

1: a) Raw Bigram probability:

$$P(\text{store} \mid \text{computer}) = \frac{\text{count}(\text{computer}, \text{store})}{\sum \text{count}(\text{computer})}$$

$$= \frac{4}{10} = 0.4$$

$$P(\text{monitor} \mid \text{computer}) = \frac{\text{count}(\text{computer}, \text{monitor})}{\sum \text{count}(\text{computer})}$$

$$= \frac{4}{10} = 0.4$$

Both the words are equally likely to occur.

⑥: Kneser - Ney Smoothing

$$d = 0.5$$

$$\hat{P}(w_i, w_{i-1}) = \frac{\max(C(w_{i-1}, w_i) - d, 0)}{C(w_{i-1})} +$$

$$\lambda(w_{i-1}) \frac{|\{v : C(v, w_i) > 0\}|}{\sum |\{v : C(v, w') > 0\}|}$$

$$\lambda(w_{i-1}) = \frac{d}{\sum c(w_{i-1})} \times |w: C(w_{i-1}, w) > 0|$$

P(store, computer) ↴

discounted bigram probability of store| computer^y

$$= \frac{c(\text{computer store}) - d}{c(\text{computer})}$$

$$= \frac{4 - 0.5}{10} = \frac{3.5}{10} = 0.35 \quad \text{--- (1)}$$

$$\lambda(w_{i-1}) = \lambda(\text{computer})$$

$$= \frac{0.5}{10} \times 2$$

$$= \frac{1}{10} \quad \text{--- (2)}$$

$$\frac{|v: C(v, w_i) > 0|}{\sum |v: C(v, w) > 0|} = \frac{\# \text{ of Diff String type preceding final word}}{\text{length of Bigram table for monitor}}$$

$$\begin{aligned} \text{computer monitor} &= 4 \\ \text{monitor monitor} &= 1 \end{aligned} \quad = \quad \frac{3}{9}$$

$$\text{computer store} = 4 \quad \text{--- (3)}$$

$$\text{Keyboard store} = 2$$

$$\text{monitor store} = 2$$

combining equation ①, ②, ③

$$\begin{aligned}
 P(\text{store, computer}) &= 0.35 + \frac{\frac{1}{10} \times \frac{1}{2}}{0.05} \\
 &= 0.35 + 0.05 \\
 &= 0.4
 \end{aligned}$$

PC monitor, computer)

discounted bigram prob { monitor, computer)

$$= \frac{4 - 0.5}{10} = \frac{3.5}{10} = 0.35 - ①$$

$$\pi(w_{i+1}) = \pi(\text{computer})$$

$$= \frac{0.5}{10} \times 3 - ②$$

$$\frac{|\{v : C(v, w_i) > 0\}|}{|\{v : C(v, w) > 0\}|} = \frac{\# \text{ of Diff string type preceding final word}}{\text{Length of Bigram}}$$

table for store

$$= \frac{2}{9} - ③$$

combining ①, ②, ③

$$P(\text{monitor, computer}) = 0.35 + \frac{0.5}{10} \times \frac{2}{g_3}$$

$$= 0.3833$$

We can see that:

$$P(\text{store, computer}) > P(\text{monitor, computer})$$

Word Computer store is more likely than computer monitor.

The result is changed because we now also considered absolute discounting and we redistribute our resulting probability using $P_{\text{continuation}}$.

$$\textcircled{C}: d = \underline{0.1}$$

$$P(\text{store, computer}) = ?$$

$$P(w_i, w_{i-1}) = \frac{\max(C(w_{i-1}, w_i) - d, 0)}{C(w_{i-1})} +$$

$$\lambda(w_{i-1}) \frac{|\{v: C(v, w_i) > 0\}|}{\sum |\{v: C(v, w') > 0\}|}$$

$$\lambda(w_{i-1}) = \frac{d}{\sum C(w_{i-1})} \times |\{w: C(w_{i-1}, w) > 0\}|$$

$$P(\text{store, computer}) = \frac{4 - 0.1}{10} +$$

$$\frac{0.1}{10} \times 3 \times \frac{3}{9}$$

$$= \frac{3.9}{10} + \frac{0.1}{10}$$
$$= 0.4$$

$$P(\text{monitor, computer}) = \frac{4 - 0.1}{10} +$$

$$\frac{0.1}{10} \times 3 \times \frac{2}{9}$$

$$= 0.39 + 0.03 + 2/9$$
$$= 0.3966$$

$$P(\text{store, computer}) > P(\text{monitor, computer})$$

"Computer store" is more likely than "computer monitor" because we build the smoothening on the basis that words that appear more context in past are likely to appear in new contexts as well, and store has appear in more context in the past.

Q:

term-Doc.	Doc1	Doc2	Doc3
car	27	4	24
insurance	3	18	0
auto	0	33	29
best	14	0	17

$$\text{Term frequency} = \log_{10}(\text{count}(t, d) + 1)$$

$$\text{Inverse Doc. freq} = \log_{10}(N / df_i)$$

car 12

insurance 6

auto 10

best 16

	Doc1		
	Tf	IDF	$Tf \times IDF$
car	1.447	0.222	0.321
insurance	0.602	0.523	0.315
auto	0	0.301	0
best	1.176	0.097	0.114

	Doc 2		
	T _I	IDF	T _I × IDF
Car	0.699	0.222	0.155
insurance	1.278	0.523	0.668
Auto	1.532	0.301	0.461
Best	0	0.097	0

	Doc 3		
	T _I	IDF	T _I × IDF
Car	1.398	0.222	0.310
insurance	0	0.523	0
Auto	1.477	0.301	0.445
Best	1.255	0.097	0.122

⑥: Cosine Similarities

$$\cos(a, b) = \frac{\sum (a_i * b_i)}{\sqrt{(\sum (a_i^2))} * \sqrt{(\sum (b_i^2))}}$$

$$\cos(\text{doc1}, \text{doc2}) =$$

$$\frac{0.321 * 0.155 + 0.315 * 0.668 + 0 + 0}{0.464 * 0.827}$$

$$\frac{0.260}{0.464 * 0.827} = 0.677$$

$$\cos(\text{doc2}, \text{doc3}) =$$

$$\frac{0.155 * 0.310 + 0.461 * 0.445}{0.464 * 0.827}$$

$$= 0.551$$

$$\cos(\text{doc1}, \text{doc3}) =$$

$$= \frac{0.321 * 0.310 + 0.114 * 0.122}{0.464 * 0.827}$$

$$= 0.440$$

③ $P(O|I)$ - Given an inside word w_I , what is the probability that and outside word w_0 lies in the context window of w_I .

$$CE(p, q) = - \sum_i p_i \log(q_i)$$

where $p_i \Rightarrow$ True Probability distribution with vector [I being correct word in the context & rest of the words as 0]

$q_i =$ our prediction
Let's suppose

$$P = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

then

$$\begin{aligned} CE(p, q) &= - \sum p_i \log q_i \\ &= - \log q_0 - \sum_{w \neq 0} 0 \cdot \log q_w \\ &= - \log q_0 \end{aligned}$$

$$= - \log P(O = w_0 | I = w_I)$$

→ Distribution we want to learn

⑥:

$$CE(p, q) = -\log P(O = \omega_0 | I = \omega_I)$$

$$= -\log \frac{\exp(U_{\omega_0}^T \cdot v_{\omega_I})}{\sum \exp(U_{\omega_i}^T \cdot v_{\omega_I})}$$

$$\log \frac{q}{b} = \log q - \log b$$

$$= -\log \exp(U_{\omega_0}^T \cdot v_{\omega_I}) + \log \sum \exp(U_{\omega_i}^T \cdot v_{\omega_I})$$

$$= -U_{\omega_0}^T \cdot v_{\omega_I} + \log \sum \exp(U_{\omega_i}^T \cdot v_{\omega_I})$$

$$\frac{\partial (E(p, q))}{\partial v_{\omega_I}} = -U_{\omega_0}^T + \frac{\sum U_{\omega_i}^T \cdot \exp(U_{\omega_i}^T \cdot v_{\omega_I})}{\sum \exp(U_{\omega_i}^T \cdot v_{\omega_I})}$$

$$= -U_0 + \frac{\sum U_{\omega_i} \exp(U_{\omega_i}^T \cdot v_{\omega_I})}{\sum \exp(U_{\omega_i}^T \cdot v_{\omega_I})}$$

$$= -U_0 + \sum_{i=1}^J P(O = \omega_0 | I = \omega_i) \cdot U_i$$

$\hat{y} = \text{predicted distribution}$

$$\frac{\partial CE(p, q)}{\partial v_{\omega_I}} = -U_0 + \sum_{i=1}^J \hat{y}_i \cdot U_i$$

$$\textcircled{C}: CE(p, q) = -\log P(O = o_i | I = w_i)$$

$$\frac{\partial}{\partial v_w} CE(p, q) = \frac{\partial}{\partial v_w} \left[-\log \frac{\exp(v_w^\top \cdot v_{w_I})}{\sum \exp(v_w^\top \cdot v_{w_I})} \right]$$

$$= \frac{\partial}{\partial v_w} \left[-\log \exp(v_w^\top \cdot v_{w_I}) + \log \sum \exp(v_w^\top \cdot v_{w_I}) \right]$$

$$= \frac{\partial}{\partial v_w} \left[-v_w^\top \cdot v_{w_I} + \log \sum \exp(v_w^\top \cdot v_{w_I}) \right]$$

$$= \underbrace{\frac{\partial}{\partial v_w} \left[-v_w^\top \cdot v_{w_I} \right]}_{\text{I}} + \underbrace{\frac{\partial}{\partial v_w} \left[\log \sum \exp(v_w^\top \cdot v_{w_I}) \right]}_{\text{II}}$$

$$\text{I: } \frac{\partial}{\partial v_w} \left[-v_w^\top \cdot v_{w_I} \right] = \begin{cases} -v_{w_I}, & \text{if } w = 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\text{II: } \frac{\partial}{\partial v_w} \left[\log \sum \exp(v_w^\top \cdot v_{w_I}) \right]$$

$$\begin{aligned}
 &= \frac{\frac{\partial}{\partial v_w} (\mathbf{v}_w^\top \cdot \mathbf{v}_{wI}) * \sum \exp(\mathbf{v}_w^\top \cdot \mathbf{v}_{wI})}{\sum \exp(\mathbf{v}_w^\top \cdot \mathbf{v}_{wI})} \\
 &= \hat{y} \cdot \mathbf{v}_{wI} \quad \left[\because \hat{y} = \frac{\sum \exp(\mathbf{v}_w^\top \cdot \mathbf{v}_{wI})}{\sum \exp \mathbf{v}_w^\top \cdot \mathbf{v}_{wI}} \right]
 \end{aligned}$$

Combining I and II

$$\frac{\partial}{\partial v_w} = \begin{cases} \hat{y}_0 \cdot \mathbf{v}_{wI} - \mathbf{v}_{wI}, & w=0 \\ \hat{y}_w \cdot \mathbf{v}_{wI}, & \text{otherwise} \end{cases}$$

①: We use negative Sampling as an alternative to Naive Softmax classifier. The goal of negative Sampling is:

②: Maximize the similarity of our target word, context word pair drawn from positive sample

⑥ Minimize the similarity of target word, context word pair drawn from negative example.

Why K if we draw K negative samples from the vocabulary, the advantage is that whereas in Naive Softmax weight update happen for the whole words and that can be computationally expensive, but in negative sampling it is $O(K)$ where K is typically 10 - 20.

The loss function for negative sampling is given as:

$$L(v_{w_1}, v_{w_0}, k) = -\log \sigma(v_{w_0}^T \cdot v_{w_1}) - \sum_{k=1}^K \log \sigma(-v_k^T v_{w_1})$$

Q4'

Part-1: Generation%

Q: what do you think? Is it as good as 1000 monkeys writing at 100 typewriters?

⇒ I think the model predicts words that doesn't exist in the vocabulary, given that it is dependent on next most likely character.

Q 4.b: They all start with f in char-level Ngram and for all word level N gram it starts with ln. Did you get such results? Explain what is going on?

Aus.: Yes, I got similar results.

⇒ we can think of this in terms of n.
if $n=0$.

↳ Then the letters are context independent and anything can be generated for starting char. The same is true for word-level N gram model. The first word being ln, $n=0$, will always give us ln as first word.

if $n=1$.

→ Then f is always chosen as the first char because it is the opening letter of the input file. Similarly ln is chosen always as ln is the first word in the word level input file.

if $n=2$

→ the character depends on the f, meaning only char that have high probability of following a f at the beginning of the sentence can appear. for word level N-gram, the word with high probability following "ln" word can appear.

and finally with $n > 3$ we always obtain a word as the opening in both char-level and word level.

Q4.C:

Shakespeare - sonnet Article

Perplexity (word)	Perplexity (char)	Order (n)
4658.7	23.45	0
12061	14.26	1
4338	10.28	2
5130	7.961	3

We can observe from the above table that perplexity improves with higher order ($n=1, 2, 3$) as than we have more context to look up, which improves the probability of the predicted word. This is consistent with our intuition as well.

However, for word level Ngram we can observe that perplexity performance reduces from $n=2$ to $n=3$, and that is perhaps because for seeing 3 previous words, it is less predictable to guess the frequency distribution in test file than the train file. Also in word level, the next word prediction is more affected by immediate previous word.

PART 2: PERPLEXITY, SMOOTHING AND INTERPOLATION.

Q4-d. In your report experiment with few different lambdas and values of K and discuss their effects for both char level and word level Ngram.

Aus: Considering higher order n gives good result (from 4-c), we fix n=2 and experiment with K, λs.

CHAR-LEVEL (Shakespeare-sonnet - txt)

→	λ_1	λ_2	λ_3	K	Perplexity
	0.9	0.05	0.05	0 1 2	7 7.9 8
	0.1	0.2	0.7	1 2	14 53 234
	0.7	0.2	0.1	0 1 2	8 13 33
	0.5	0.4	0.1	0 1 2	8.9 15 46
	0.3	0.4	0.3	0 1 2	10 21 78

WORD - LEVEL

(val-e.txt)

λ_1	λ_2	λ_3	κ	Perplexity
0.1	0.2	0.7	0.1 1 2 0.1 1 2 0.1 1 2 0.1 1 2	3201 5744 12855 1750 2304 2100 1651 2472 2707 3850 6559 7516
0.7	0.2	0.1		
0.5	0.4	0.1		
0.1	0.1	0.8		

Based on the above 2 tables, I observed that for char level n-gram the perplexity performance improves if we give more weight towards Lambda-1, and that is reasonable as in char level n-gram the higher the context length we have the better the prediction in words.

In Word -level n-gram the observations were opposite. Higher weights towards Lambda - n gives better performance in perplexity because in word level n-gram the next word prediction is dependent on few immediate previous word.

for K, I observed lower value of K favours better perplexity.
finally the best hyperparameters that I was able to get were.

char-level

(Shakespeare-sonnet.txt)

$$K = 0.001$$

$$n = 4$$

$$\lambda_1 = 0.9$$

$$\lambda_L = 0.05$$

$$\lambda_3 = 0.02$$

$$\lambda_4 = 0.02$$

$$\lambda_5 = 0.01$$

$$\text{perplexity} = 5.2$$

Word level

(val-e.txt)

$$K = 0.001$$

$$n = 2$$

$$\lambda_1 = 0.1$$

$$\lambda_2 = 0.2$$

$$\lambda_3 = 0.7$$

$$\text{perplexity} = 9.75$$

Ques: In the results of RNN and char level n-gram model we see that perplexity (char level n-gram is 5.2 (with hyperparameter tuning) whereas in RNN the perplexity comes around 2.2 which is better than char level n-gram, because in RNN the model accounts dynamically previous n-char, so we no need to manually look for best set of hyperparameters in RNN, and that is why Deep learning models are efficient.