

A SHORT STORY OF MY APPROACH

I have preprocessed the dataset using PyCaret. Before choosing PyCaret, I employed several methods for Data Preprocessing. Following is a short Fairytale story of my approach.

The Story:

Once upon a time, there was a dataset filled with tweets that needed to be preprocessed. The first step was to load the dataset, which was stored in a CSV file. The unnecessary 'id' column was removed, as it didn't contribute to the analysis.

To prepare the tweet column for further analysis, several preprocessing methods were applied. The adventure began with removing URLs from the tweets using regular expressions, ensuring that website links didn't interfere with the text analysis.

Next, special characters and numbers were eliminated, as they didn't provide valuable insights into the emotions expressed in the tweets. The objective was to focus on the textual content itself. All the letters were converted to lowercase to achieve consistency in the subsequent analysis.

The journey continued with tokenization, splitting the tweets into individual words. However, not all words were equally important for the analysis. Therefore, stop words such as "and," "the," and "but" were removed to eliminate noise and concentrate on meaningful words.

To enhance the understanding of the emotional intensity, the words were lemmatized. This process reduced each word to its base or root form, ensuring that different variations of a word were considered as one entity.

With the tweets now preprocessed and ready for analysis, it was time to delve into sentiment analysis. The first method employed was the AFINN-111 lexicon, which assigns a sentiment score to each word based on pre-determined values.

By summing up the sentiment scores of the words in a tweet, an overall sentiment score for the tweet was obtained.

Another approach utilized VaderSentiment from the Natural Language Toolkit (NLTK). This lexicon-based sentiment analysis tool provided a compound sentiment score for each tweet, taking into account the intensity and polarity of the expressed emotions.

In search of additional insights, the TextBlob library was called upon. TextBlob implemented a sentiment analysis algorithm, determining the polarity (positive, negative, or neutral) of each tweet. The resulting sentiment score provided a different perspective on the emotional intensity.

To unlock the power of word embeddings, the Word2Vec model from the Gensim library was employed. This model learned vector representations of words by considering their context in the tweets. Each word was transformed into a high-dimensional vector, capturing its meaning and semantic relationships with other words. The tweet-level word embeddings were then aggregated by calculating the average value.

However, after applying these methods, the data explorer faced two challenges. First, there was skewness in the data, indicating an imbalance in the distribution of the features. This could potentially affect the performance of machine learning models. Second, there was low correlation between the features, which could hinder the model's ability to capture meaningful patterns.

To address these challenges and overcome the limitations of manual preprocessing, the data explorer decided to leverage the power of PyCaret. PyCaret is a Python library specifically designed for automating the end-to-end machine learning workflow, including data preprocessing, feature selection, model training, and evaluation.

By using PyCaret, the data explorer could efficiently handle the skewness in the data using techniques such as data balancing and feature transformations.

PyCaret's built-in functionality for feature engineering and selection allowed for automatic generation of new features and identification of the most relevant ones.

Additionally, PyCaret provided a comprehensive suite of machine learning algorithms, enabling the data explorer to compare multiple models and select the best-performing one for the given task. The library also facilitated hyperparameter tuning, cross-validation, and performance evaluation, streamlining the model development process.

With the decision to use PyCaret, the data explorer embraced a more systematic and automated approach to data preprocessing, empowering them to overcome the challenges of skew.