# Ford

# Kiwi: Python Tool for Text Processing and Classification

**Neelima Pulagam, Sai Marasani, Brian Sass**
**Ford Motor Company**

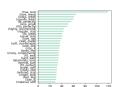## An NLP tool for beginners and experts alike…

## Introduction

Text analysis is a ubiquitous and important task as most unstructured data is textual. Extracting and classifying text can be extremely useful and help automate manual, time consuming tasks. On the other hand, machine learning can be a difficult field to breach. We propose a user-friendly text classification tool that will organize and categorize unstructured text data. This tool will allow users within the field to avoid creating boilerplate code for basic natural language processing (NLP) tasks and new to machine learning and NLP plug and play with various models and methods. Our main goal is to make natural language processing accessible and easy.

## Uploading Files and Data Visualization

The user begins by choosing a CSV file of text data that has already been cleaned. This includes removing null values and special characters. The chosen file will be displayed in the data upload frame.

Once the user has selected the CSV file with text data, the corresponding input text column and output label column, they can visualize the data. They can get a visual of the class distribution as a pie chart, word cloud and histograms for n-grams.



## Machine Learning Models

Once the input text and output columns have been chosen, the user can then run and test a few different types of models. To process the data, the user can choose between three different vectorizers.
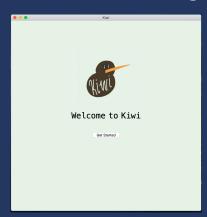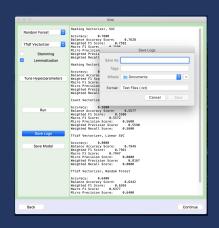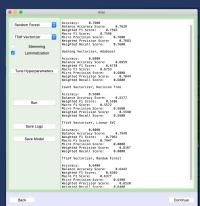
Select Classifier
- Adaboost
- Random Forest
- Extra Trees
- Gradient Boosting
- Extra Tree
- Decision Tree
- K-Neighbors
- Logistic Regression
- SGD
- SVC
- Linear SVC
- Bernoulli NB
- Passive Aggressive
- Ridge
- Run All Classifiers

Select Vectorizer
- Tfidf Vectorizer
- Hashing Vectorizer
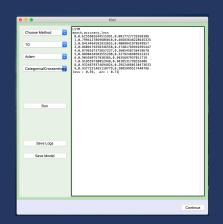- Count Vectorizer
- Run All Vectorizers

## Machine Learning (Cont.)

Additionally, there are two options for users to preprocess their text data by utilizing stemming and lemmatization. Stemming, the process of reducing a word to its stem (i.e. cats -> cat or going -> go) and lemmatization, a technique that utilizes part of speech tagging to reduce a word to its root (i.e. was -> be or has -> have) sometimes offer better performance.

## Hyper-parameter Tuning

Hyper-parameter tuning is another useful technique that searches for through different combinations of parameters to find the ideal model architecture with the best performance. Our tool currently allows users to simply run randomized search on a model of their choosing. The parameters were manually hardcoded into a dictionary and the parameters for the selected model are retrieved when the the 'Tune Hyperparameters' button is pressed.

## Deep Learning Methods

Currently there are two architectures implemented: LSTMs and RNNs. These models were chosen as they had the shortest training time among the several available architectures. Also, these architectures do not need GPU capabilities to run efficiently compared to other models. The users can currently choose from a set of options the number of epochs, optimizer and loss function for each model. The users will have the accuracy and loss function for each epoch.

## Future Work

We see a lot of potential and growth for this tool. This includes giving the user more granular control of choosing the model architecture. This includes which hyperparameters they wish to change or building models that utilize several different types of input. Additionally, there are a lot of deep learning architectures that can be added. Next steps for the tool also include more text processing methods and text embeddings. We also plan on integrating existing frameworks such as ML Flow, Optuna, etc. to provide a more user friendly experience. This project holds much potential that we hope to exploit in the near future.

## References

Agarap, Abien Fred. (2018). Statistical Analysis on E-Commerce Reviews, with Sentiment Classification using Bidirectional Recurrent Neural Network.

Women's E-Commerce Clothing Reviews, 2018, June 30, 2021. https://www.kaggle.com/datasets/nicapotato/womens-ecommerce-clothing-reviews