

Project Name: Natural Language Processing Project

Project Type: Required partner project

Project Goals:

Can the programming language of a GitHub repository be predicted from the text of the README document?
Create a dataset of GitHub readme documents, Explore using NLP techniques, Build a model that predicts the programming language using NLP techniques to extract features and classification algorithms for prediction.

Deliverables: Well documented Jupyter Notebook, slide deck, and 5 min presentation

Stage	Tools	Brief Description of Process	Challenge Resolution
Plan	<ul style="list-style-type: none">• Visual Studio	<ul style="list-style-type: none">• Obtained requirements from curriculum• Decided on subject for dataset gathering• With partner determined MVP questions to answer for general audience• Established workload timeline and split up tasks for each of us to address	<ul style="list-style-type: none">• No unusual challenges in this section
Acquire	<ul style="list-style-type: none">• Visual Studio• .py script	<ul style="list-style-type: none">• Built function to get urls from 1st search page and then get the readme text and language (if present) and return dataframe.• Broadened function to get urls from multiple search pages• Scraped 100 search pages with 10 urls each with filtering for only those with readme and language returned 800 observations	<ul style="list-style-type: none">• First time using web scraping to obtain data• Error getting data from more than 10 pages - used sleep function - could have/should use GitHub API for this
Prepare	<ul style="list-style-type: none">• Visual Studio• .py script• Jupyter Notebook	<ul style="list-style-type: none">• Identified the 4 most common languages and filter data to include only those observations (reduced dataset to approx 400 rows)• Cleaned data (lowercase, convert to only	<ul style="list-style-type: none">• First time using regex

		<p>ascii characters and recode to utf-8), use regex to remove anything not a-z, 0-9, or a space</p> <ul style="list-style-type: none"> • Tokenized using NLTK Toktoktokenizer, removed stopwords, then created column of stemmed vs lemmatized words • Split dataset into train, validate, test, and returned additional train_explore dataset 	
Explore	<ul style="list-style-type: none"> • Jupyter Notebook • Visual Studio 	<ul style="list-style-type: none"> • Checked for imbalance in the dataset language distribution • Found top 5 words for each language type and discovered 3 words were common to all, returned to prepare and added 'file', 'data', and 'environmental' to the stopwords list • Created visualization of distribution of top 20 words by language • Performed hypothesis test to determine if the length of the readme text was significantly different from the overall mean length. It was not for 3 of the languages, only HTML showed significance • Functions are found in explore.py 	<ul style="list-style-type: none"> • Did not use word cloud with question mark mask in final notebook, but included the visualization in the slide deck
Model	<ul style="list-style-type: none"> • Jupyter Notebook • Visual Studio 	<ul style="list-style-type: none"> • Partner extracted features using both Bag of Words (BOW) and TF-IDF (Term Frequency-Inverse Document Frequency). • Established a baseline of 35% by predicting the most frequent language (Python) for all observations • Then used multiple classification models to predict the language including Logistic 	<ul style="list-style-type: none"> • Model results accuracy were disappointing and under other class teams. Scheduled office time with instructor to review process to ensure I implemented everything correctly

		<p>Regression, Decision Tree, Random Forest, K Nearest Neighbors (KNN), and both multinomial and complement Naive Bayes</p> <ul style="list-style-type: none"> • Functions are found in model.py 	
Evaluate	<ul style="list-style-type: none"> • Jupyter Notebook • Visual Studio 	<ul style="list-style-type: none"> • Top 3 models based on accuracy were run on the validate data and the best performing, the Logistic Regression using TF-IDF was run on the test set. The final result was an average of 47% accuracy on unseen data. • Functions are found in model.py 	<ul style="list-style-type: none"> • Partner created the functions to return evaluations
Model Explanation	How does your algorithm work?	<ul style="list-style-type: none"> • <i>From curriculum</i> • Technically a regression algorithm (goal is to find the values for the coefficients that weight each input variable) • Used for predicting discrete outcomes (binomial and multinomial) • Because the prediction for the output is transformed using the logistic function, a non-linear function, it is a classification algorithm. • The output is a value between 0 and 1 that represents the probability of one class over the other. • Like linear regression, logistic regression works better when you remove attributes that are either unrelated to the output variable or correlated to other attributes. 	

Delivery	<ul style="list-style-type: none">• Jupyter Notebook• Canva	<ul style="list-style-type: none">• Partner create slide deck using Canva• I added Summary, Conclusions, and Next Steps to the Final Jupyter Notebook	<ul style="list-style-type: none">• No unusual challenges in this section
-----------------	--	--	---