**Response to Justin's Questions**

*Which similarity metric seemed to work better? Why do you think that is?*

Manhattan Distance barely worked. It would often find close to zero correct answers for a given file. Cosine distance and Euclidean distance had the same result and were both quite effective. There were no instances where at least a few matches could not be found. With a better vector model, the results could likely be even greater.

*What input files did better than others?*

Past tense (gram4) was far by the most effectively operating on. City and state (city-and-state) also yielded accurate solutions. Nationality adjective (gram6) also proved to be easier to analyze for the program. Familial Connections and superlatives tended to be harder for word_analogy to solve.

*Did normalization help? In what cases?*

The vectors were already normalized, therefore, normalization never made a difference in the output. There was no difference in content between output1x/eval.txt and output0x/eval.txt.