# Word2Vec Final

## Jean Mark Gawron

## May 7, 2020

# 1 PMI

| w\c | professional_JJ | college_JJ | olympic_JJ | equine_JJ |
|---|---|---|---|---|
| golf_NN | 3 | 1 | 2 | 0 |
| pro_NN | 0 | 0 | 0 | 1 |
| amateur_NN | 1 | 1 | 4 | 2 |
| champion_NN | 2 | 3 | 6 | 4 |
| tennis_NN | 8 | 3 | 1 | 0 |

Find the following:, using the formulae in vectors1.pdf, slides slides 20-24.

1.1. $p(w = \text{golf\_NN}, c = \text{professional\_JJ})$

1.2. $p(w = \text{golf\_NN})$

1.3. $p(c = \text{professional\_JJ})$

1.4. Compute the PPMI score for word *golf_NN* and context *professional_JJ*.

1.5. Suppose P(w = i | c = j) is equal to P(w = i). What can we say about PPMI(w=i,c=j)? If you don't remember the discussion of this case in class, use the chain rule to turn this into a fact about P(i,j).

1.6. Suppose P(i | j) = 2 * P(i) and suppose neither P(i) nor P(j) is equal to 0. Can the PPMI value of target word $i$ and context word $j$ be 0? Explain. Use an example with made-up counts, if it helps.

1.7. Is it possible for the PPMI value of target word $i$ and context word $j$ to be negative? Why or Why not?

1.8. Is it possible for the PMI value of target word $i$ and context word $j$ to be negative?

1.9. When is PPMI undefined?

## 2   Cosine similarity

Use slides 33-37 to help with the following.

2.1. Using the **counts** (rather than the PPMI values), compute the cosine similarity of target words *champion* and *tennis*. Note: It should be a number between 0 and 1.

2.2. Same two words: Now compute the cosine similarity using vectors with PPMI values. Note: For this problem, think of the log of a probability of 0 as a negative number with a **very** high absolute value. So for the purposes of PPMI any 0 probabilities are going to yield a PPMI of 0. Show at least this much of your work: What are the two PPMI vectors? What are the vectors after they are divided by their length? (We say the have been **normalized**).