# Data Analysis of Pooled Dataset

We carried out dataset analysis of pooled dataset as follows.
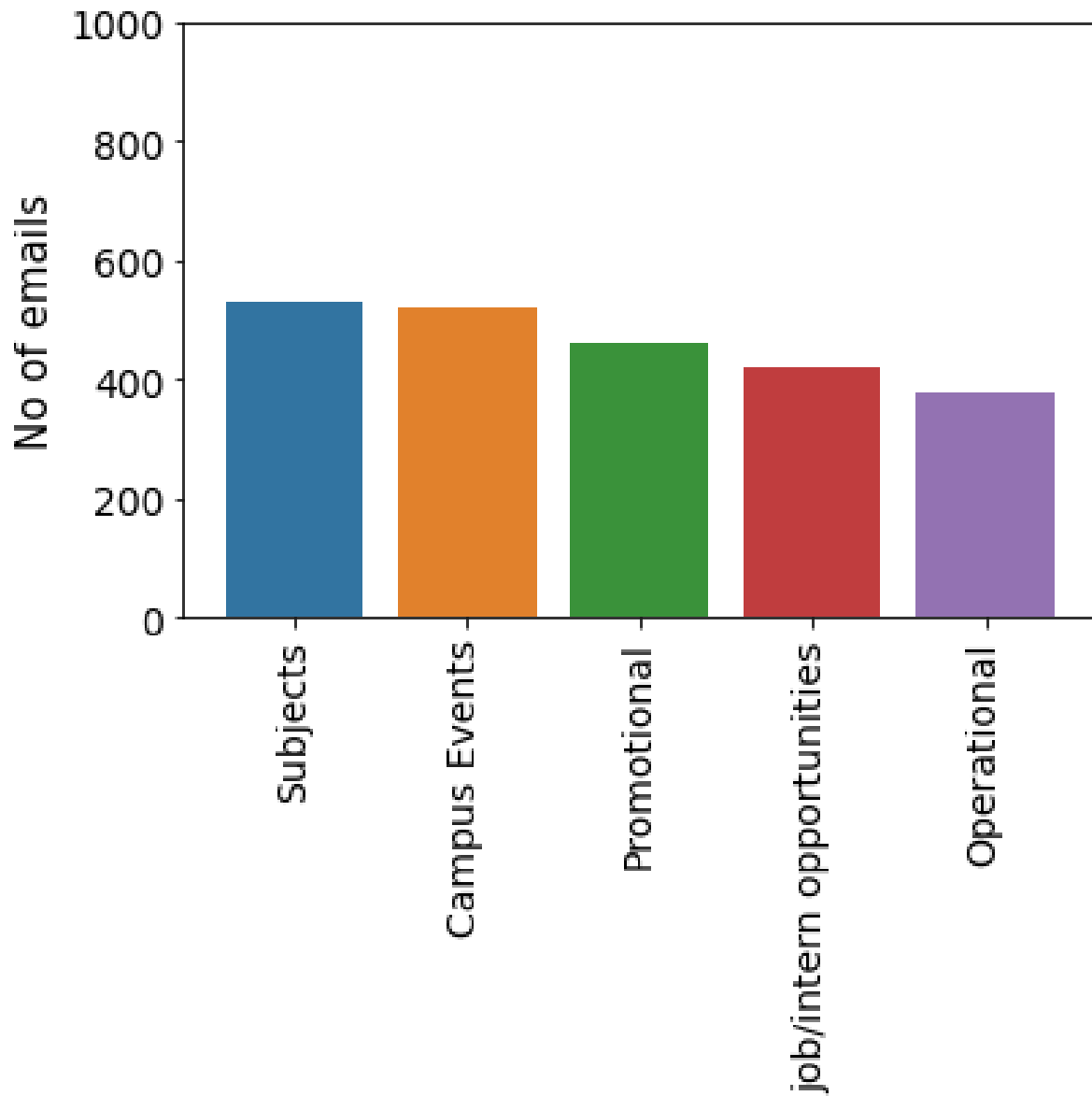


*Fig. 1 Number of Emails for each label*
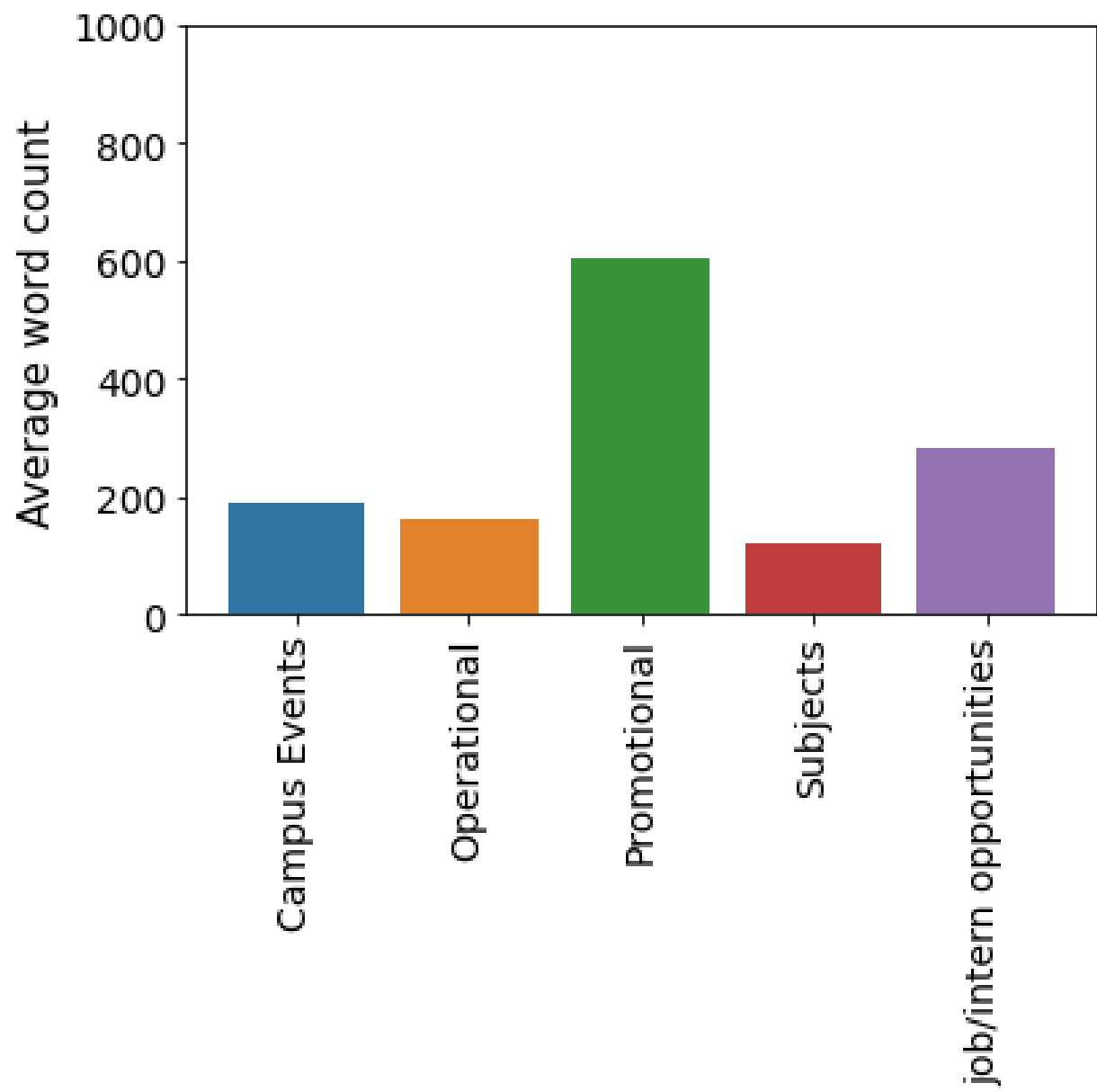
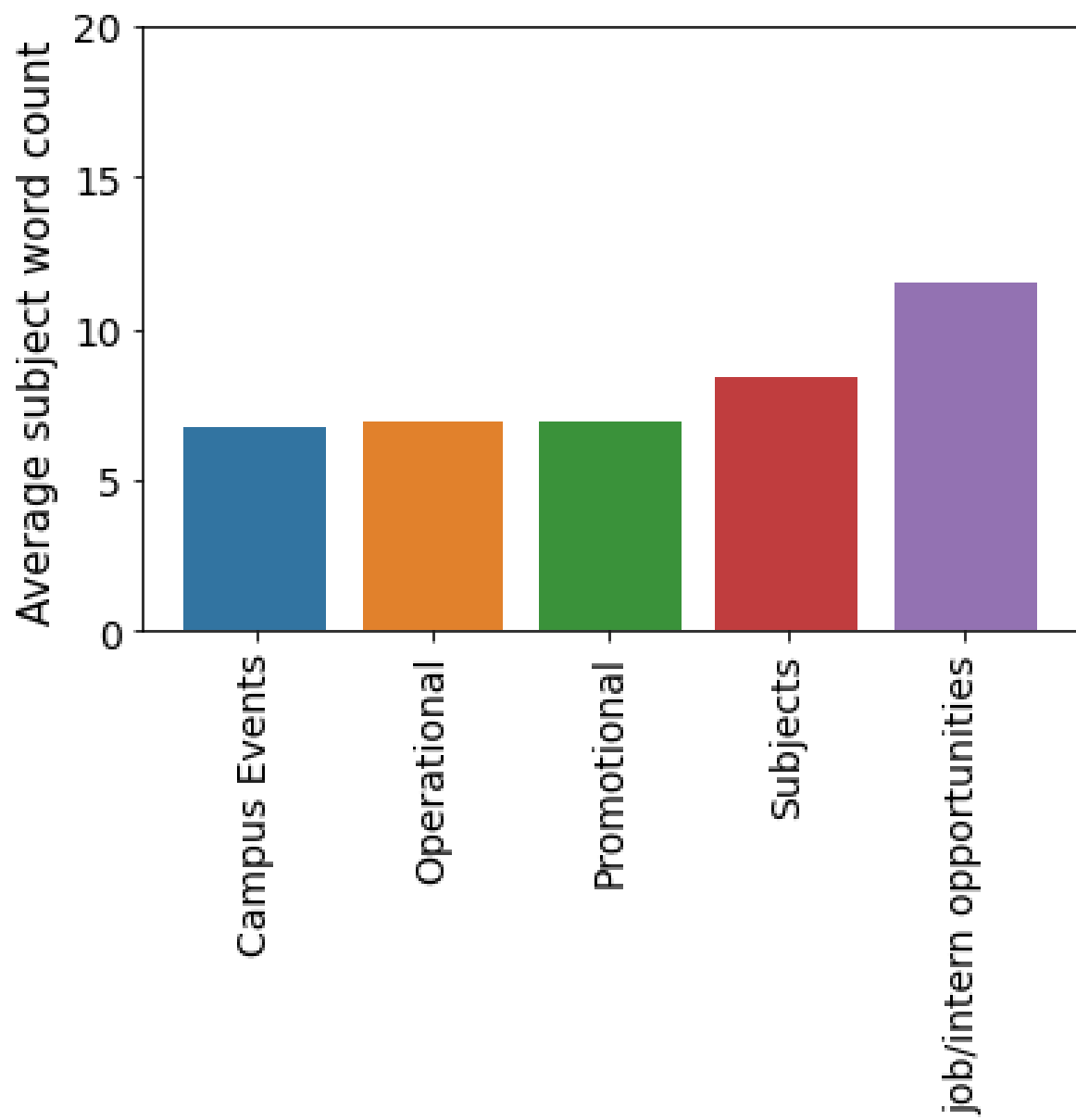*Fig. 2 Average word count for each label*

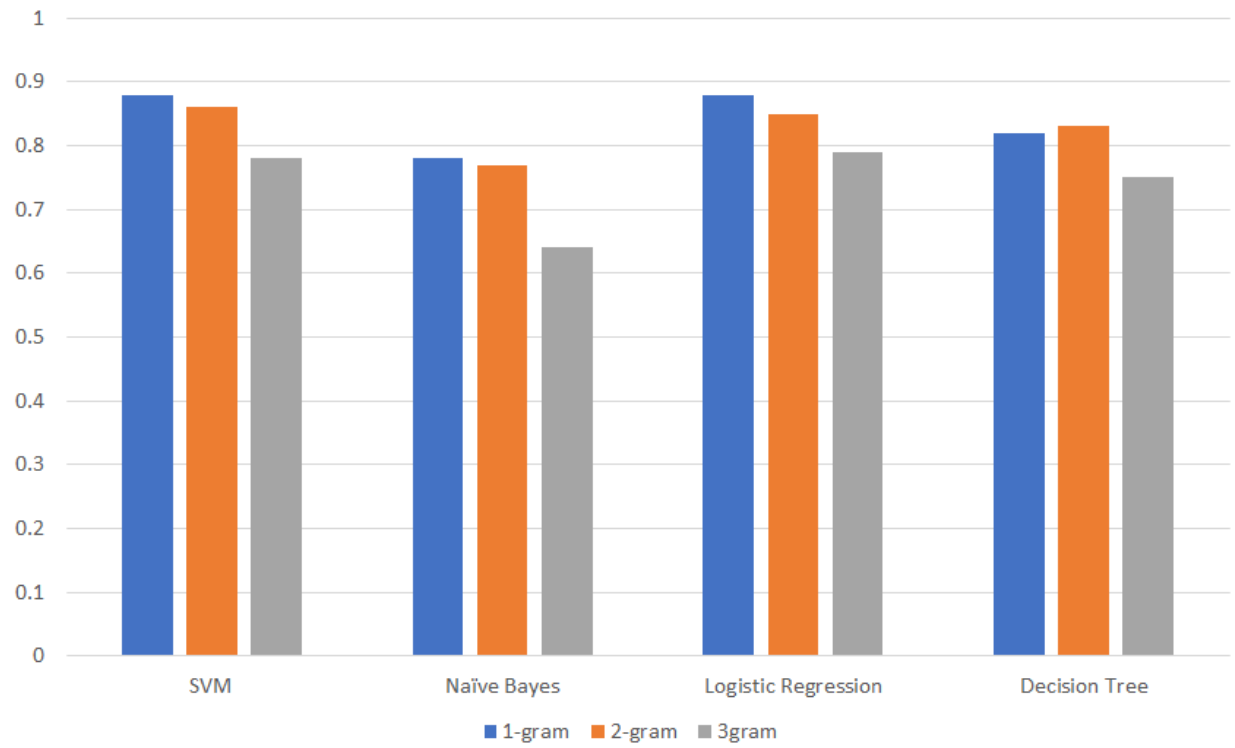*Fig. 3 Average subject word count for each label*

# Comparison Study

Here we compare different Machine Learning algorithms using 3 different N-gram language models.

Point to note regarding LDA, we tried it but the system crashed for 2-gram and 3-gram dataset on Google Colab.
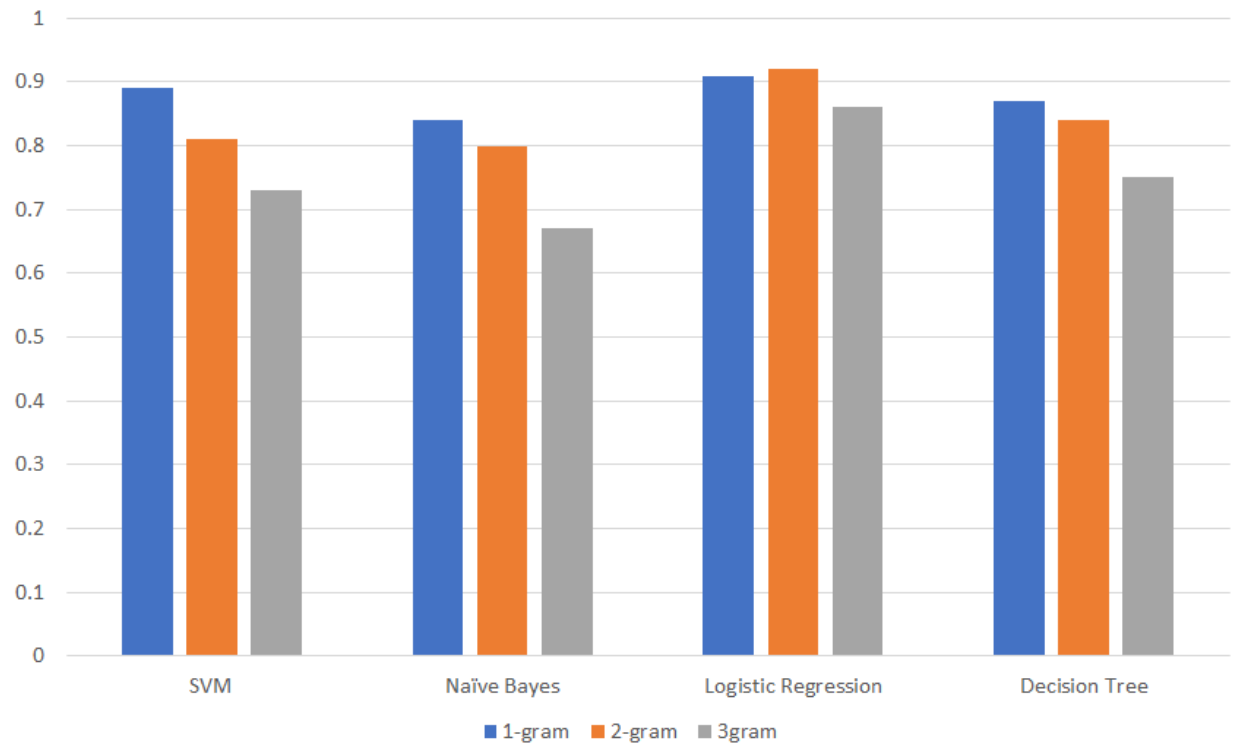
The values that follow are weighted average F1 scores.

|  | 1-gram | 2-gram | 3-gram |
| --- | --- | --- | --- |
| **Ayush Sharma** | | | |
| MultinomialNB | 0.84 | 0.83 | 0.67 |
| SVM | 0.87 | 0.87 | 0.80 |
| Logistic Regression | 0.89 | 0.86 | 0.81 |
| Decision Tree | 0.82 | 0.84 | 0.75 |
| **Naman Goenka** | | | |
| MultinomialNB | 0.84 | 0.80 | 0.67 |
| SVM | 0.89 | 0.81 | 0.73 |
| Logistic Regression | 0.91 | 0.92 | 0.86 |
| Decision Tree | 0.87 | 0.84 | 0.79 |
| **Mohul Maheshwari** | | | |
| MultinomialNB | 0.78 | 0.77 | 0.64 |
| SVM | 0.88 | 0.86 | 0.78 |
| Logistic Regression | 0.88 | 0.85 | 0.79 |
| Decision Tree | 0.82 | 0.83 | 0.75 |

Mohul Maheshwari

Legend: 1-gram, 2-gram, 3gram

Naman Goenka

Legend: 1-gram, 2-gram, 3gram

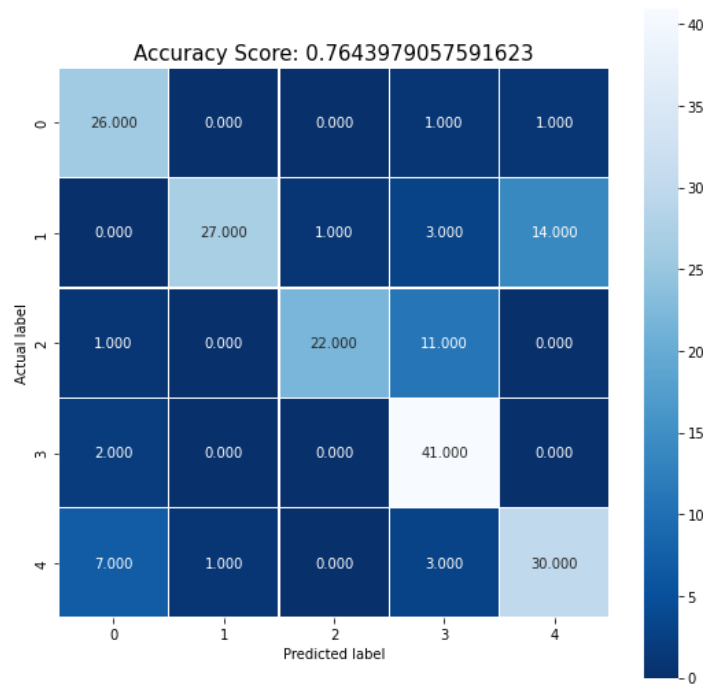SVM — Naïve Bayes — Logistic Regression — Decision Tree

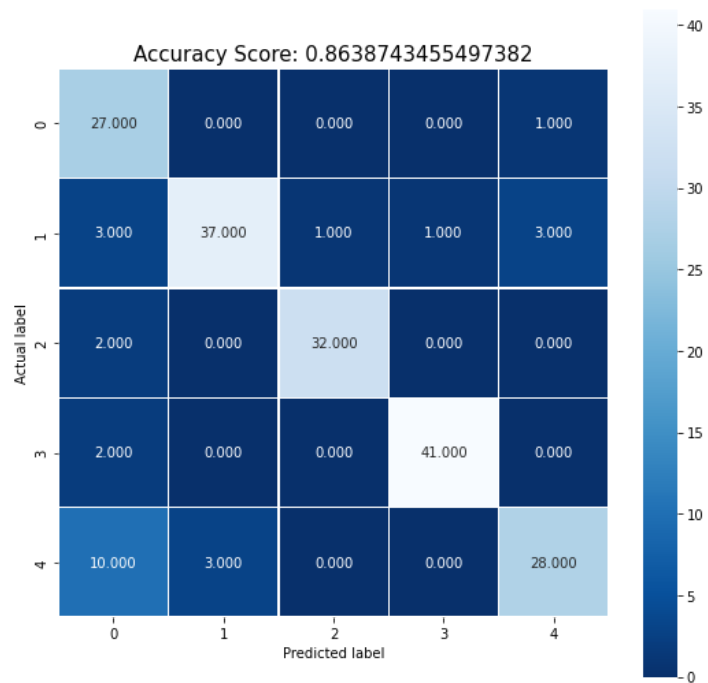Ayush Sharma

# Data_1

**MultinomialNB**

Unigram



Bigram

## Trigram



## **SVM**

## Unigram

## Bigram



Accuracy Score: 0.8638743455497382

## Trigram



Accuracy Score: 0.7696335078534031

## Logistic Regression

### Unigram



### Bigram

## Trigram



Accuracy Score: 0.7801047120418848

## **Decision Tree**

## Unigram



Accuracy Score: 0.8324607329842932

## Bigram



Accuracy Score: 0.8324607329842932

## Trigram



Accuracy Score: 0.7539267015706806

# Data_2

**MultinomialNB**

Unigram



Bigram

# Trigram



Accuracy Score: 0.6878306878306878

# SVM

## Unigram



Accuracy Score: 0.8888888888888888

## Bigram

### Accuracy Score: 0.8306878306878307

|  | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| **0** | 14.000 | 0.000 | 0.000 | 16.000 | 0.000 |
| **1** | 1.000 | 45.000 | 0.000 | 1.000 | 0.000 |
| **2** | 2.000 | 0.000 | 27.000 | 3.000 | 0.000 |
| **3** | 0.000 | 0.000 | 0.000 | 40.000 | 0.000 |
| **4** | 4.000 | 1.000 | 0.000 | 4.000 | 31.000 |

Actual label / Predicted label

## Trigram

### Accuracy Score: 0.746031746031746

|  | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| **0** | 10.000 | 0.000 | 0.000 | 20.000 | 0.000 |
| **1** | 0.000 | 41.000 | 0.000 | 6.000 | 0.000 |
| **2** | 1.000 | 0.000 | 24.000 | 7.000 | 0.000 |
| **3** | 0.000 | 0.000 | 0.000 | 40.000 | 0.000 |
| **4** | 4.000 | 2.000 | 0.000 | 8.000 | 26.000 |

Actual label / Predicted label

## Logistic Regression

### Unigram

**Accuracy Score: 0.91005291005291**



### Bigram

**Accuracy Score: 0.9259259259259259**

## Trigram



Accuracy Score: 0.8624338624338624

## **Decision Tree**

## Unigram



Accuracy Score: 0.8783068783068783

## Bigram



Accuracy Score: 0.8571428571428571

## Trigram



Accuracy Score: 0.798941798941799

# Data 3

**MultinomialNB**

Unigram
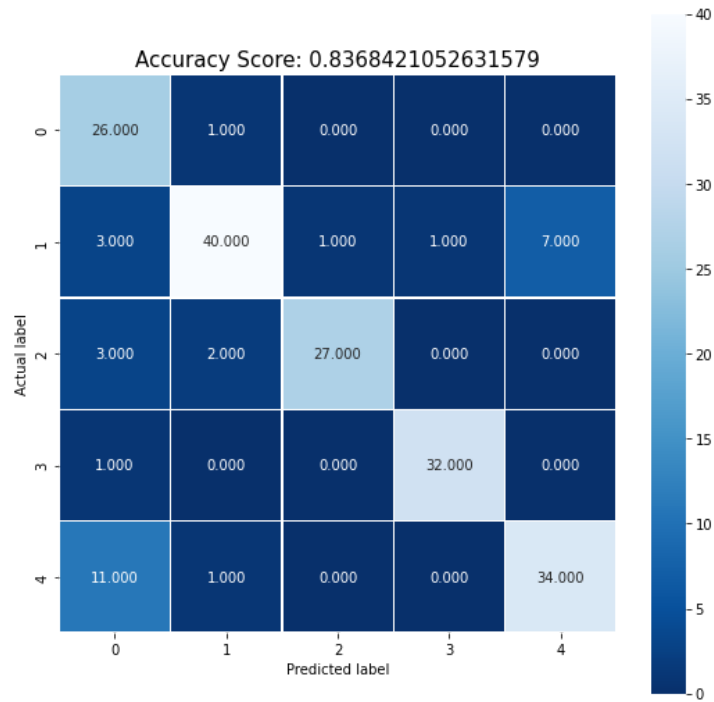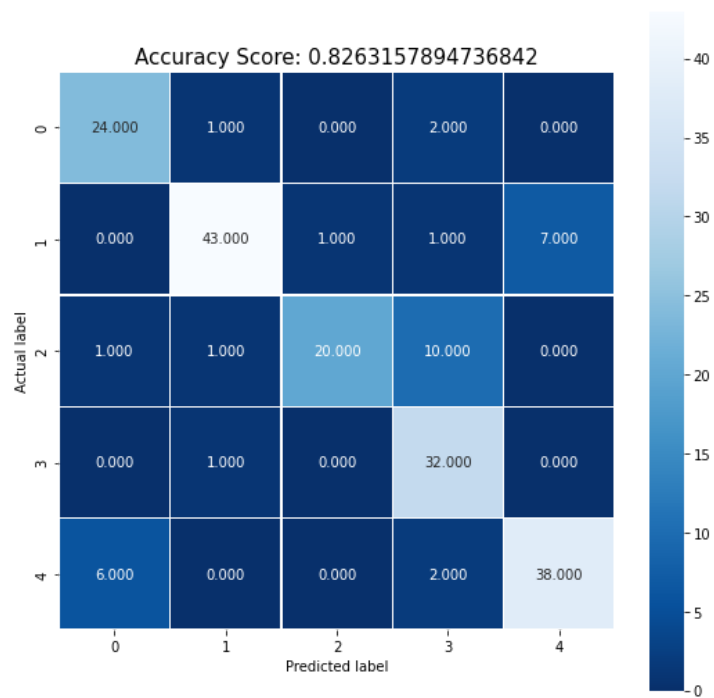


Bigram

## Trigram

### Accuracy Score: 0.6578947368421053

|   | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| 0 | 10.000 | 0.000 | 0.000 | 17.000 | 0.000 |
| 1 | 0.000 | 39.000 | 1.000 | 5.000 | 7.000 |
| 2 | 1.000 | 1.000 | 14.000 | 16.000 | 0.000 |
| 3 | 0.000 | 2.000 | 0.000 | 31.000 | 0.000 |
| 4 | 1.000 | 0.000 | 1.000 | 13.000 | 31.000 |

Actual label / Predicted label

## SVM

## Unigram

### Accuracy Score: 0.868421052631579

|   | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| 0 | 25.000 | 0.000 | 1.000 | 1.000 | 0.000 |
| 1 | 3.000 | 45.000 | 2.000 | 0.000 | 2.000 |
| 2 | 3.000 | 2.000 | 27.000 | 0.000 | 0.000 |
| 3 | 0.000 | 2.000 | 0.000 | 31.000 | 0.000 |
| 4 | 5.000 | 2.000 | 2.000 | 0.000 | 37.000 |

Actual label / Predicted label

## Bigram



Accuracy Score: 0.8631578947368421

## Trigram



Accuracy Score: 0.7842105263157895

**Linear regression**

## Unigram


Accuracy Score: 0.8842105263157894

## Bigram


Accuracy Score: 0.8526315789473684

## Trigram

Accuracy Score: 0.7947368421052632



## **Decision Tree**

## Unigram

Accuracy Score: 0.8210526315789474

## Bigram



Accuracy Score: 0.8368421052631579

## Trigram



Accuracy Score: 0.7473684210526316