

Student ID:

CS 442b-542b

Name:

Short Exam 2

Instructions: Show all the work you do. Use the back of the page, if necessary. Calculators are allowed, laptops are not allowed.

Problem 1 (20%): The training data consists of the text: “I like to play with cats, and cats like to chase my toy mouse. Cats like me, and I like cats.” Ignore the sentence punctuation for this problem, that is treat the training text as one continuous sentence. Also ignore capitalization, that is “cats” and “Cats” should be treated as the same word. Suppose that we use trigram language model with maximum likelihood estimation of parameters.

(a) Compute the probability of a sentence “Cats like to play.” Make sure to show all the work.

Solution:

$$P(\text{Cats like to play}) = P(\text{play} \mid \text{cats like to}) * P(\text{to} \mid \text{cats like}) * P(\text{like} \mid \text{cats}) * P(\text{cats}) \\ \approx P(\text{play} \mid \text{like to}) * P(\text{to} \mid \text{cats like}) * P(\text{like} \mid \text{cats}) * P(\text{cats});$$

Let $C(\text{like to})$ be the count of how many times where “cats like” occurs in training data:

$$P(\text{play} \mid \text{like to}) = P(\text{like to play}) / P(\text{like to}) = C(\text{like to play}) / C(\text{like to}) = 1/2 = 0.5;$$

$$P(\text{to} \mid \text{cats like}) = P(\text{cats like to}) / P(\text{cats like}) = C(\text{cats like to}) / C(\text{cats like}) = 1/2 = 0.5;$$

$$P(\text{like} \mid \text{cats}) = P(\text{cats like}) / P(\text{cats}) = C(\text{cats like}) / C(\text{cats}) = 2/4 = 0.5$$

$$P(\text{cats}) = C(\text{cats}) / N = 4/21 = 0.1905$$

$$P(\text{Cats like to play}) = 0.5 * 0.5 * 0.5 * 0.1905 = 0.0238;$$

(b) Compute the probability of a sentence “Black cat jumped higher than any other cat”

Solution: $P(\text{Black cat jumped higher than any other cat})$

$$\approx P(\text{cat} \mid \text{any other}) * P(\text{other} \mid \text{than any}) * P(\text{any} \mid \text{higher than}) \dots * P(\text{black});$$

$$P(\text{black}) = C(\text{black}) / N = 0, \text{ therefore, } P(\text{black cat jumped higher than any other cat}) = 0$$

Problem 2 (10%): Suppose you have the following training data: “I like to play with cats, and cats like me”. You decide to use add 1 smoothing (Laplace’s law). What is the probability of nGram “a black cat jumps”. Assume that vocabulary size is 10.

$$\text{Solution: } P(\text{a black cat jumps}) = [C(\text{a black cat jumps}) + 1] / (N + B);$$

$C(\text{a black cat jumps}) = 0$; “a black cat jumps” is a 4-gram, and there are seven 4-grams in the training data. So $N = 7$;

B is the number of faked 4-grams, $B = 10^4$

$$P(\text{a black cat jumps}) = [C(\text{a black cat jumps}) + 1] / (N+B) = 1 / 10007 = 0.00009993$$

Problem 3 (10%): We have the following training text tagged with parts of speech: “The\DET time\NOUN flies\VERB very\ADV fast\ADV. Fast\ADJ flies\NOUN move\VERB. Time\VERB the\DET flies\NOUN fast\ADV!. Flies\NOUN do\VERB not\ADV fast\NOUN.” Tag the phrase “Fast flies” using Charniak’s tagger.

Solution: Let $C(\text{fast}\backslash\text{NOUN})$ be the number of times that “fast” is tagged as a noun:
 $C(\text{fast}\backslash\text{NOUN}) = 1$; $C(\text{fast}\backslash\text{VERB}) = 0$; $C(\text{fast}\backslash\text{DET}) = 0$; $C(\text{fast}\backslash\text{ADJ}) = 1$; $C(\text{fast}\backslash\text{ADV}) = 2$,
therefore the tag for “fast” is “\ADV”;
 $C(\text{files}\backslash\text{NOUN}) = 3$; $C(\text{files}\backslash\text{VERB}) = 1$; $C(\text{files}\backslash\text{DET}) = 0$; $C(\text{files}\backslash\text{ADJ}) = 1$; $C(\text{files}\backslash\text{ADV}) = 2$,
therefore the tag for “files” is “\NOUN”, so the final tagging is **ADV NOUN**

Problem 4(10%): Using the same training text as in problem 3, estimate $P(\text{VERB} \mid \text{NOUN})$ using MLE.

Solution: Let $C(\text{NOUN})$ be the number of times tag “NOUN” occurs in training data, $C(\text{NOUN VERB})$ be the number of times that a “VERB” follows a “NOUN”:

$$P(\text{VERB} \mid \text{NOUN}) = P(\text{NOUN VERB})/P(\text{NOUN}) = C(\text{NOUN VERB})/C(\text{NOUN}) = 3/5 = 0.6$$

Problem 5 (25%): Suppose you are using Markov model POS tagger, and for the i th word in the sentence you are tagging, the possible part of speech labels are “ADV”, “ADJ”, “NOUN”. For the $(i+1)$ th word in the sentence, the possible POS tags are “NOUN” and “VERB”. You have run dynamic programming for i iterations, and at the i th iteration, you have computed:

$$C(\text{word } i, \text{ADJ}) = 4, C(\text{word } i, \text{ADV}) = 3, C(\text{word } i, \text{NOUN}) = 9. \text{ Furthermore,}$$

$L(\text{NOUN} \text{ADJ}) = 1$	$L(\text{VERB} \text{ADJ}) = 10$	$L(\text{word } i+1 \text{VERB}) = 1$
$L(\text{NOUN} \text{ADV}) = 10$	$L(\text{VERB} \text{ADV}) = 2$	$L(\text{word } i+1 \text{NOUN}) = 3$
$L(\text{NOUN} \text{NOUN}) = 3$	$L(\text{VERB} \text{NOUN}) = 1$	

(a) (20%) After iteration $i+1$, what are the values of:

$C(\text{word } i+1, \text{VERB})$

$C(\text{word } i+1, \text{NOUN})$

$P(\text{word } i+1, \text{VERB})$

$P(\text{word } i+1, \text{NOUN})$

Solution:

$$C(\text{word } i+1, \text{VERB}) = \min\{C(\text{word } i, t) + L(\text{VERB}|t)\} + L(\text{word } i+1 | \text{VERB})$$

$$C(\text{word } i, \text{ADJ}) + L(\text{VERB} | \text{ADJ}) = 4 + 10 = 14;$$

$$C(\text{word } i, \text{ADV}) + L(\text{VERB} | \text{ADV}) = 3 + 2 = 5;$$

$$C(\text{word } i, \text{NOUN}) + L(\text{VERB} | \text{NOUN}) = 9 + 1 = 10;$$

$$\text{Therefore, } C(\text{word } i+1, \text{VERB}) = C(\text{word } i, \text{ADV}) + L(\text{VERB} | \text{ADV}) + L(\text{word } i+1 | \text{VERB}) = 5+1 = 6, P(\text{word } i+1, \text{VERB}) = \text{ADV};$$

$$C(\text{word } i+1, \text{NOUN}) = \min\{C(\text{word } i, t) + L(\text{NOUN}|t)\} + L(\text{word } i+1 | \text{NOUN})$$

$$C(\text{word } i, \text{ADJ}) + L(\text{NOUN} | \text{ADJ}) = 4 + 1 = 5;$$

$$C(\text{word } i, \text{ADV}) + L(\text{NOUN} | \text{ADV}) = 3 + 10 = 13;$$

$$C(\text{word } i, \text{NOUN}) + L(\text{NOUN} | \text{NOUN}) = 9 + 3 = 12;$$

$$\text{Therefore, } C(\text{word } i+1, \text{NOUN}) = C(\text{word } i, \text{ADJ}) + L(\text{NOUN} | \text{ADJ}) + L(\text{word } i+1, \text{NOUN}) = 5+3 = 8, P(\text{word } i+1, \text{NOUN}) = \text{ADJ};$$

(b) (5%) Suppose now that the length of the sentence is $k > i+1$. After $(i+1)$ th iteration, can you tell which part of speech the $(i+1)$ th word will get tagged with by the Markov tagger? If yes, which one, if no, why not?

Solution: We can not know any of the tags. Let’s consider word j , where $j < k$. If we want to know the tag of word j , we should first know the tag of word $j+1$, say TAG $J+1$, then $P(\text{word } j+1, \text{TAG } J+1)$ is the tag of word j . If we want know the tag of word $j+1$, we must know the tag of word $j+2$, therefore, we can not know any of the tags until we find the tag for the last word of the sentence.

Problem 6(25%) : Suppose inverted index has 4 terms and 2 documents, and the term/document counts are as follows:

$$document1 = \begin{bmatrix} term1 & term2 & term3 & term4 \\ 0 & 2 & 0 & 9 \end{bmatrix} \quad document2 = \begin{bmatrix} term1 & term2 & term3 & term4 \\ 3 & 0 & 4 & 6 \end{bmatrix}$$

(a) (15%) Suppose that we use term counts for the weights. User query is “term3 term1 term2”, that is it has one occurrence of term3, 1 occurrence of term1, and one occurrence of term2. Using the cosine similarity, if we are only asked to retrieve one document, should it be document 1 or 2? To simplify calculations, **do not normalize** vectors to length 1 for this question. Show all the work you do.

Solution: Let d_1 be the vector for document 1, d_2 be the vector for document 2, q be the query vector. Now we have:

$$d_1 = [0 \ 2 \ 0 \ 9], \quad d_2 = [3 \ 0 \ 4 \ 6], \quad q = [1 \ 1 \ 1 \ 0]$$

$$\begin{aligned} \cos(d_1, q) &= (d_1 \cdot q) / (\|d_1\| \times \|q\|) = (0 \cdot 1 + 2 \cdot 1 + 0 \cdot 1 + 9 \cdot 0) / (\sqrt{4+81} \cdot \sqrt{1+1+1}) \\ &= (2) / (15.97) = 0.1252 \end{aligned}$$

$$\begin{aligned} \cos(d_2, q) &= (d_2 \cdot q) / (\|d_2\| \times \|q\|) = (3+4) / (\sqrt{9+16+36} \cdot \sqrt{1+1+1}) \\ &= 0.52 \end{aligned}$$

So document 2 is closer to the query vector. Return document 2;

(b) (10%) Which term will have the highest weight in document 2 if instead of frequencies we use $tf \cdot idf$ for the weights? What is the value of this highest weight?

Solution:

Let $w(1,2)$, $w(2,2)$, $w(3,2)$, $w(4,2)$ be the weight of term 1, 2, 3 and 4 in document 2. Let $tf(1,2)$, $tf(2,2)$, $tf(3,2)$ and $tf(4,2)$ be the frequency of term 1, 2, 3 and 4 in document 2. M is the number of document, that is 2. Let $df(1)$, $df(2)$, $df(3)$ and $df(4)$ be the numbers of documents which contain term 1, 2, 3 and 4.

$$w(1,2) = tf(1,2) * idf(M/df(1)) = 3 * \log(2/1) = 3 * \log 2$$

$$w(2,2) = tf(2,2) * idf(M/df(2)) = 0 * \log(2/1) = 0 ;$$

$$w(3,2) = tf(3,2) * idf(M/df(3)) = 4 * \log(2/1) = 4 * \log 2$$

$$w(4,2) = tf(4,2) * idf(M/df(4)) = 6 * \log(2/2) = 0;$$

The largest weight is $w(3,2) = 4 * \log 2$