

Name: _____

CS4442/9542: Artificial Intelligence II
Winter 2013: Quiz 2 Solution

Instructions:

Show all the work you do. Use the back of the page, if necessary. Calculators are allowed, laptops, cell phones, or any other communication devices are not allowed. This is an open notes exam. However, the sample quiz solution is not allowed.

1. (10%) Suppose our corpus is the following text:

I liked Herman. Herman was a hippopotamus. And Herman thought about his friend the Mouse. I thought it made a lot of sense, for Herman the hippo to jump our fence. Now, Herman followed me to my house, where he stayed about the pool, and made no noise, like a mouse.

Ignore punctuation, capitalization (i.e. “To” and “to” should be treated as the same word) and treat the corpus as one continuous sentence.

- (a) (5%) Using unigram model, which of the following sentences is more likely under ML estimation and why? You do not have to compute the exact probabilities of each sentence, but may do so, if you wish.
- Herman liked to jump
 - I liked the pool

Solution:

$C(\text{Herman}) = 5$, $C(\text{liked}) = 1$, $C(\text{to}) = 2$, $C(\text{jump}) = 1$.

$C(\text{I}) = 2$, $C(\text{the}) = 1$, $C(\text{pool}) = 1$.

The combined unigram counts for the first sentence are higher, so the probability of the first sentence is higher.

- (b) (5%) Now use the bigram model. Which of the following two sentences are more likely under the ML estimation and why? Again, you do not have to compute the exact probabilities.
- Herman was a pool mouse
 - I thought about the mouse

Solution:

The bigram “a pool” is missing in the training set, so probability of the first sentence is zero. Second sentence has all bigrams present in the training corpus, so the probability of the second sentence is higher.

2. (30%) Our corpus is the following text:

Herman was a mouse. Herman lives in a house with a mouse.

Assume the vocabulary size $V = 10$. Use a bigram model and Good-Turing method for estimating the probability of the unseen bigrams. Ignore punctuation. For the observed bigrams, use probabilities based on ML estimation. You do not have to renormalize the numbers to get the true probabilities, but if you want to, you may do so. What is the probability of the following bigrams. In your answer, you can leave expressions not calculated out completely, that is something like $2/30 - 1/6$ is fine.

- *mouse was*
- *Herman lives*
- *a mouse*

Solution: Let us count the bigram first:

$C(\text{Herman was}) = 1$, $C(\text{was a}) = 1$, $C(\text{a mouse}) = 2$, $C(\text{mouse Herman}) = 1$, $C(\text{Herman lives}) = 1$, $C(\text{lives in}) = 1$, $C(\text{in a}) = 1$, $C(\text{a house}) = 1$, $C(\text{house with}) = 1$, $C(\text{with a}) = 1$.

From the above counts, $N_1 = 9$ and $N_2 = 1$. There are total of 10^2 bigrams possible, but only 10 are observed. So $N_0 = 10^2 - 10 = 90$. Remember that we take the training data set size to be the number of words, so $N = 12$.

We estimate probability of unseen bigrams according to Good-Turing. The bigram “mouse was” is unobserved, so

$$P_{GT}(\text{mouse was}) = (0 + 1) \frac{N_1}{N \cdot N_0} = \frac{9}{12 \cdot 90} = \frac{1}{120}$$

The other two bigrams are observed, so we use ML estimation for them.

$$P(\text{Herman lives}) = \frac{C(\text{Herman lives})}{N} = \frac{1}{12}$$

$$P(\text{a mouse}) = \frac{C(\text{a mouse})}{N} = \frac{2}{12} = \frac{1}{6}$$

3. (30%)

Suppose we have a sentence *To play pay*, and possible POS tags are as follows:

To: PREP CONJ	Play: NOUN VERB	pay: NOUN VERB
---------------	-----------------	----------------

Also, the log-likelihoods are as follows:

$$L(\text{to}|\text{PREP}) = 2 \quad L(\text{play}|\text{NOUN}) = 1 \quad L(\text{pay}|\text{NOUN}) = 4$$

$$L(\text{to}|\text{CONJ}) = 1 \quad L(\text{play}|\text{VERB}) = 2 \quad L(\text{pay}|\text{VERB}) = 2$$

$$L(\text{VERB}|\text{PREP}) = 10 \quad L(\text{NOUN}|\text{NOUN}) = 13$$

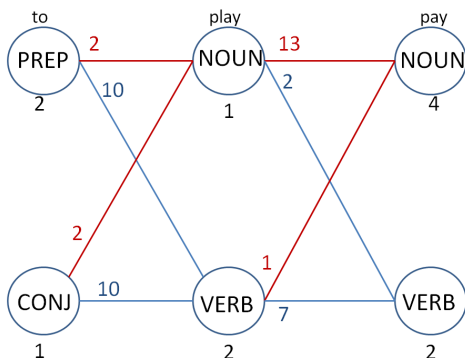
$$L(\text{NOUN}|\text{PREP}) = 2 \quad L(\text{VERB}|\text{NOUN}) = 2$$

$$L(\text{VERB}|\text{CONJ}) = 10 \quad L(\text{NOUN}|\text{VERB}) = 1$$

$$L(\text{NOUN}|\text{CONJ}) = 2 \quad L(\text{VERB}|\text{VERB}) = 7$$

- (a) (20%) For the given sentence, build the graph for solving the POS tagging problem via dynamic programming, like the one used in lecture notes. Show all the connections and the weight of all connections and nodes.

Solution:



- (b) (10%) Tag the sentence. You can do it by observation or by running dynamic programming. If you do it by observation, justify with a couple of words. If you do dynamic programming, show all the details.

Solution: It can be easily done by observation for this graph, but I will run the full DP algorithm:

Initialization:

$$C(\text{to}, \text{prep}) = 2, P(\text{to}, \text{prep}) = \text{null}; C(\text{to}, \text{conj}) = 1, P(\text{to}, \text{conj}) = \text{null}$$

Iteration 1:

$$C(\text{play}, \text{noun}) = 1 + \min\{C(\text{to}, \text{prep}) + 2, C(\text{to}, \text{conj}) + 2\} = 4, P(\text{play}, \text{noun}) = (\text{to}, \text{conj})$$

$$C(\text{play}, \text{verb}) = 2 + \min\{C(\text{to}, \text{prep}) + 10, C(\text{to}, \text{conj}) + 10\} = 13, P(\text{play}, \text{verb}) = (\text{to}, \text{conj})$$

Iteration 2:

$$C(\text{pay}, \text{noun}) = 4 + \min\{C(\text{play}, \text{noun}) + 13, C(\text{play}, \text{verb}) + 1\} = 18, P(\text{pay}, \text{noun}) = (\text{play}, \text{verb})$$

$$C(\text{pay}, \text{verb}) = 2 + \min\{C(\text{play}, \text{noun}) + 2, C(\text{play}, \text{verb}) + 7\} = 8, P(\text{pay}, \text{verb}) = (\text{play}, \text{noun})$$

The smallest cost path is 8, going through (pay,verb). Tracing back from (pay,verb) through the parent array “P” the tagging is “conj noun verb”

4. (30%)

(a) (20%)

Suppose inverted index has 3 terms and 3 documents, and the term/document counts are as follows:

$$doc_1 = \begin{bmatrix} term_1 & term_2 & term_3 \\ & 3 & 1 & 2 \end{bmatrix}, doc_2 = \begin{bmatrix} term_1 & term_2 & term_3 \\ & 4 & 3 & 1 \end{bmatrix},$$
$$doc_3 = \begin{bmatrix} term_1 & term_2 & term_3 \\ & 1 & 2 & 1 \end{bmatrix}.$$

Use term counts for the weights. User query is “ $term_1 term_2$ ”, i.e. it has one occurrence of $term_1$, and one occurrence of $term_2$. Using the cosine similarity, rank all the documents according to their similarity to the query, where rank 1 means the highest similarity, and rank 3 means the lowest similarity. Show all work.

Solution:

Let us first normalize the documents:

$$d_1 = \left[\frac{3}{\sqrt{14}}; \frac{1}{\sqrt{14}}; \frac{2}{\sqrt{14}} \right]$$
$$d_2 = \left[\frac{4}{\sqrt{26}}; \frac{3}{\sqrt{26}}; \frac{1}{\sqrt{26}} \right]$$
$$d_3 = \left[\frac{1}{\sqrt{6}}; \frac{2}{\sqrt{6}}; \frac{1}{\sqrt{6}} \right]$$

The query is $q = [1 \ 1 \ 0]$. Normalizing it for length (it is actually not necessary for relative ranking, but still let's do it):

$$q = \left[\frac{1}{\sqrt{2}}; \frac{1}{\sqrt{2}}; 0 \right]$$

Now, cosine similarity is:

$$\cos(d_1, q) = d_1 \cdot q = \frac{3}{\sqrt{14}} \cdot \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{14}} \cdot \frac{1}{\sqrt{2}} = \frac{4}{\sqrt{28}} \approx 0.76$$

$$\cos(d_2, q) = d_2 \cdot q = \frac{4}{\sqrt{26}} \cdot \frac{1}{\sqrt{2}} + \frac{3}{\sqrt{26}} \cdot \frac{1}{\sqrt{2}} = \frac{7}{\sqrt{52}} \approx 0.97$$

$$\cos(d_3, q) = d_3 \cdot q = \frac{1}{\sqrt{6}} \cdot \frac{1}{\sqrt{2}} + \frac{2}{\sqrt{6}} \cdot \frac{1}{\sqrt{2}} = \frac{3}{\sqrt{12}} \approx 0.89$$

Therefore the ranking is d_2, d_3, d_1 .

(b) (10%) Suppose inverted index has 3 terms and 3 documents, and the term/document counts are as follows:

$$doc_1 = \begin{bmatrix} term_1 & term_2 & term_3 \\ 100 & 10 & 1000 \end{bmatrix}, doc_2 = \begin{bmatrix} term_1 & term_2 & term_3 \\ 10 & 0 & 100 \end{bmatrix},$$
$$doc_3 = \begin{bmatrix} term_1 & term_2 & term_3 \\ 10000 & 10 & 0 \end{bmatrix}.$$

Suppose you decide to use tfidf term weighting.

- What is the weight for term 2 in document 2?
- What is the weight for term 3 in document 2?

You don't have to work out the expressions completely, you can leave answer in the format $\frac{40}{3} \log(4/5)$.

Solution:

In case $tf_{td} = 0$, the weight of term t in document d is defined as 0. Otherwise, it is defined as:

$$w_{td} = (1 + \log tf_{td} \times \log (N/df_t))$$

Since the count of term 2 in document 2 is 0, $w_{22} = 0$.

The count of term 3 in document 2 is 100. The document frequency of term 3 is 2. The total number of documents is 3. Therefore,

$$w_{32} = (1 + \log 100 \times \log (3/2)) = 3 \log (3/2)$$