

Due Date: *Friday, February 17, 11:55PM Eastern*

Let $P(w)$ be the probability of w :

As an example, given the following corpus: “The Cat in the Hat is the best cat in the hat”, a **unigram** LM language model would be: (using fractions for clarity)

$$P(\text{the}) = 4/12$$

$$P(\text{cat}) = 2/12$$

$$P(\text{in}) = 2/12$$

$$P(\text{hat}) = 2/12$$

$$P(\text{is}) = 1/12$$

$$P(\text{best}) = 1/12$$

For unigrams, the probability of ‘cat’ appearing anywhere in the corpus is 2/12 using maximum likelihood estimation MLE (a.k.a. **word count**) - note this is a very simplistic model – the closed universe model.

A bigram (n-gram, n=2) LM:

$$P(\text{the cat}) = 1/8$$

$$P(\text{cat in}) = 2/8$$

$$P(\text{in the}) = 2/8$$

$$P(\text{the hat}) = 2/8$$

$$P(\text{hat is}) = 1/8$$

$$P(\text{is the}) = 1/8$$

$$P(\text{the best}) = 1/8$$

$$P(\text{best cat}) = 1/8$$

However, most likely is that we are not interested in the probability of the phrase, but in the conditional probability of ‘cat’ given that the word ‘the’ has been seen. In our example, the probability of ‘cat’ given ‘the’, $P(\text{cat}|\text{the})$, is given by Bayes theorem:

$$P(B \text{ given } A) = P(A \text{ and } B) / P(A)$$

In this homework, let’s approximate this using the closed-corpus assumption (no unseen words exist, so no smoothing for those statisticians in class):

$$P(\text{cat}|\text{the}) = P(\text{the cat}) / P(\text{the}) = (1/8) / (4/12) = 0.375$$

Homework Rules:

- punctuation does NOT count; so the words is ‘(1991)’ and ‘1991’ are the same.
You must parse your input: replace **all characters not in this set: [a-z, A-Z, 0-9]** with spaces.
- all text should be normalized to lowercase
- Ignore lines with less than 3 words.
- Input should be lines of text (separated by new line and/or carriage return)

Input for this problem: **hw1dir1.zip** (provided in class website)

Write your own code in your language of choice, but your code MUST be Hadoop MapReduce. For Python, use Hadoop streaming. Submit the result and code.

Solve: Compute the Bigram LM, of the form “word1, word2”, $P(\text{word2} \mid \text{word1})$

$$P(w_2 \mid w_1) = P(w_1 w_2) / P(w_1)$$
$$= \frac{\frac{\text{count}(w_1 w_2)}{\text{total}(\text{bigrams})}}{\frac{\text{count}(w_1)}{\text{total}(\text{unigrams})}}$$

- Your input is lines of text.
- Hadoop programs can individual files, directories (which are recursed) or simple wildcard patterns to match multiple files.
- You can/could use multiple jobs!
-

IMPORTANT: Solving this by *a single map/reduce* program is very difficult. Why? Recall the discussion in class: mapper and reducers cannot hold/keep state across object instances. Your mapper is not guaranteed to exist past a single *input split*, and your reducer is only guaranteed to exist for one 1 (unless you use a single reducer).

Hints:

- Compute 4 things: unigrams counts, bigrams counts, total unigrams, total bigrams
- State does not save: you cannot count totals unless you use Hadoop global counters...
- You define what the keys and values of the mapper.....and the values.....
- you can daisy chain jobs

3. EXTRA CREDIT: 25 points

Using the solution of problem 2, find the x below (print the word in this sequence with the highest probability)

united states ____x____

For $P(x \mid \text{united states}) = p$, find the x with the highest p