

Prediction Competition 6: Text Analysis, Economic News and ML Research

Anonymized name: KentoNanami

Q1

The assignment aims to classify which of the 15,578 news article text snippets contain words related to “economic”, “economy”, “economics”, or the character sequence “econom”. Prior to importing the dataset, adjustments were made including adding a 'text' header, resolving duplicate issues caused by text starting with “=”, and removing problematic symbols like '<' which affected HTML interpretation. After ensuring data integrity, NLP functions were tested and a cleaning function was developed. The Bag of Words (BoW) method was then used to represent text as word occurrences, followed by TF-IDF to evaluate word importance relative to a corpus. Subsequently, a neural network with batch normalization and two hidden layers was constructed, utilizing 'relu' activation function. Output layer with 'softmax' activation was employed for binary classification. Assessing predictions involved creating a confusion matrix and monitoring loss graphs. To address misclassification issues, a threshold based on neural network output probability was implemented and fine-tuned.

Q2



Figure 1: A wordcloud of most common words scaled by size in the text data with a word containing 'econom'



Figure 2: A wordcloud of most common words scaled by size in the text data without a word containing 'econom'

Q3

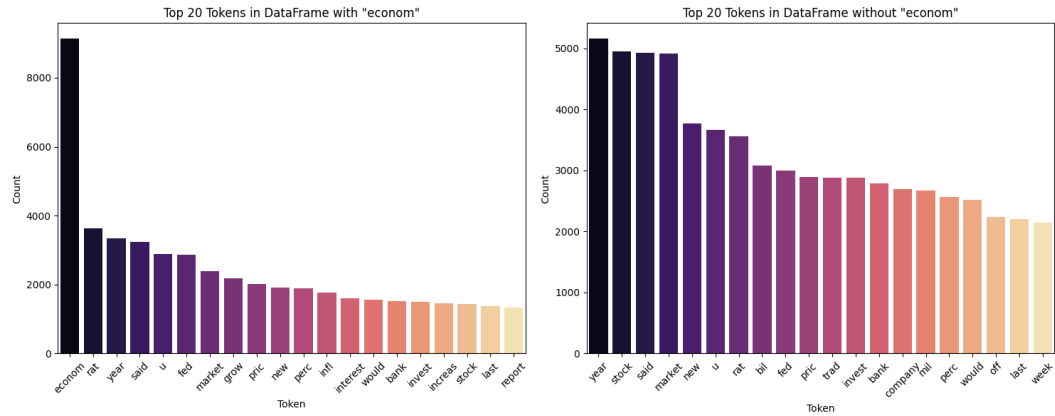


Figure 3: A multigraph plot of the counts of 20 most common words in the text data with a word containing 'econom' in it vs text data without 'econom' in it.

Q4

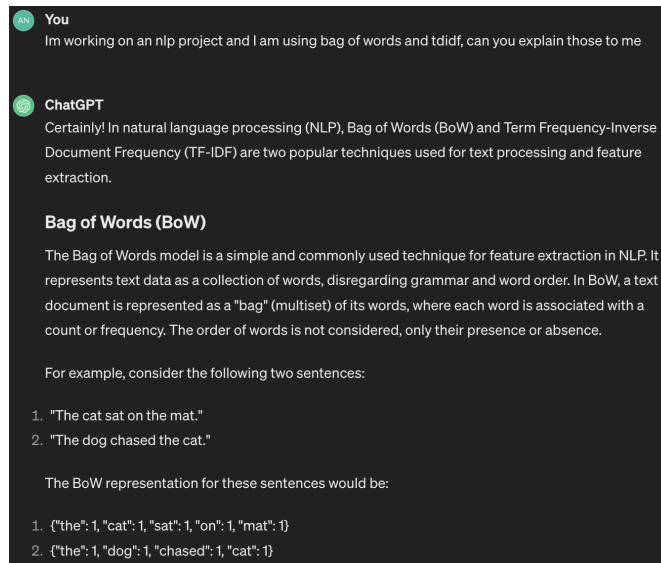


Figure 4: Conversation with Chat GPT

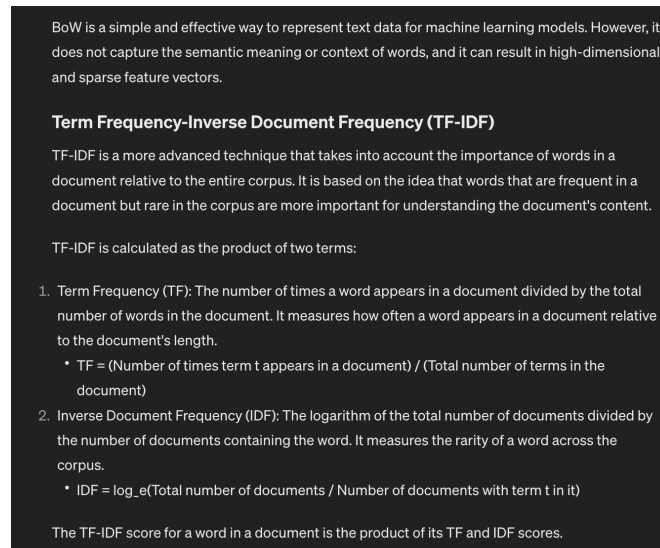


Figure 5: Conversation with Chat GPT

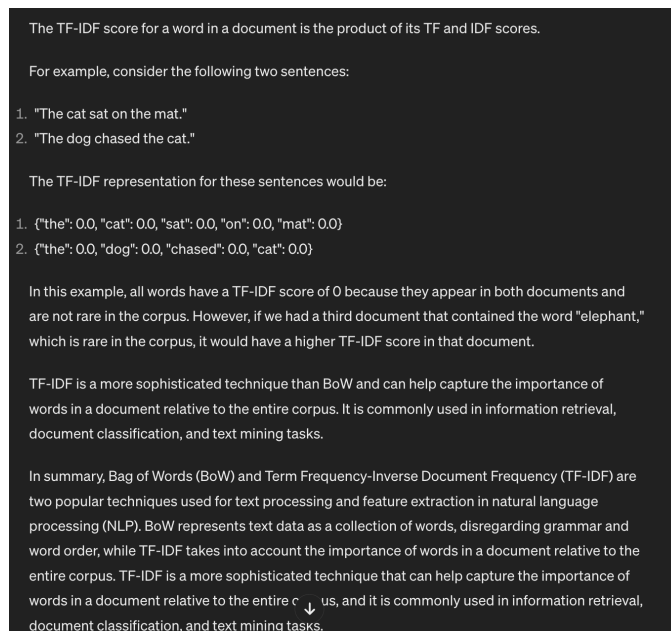


Figure 6: Conversation with Chat GPT