

Kaggle Competition

Predicting word complexity

Machine Learning Course

2st Semester of 2023-2024

Turis Gavriil-Vlad

`gavriil-vlad.turis@s.unibuc.ro`

Ilie Octavian Tudor

`octavian-tudor.ilie@s.unibuc.ro`

Abstract

This document is a report of a Kaggle competition that is a project for the machine learning course. Our team is formed by: Turis Gavriil-Vlad and Ilie Octavian Tudor. We needed to predict the word complexity for this competition using a regression model. We had 8633 words for the train data and they were in different languages : catalan, english, filipino, french, german, italian, japanese, portuguese, sinhala, spanish. The test data was formed by 5623 words and we needed to predict the complexity column. To evaluate our progression we used: Coefficient of Determination, and Person Correlation Coefficient.

1 Introduction

- In this project, we aim to develop a machine-learning model to predict the complexity of a given word within a sentence
- Turis Sef: Researched and identified a comprehensive set of features related to word complexity, from the most general to the most detailed and performed grid search to find the best parameters for the model.

Tudor Sef: Analyzed and selected the best features for predicting word complexity by testing different combinations of features.
- Evaluated the impact of different feature sets on model performance.
 - Evaluated the impact of different parameters for the model performance
 - Evaluated the performance of K-Nearest Neighbors (KNN) and Random Forest Regressor models .
 - Analyzed the importance of each feature in the obtained score.
- Tested the effects of including too many features versus too few features.

- Searching for general features Testing KNN vs Linear Regressor vs MLP vs RandomForest+Regressor with the general obtained features. Searching for more complex features (which use libraries like wordnet , spacy , BERT). Testing the Random Forest Regressor model with all obtained features. Testing the Random Forest Regressor model with an optimal combination of features: Using Grid-Search to find the best parameters
- -Discover all the steps involved in creating an artificial intelligence model.
 - Gain familiarity with the Kaggle platform and its tools for data science competitions and project collaboration.
 - Train our analysis and research skills in solving a given task.
 - Evaluate our interest in the field of artificial intelligence to see if it aligns with our career goals.
- Studies on how Random Forest works and its effectiveness
 - Research on the capabilities and applications of SpaCy, WordNet, BERT and RoBERT libraries
 - Exploration the embedding technique, which transforms words into numerical vectors to capture semantic meaning
 - Research on the parameters of RandomForestRegressor, such as the number of trees, maximum depth, and splitting criteria
- 2 Approach
- Github Link: <https://github.com/Tudorr02/IA>
- For coding we used Python and the IDE used was Visual Studio Code and also the Kaggle platform.

077	• The training process takes now 5 seconds, but	3	Conclusions and Future Work	126
078	in the past when we tried various features we			
079	had waited several minutes.			
080	• We used a regression model named Random		• From our point of view, we could have done	127
081	Forrest Regression, Linear Regression, Spacy		differently plenty of things. For instance, the	128
082	library, and WordNet. Also, we tried K-		communication in the team could have been	129
083	Neibgbour Regression, MLP Regressor and		better. We certainly, get along very well, but	130
084	some transformer models like BERT and		in the heat of the moment and with other con-	131
085	RoBERTa, for word embeddings but with no		siderable tasks in this period on our minds,	132
086	success.		both at work and as well as at university with	133
087	• We added features and we selected them based		the examination session and other projects.	134
088	on their importance and their effectiveness,		I tend to think that, we could enhance our	135
089	and also we tried to not overfit the model. An-		communication overall and maintain a calmer	136
090	other trick we used was preloading the spaCy		environment.	137
091	language models and storing them into a dic-		• Definitely, by adding more effective features	138
092	tionary. This approach significantly optimize		and doing a more elaborate grid search for the	139
093	data processing.		newly added features.	140
094	• After testing different combinations of		• The project itself, was extremely engaging,	141
095	features types of features, we finally choose		however the fact that we had only five sub-	142
096	12 of them. We will list them with a brief		missions daily and other projects and exams	143
097	description: Word Length :		on our head, made it less enjoyable. On top	144
098	Calculates the length of the word.		of that, the concept of the project was very	145
099	Vowel Count : Counts the number of vowels		interesting and highly relevant nowadays.	146
100	in the word.		• By doing this project, we learned considerably	147
101	Consonant Count : Counts the number of		new pieces of information about how a proper	148
102	consonants in the word.		machine-learning project is done and about.	149
103	Part of Speech : Retrieves the part of speech		Moreover, we learned about relevant aspects	150
104	tag for the word.		that make a word more complex, and how to	151
105	Word Frequency : Computes the frequency of		use artificial intelligence concepts to solve a	152
106	the word from an external corpus.		given task.	153
107	Synsets Count : Counts the number of synsets		• From our point of view, an Image classifier	154
108	(sets of cognitive synonyms) for the word in		and a football match score prediction would	155
109	WordNet.		be nice.	156
110	Synsets Depth : Determines the maximum		• More daily submissions and more train data.	157
111	depth of the synsets in the WordNet hierarchy.			
112	Hypernyms Count : Counts the number of			
113	hypernyms (more general terms) for the word			
114	in WordNet.			
115	Hyponyms Count : Counts the number of			
116	hyponyms (more specific terms) for the word			
117	in WordNet.			
118	Root Distance : Calculates the distance of the			
119	word from the root of its dependency tree in			
120	the sentence.			
121	Syllable Count : Determines the number of			
122	syllables in the word.			
123	Capitalization Ratio : Calculates the ratio			
124	of capital letters to the total letters in the word.			
125				