

Weight

For the word “theater” and “magic”, we see that TF and TFIDF produce results that are somewhat similar. There appears to be no pattern to the words that are different. However, PMI produces very specific words. For example, “theater” is most similar to “vishnevskaya” and “magic” is very similar to “azkaban” and “psilocybin.” This makes sense, as when PMI is building a vector for these uncommon words, it takes into account the low probability of these words in its calculation.

“theater” and “magic” were both common words, but “termites” is a rare word, only appearing 16 times. In this case all three weighting measures produced fairly similar results.

We conclude that TF and TFIDF are fairly similar, while PMI produces more uncommon words. However, if the original word is uncommon the three methods perform similarly.

Distance

The Euclidean and Cosine measures produce the exact same list of words, in the same order. However, L1 produces a list of bizarre words that are not related to magic. It thus seems that L1 is not a good metric.