

CPSC 532P / LING 530A: Deep Learning for Natural Language Processing (DL-NLP)

Muhammad Abdul-Mageed

muhammad.mageed@ubc.ca

Natural Language Processing Lab

The University of British Columbia

Table of Contents

1 Language Models

- Motivation
- N-gram Language Models
- Evaluation
- Words as Vectors

Motivations

What is more probable?

- Brewing a great cup of **coffee/tea/espresso** ...
- Brewing a great cup of **entropy** ...

Language Models

- **Assign probabilities to sequences of words** (or characters, etc.).
- Play a **very significant role in NLP**.
- **Classically**, they are motivated by usage in tasks like **handwriting recognition, spelling correction, speech recognition, and machine translation**.
- **Recently**, their applications are exploding, including because of being an important block in learning (contextualized) word vectors (**referred to as unsupervised pre-training or generative pre-training**) ...

Statistical Language Models

- **A statistical LM:** The conditional probability of the next word given all the previous words.
- **Note:** We can also train **bidirectionally** (as in ELMo), or **mask words** (as in BERT). More later...

1: Statistical LM

$$\hat{P}(W_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1})$$

Brewing a great cup of coffee

- coffee (w_t)
- Brewing a great cup of (w_1^{t-1})
- Brewing a great cup of coffee (W_1^T)

- **Markov assumption:** We do not have to look too far in the past

2: n-gram LM

$$\hat{P}(w_t | w_1^{t-1}) \approx \hat{P}(w_t | w_{t-N+1}^{t-1})$$

3: Bigram LM (one word in the past)

$$\hat{P}(w_t | w_{t-1}) = \hat{P}(w_t | w_{t-1})$$

A Simple LM

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

Here are the calculations for some of the bigram probabilities from this corpus

$$P(I|<s>) = \frac{2}{3} = .67 \quad P(\text{Sam}|<s>) = \frac{1}{3} = .33 \quad P(\text{am}|I) = \frac{2}{3} = .67$$

$$P(</s>|\text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam}|\text{am}) = \frac{1}{2} = .5 \quad P(\text{do}|I) = \frac{1}{3} = .33$$

Figure: A simple LM on a mini-corpus. [From J&M, 2017, ch. 3]

Computing $P(\text{sequence})$

$$\begin{aligned} P(i|\langle s \rangle) &= 0.25 & P(\text{english}|\text{want}) &= 0.0011 \\ P(\text{food}|\text{english}) &= 0.5 & P(\langle /s \rangle|\text{food}) &= 0.68 \end{aligned}$$

Now we can compute the probability of sentences like *I want English food* or *I want Chinese food* by simply multiplying the appropriate bigram probabilities together, as follows:

$$\begin{aligned} P(\langle s \rangle \text{ i want english food } \langle /s \rangle) \\ &= P(i|\langle s \rangle)P(\text{want}|i)P(\text{english}|\text{want}) \\ &\quad P(\text{food}|\text{english})P(\langle /s \rangle|\text{food}) \\ &= .25 \times .33 \times .0011 \times 0.5 \times 0.68 \\ &= .000031 \end{aligned}$$

Figure: Computing $P(\text{sequence})$ with an LM. [See details in J&M, 2017, ch. 3]

Notes on Computing Probs from LM

- We represent LM probabilities as **log probabilities**
- This **avoids numerical underflow** (if we were to use raw format)
- **Adding in log space = multiplying in linear space**
- To convert back, we exponentiate:

4: Log Probabilities

$$p_1 * p_2 * p_3 = \exp(\log_2(p_1) + \log_2(p_2) + \log_2(p_3))$$

Intrinsic vs. Extrinsic Evaluation

- Many models in NLP can be evaluated **intrinsically** and **extrinsically**
- **Intrinsic evaluation:** We use **perplexity**
- **Extrinsic evaluation:** Plugging LMs in an application like a **speech recognizer** or **MT system** is best method
- Again, modern language models are everywhere. Think about **ELMo** and **BERT**, for example.
- We need to split our data, possibly into **80% train**, **10% dev**, and **10% test**.
- So, **don't train your LM on data that is part of your downstream task**. (Some popular papers unfortunately trained an LM on Wikipedia and tested on tasks whose data come from Wikipedia.)

Perplexity as a Branching Factor

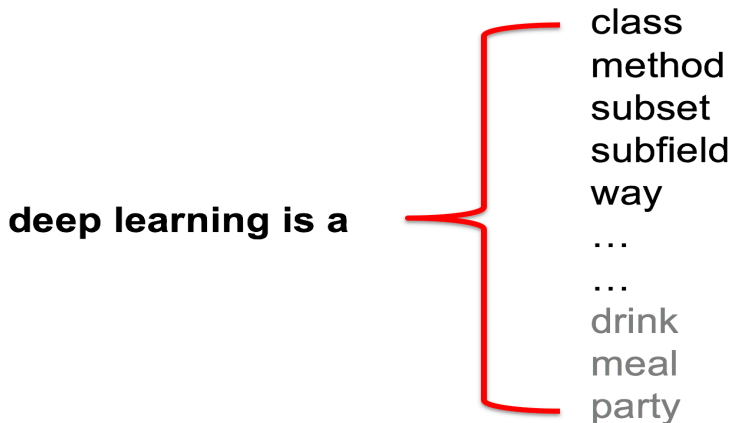


Figure: Perplexity is a branching factor. It will tell us how many words can follow a given word. Lower is better. Why?

Perplexity

- **Perplexity** is the **weighted average branching factor**
- It is the **inverse probability on unseen data**, normalized by the number of words (for word-level perplexity).
- **Our goal is to maximize probability of unseen data.**
- Minimizing perplexity = maximizing probability.
- Perplexity is closely related to **entropy** (**recall: entropy is the avg amount of info.**)

What's a good LM?

One that best predicts unseen data (test set). It's one that gives **highest $P(S)$** for each **sentence S** in test.

5: Perplexity

$$\begin{aligned} PP(W) &= \hat{P}(w_1, w_2 \dots w_T)^{-\frac{1}{T}} \\ &= \sqrt[T]{\frac{1}{\hat{P}(w_1, w_2 \dots w_T)}} \end{aligned}$$

6: Chain Rule

$$PP(W) = \sqrt[T]{\prod_{i=1}^T \frac{1}{\hat{P}(w_i | w_1 \dots w_{i-1})}}$$

7: Perplexity for Bigrams

$$PP(W) = \sqrt[T]{\prod_{i=1}^T \frac{1}{\hat{P}(w_i|w_{i-1})}}$$

Sample Generation From Shakespeare

1 gram	-To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have
2 gram	-Hill he late speaks; or! a more to leg less first you enter
3 gram	-Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.
4 gram	-What means, sir. I confess she? then all sorts, he is trim, captain.
5 gram	-Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.
6 gram	-This shall forbid it should be branded, if renown made it empty.
7 gram	-King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;
8 gram	-It cannot be but so.

Figure 3.3 Eight sentences randomly generated from four n-grams computed from Shakespeare's works. All characters were mapped to lower-case and punctuation marks were treated as words. Output is hand-corrected for capitalization to improve readability.

Figure: [From J&M, 2017, ch. 3]

Sample Generation From WSJ

1 gram	Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives
2 gram	Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her
3 gram	They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

Figure 3.4 Three sentences randomly generated from three n-gram models computed from 40 million words of the *Wall Street Journal*, lower-casing all characters and treating punctuation as words. Output was then hand-corrected for capitalization to improve readability.

Figure: [From J&M, 2017, ch. 3]

Smoothing

- We will see words at test time that we have not seen in **train**
- Can't assign these **zero probability** (otherwise we'll have **division by zero!**)
- Called **out-of-vocabulary (OOV)**
- We have to do some **smoothing**, e.g.:
 - Laplace smoothing (add-one)
 - Add-K smoothing
 - back-off and interpolation smoothing
 - Kneser-Ney smoothing
 - ...
- For details, see J&M (2017, ch. 03) ...

Problems with n-gram LM

- Limited to a short sub-sequence (e.g., 2 words, for $n=3$).

Data scarcity

LMs for larger n-gram suffer from data scarcity.

- Does not easily capture word "similarity".

Challenge with word similarity

Consider the following two sentences where "cat" and "dog" have similar semantic and grammatical roles. [Bengio et al., 2003].

- 1 "The cat is walking in the bedroom"
- 2 "A dog was running in a room"

- Bengio et al. (2003, p. 1139. JML paper):

LM Via A Neural Net

- Express each **word** as a **feature vector** (real-valued).
- Express the **joint probability function** of word sequences in terms of these word vectors.
- **Learn the word vectors and the joint probability function simultaneously** (using a **neural network**).

Why Does it Work?

- It is possible to naturally generalize (i.e., transfer probability mass) from (1) to (2-4) such that **dog** and **cat** end up with similar vectors:

Example

- 1 The **cat** is walking in the bedroom
- 2 A **dog** was running in a room
- 3 The **cat** is running in a room
- 4 A **dog** is walking in a bedroom
- 5 The **dog** was walking in the room
- 6 ...

Bengio et al., 2003 Neural Model

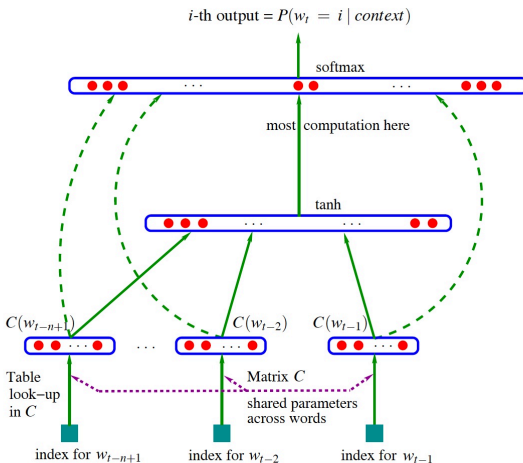


Figure 1: Neural architecture: $f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$ where g is the neural network and $C(i)$ is the i -th word feature vector.

It Still Doesn't Scale...

It is such a great idea, **but there is a bottleneck**...

- **Computationally costly** to obtaining the output probabilities (since softmax operates on **all** vocabulary).
- **Bottleneck** in the computation of the activations of the output layer.
- An n-gram model **does not require the computation of the probabilities for all the words in the vocabulary**.
- **Mikolov et al. (2013)** Scale learning word vectors with a large vocabulary (V).