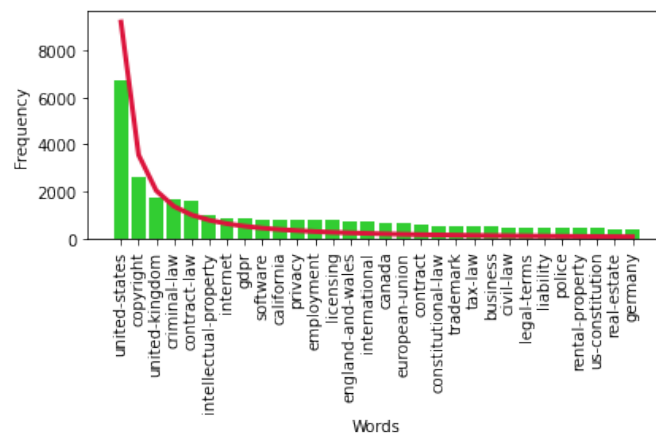




The most interesting tokens to notice are "est", "wal", "busi", "us", and "st". For "est", "us", and "st" we can presume that these come from the ends of words, such as "best", "most", "thus" or for certain Latin words such as "corpus" and "mandamus". The word token "busi" could be derived from "business" or compound words such as "agribusiness", but may also be influenced by words such as "abusive". The token "wal" is interesting, as it is less clear what tokens other than "wales" this may have been derived from; some potential possibilities are "walk", "firewall", or "withdrawal".

4 Step 4

For our third graph, we are observing the frequency of the words in our dataset to determine whether our words maintain Zipf's Law, where the frequency of each word should remain proportional to its rank.



In this dataset, we can see the frequency of each word represented by the green bars, and the red line representing the projected frequency of each word based on Zipf's Law. While we see some level of fluctuation with the words in the center of our graph, the general trend of the words seems to follow that of Zipf's Law. From this we can conclude that we have a valid collection of words for our dataset.