

LANGUAGE MODELS

Francois Role - francois.role@parisdescartes.

October 2018

1 Language Models - Co-occurrence

1.1 N-gram Language Models

A **language model** is designed to measure how likely it is that a sequence would be uttered. This is useful in many settings. One are where language model are widely is speech recognition.

More precisely, a language model is a function that takes a word sequence and returns a probability that it would be uttered.

A good language model will assign a greater probability to the sequence "le ciel est bleu" than to the sequence "est le bleu ciel".

In theory, given a sequence of words $W = w_1 w_2 \cdots w_n$ its probability can be estimated as :

$$p(w_1 w_2 \cdots w_n) = p(w_1) p(w_2 | w_1), \cdots, p(w_n | w_1 w_2 \cdots w_{n-1}) \quad (1)$$

$$= \prod_{i=1}^n p(w_i | w_1^{i-1}) \quad (2)$$

In practice, because of **data sparsity**, the history of previous words has to be limited to m words using the **Markov assumption** that considering only a small number of previous m words is enough to determine the probability of the next words. So, we have :

$$p(w_1 w_2, \cdots, w_n) = p(w_1) p(w_2 | w_1), \cdots, p(w_n | w_{n-m} \cdots w_{n-1}) \quad (3)$$

$$= \prod_{i=1}^n p(w_i | w_{i-m}^{i-1}) \quad (4)$$

A model based on a two-word history is called a **trigram model**, where statistics are collected over sequences of three words called **trigrams**: a two-word history is used to predict a third word. A model based on a one-word history is a **bigram model**. If a model relies on single words the language model is a **unigram model**.

So, for a N-gram model the probability of a sequence is :

$$p(w_1 w_2, \cdots, w_n) = \prod_{i=1}^n p(w_i | w_{i-N+1}^{i-1})$$

1.2 Estimating the conditional probabilities

N-gram probabilities are estimated by dividing the observed frequency of a particular sequence by the observed frequency of a prefix of this sequence. The ratio between these frequency counts is called a **relative frequency**.

Here is an example of how to estimate a conditional probability for a trigram model:

$$p(w_3|w_1w_2) = \frac{\text{count}(w_1w_2w_3)}{\text{count}(w_1w_2)} \quad (5)$$

$$= \frac{\text{count}(w_1w_2w_3)}{\sum_w \text{count}(w_1w_2w)} \quad (6)$$

More generally, the formula for estimating N-gram probabilities (that is, using an history of $m = N - 1$ words is:

$$p(w_i|w_{i-N+1}^{i-1}) = \frac{\text{count}(w_{i-N+1}^i)}{\text{count}(w_{i-N+1}^{i-1})} \quad (7)$$

$$= \frac{\text{count}(w_{i-N+1}^i)}{\sum_w \text{count}(w_{i-N+1}^{i-1}w)} \quad (8)$$

However, basing the estimations on raw frequency counts has severe limitations.

Suppose we have to estimate the probability of the sequence A B C D E using a bigram language model and the corpus from which the model was built did not include the sequence B C. So our estimation for $p(C|B)$ was $\frac{\text{count}(B,C)}{\sum_w \text{count}(B,w)} = \frac{0}{\sum_w \text{count}(B,w)} = 0$ and so the predicted probability for A B C D E will be 0.

More generally, all unseen n-grams will be assigned a probability of 0 which is a very hard decision. We want to be able to generalize in presence of unseen data.

1.3 Smoothing Techniques

Add-one (Laplace) smoothing To avoid zero counts and generalize better, we pretend that there exists an additional document where each sequence of length N appears once.

$$p(w_i|w_{i-N+1}^{i-1}) = \frac{1 + \text{count}(w_{i-N+1}^i)}{\sum_w [1 + \text{count}(w_{i-N+1}^{i-1}w)]} = \frac{1 + \text{count}(w_{i-N+1}^i)}{|V| + \sum_w \text{count}(w_{i-N+1}^{i-1}w)}$$

For a bigram model, we have:

$$p(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1}w_i) + 1}{\sum_w [1 + \text{count}(w_{i-1}w)]} = \frac{\text{count}(w_{i-1}w_i) + 1}{|V| + \sum_w \text{count}(w_{i-1}w)}$$

where V is the vocabulary.

The limit of this additive smoothing is that it treats unseen bigrams involving frequent words and unseen bigrams involving rare words on the same footing.

More Advanced Smoothing Techniques Kneser-Ney and Good-Turing smoothing are more complex methods that allow to estimate the count of things we've never seen in a more principled way.

Exercise

(a) According to a bigram model trained on the following sentences (and without using any smoothing technique):

I saw the boy the girl is in the garden the boy is funny the girl is funny the baby is smiling the girl is smiling

1. What is the probability $p(\text{funny}|\text{is})$?
2. What is the most probable word after the sequence "the boy is" ?
3. Rank the following sequences based on their probabilities ? I saw the smiling the bird
I saw the baby is funny
4. What is the probability of I saw the smiling, this time using Add-one smoothing.

1.4 Measuring the Quality of a Language Model

Task-based Evaluation To compare two models A and B , a way is to assess (e.g. accuracy) how each performs on a given task. It requires designing, setting up and running a task. A cheaper way is to compute some quality measures directly from the model.

Cross-entropy and Perplexity Let L be a language. Using a training set T drawn from L , we estimate the unknown probability distribution p associated with L . This results in an estimated probability distribution q . To sum up p is the true (unknown) distribution of the words in L and q is the distribution of words estimated from the training corpus.

To evaluate q we assess how well it predicts (what probability it gives) to a separate test sample x_1, x_2, \dots, x_n also drawn from L . The better q the higher probabilities it should assign to the test events.

The perplexity of the model q is defined to be:

$$2^{-\sum_{i=1}^n \frac{1}{N} \log_2 q(x_i)}$$

The better q , the higher probabilities it should assign to the test events, so the lower perplexity. How is this quantity derived? We will show it in the case of a bigram model.

Let T be a (improperly named) **test set**, that is a sequence of tokens w_1, w_2, \dots, w_n , and let q be a probability distribution trained from a training set. For a bigram model, the likelihood of the test set T according to q is:

$$L = \prod_{i=1}^n q(w_i | w_{i-1}) \quad (9)$$

We want to minimize the opposite of the log-likelihood $\sum_{i=1}^n -\log q(w_i | w_{i-1})$, or, equivalently, the mean of this value :

$$\frac{1}{n} \sum_{i=1}^n -\log q(w_i|w_{i-1}) = \frac{1}{n} \left(\sum_{i=1}^n \log \frac{1}{q(w_i|w_{i-1})} \right) \quad (10)$$

This quantity is precisely the exponent in the perplexity formula. It also turns out that it equals the geometric mean of the inverses of the probabilities) is:

$$\log \sqrt[n]{\prod_{i=1}^n \frac{1}{q(w_i|w_{i-1})}} = \log \left(\prod_{i=1}^n \frac{1}{q(w_i|w_{i-1})} \right)^{1/n} \quad (11)$$

$$= \frac{1}{n} \sum_{i=1}^n \log \frac{1}{q(w_i|w_{i-1})} \quad (12)$$

$$= \frac{1}{n} \sum_{i=1}^n -\log q(w_i|w_{i-1}) \quad (13)$$

So, the exponent $-\sum_{i=1}^n \frac{1}{N} \log_2 q(x_i)$ found in the perplexity formula is the log of the geometric mean of the **word perplexities** (inverse predicted probabilities for the words). This exponent $-\sum_{i=1}^n \frac{1}{N} \log_2 q(x_i)$ is in fact an approximation of the cross-entropy of a (long-enough) sequence x_1, \dots, x_n according to the Shannon-McMillan-Breiman theorem. This is why some say that "the cross-entropy is the logarithm of perplexity".

To summarize, the perplexity of a sequence x_1, \dots, x_n can be computed as :

$$2^{-\sum_{i=1}^n \frac{1}{N} \log_2 q(x_i)} = 2^{\log \sqrt[n]{\prod_{i=1}^n \frac{1}{q(w_i|w_{i-1})}}} = \sqrt[n]{\prod_{i=1}^n \frac{1}{q(w_i|w_{i-1})}}$$

Exercise

(b)

Using the model of exercise (a), compute the perplexity of the sequence:

I saw the baby

Aside Note:

The fact that the cross-entropy is the logarithm of perplexity can also be demonstrated in the context of logistic regression.

For one instance and K classes the cross-entropy error is:

$$-\sum_{k=1}^K t_k \log y_k = -t_k^* \log y_k^* = -\log y_k^*$$

where y_k^* is the probability computed by the model for the correct word.

Perplexity is the inverse predicted probability for the correct word:

$$\frac{1}{y_k^*}$$

Since $-\log y_k^* = \log \frac{1}{y_k^*}$ the cross-entropy is the logarithm of perplexity.

Minimizing the arithmetic mean of the cross-entropy is equivalent to minimizing the log of the geometric mean of the perplexity.

$$\begin{aligned}
& \frac{(-\log y_{1*}) + \dots + (-\log y_{n*})}{n} = \\
& \frac{\log \frac{1}{y_{1*}} + \dots + \log \frac{1}{y_{n*}}}{n} \\
& \frac{\log(\frac{1}{y_{1*}} \times \dots \times \frac{1}{y_{n*}})}{n} \\
& \frac{1}{n} \log(\frac{1}{y_{1*}} \times \dots \times \frac{1}{y_{n*}}) \\
& \log \sqrt[n]{\frac{1}{y_{1*}} \times \dots \times \frac{1}{y_{n*}}}
\end{aligned}$$

2 Statistical Measures of Word Association

**** Positive Pointwise Mutual Information**** (PMI) has been initially designed to help find interesting **collocations** in a corpus.

2.1 PMI for Finding if a Bigram is a collocation

For two word x and y , PMI is defined as:

$$pmi(x, y) = \log \frac{p(x, y)}{p(x)p(y)}$$

and is estimated as:

$$\hat{pmi}(x, y) = \frac{count(x, y)N}{count(x)count(y)}$$

where N is the total number of tokens in the training corpus, and $count(x, y)$, $count(x)$, $count(y)$ are the number of occurrences of the bigram xy , of the word x and of the word y resp.

Note that in this case PMI is not symmetric.

2.2 PMI for Finding General Syntagmatic Relations

By **general syntagmatic relation between two words**, we mean that if we see some word occur in some context we tend to see the other word.

The context can be a sentence, a paragraph, a document, etc. Here $count(x)$, $count(y)$, and $count(x, y)$ represent the number of contexts that contain x , y and both x and y resp. and N is the total number of contexts.

In fact, today, PMI is not only used to spot exact collocations but, more generally to discover pairs of words tending to occur within a **window** of n words of each other. The counts are then retrieved from a **word-word matrix** where the cell i, j represents the number of times word i occurs in a window of $\pm n$ words around word j .

Today, PMI is widely applied in NLP to produce semantic representations for words and documents. Recently, connections have been made between PMI and **word embeddings**.

Other techniques, inspired from the hypothesis testing field are also used (see exercise 2 at the end of this section).

2.3 Limitations of PMI

PMI has several shortcomings:

- Its values are not bounded: they range from negative to positive infinity.
- Negative PMI values are difficult to interpret.
- it is biased towards bigrams involving rare words.

These problems are partially dealt with as follows. First, to get rid of negative values, one can use a variant called **Positive PMI** (PPMI) instead of PMI. PPMI is defined as:

$$PPMI(x, y) = \max(\log \frac{p(x, y)}{p(x)p(y)}, 0)$$

To have values in 0, 1 one can use the NPMI defined as :

$$\frac{PMI(x, y)}{-\log p(x, y)}$$

To overcome the bias problem, one can either raise the dividend $p(x, y)$ to a power greater than 0 (typically 2 or 3). This has the effect of boosting frequent pairs. Alternatively, one can raise the terms in the denominator to a power slightly below 1 (typically 0.75). This will increase the probability assigned to rare words compared to frequent words, and thus diminish the PMI values of pairs involving rare words.

Exercise

(a) Given the following message recently received from the Z μ 125 α planet
'aaa bbb aaa ccc eee ddd eee fff aaa eee fff eee aaa ccc fff eee bbb aaa

1. In order to see if *aaa bbb* or *aaa ccc* are collocations, compute $pmi(aaa, bbb)$ and $pmi(aaa, ccc)$ and decide which of the two bigrams is more likely to be a collocation.
2. Could the bigram *bbb aaa* play a role if we were interested not in collocations but general syntagmatic relations between *aaa* and *bbb* ?
3. (bonus question) What is the meaning of the message?

(b) Given the following table showing the dependence of occurrences between the words *text* and *information*, The chi-square test can be used to decide if *text information* is likely to be a collocation.

	$w_1 = text$	$w_1 \neq text$
$w_2 = information$	8	4667
$w_2 \neq information$	15820	14287173

There are 8 occurrences of *text information*, 4667 bigrams ending with *information* but not beginning with *text*, 15820 bigrams beginning with *text* and ending with a word different from *information*, and 14287173 bigrams that contain neither words.

1. Since the χ^2 statistics is $\sum_{i,j} \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$. For example what is the value of the expected frequency $e_{1,1}$?
2. The test statistics has a simple form for 2-by-2 tables :

The number of tokens multiplied by the squared determinant of the table, divided by the product of the sums of all lines and columns.