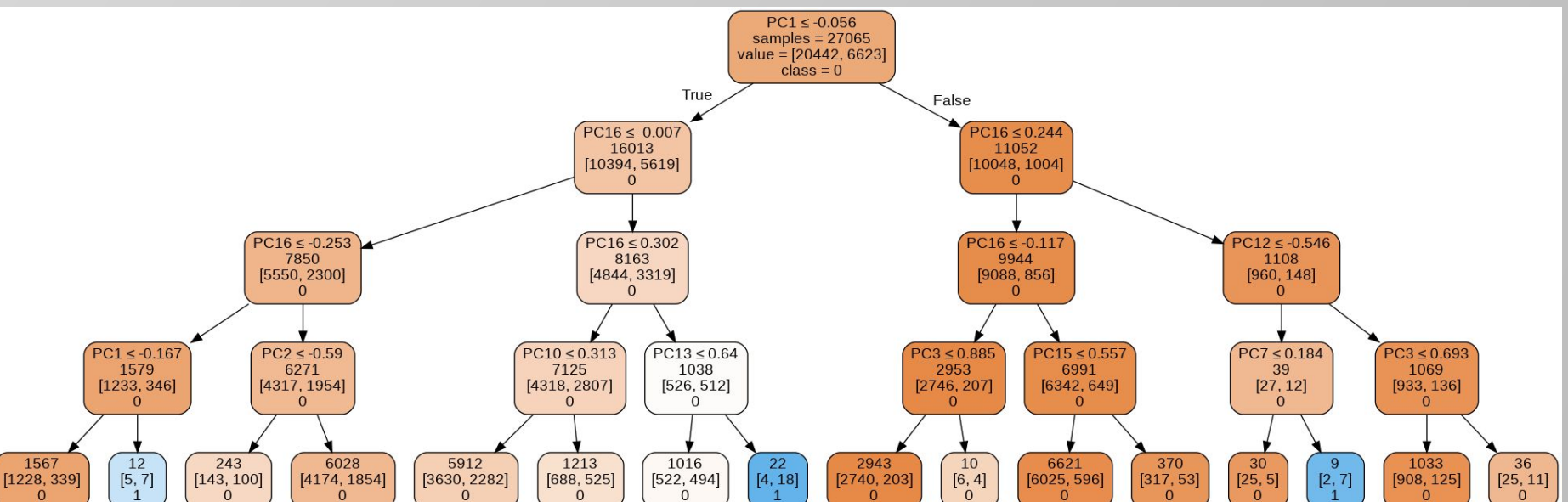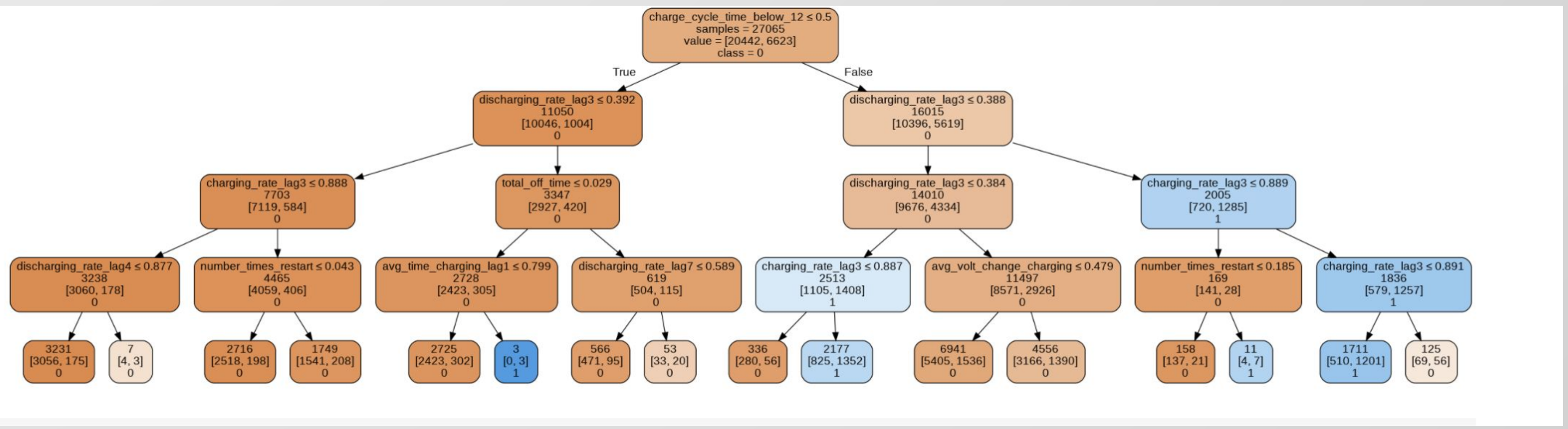

PCA

```
Confusion Matrix (Accuracy 0.7546)

        Prediction
Actual    0      1
  0    20054   388
  1     6253   370

Accuracy: 0.7591251662479681
Precision score: 0.4819277108433735
Recall score: 0.09888751545117429
F1 score: 0.1641025641025641
```

## Fit a Decision Tree using the normalized data(left) and PCA's(right)



## Summary Table for Evaluation of Models

| Model | | KNN with K = 8 | | Decision Tree | | Logistic Regression | | Naïve Bayes | |
|---|---|---|---|---|---|---|---|---|---|
| | Dataset | Normalized Dataset | PCA | Normalized Dataset | PCA | Normalized Dataset | PCA | Normalized Dataset | PCA |
| Metric | Accuracy_Score | 0.7477 | 0.7328 | 0.7544 | 0.7603 | 0.3829 | 0.7616 | 0.7609 | 0.7591 |
| | Precision Score | 0.4000 | 0.3304 | 0.3750 | 0.4474 | 0.2727 | 0.5243 | 0.3591 | 0.4819 |
| | Recall Score | 0.1100 | 0.1143 | 0.0408 | 0.0105 | 0.9487 | 0.0334 | 0.7053 | 0.0989 |
| | F1-score | 0.1726 | 0.1699 | 0.0736 | 0.0205 | 0.4237 | 0.0628 | 0.4759 | 0.1641 |