

i.) Features

All of the features designed when building the model, including features not used in the final model.

Lengths: The length of sentence 1 and the length of sentence 2 as two separate features.

Length Difference: The absolute value of the difference between the two sentences.

Overlap: A ratio of shared words to unique words used to get a proportion of overlap between the two sentences.

Sentence Length Difference: The length difference as proportion. Calculated by the length of sentence 1 minus the length of sentence 2 divided by the length of sentence 1.

Sentence Length Difference 2: An alternate way of finding sentence length difference.

Calculated by the reciprocal of d raised to the power of the length of sentence 1 minus the length of sentence 2. d is a constant, 0.8 was used for the value of d .

Levenshtein Distance: Used to measure the difference between the two sentences. Levenshtein distance is the minimum number of single-character edits used to change sentence 1 to sentence 2. I used the levenshtein method from the pylev library, which is an implementation that uses an iterative version of the Wagner-Fischer algorithm.

Cosine Similarity: Used to measure the similarity between two non-zero vectors of an inner product space that measures the cosine between them. First, I vectorized the sentences and then used the cosine similarity formula on the two sentence vectors. I experimented with removing stopwords from the sentences, but I observed better accuracy without removing the stopwords.

Shared Words: The number of shared words between sentence 1 and sentence 2. Found using the intersection of the set of sentence 1 and the set of sentence 2.

BLEU Score: BLEU (bilingual evaluation understudy) score is a number between 0 and 1 that measure the similarity of machine-translated text to reference translations. The formula consists of brevity penalty and n-gram overlap. I used 3 different BLEU scores, for unigrams, bigrams, and trigrams. Experimenting with higher n-values, I found that it did not help the model because most of the BLEU scores were 0. I used the bleu_score method from the nltk library.

Meteor Score: Meteor score is a metric based on the harmonic mean of unigram precision and recall. I used the meteor_score method from the nltk library.

Jaccard Similarity: Jaccard similarity measures the similarity between the 2 sentences using the shared values and the unique values. Jaccard similarity was calculated by the intersection of the set of sentence 1 and the set of sentence 2 divided by the union of the set of sentence 1 and the set of sentence 2.

NIST Score: NIST evaluates the quality of the text using machine translation. This is based off of BLEU score, but it differs because NIST score gives more weight to rarer n-grams. NIST also uses a different calculation for brevity penalty. I used the nist_score method from the nltk library.

I mostly used trial and error to select the features for the model. The features I selected were BLEU score, overlap, cosine similarity, levenshtein distance, length difference, meteor score, and NIST score. Using this combination of features, I observed the highest accuracy on the dev set.

ii.) Data Preprocessing and Feature Preprocessing

Methods used to preprocess the data before extracting the features, including preprocessing methods not used in the final model.

Remove spaces before punctuation: Because I planned to expand all of the contractions and remove the punctuation, I had to fix the spacing because there was an additional space before any punctuation. I used regex to remove all of these spaces.

Expand contractions: I expanded all of the contractions using `contractions.fix`. Because a contraction and the expanded version have the same meaning, expanding the contractions can improve results when finding common words or jaccard similarity.

Lowercase: Strings are case-sensitive, so it makes sense to change all of the letters to lowercase, as this does not change the meaning of the sentence.

Remove punctuation: In most cases, punctuation does not change the meaning of the sentences, so it makes sense to remove the punctuation.

Remove digits: From doing research on the best preprocessing methods, I found that some sources stated that removing digits from the sentences can improve accuracy. But, I found that this significantly lowered the accuracy, so I did not use this preprocessing method in the final model.

Remove stopwords: Stop words are filler words that do not change the meaning of the sentence, so it initially made sense to me to remove these. However, removing the stop words significantly decreased the accuracy of the model, so I kept all of the stop words.

Lemmatization: Lemmatization changes any kind of word to its base root model. This can make it easier to find synonyms. I observed better results with lemmatization rather than stemming.

Stemming: Stemming reduces a word to its word stem. I experimented with stemming, but did not use it in the final model, as I found better results by lemmatizing.

Remove white space: I removed all additional white space that was not needed.

Split: Because most features I was planning on using required lists of words as input, I split all of the sentences by white space, resulting in lists of words.

Similar to features, I also used trial and error to determine which preprocessing methods to use. I found that I achieved better results using more simple methods of preprocessing. Otherwise, preprocessing can remove too much information from the sentences. In the final model, the preprocessing I used included removing spaces before punctuation, expanding contractions, converting to lowercase, removing punctuation, lemmatizing, removing white spaces, and splitting.

For feature preprocessing, I used `StandardScaler()` from the `sklearn` library to normalize all of the features.

iii.) Algorithms and Libraries

pandas:

- I used pandas to read in the csv files and store the data in data frames.

nltk:

- stopwords: used to remove stop words when needed. I did not end up using this in the final model
- WordNetLemmatizer: used during preprocessing to lemmatize the inputs.
- PorterStemmer: I experimented with stemming during preprocessing, but did not end up using this in the final model.

sklearn:

- pipeline.make_pipeline(): Used to construct a pipeline from the given estimators.
- svm: SVM was the model that I fit the data to, using a Gaussian kernel.
- StandardScaler: used to scale and normalize all of the features.
- LogisticRegression: experimented with fitting the data to a logistic regression model, but did not achieve good results.

contractions: Used to remove contractions when preprocessing the data.

regex: Used during preprocessing to easily remove all punctuation and white space.

string: Used to easily detect punctuation so I could remove it during preprocessing.

pylev: Used the pylev library to calculate the levenshtein distance.

warnings: Used to suppress the warnings from calculating the BLEU score.

iv.) Experience and Lessons

This was my first time creating a machine learning model, so overall I learned a lot. Originally, my approach was to use a lot of preprocessing techniques and add a lot of features, but I realized that this was not the best tactic. During preprocessing, I initially stripped too much information from the data, causing my accuracy to be low. I also originally created too many features that did not work well together. Most of my time on this project was spent implementing features and using trial and error to decide which features helped the model. I also spent a lot of time experimenting with different models, SVM vs. logistic regression and deciding which had the best performance. Then, I used the hyperparameters to improve the accuracy, such as changing the values of C and gamma. Although there is much room for improvement with this model, I still feel that I learned a lot and gained valuable experience with machine learning.