

Predicting Coding Languages in Doctor Who Readme Files

By Caitlyn Carney



Executive Summary

What should you expect going into this?



Explore

What was learned in exploration?



Model

What models were ran and which one was the best.



Conclusion

What was learned overall.



Questions

Are there specific words used in different programming languages that can they help us predict the language being used?

Goal

Create a NLP model to predict the programming language used in a github repository based on the words and word combinations.

The Process


Acquire through web-scraping, prepare the content, explore the data, create models, and pick the most accurate one to use.

Key Findings

Each language has different frequent words from each other, Java Script is the most commonly used, and the top 10 most frequent words of all are all Java Script terms.

Conclusion

There are specific words used in readme files that key into what language is being used. The SDG Classifier model performed the best, and beat the base line by 14.3%.





Exploration

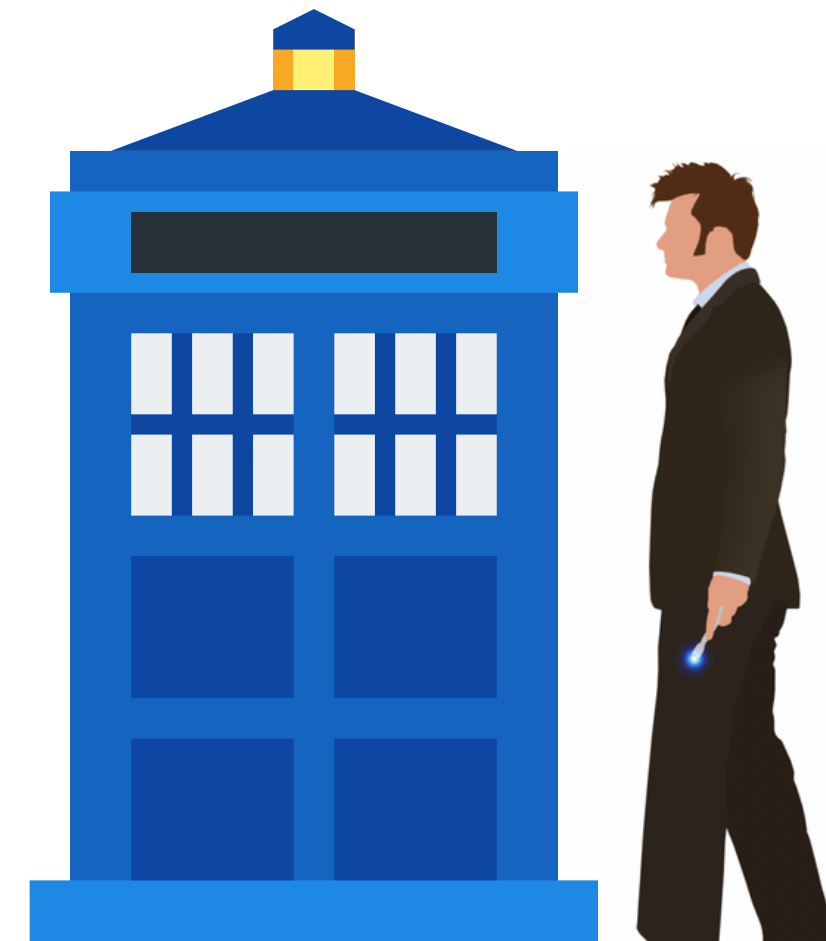
A silhouette of a person wearing a backpack, filled with a word cloud of programming and web-related terms. The words are in various sizes and colors (blue, black, white). The silhouette is facing right. The word cloud includes terms like 'spring', 'data', 'rest', 'self', 'href', 'lastname', 'baggins', 'curl', 'x', 'person', 'objects', 'links', 'self', 'baggins', 'links', 'type', 'application', 'json', 'href', 'localhost', 'people', 'pages', 'size', 'or', 'true', 'templated', 'true', 'firstname', 'frodo', 'page', 'size', 'data', 'rest', 'self', 'href', 'spring', 'boot', 'methode', 'retourne', 'value', '20', 'curl', 'localhost', 'people', 'localhost', 'people', 'firstname', 'spring', 'data', 'binocraft', 'mod', 'number', '0', 'localhost', 'people', 'curl', 'methode', 'lastname', 'baggins', 'href', 'localhost', 'people', 'spring', 'mvc'.

baker 1975
baker 1980
davison 1984
baker 1981
pertwee 1971
mccooy 1988
baker 1985
baker 1978
pertwee 1974
pertwee 1972
sylvester mccooy
peter davison
patrickroughton
netflixamazonhulu jon
jon pertwee
mccooy 1989
trougton 1968
mccooy 1981
hartnell 1964
baker 1979
davison 1982
baker 1977
pertwee 1970
davison 1983
tom baker
netflixamazonhulu tom
hartnell 1965
pertwee 1973
william hartnell
colin baker

I learned that there are specific phrases/words that are used more often in each coding language.

I also learned that there were no commonalities between the languages and their top 20 words/phrases.

Java script is the most commonly used coding language
and makes up the top 10 most frequent words of all.



Modeling

I made my baseline:

29.5%

***Tested train on 6 models. My
top models were:***

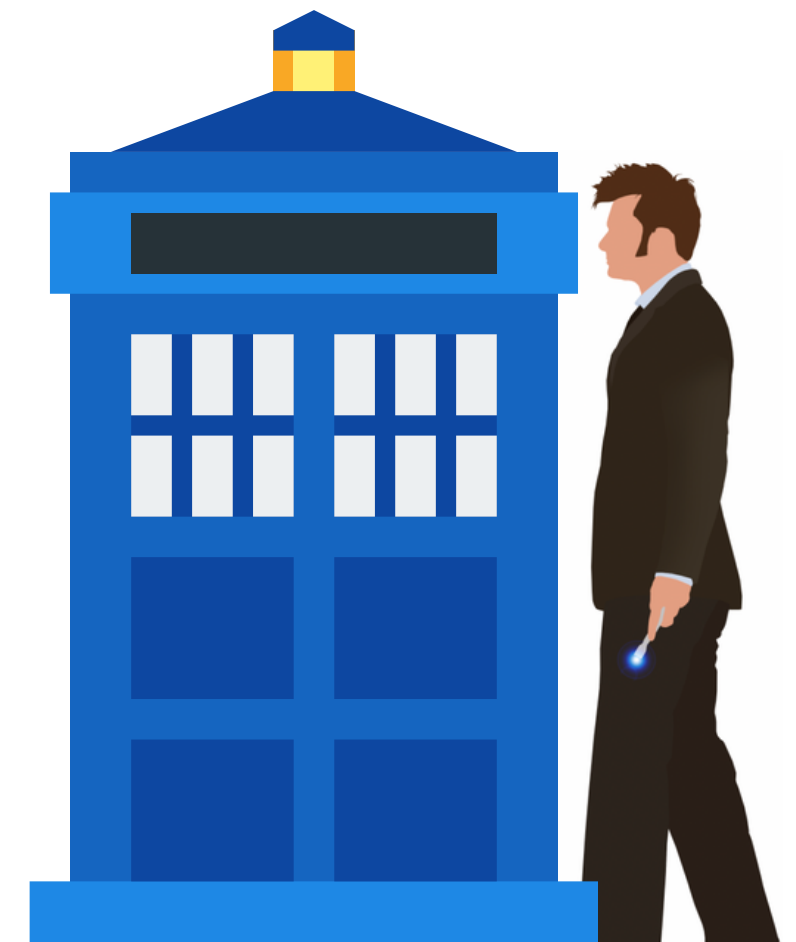
Logistic Regression Train	95.5%
Ridge Classifier Train	97.7%
SDG Classifier Train	97.7%

I then ran them on validate:

Logistic Regression Validate	42.1%
Ridge Classifier Validate	36.8%
SDG Classifier Validate	57.9%

The best model is:

SDG Classifier Test	43.8%
---------------------	-------



Conclusion

Exploration

There are specific words and phrases that feed into each coding language individually, with no overlap within the top 20 of each language.

Java Script is the most prominently used coding language out of them all. This coding language makes up for the top 10 most frequently used words across all languages.

Modeling

Each model beat the baseline in train, and validate.

However, the best model overall was the SDG Classification model which beat my the baseline of 29.5% by 14.3%. Putting it at a 43.8% accuracy reading.





Questions?



[*https://github.com/CaitlynCarney/coding_language_prediction*](https://github.com/CaitlynCarney/coding_language_prediction)



[*https://www.linkedin.com/in/caitlyn-carney/*](https://www.linkedin.com/in/caitlyn-carney/)

