

Grading Rubric, NLP

Homework 1: Corpus Statistics

1. (10 pts) Choose or Collect Appropriate Data: the two documents should be sufficiently different to yield good questions and of sufficient length for the word frequency and bigram lists to be useful.
2. (30 pts) Process each document and produce the frequency, bigram frequency and bigram PMI score lists, with processing steps chosen to produce lists suitable for analysis of your question. Discuss any issues with the lists.
 - a. (10 pts) Description of processing steps: tokenization, lower case, stopwords or lemmatization, word frequencies, bigram frequencies and bigram PMI with frequency filter of 5 or greater, and state why you chose those options.
 - b. (10 pts) Describe how the bigrams scored by frequency are different that the bigrams scored by PMI, using the definitions or an overall characterization of the words that score highly on the lists.
- 3a. (10 pts) Define a comparison question between the two documents.
- b. (20 pts) Answer the question by picking examples from the lists and discussing how they show that the documents are different, not just reporting numbers. Discussion may include collection steps if significant, which will count towards discussion of differences.

Additional merit (5 pts): Choose some aspect of the processing or analysis that requires additional thought or work. Some options are

- If you collect your own data, describe that work or steps necessary to obtain documents ready for processing (add section to part 1)
- During processing, describe additional steps, for example, if you define or modify stopword lists to suit your documents or analysis question (from part 2a)
- Make trigram lists and include in your discussion
- Expand the question or analysis (from part 3).

Interpretation of numeric grades as letter grades:

90 – 100 A
85 – 89.9 A-
80 – 84.9 B+
75 – 79.9 B
70 – 74.9 B-

below 70 has similar interpretation in the C and lower range.

Late assignment submissions will be accepted, but will be penalized:

1 Week late- 10 points taken off (2/3 letter grade)

2 Weeks or more late- 20 points taken off (1 1/3 letter grade)

Late assignments may possibly be excused by emailing the instructor with the appropriate excuse.