

Smoothing

1. Applied to the training corpus, the log likelihood will be the largest when $\alpha = 0$, and smallest when $\alpha \rightarrow \infty$. Proof:

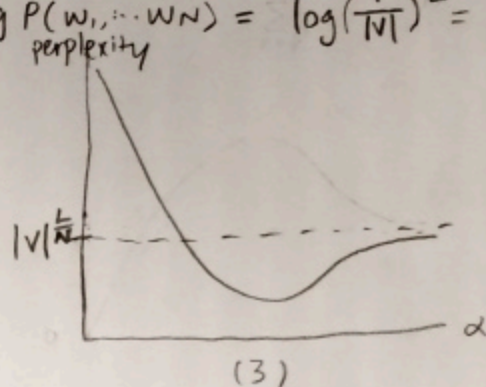
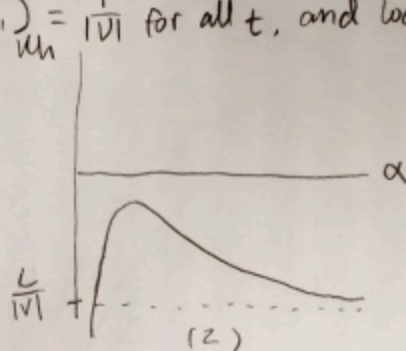
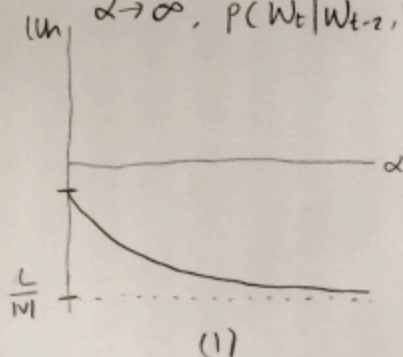
$$\log(P(w_1, w_2, \dots, w_N)) = \log \prod_{n=1}^N \frac{\text{Count}(w_t, w_{t-1}, w_{t-2}) + \alpha}{\text{Count}(w_{t-1}, w_{t-2}) + \alpha |V|}$$

Suppose $C_1 = \text{Count}(w_t, w_{t-1}, w_{t-2})$ and $C_2 = \text{Count}(w_{t-1}, w_{t-2})$.

$$\frac{C_1 + \alpha}{C_2 + \alpha |V|} - \frac{C_1}{C_2} = \frac{C_2(C_1 + \alpha) - C_1(C_2 + \alpha |V|)}{C_2(C_2 + \alpha |V|)} = \frac{\alpha(C_2 - C_1 |V|)}{C_2(C_2 + \alpha |V|)} < 0,$$

because $C_1 \geq 1$ and $|V| > C_2$. [$|V| > C_2$ is generally true, unless with extremely large corpus and small vocabulary].

Intuitively, probabilities are allocated to trigrams "most efficiently". In the worst case, when $\alpha \rightarrow \infty$, $P(w_t | w_{t-2}, w_{t-1}) = \frac{1}{|V|}$ for all t , and $\log P(w_1, \dots, w_N) = \log \left(\frac{1}{|V|}\right)^L = \frac{L}{|V|}$.



2. Now using an unseen test corpus, the log likelihood is $-\infty$ if $\alpha = 0$, because it will assign $-\infty$ to unseen words. When $\alpha \rightarrow \infty$, each word/trigram is again assigned equal probability, similar to the case discussed above.

However, there exists an α that maximizes the test corpus, which balances between the frequencies observed in the training corpus and the possibilities of unseen words.

3. As an extension of question (2), $\text{perplexity} = 2^{-\frac{1}{N} \log \text{likelihood}}$.

4. When $\alpha \rightarrow \infty$, $\text{perplexity} \rightarrow 2^{-\frac{1}{N} \log \left(\frac{1}{|V|}\right)^L} = \left(\frac{1}{|V|}\right)^{-\frac{L}{N}} = |V|^{\frac{L}{N}}$
When $\alpha = 0$, average likelihood $\rightarrow -\infty$, $\text{perplexity} \rightarrow \infty$.

Neural Language Models

	Serial	Parallel	# units
1. ① read	$(n-1)d$	1	$(n-1)d$
② $x \rightarrow z^n$	$[(n-1)d + 1] \cdot m$	$(n-1)d + 1$	m
③ $z^n \rightarrow f^n$	m	1	m
④ $f^n \rightarrow z^0$	$(1+m) \cdot V $	$1+m$	$ V $
⑤ $z^0 \rightarrow p$	$ V $	1	$ V $

Therefore, the computation needed is capped by

$$(n-1)d + [(n-1)d+1] \cdot m + m + (1+m) \cdot |V| + |V| \rightarrow O(ndm + mv) \\ \rightarrow O((nd+M) \cdot m)$$

Shown in the table above, with unlimited computational resource,

$$O(1+(n-1)d+1+1+m+1) \rightarrow O(nd+m).$$

$$\begin{aligned} \delta_k^o &= \frac{\partial \log P(W_t | W_{t-1}, W_{t-2})}{\partial z_k^o} = \frac{\partial}{\partial z_k^o} \left(\log \frac{\exp(z_w)}{\sum_l \exp(z_l)} \right) = \frac{\partial}{\partial z_k^o} (\log(\exp(z_w)) - \log \sum_l \exp(z_l)) \\ &= \frac{\partial z_w}{\partial z_k^o} - \frac{\partial \log(\sum_{l \neq k} \exp(z_l) + \exp(z_k))}{\partial z_k^o} \rightarrow \frac{\exp(z_k)}{\sum \exp(z_l) + \exp(z_k)} = p_k \\ &= \begin{cases} 1 - p_k & \text{if } k=w \\ -p_k & \text{if } k \neq w \end{cases} \end{aligned}$$

$$\begin{aligned} \delta_j^h &= \frac{\partial \log P_w}{\partial z_j^h} \\ &= \frac{\partial \log P_w}{\partial f_j^h} \cdot \frac{\partial f_j^h}{\partial z_j^h} = \frac{\partial f_j^h}{\partial z_j^h} \cdot \sum_k \frac{\partial \log P_w}{\partial z_k^o} \cdot \frac{\partial z_k^o}{\partial z_j^h} \\ &= (1 - (f_j^h)^2) \sum_{k=1}^{|V|} \delta_k^o \cdot w_{jk} \end{aligned}$$

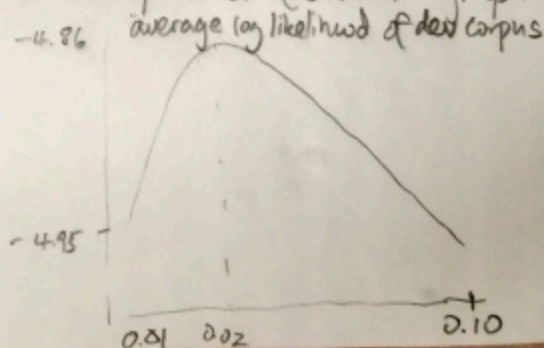
$$\begin{aligned} \delta_i^x &= \frac{\partial \log P_w}{\partial x_i} \\ &= \sum_{j=1}^m \frac{\partial \log P_w}{\partial z_j^h} \cdot \frac{\partial z_j^h}{\partial x_i} = \sum_{j=1}^m \delta_j^h \cdot w_{ij}^h \end{aligned}$$

Step to update $V(W_{t-2})$

- take a trigram randomly or in serial: W_t, W_{t-1}, W_{t-2}
- select the corresponding transformation/lookup matrix to calculate the input vector and then activation function and probability
- calculate the gradient using the formulae above, and use the gradient to update the parameters.

For $V(W_{t-2})$ specifically, $[V(W_{t-2}) \leftarrow V(W_{t-1})] \leftarrow V_{old} + \underset{\text{learning rate}}{\eta} \nabla_V \log P_w$
and in this case $V(W_{t-2})$ and $V(W_{t-1})$ are updated individually.

(a) Looping through different α 's, I found $\alpha=0.02$ to be optimal. (Shown in graph below)
With $\alpha=0.02$, the test likelihood is -4.8457 .



(b) My best configuration is: $n=3$, $d=75$, $m=60$. The test likelihood is -4.5873 .
First, I conducted a "coarse" grid search;

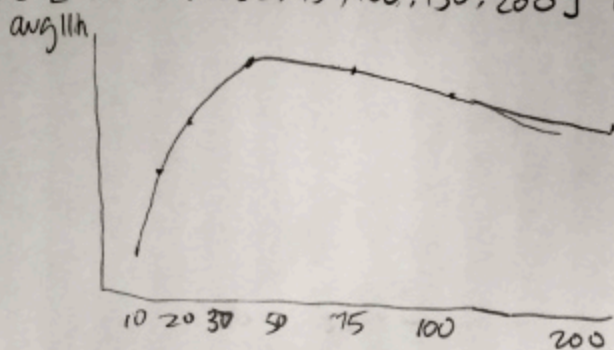
$n \in [2, 3, 4, 5]$, $d \in [5, 10, 20, 50]$, $m \in [10, 30, 50, 100]$

and found that tri-gram has the best generalization, whereas larger d and m give better performance.

Next, I conducted a 2nd coarse search, and found that performance decreases when d and m exceed some large values, around 100.

Then, I conducted finer searches, and found $d=75$ and $m=60$ to be optimal.

(c) As discussed above, the performance decreases after m reaches some level. To be more exact and also time-saving, I picked $n=3$, $d=10$ and experimented $m \in [10, 20, 30, 50, 75, 100, 150, 200]$ and obtained:



5. (a) No, it usually won't be zero.

I noticed that when indexing, all low frequency words are indexed to 0, and thus all bundled together. Therefore if a word in test is unseen in training, it is likely that it got bundled with other low-frequency words in the training set and assigned a probability greater than 0. This probability may not make sense though.

(b) Yes it can assign somewhat reasonable probability.

This is mainly due to how input vectors are constructed - preceding combinations don't have to be exact, but being similar is good enough.

For example, if "I am cool" and "You are cool" are in training and we encounter "We are cool" in test, then we expect the activation and softmax layers give high probabilities to cool because "we are" should be spatially close to "I am" and "You are" in the vector space.

6. The first sentence is more likely.

Average LLH (1) = -5.2679

Average LLH (2) = -6.0529 . with $n=3$, $d=75$ and $m=60$.

Recurrent Neural Network

1. Taking away h^{t-1} , the model acts as a bigram. Normally, we can find parameters that "overfit" the sentence to make probability close to 1, but we CANNOT in this case. Note that "very" appears twice, and therefore $P(\text{very}|\text{very}) + P(\text{doll}|\text{very}) \leq 1$. The maximum possible product of these two terms alone is much smaller than 1 (≤ 0.25).

2. Yes.

Removing W_{t-1} makes the network able to be treated as any n-gram except that the number of hidden unit is smaller than vocabulary. Parameters could be manipulated to span all existing n-grams and assign probabilities of 1 to each.