# reddit Classification

# Hello!

## I am Alex Fioto

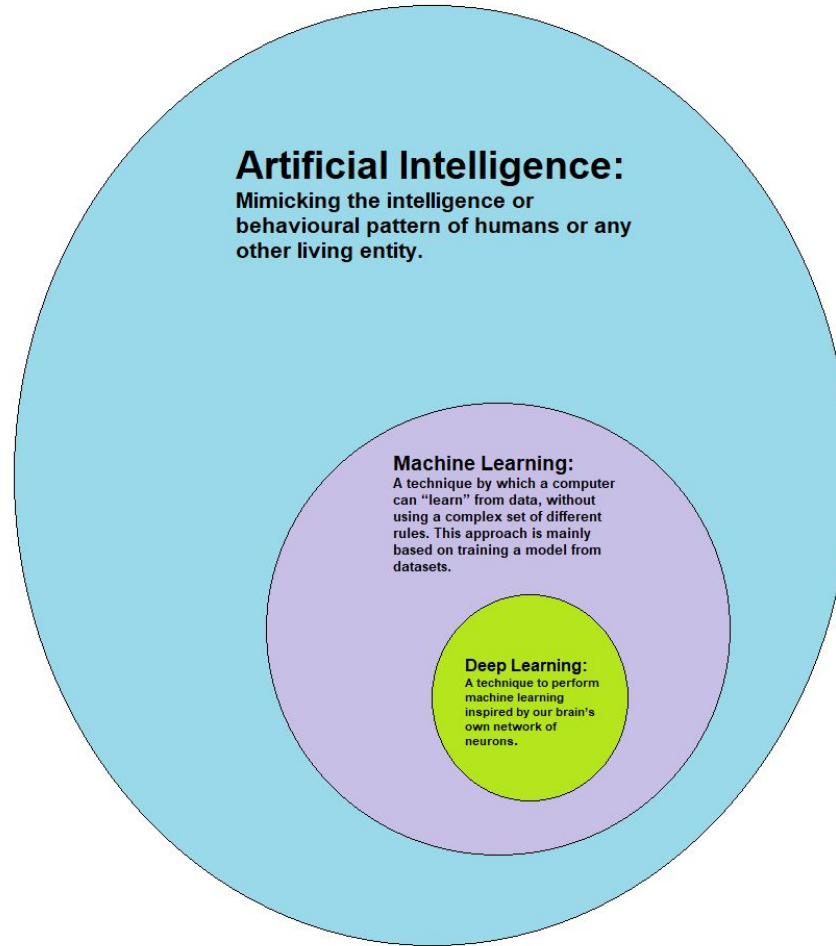I am here to talk about classifying subreddit posts using scikit-learn classifier models

# Problem:

Trained on data from two similar subreddits, how well can a model classify a Reddit post?

# Machine Learning vs Artificial Intelligence

# How similar?

**Artificial Intelligence:**
**Mimicking the intelligence or behavioural pattern of humans or any other living entity.**

**Machine Learning:**
**A technique by which a computer can "learn" from data, without using a complex set of different rules. This approach is mainly based on training a model from datasets.**

**Deep Learning:**
**A technique to perform machine learning inspired by our brain's own network of neurons.**

# Data Collection

| | Subreddit | Rows |
|---|---|---|
| Machine Learning | **r/MachineLearning/** | **9743** |
| Artificial Intelligence | **r/ArtificialInteligence/** | **9998** |

# Text Preprocessing

**Feature Creation**

Combining all text from each post to create an "all_text" feature

**Lemmatizing**

Create custom lemmating function and apply to "all_text"
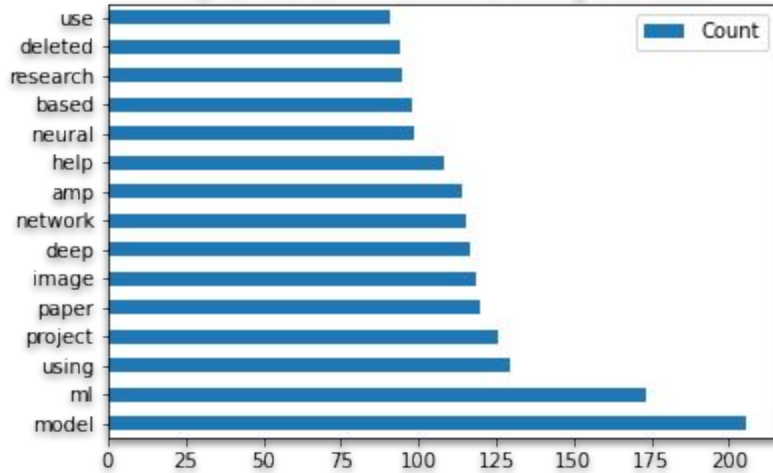
**Stemming**
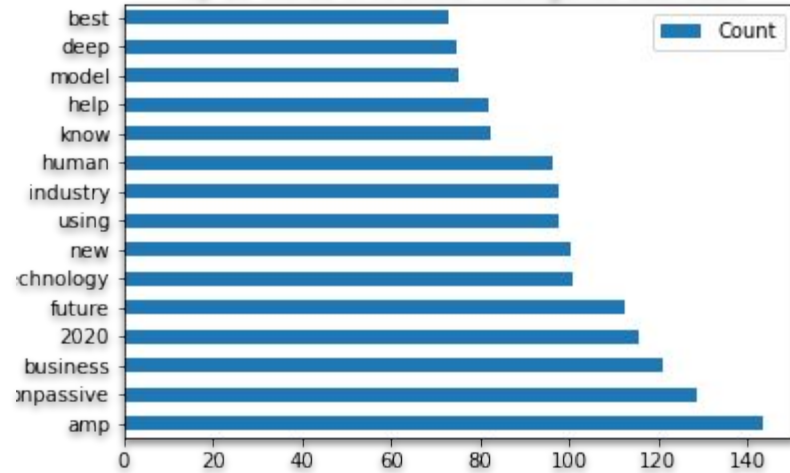
Create custom stemming function and apply to "all_text"

# Common words



Top 15 Words in Machine Learning Subreddit



Top 15 Words in Artificial Intelligence Subreddit

[www.onpassive.com](http://www.onpassive.com)

# Classifiers Used:

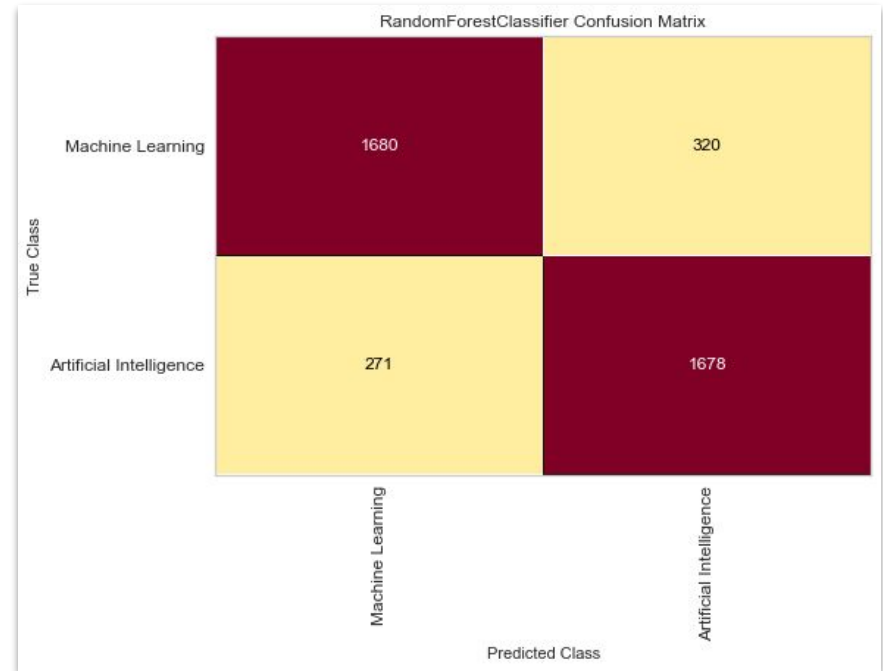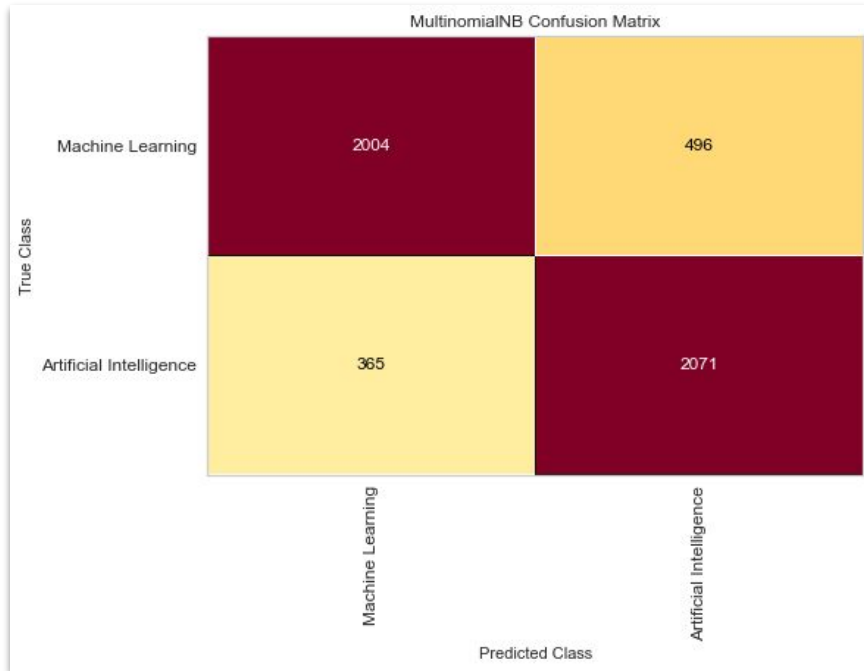**Naive-Bayes**

▷ Simple
▷ Quick to train
▷ Interpretable

**Random Forest**

▷ Powerful
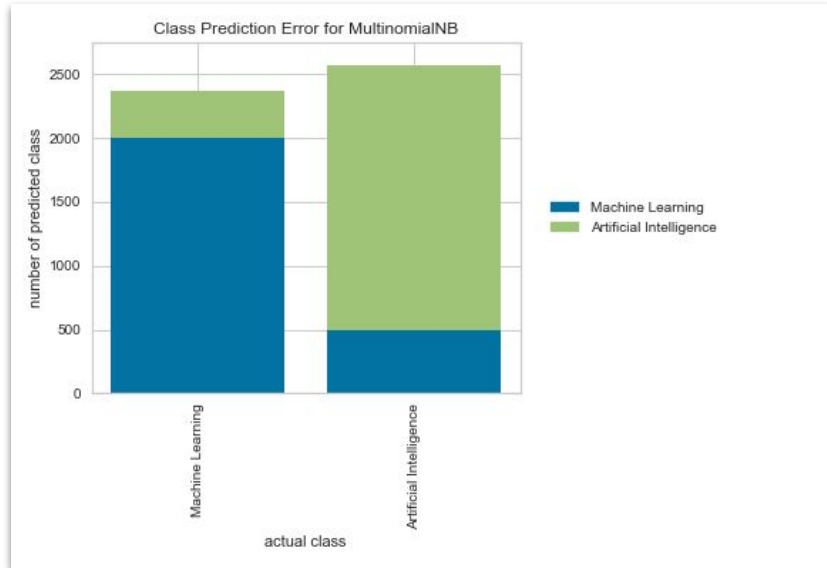▷ Interpretable
▷ Accurate

# Classification Metrics
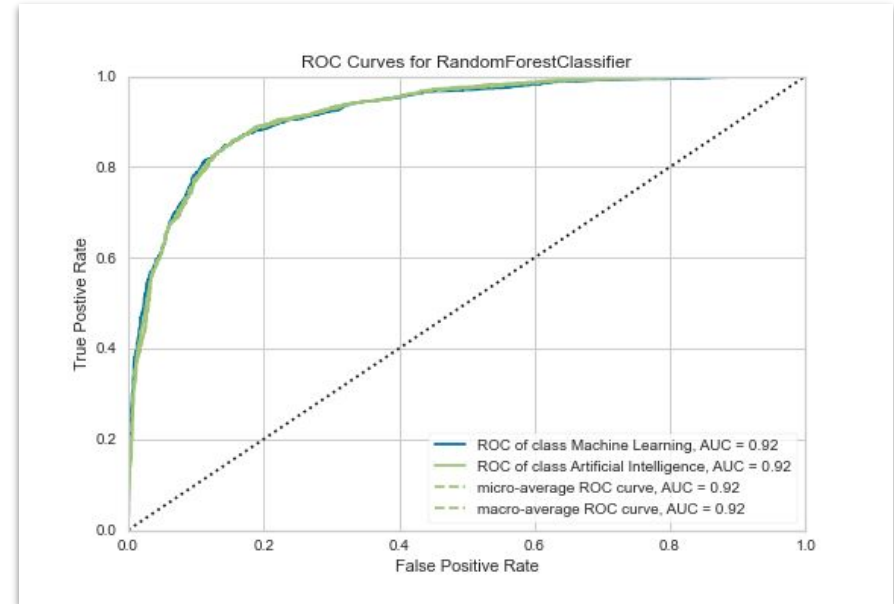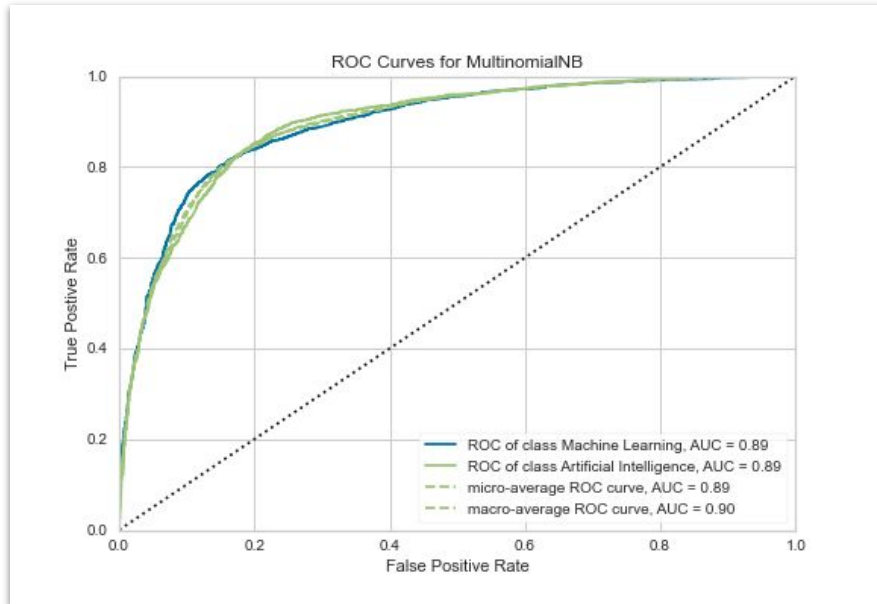
▷ Accuracy

▷ F-1 Score

▷ Confusion Matrix

# Confusion Matrices

# Class Prediction Error

# ROC Curves

# Model Results

| | Test Accuracy | Train Accuracy | Test F-1 Score | Train F-1 Score |
|---|---|---|---|---|
| Naive-Bayes | **0.8256** | **0.8321** | **0.8279** | **0.8353** |
| Random Forest | **0.8519** | **0.9851** | **0.8513** | **0.9850** |

# Hyperparameters Used:

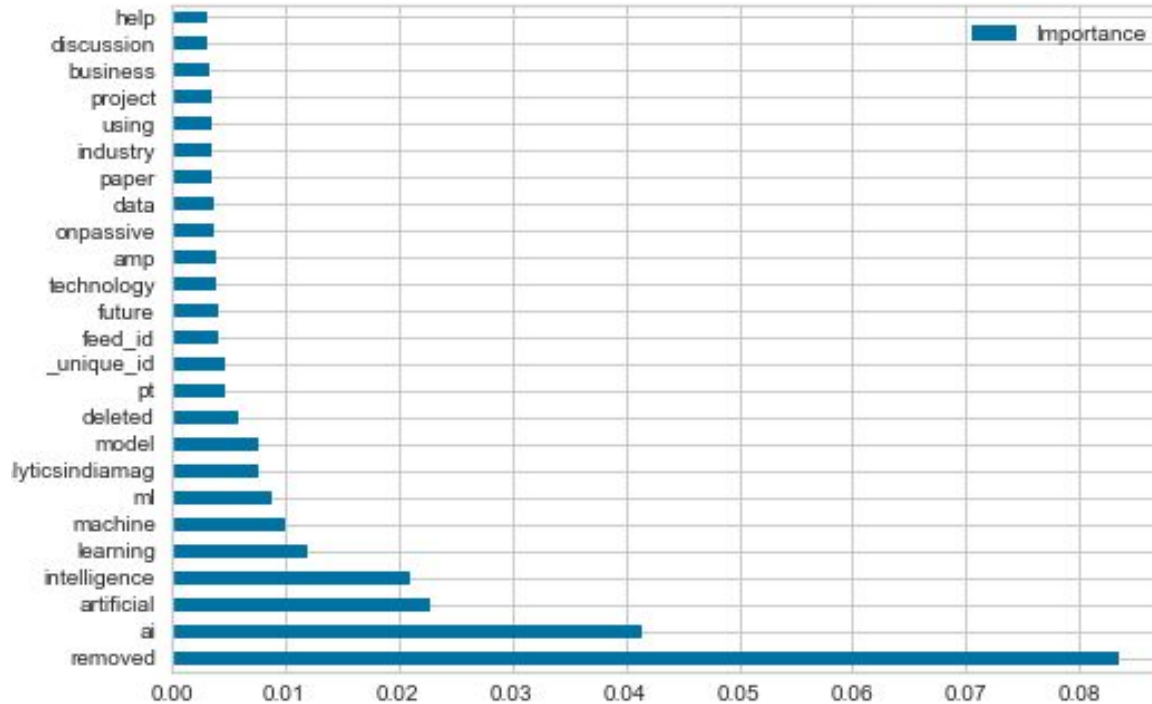## Naive-Bayes

▷   Tfidfvectorizer

▷   alpha=250

▷   fit_prior=True

## Random Forest

▷   Tfidfvectorizer

▷   n_estimators=500

▷   criterion='entropy'

▷   max_depth=None

▷   min_samples_split=2

▷   min_samples_leaf=1

▷   max_features='auto'

▷   max_leaf_nodes=None

▷   min_impurity_decrease=0.0

# Random Forest Feature Importance

# Conclusion

**Summary**

- Random Forest
- Many features
- 85% Accuracy

**Limitations**

- Narrow view
- "Removed" Dependent

**Future**

- Collect more data
- Logistic Regression
- Neural Net

# Thanks!

## Any questions?