

# Predicting Subreddits

Wesley Osborne

# Introduction



# Problem Statement

Maker Faire Conference wants to do a fun interactive app for their attendees at the upcoming conference. For the attendees who are fans of Arduinos and Raspberry Pi they want to build a classification model that will identify which person is a fan of either device based on the text they enter into the app. The goal of the model is to be as accurate as possible. The hope for this project is to delight the attendees and retain attendance for future conferences.



# Data Source

r/raspberry\_pi

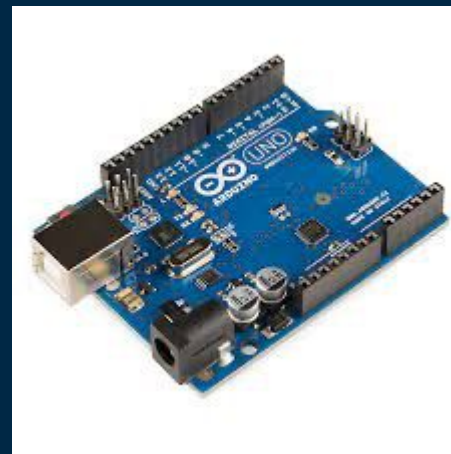


[image source](#)



[image source](#)

r/arduino

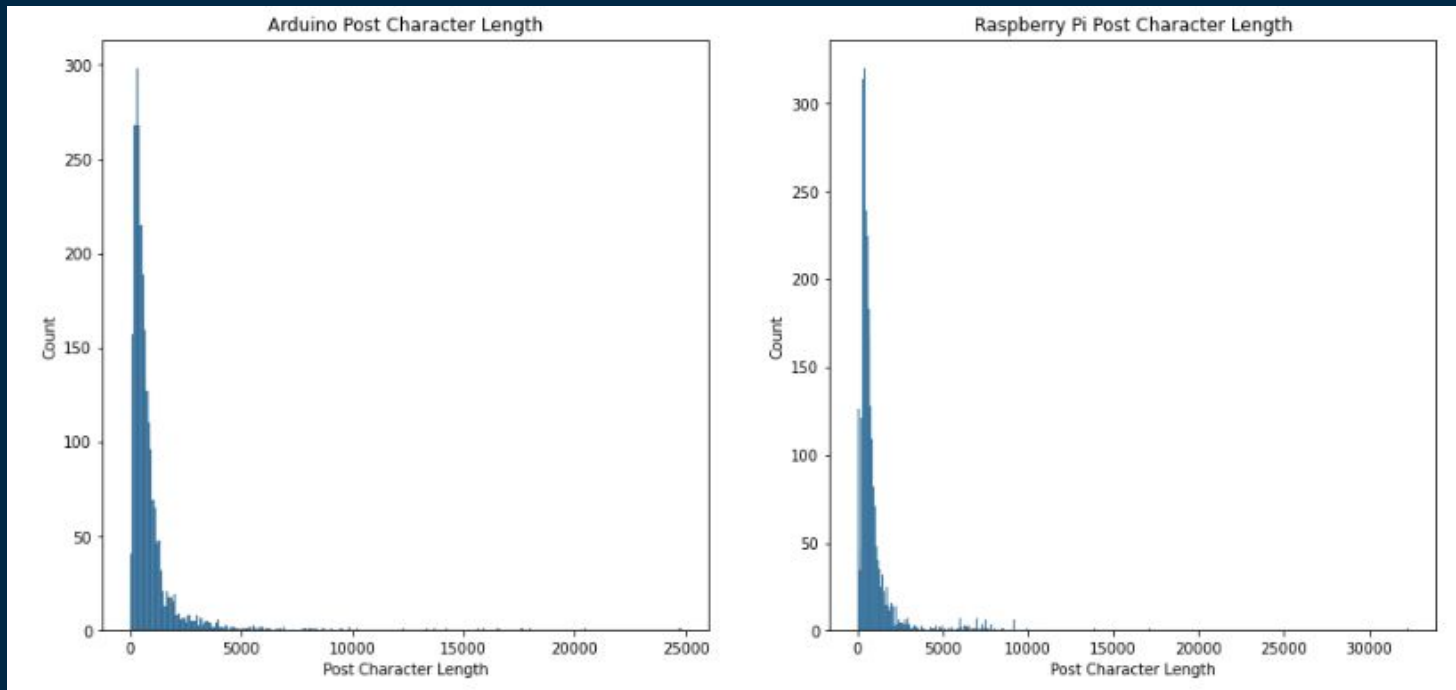


[image source](#)

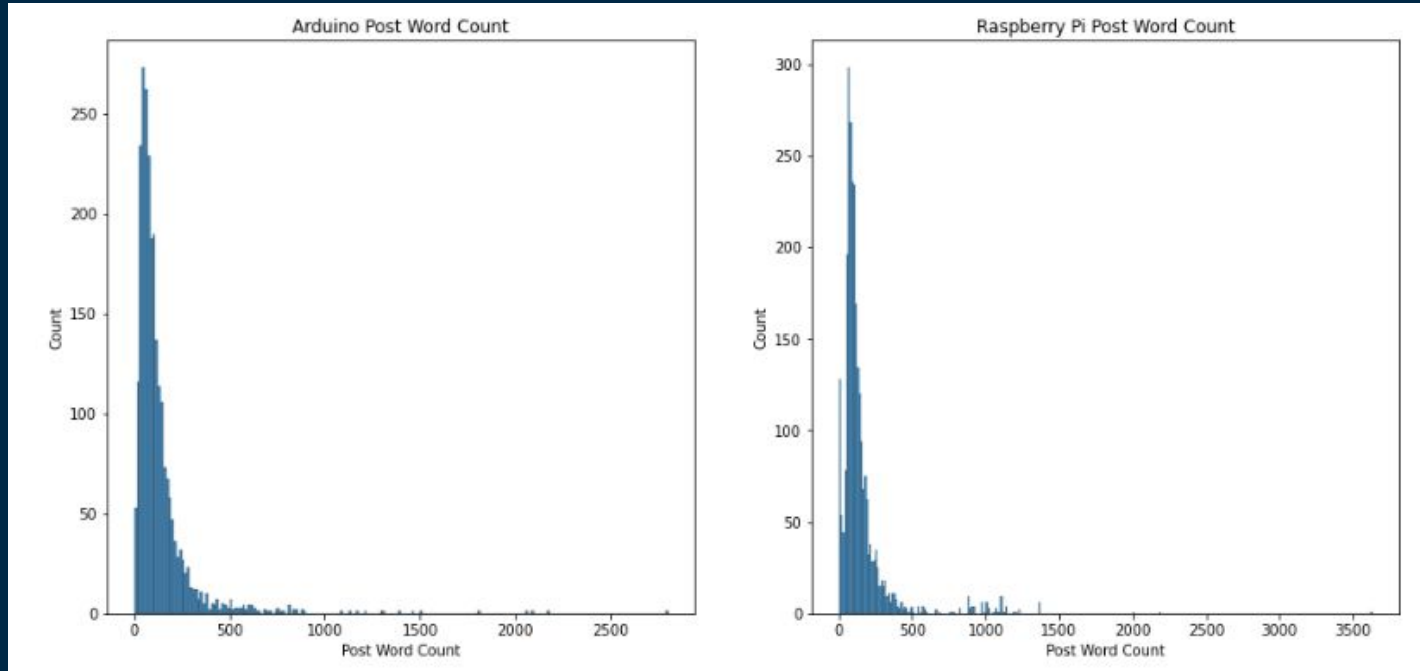
The background is a dark blue gradient. It is decorated with various geometric elements: thin white vertical lines of varying lengths, small solid squares in teal, orange, and pink, and thin white square outlines. These elements are scattered across the slide, with some appearing to be connected by thin lines, creating a modern, abstract aesthetic.

# Exploratory Data Analysis

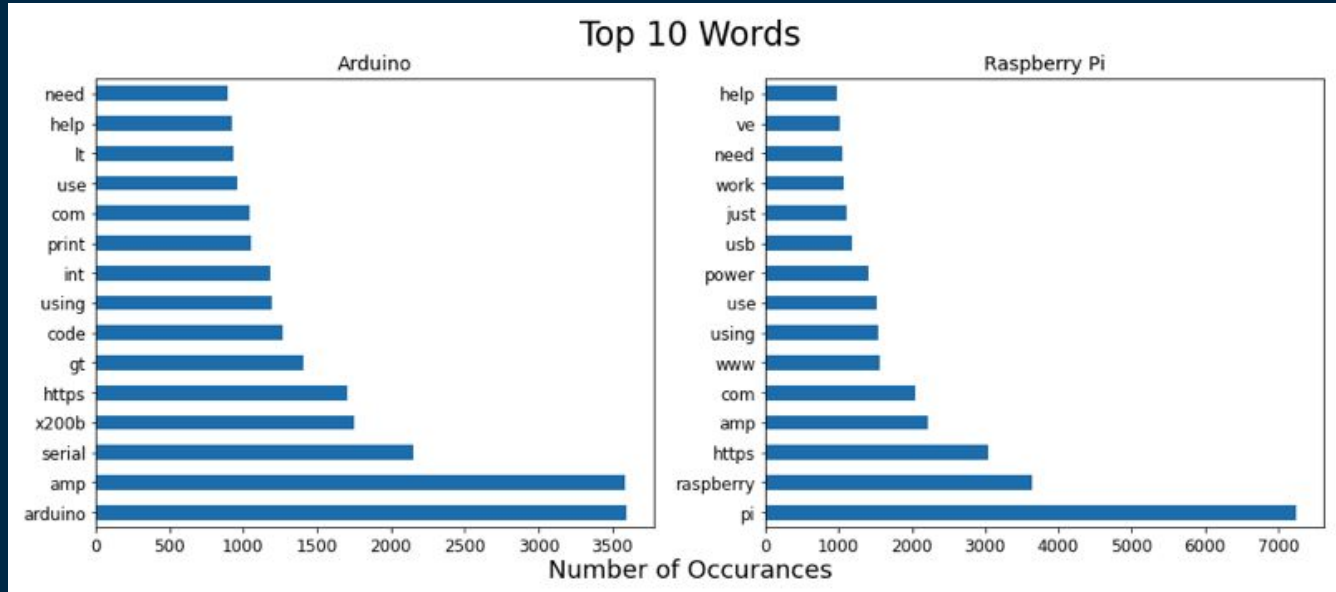
# Insights



# Insights



# Insights





# Model Workflow

The background is a dark blue gradient. It is decorated with various geometric elements: thin white vertical lines of varying lengths, small squares in teal, orange, and pink, and larger squares with thin orange outlines. These elements are scattered across the slide, creating a modern, abstract aesthetic.

# 51%

Baseline Accuracy



# Model 1: Count Vectorizer

## Random Forest Classifier

### Best Parameters

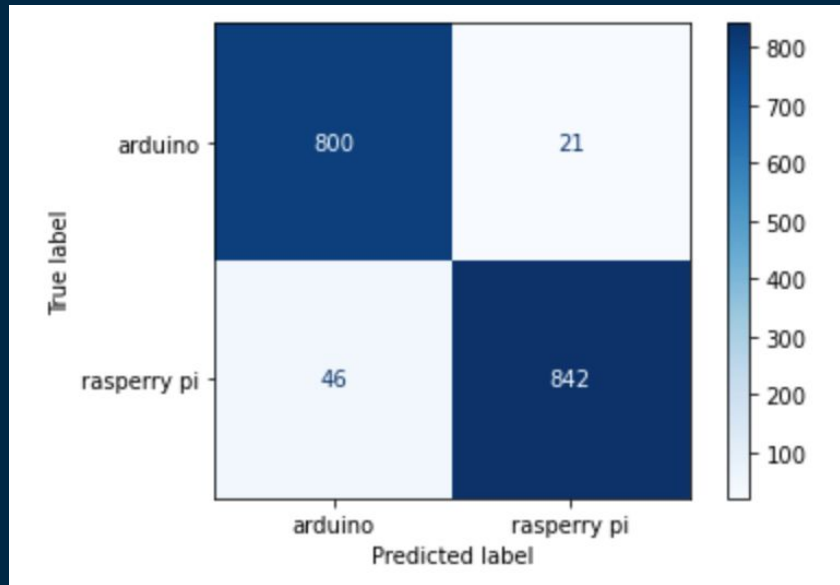
- Cvec max df: 0.9
- Cvec max features: 2000
- Cvec min df: 3
- Cvec ngram range: (1, 1)
- Rf max depth: None
- Rf max features: 'sqrt'
- Rf n estimators: 100
- Best Score: 95.39%

### Metrics

#### Accuracy Score

96.08%

# Confusion Matrix



Total Errors

67

# Model 2: Count Vectorizer

## K-Nearest Neighbors Classifier

### Best Parameters

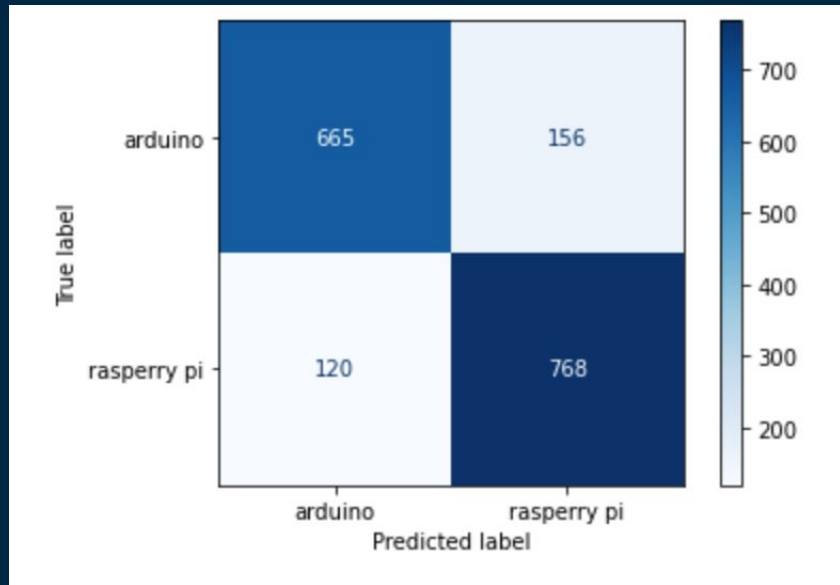
- Cvec max df: 0.9
- Cvec max features: 2000
- Cvec min df: 2
- Cvec ngram range: (1, 2)
- Knn n neighbors: 5
- Knn p: 2
- Knn weights: distance
- Best Score: 81.43%

### Metrics

#### Accuracy Score

83.85%

# Confusion Matrix



Total Errors

276

# Model 3: TF-IDF Vectorizer

## Random Forest Classifier

### Best Parameters

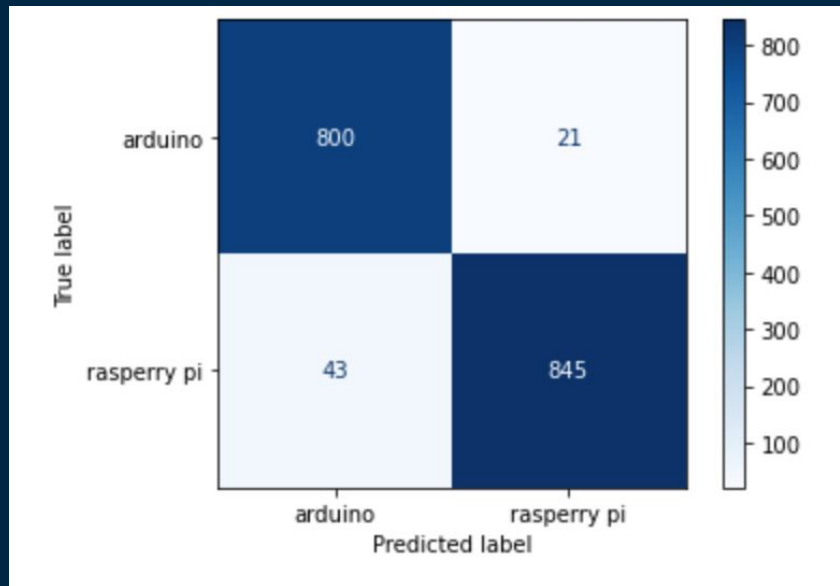
- Tvec max df: 0.9
- Tvec max features: 3000
- Tvec min df: 3
- Tvec ngram range: (1, 1)
- Rf max depth: None
- Rf max features: 'sqrt'
- Rf n estimators: 100
- Best Score: 95.27%

### Metrics

#### Accuracy Score

96.26%

# Confusion Matrix



Total Errors

64



# Model 4: TF-IDF Vectorizer

## K-Nearest Neighbors Classifier

### Best Parameters

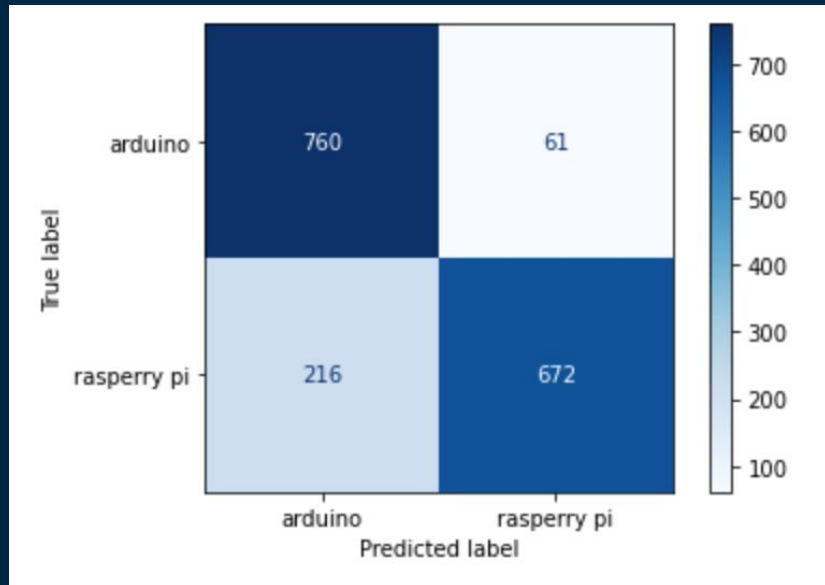
- Tvec max df: 0.9
- Tvec max features: 3000
- Tvec min df: 2
- Tvec ngram range: (1, 1)
- Knn n neighbors: 5
- Knn p: 2
- Knn weights: distance
- Best Score: 84.16%

### Metrics

#### Accuracy Score

83.80%

# Confusion Matrix



Total Errors

277

# Recommendation

## Random Forest Classifier: TF-IDF Vectorizer

### Best Parameters

- Tvec max df: 0.9
- Tvec max features: 3000
- Tvec min df: 3
- Tvec ngram range: (1, 1)
- Rf max depth: None
- Rf max features: 'sqrt'
- Rf n estimators: 100
- Best Score: 95.27%

### Metrics

#### Accuracy Score

96.26%

# Further Investigations

- AdaBoost
- Gradient Boost
- XGBoost

The background is a dark blue gradient. It is decorated with various geometric elements: thin white vertical lines of varying lengths, small squares in teal, orange, and pink, and larger squares in teal and orange. The word "Questions?" is centered in a white, sans-serif font, with the question mark being a larger, orange-colored character.

# Questions?