David M. Arnold

Data Science – Winter 2021 Cohort

04/06/2021

# USING MACHINE LEARNING TO DETERMINE WHAT MAKES A REDDIT MEME POPULAR

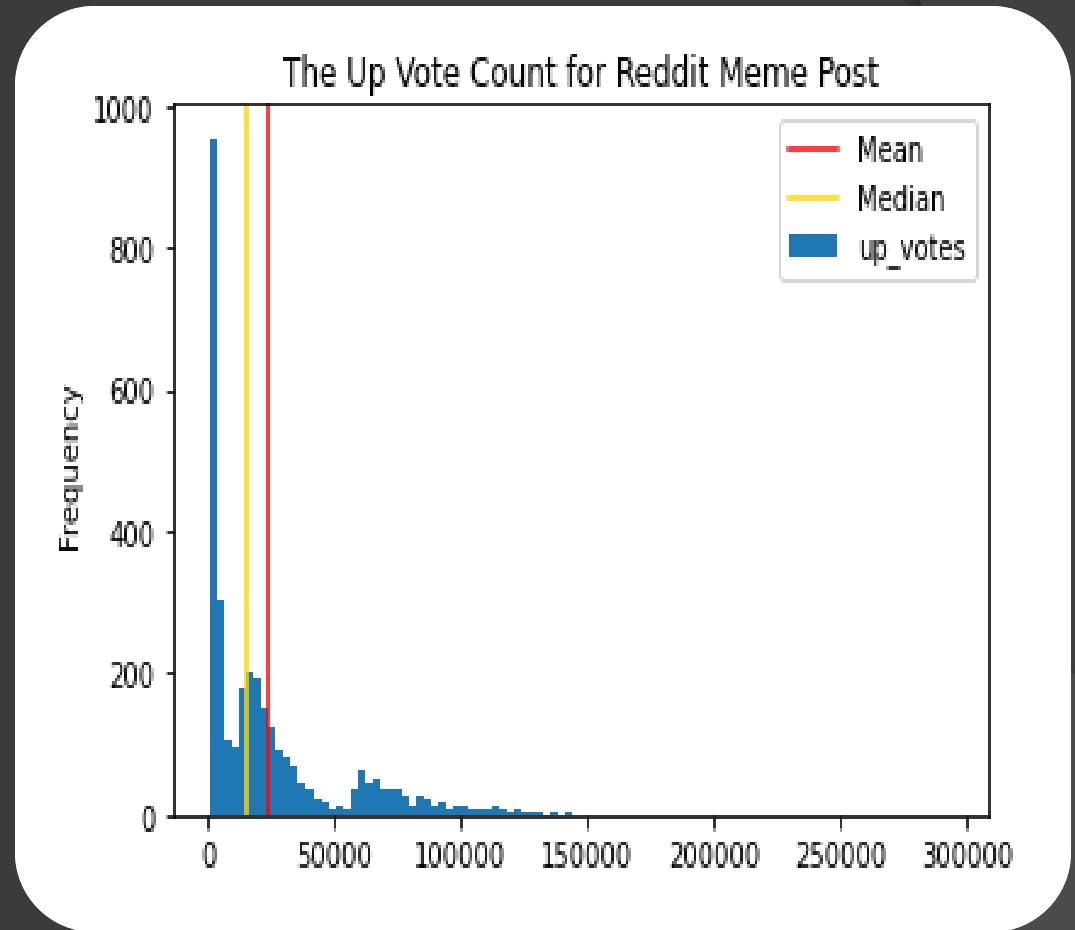# Overview

- Project Introduction / Background
- Target / Features
- Machine Learning Model Results
- Image Classification Results
- Summary & Recommendations

# Project Introduction /Background

- What Makes a Popular Meme On Reddit?
  - Content
    - Sentiment / Feeling
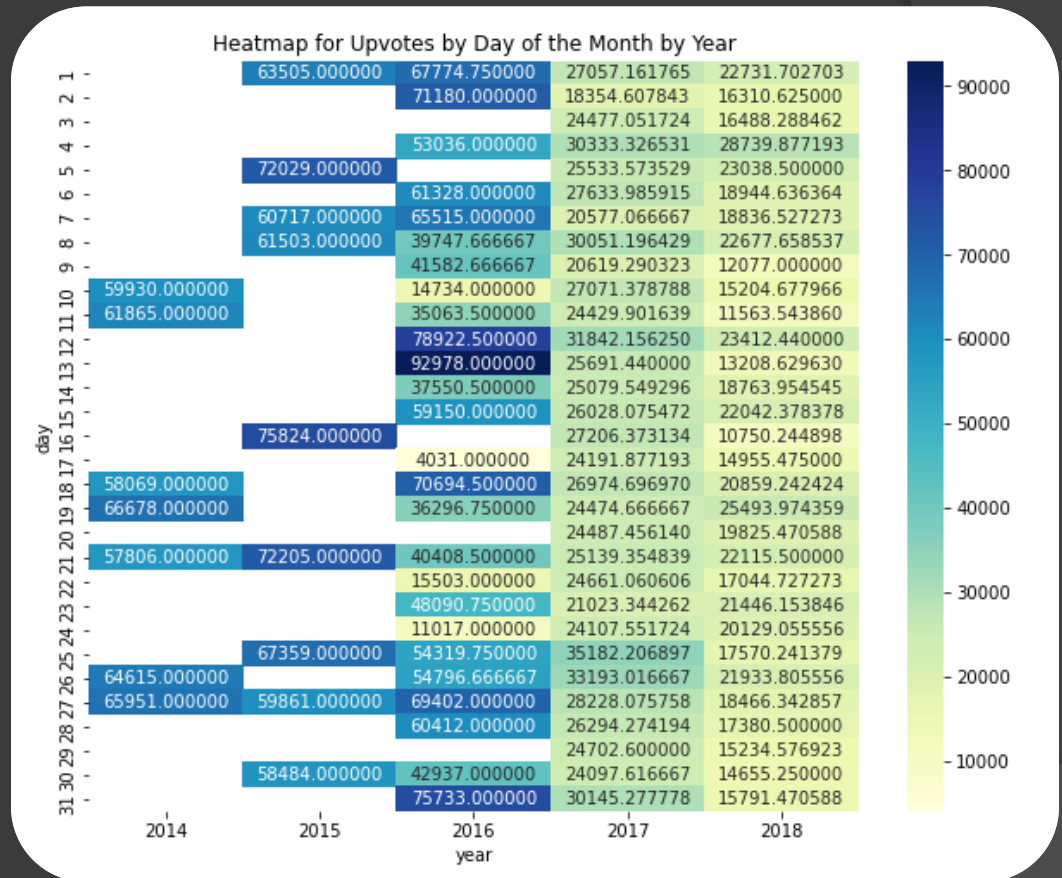  - User Submitting
  - Date Posted
  - Title / Description

# Target – Upvotes

- Upvotes
- Binary Classification
  - Good Meme
  - Over 24,000 votes
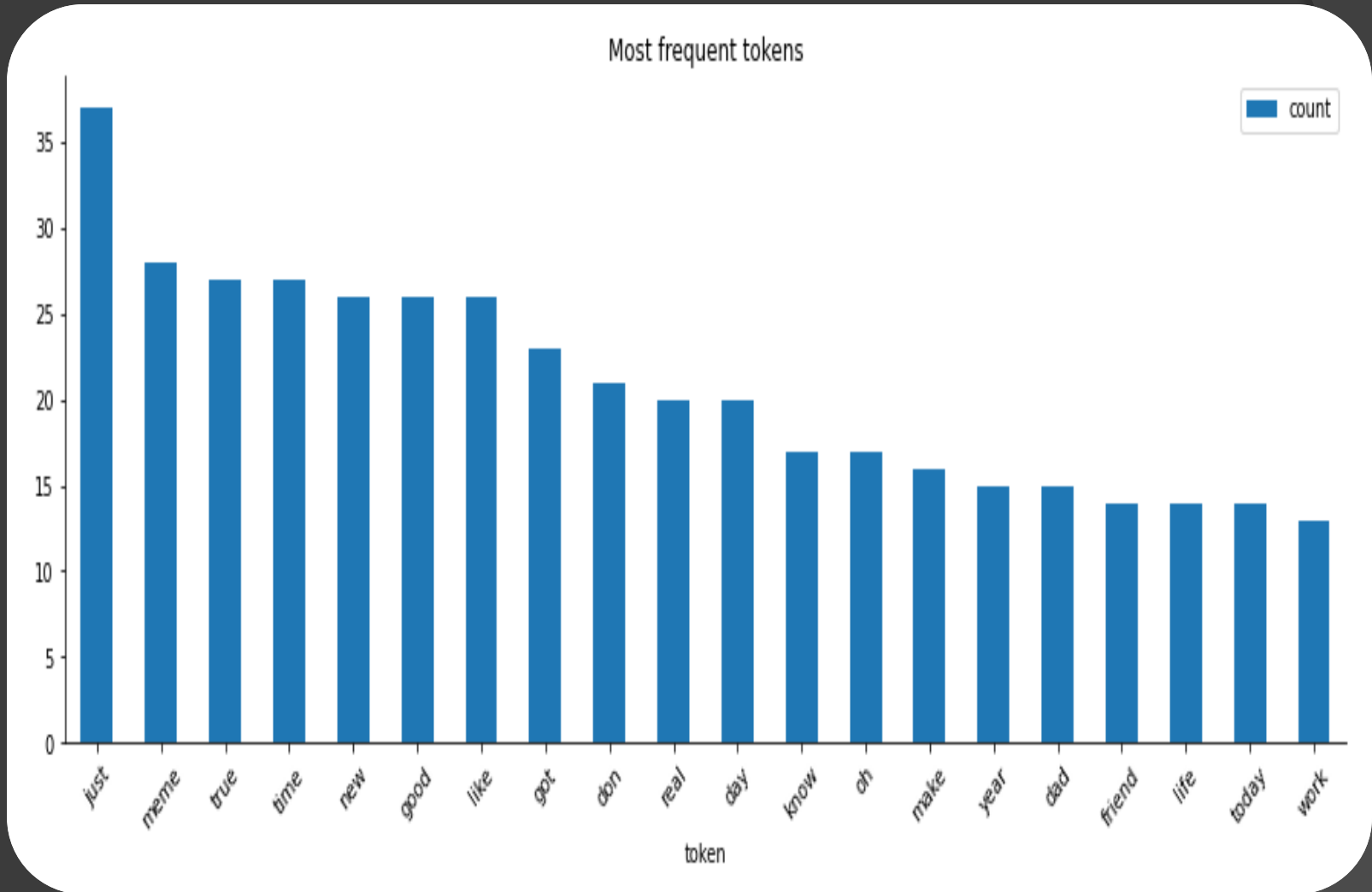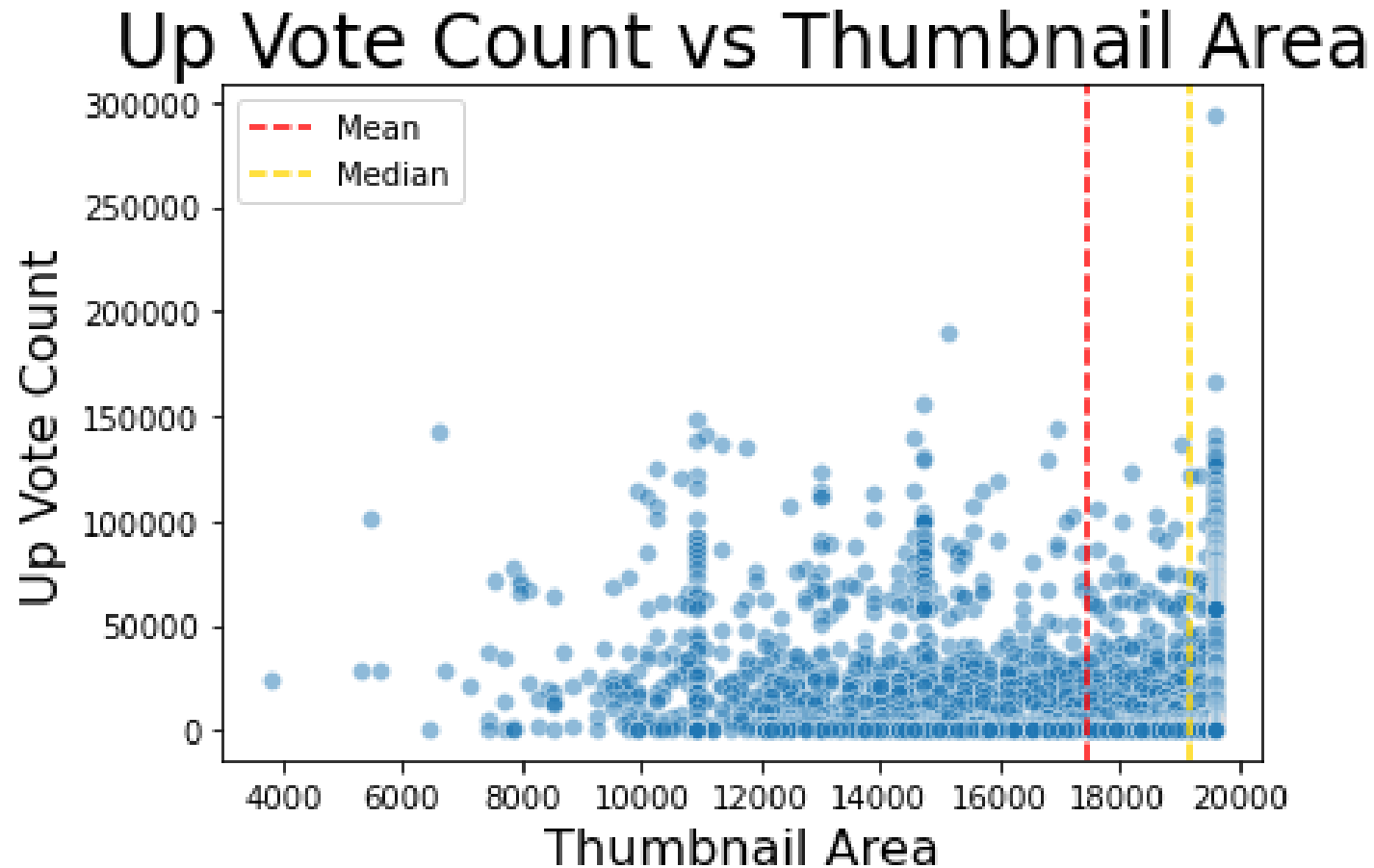  - Bad Meme
  - Under 24,000 votes

# Features – Date Posted

- ◉ Fewer Posts / More Votes
- ◉ Separate into
  - Year
  - Months
  - Weekdays
  - Weekends



Heatmap for Upvotes by Day of the Month by Year

# Features – Title Column



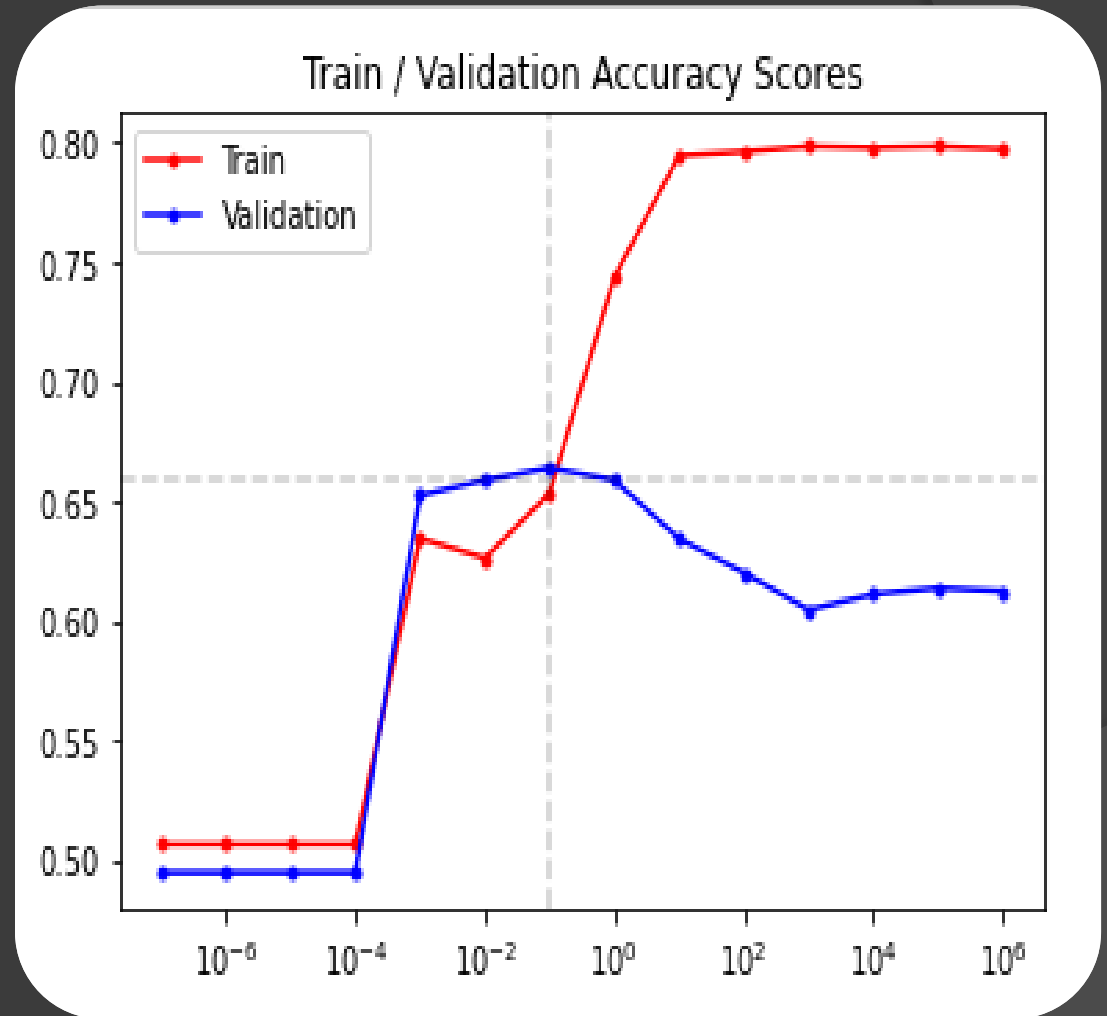Most frequent tokens

# Features – Thumbnail Area
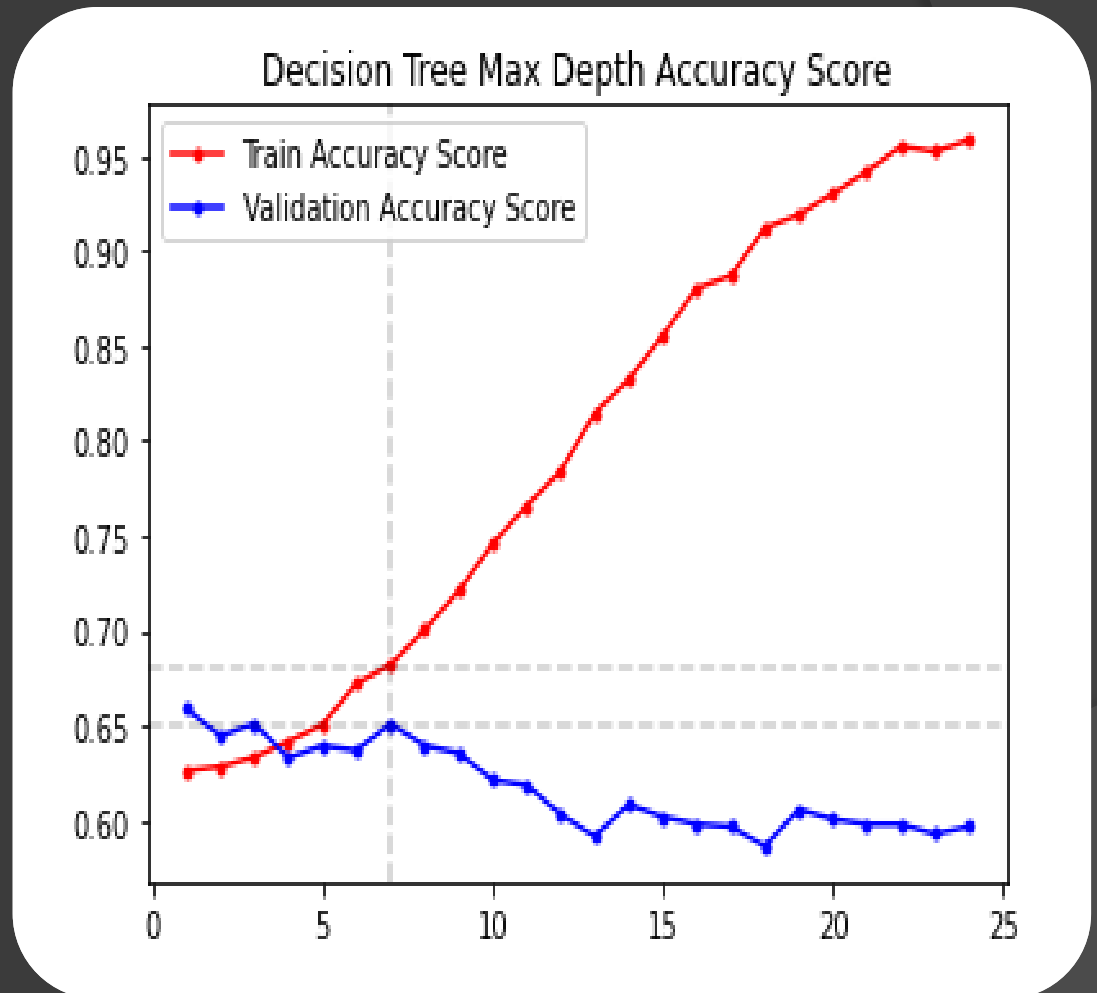


Up Vote Count vs Thumbnail Area

# ML Modeling – Logistic Regression

- "C" value – 0.1
- MinMax Scaler
- PCA 10 features
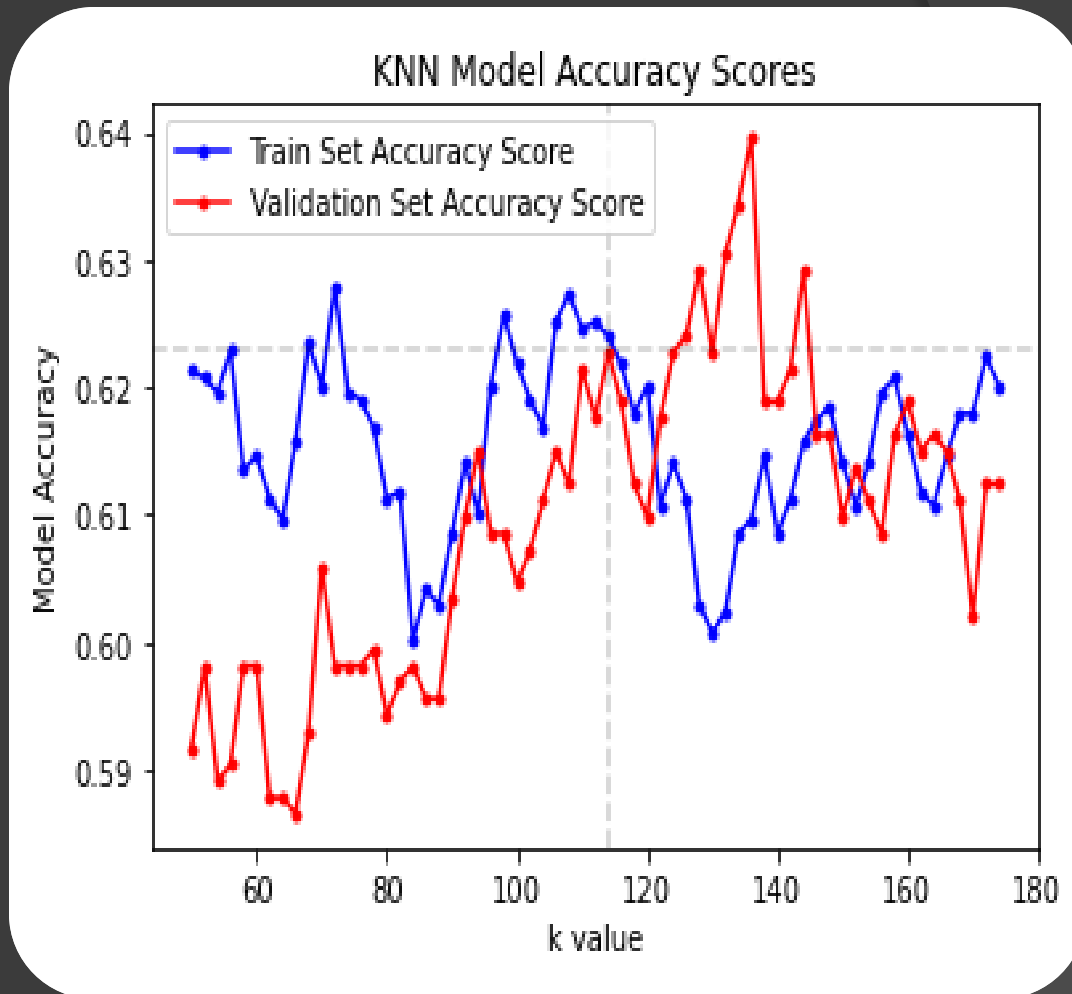- Accuracy Score
  - Train 65.3%
  - Validation 65%
  - Test 66.4%



Train / Validation Accuracy Scores

# ML Modeling – Decision Trees

- ◎ Max Depth = 7
- ◎ No scaling/PCA
- ◎ Accuracy Score
  - Train 68.2%
  - Validation 65%
  - Test 65.2%

Decision Tree Max Depth Accuracy Score

- Train Accuracy Score
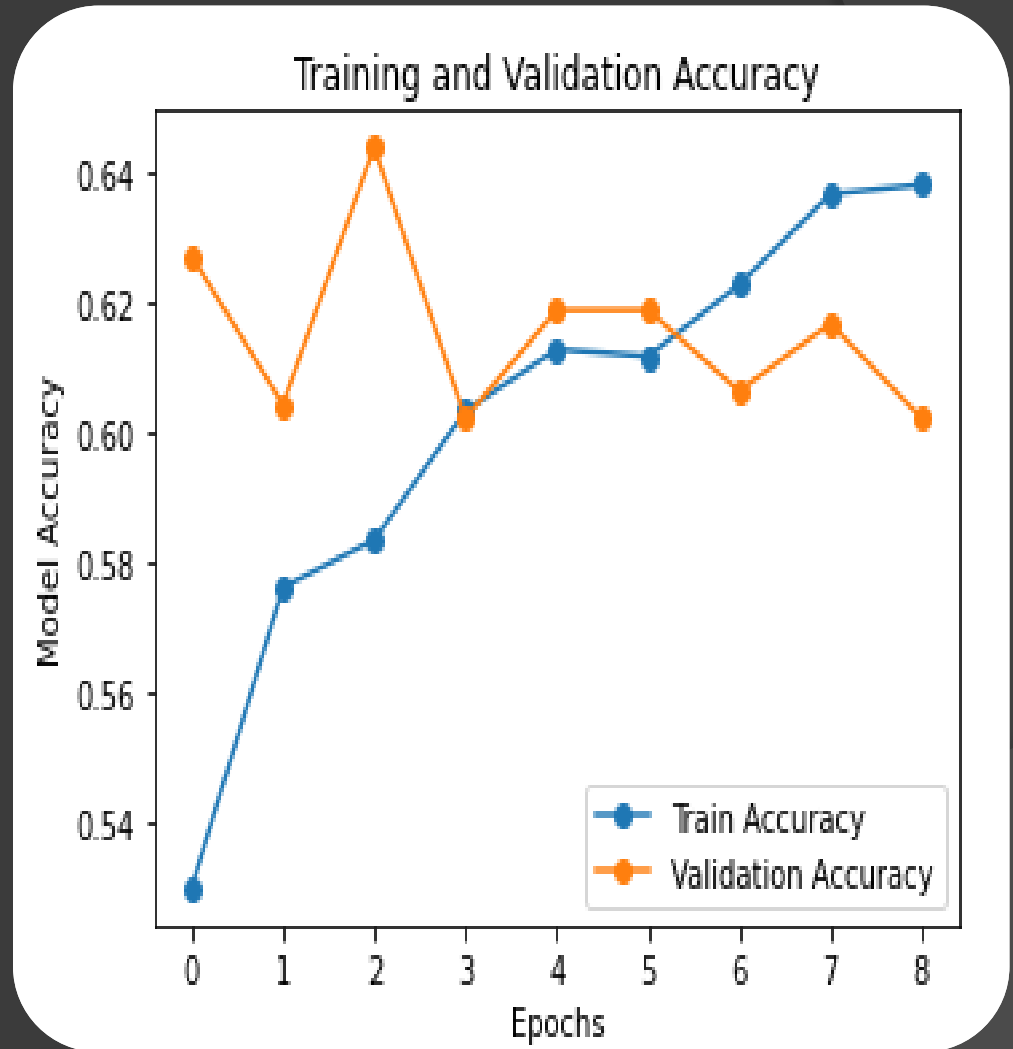- Validation Accuracy Score

# ML Modeling – K-Nearest Neighbors

- "k" value – 66
- MinMax Scaler
- Accuracy Score
  - Train 61.6%
  - Validation 62%
  - Test 59.1%

# ML Modeling – Convolutional Neural Network (CNN)

- ⦿ Classification
  - Good / Bad
- ⦿ Pre-Trained Model
  - VGG16
- ⦿ Pixel Size
  - 64 x 64
  - 256 x 256
- ⦿ Accuracy
  - 62-65%



Training and Validation Accuracy

# Summary

- ML / Neural Networks
  - 59 – 68% Test Accuracy
  - Decision Trees Best Results
- Need Confusion Matrices
- Sample / Predict Images
- Add More Datapoints
- Prediction on Newer Memes
- Share on Kaggle