

# Comparison of two-side MI

As mentioned in [Methods for building twitter-specific sentiment lexicon](#) there are two general ways that I tried to build a twitter specific sentiment lexicon. The first was to calculate the mutual information associated with the positive class and subtract from that the mutual information associated with the negative class. The other option was to take whichever had the higher value as the mutual information score, multiplying the negative class by -1.

However, upon inspection of the results, the winner-take-all method is producing a much more sensible list of vocabulary.

## The raw files can be found here

### Binary-Based

- [Winner-take-all](#)
- [Relative](#)

### Count-Based

- [Winner-take-all](#)
- [Relative](#)

## Determining a cutoff

There is some imbalance in how many terms are given higher mutual information for the positive class and the negative class.

For example, the 0 value for the winner take all binary setup occurs about two thirds of the way through. This imbalance would be problematic if all words were used to compute shifts in sentiment for the sarcasm detection part. The best solution seems to be to make the threshold some number of words from the ends (e.g. we're usign a ranking scheme to determine which words are associated strongly enough with each class to be representatives of that class).

My next steps are to determine how much overlap with the content of the sarcasm dataset there is.

# Adding a minimum count cutoff

Commit [27dd9e4300](#) adds a cutoff to how low in frequency a given token can occur in order to be considered in the mutual information calculations. The entry is still present in the results array in the program, the mutual information is just automatically set to 0 if there are less than `x` instances of a feature.

Currently the behavior is not special for counts. E.g. when a binary feature matrix has been computed, the minimum cutoff is effectively how many tweets it occurred in, The counts do not try to emulate this and instead just count the frequency of usage including multiple usages in a single tweet.

---

Revision #9

Created Thu, Dec 20, 2018 5:33 AM by [kenneth](#)

Updated Thu, Feb 6, 2020 4:53 PM by [kenneth](#)