

N-Gram Models

N-gram Model

An **n-gram** is a contiguous sequence of **n** items from a given sample of text or speech. The items can be phonemes, syllables, letters, words or base pairs according to the application. The **n-grams** typically are collected from a text or speech corpus.

Conditional Probability: $P(B | A) = \frac{P(A, B)}{P(A)}$ $P(A, B) = P(A)P(B | A)$

More variables: $P(A, B, C, D) = P(A)P(B | A)P(C | A, B)P(D | A, B, C)$

Chain Rule:

$$P(x_1, x_2, \dots, x_n) = P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2) \dots P(x_n | x_1, \dots, x_{n-1})$$

$P(\text{"about five minutes from"}) = P(\text{about}) \times P(\text{five} | \text{about}) \times P(\text{minutes} | \text{about five}) \times P(\text{from} | \text{about five minutes})$

Probability of words in sentences:

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i | w_1, w_2, w_3, \dots, w_{i-1})$$

Unigram(1-gram): **No history is used.**

Bi-gram(2-gram): **One word history**

Tri-gram(3-gram): **Two words history**

Four-gram(4-gram): **Three words history**

Five-gram(5-gram): **Four words history**

Generally in practical applications, Bi-gram(previous one word), Tri-gram(previous two word, Four-gram (previous three word) are used.

Unigram(1-gram): No history is used.

“about five minutes from....”

Assume in corpus dinner word is present with highest probability.

Unigram doesn't take into account probabilities with previous words like from, minutes.

Unigram will predict dinner.

“about five minutes from **dinner**”



Bi-gram(2-gram): **One word history**

$$P(w_1, w_2) = \prod_{i=2} P(w_i | w_{i-1}) \quad P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

“about five minutes from.....”

Assumption: Next word may be college, class

$$P(\text{college} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from college})}{\text{count}(\text{about five minutes from})}$$

$$P(\text{class} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from class})}{\text{count}(\text{about five minutes from})}$$

“about five minutes from.....”

$$\text{Count}(\text{about five minutes from}) = P(\text{about} | <S>) \times P(\text{five} | \text{about}) \times P(\text{minutes} | \text{five}) \\ \times P(\text{from} | \text{minutes})$$

$$\text{Count}(\text{about five minutes from college}) = P(\text{about} | <S>) \times P(\text{five} | \text{about}) \times P(\text{minutes} | \text{five}) \\ \times P(\text{from} | \text{minutes}) \times P(\text{college} | \text{from})$$

$$\text{Count}(\text{about five minutes from class}) = P(\text{about} | <S>) \times P(\text{five} | \text{about}) \times P(\text{minutes} | \text{five}) \\ \times P(\text{from} | \text{minutes}) \times P(\text{class} | \text{from})$$

$$P(\text{college} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from college})}{\text{count}(\text{about five minutes from})}$$

$$= P(\text{college} | \text{from})$$

$$P(\text{class} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from class})}{\text{count}(\text{about five minutes from})}$$

$$= P(\text{class} | \text{from})$$

Tri-gram(2-gram): **Two words history**

$$P(w_1, w_2, w_3) = \prod_{i=3} P(w_i | w_1, w_2) \quad P(w_i | w_{i-1}, w_{i-2}) = \frac{\text{count}(w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-2}, w_{i-1})}$$

$$\text{Count}(\text{about five minutes from}) = P(\text{five} | \langle S \rangle, \text{about}) \times P(\text{minutes} | \text{about}, \text{five}) \times P(\text{from} | \text{five}, \text{minutes})$$

$$\text{Count}(\text{about five minutes from } \mathbf{college}) = P(\text{five} | \langle S \rangle, \text{about}) \times P(\text{minutes} | \text{about}, \text{five}) \times P(\text{from} | \text{five}, \text{minutes}) \times P(\mathbf{college} | \text{minutes from})$$

$$\text{Count}(\text{about five minutes from } \mathbf{class}) = P(\text{five} | \langle S \rangle, \text{about}) \times P(\text{minutes} | \text{about}, \text{five}) \times P(\text{from} | \text{five}, \text{minutes}) \times P(\mathbf{class} | \text{minutes from})$$

$$P(\mathbf{college} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from college})}{\text{count}(\text{about five minutes from})}$$

$$= P(\mathbf{college} | \text{minutes from})$$

$$P(\mathbf{class} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from class})}{\text{count}(\text{about five minutes from})}$$

$$= P(\mathbf{class} | \text{minutes from})$$

Four-gram(4-gram): **Three words history**

$$P(w_1, w_2, w_3, w_4) = \prod_{i=4} P(w_i | w_1, w_2, w_3)$$

$$P(w_i | w_{i-1}, w_{i-2}, w_{i-3}) = \frac{\text{count}(w_{i-3}, w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-3}, w_{i-2}, w_{i-1})}$$

Count(about five minutes from)= $P(\text{minutes} | \langle S \rangle, \text{about}, \text{five}) \times P(\text{from} | \text{about}, \text{five}, \text{minutes})$

Count(about five minutes from **college**)= $P(\text{minutes} | \langle S \rangle, \text{about}, \text{five}) \times$
 $P(\text{from} | \text{about}, \text{five}, \text{minutes}) \times$
 $P(\text{college} | \text{five}, \text{minutes}, \text{from})$

Count(about five minutes from **class**)= $P(\text{minutes} | \langle S \rangle, \text{about}, \text{five}) \times$
 $P(\text{from} | \text{about}, \text{five}, \text{minutes}) \times$
 $P(\text{class} | \text{five}, \text{minutes}, \text{from})$

$$P(\text{college} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from college})}{\text{count}(\text{about five minutes from})}$$

$$= P(\text{college} | \text{five minutes from})$$

$$P(\text{class} | \text{about five minutes from}) = \frac{\text{count}(\text{about five minutes from class})}{\text{count}(\text{about five minutes from})}$$

$$= P(\text{روے | college} | \text{five minutes from})$$

As no. of previous state (history) increases, it is very difficult to match that set of words in corpus.

Probabilities of **larger collection of word is minimum**. To overcome this problem, Bi-gram model is used

Exercise 1: Estimating Bi-gram probabilities

What is the most probable next word predicted by the model for the following word sequence?

Given Corpus

<S> I am Henry </S>
⌨ <S> I like college </S>
<S> Do Henry like college </S>
<S> Henry I am </S>
<S> Do I like Henry </S>
<S> Do I like college </S>
<S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

1) <S> Do ?

<S> I am Henry </S>
 <S> I like college </S>
 <S> Do Henry like college </S>
 <S> Henry I am </S>
 <S> Do I like Henry </S>
 <S> Do I like college </S>
 <S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

Next word prediction probability $W_{i-1} = \text{do}$

Next word	Probability $\frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$
$P(</S> \text{do})$	0/4
$P(<I> \text{do})$	2/4
$P(<\text{am}> \text{do})$	0/4
$P(<\text{Henry}> \text{do})$	1/4
$P(<\text{like} \text{do})$	1/4
$P(<\text{college} \text{do})$	0/4
$P(\text{do} \text{do})$	0/4

I is more probable

2) <S> I like Henry ?

<S> I am Henry </S>

<S> I like college </S>

<S> Do Henry like college </S>

<S> Henry I am </S>

<S> Do I like Henry </S>

<S> Do I like college </S>

<S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

Next word prediction probability $W_{i-1} = \text{Henry}$

Next word	Probability Next Word = $\frac{N}{D} = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$
P(</S> Henry)	3/5
P(<I> Henry)	1/5
P(<am> Henry)	0
P(<Henry> Henry)	0
P(<like Henry)	1/5
P(<college Henry)	0
P(<do Henry)	0

</S> is more probable

3) <S> Do I like ?

Use Tri-gram

$P(\text{I like}) = 3$

<S> I am Henry </S>

<S> I like college </S>

<S> Do Henry like college </S>

<S> Henry I am </S>

<S> Do I like Henry </S>

<S> Do I like college </S>

<S> I do like Henry </S>

Next word prediction probability

$W_{i-2} = \text{I}$ and $W_{i-1} = \text{like}$

Next word	Probability Next Word = $\frac{\text{count}(w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-2}, w_{i-1})}$
$P(\text{</S>} \mid \text{I like})$	0/3
$P(\text{<I>} \mid \text{I like})$	0/3
$P(\text{<am>} \mid \text{I like})$	0/3
$P(\text{<Henry>} \mid \text{I like})$	1/3
$P(\text{<like>} \mid \text{I like})$	0/3
$P(\text{<college>} \mid \text{I like})$	2/3
$P(\text{do} \mid \text{I like})$	0/3

College is probable

4) <S> Do I like college ?

Use Four-gram

<S> I am Henry </S>

<S> I like college </S>

<S> Do Henry like college </S>

<S> Henry I am </S>

<S> Do I like Henry </S>

<S> Do I like college </S>

<S> I do like Henry </S>

Next word prediction probability

$W_{i-3}=I, W_{i-2}=like, W_{i-1}=college$

Next word	Probability Next Word = $\frac{\text{count}(w_{i-3}, w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-3}, w_{i-2}, w_{i-1})}$
$P(</S> I like college)$	2/2
$P(<I> I like college)$	0/2
$P(<am> I like college)$	0/2
$P(<Henry> I like college)$	0/2
$P(<like> I like college)$	0/2
$P(<college> I like college)$	0/2
$P(<do> I like college)$	0/2

</S> is more probable

Which of the following sentence is better. i.e. Gets a higher probability with this model.
Use Bi-gram

<S> I am Henry </S>
<S> I like college </S>
<S> Do Henry like college </S>
<S> Henry I am </S>
<S> Do I like Henry </S>
<S> Do I like college </S>
<S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

1. <S> I like college </S>

<S> like college </S>=?

$$\begin{aligned} &= P(I | <S>) \times P(\text{like} | I) \times P(\text{college} | \text{like}) \times P(</S> | \text{college}) \\ &= 3/7 \times 3/6 \times 3/5 \times 3/3 = 9/70 = 0.13 \end{aligned}$$

2. <S> Do I like Henry </S>

$$\begin{aligned} &= P(\text{do} | <S>) \times P(I | \text{do}) \times P(\text{like} | I) \times P(\text{Henry} | \text{like}) \times P(</S> | \text{Henry}) \\ &= 3/7 \times 2/4 \times 3/6 \times 2/5 \times 3/5 = 9/350 = 0.0257 \end{aligned}$$

ANS: First statement is more probable

Advantages:

- Easy to understand, implement
- Can be easily converted to any grammar

Disadvantages:

- Underflow due to multiplication of probabilities
- **Solution:** Use log. Add probabilities.
- Zero probability problem
- **Solution:** Use Laplace smoothing

Which of the following sentence is better. i.e. Gets a higher probability with Bi-gram model.

<S> I am Henry </S>
<S> I like college </S>
<S> Do Henry like college </S>
<S> Henry I am </S>
<S> Do I like Henry </S>
<S> Do I like college </S>
<S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

First statement is more probable

1. <S> I like college </S>

$$=P(I | <S>) \times P(\text{like} | I) \times P(\text{college} | \text{like}) \times P(</S> | \text{college})$$

$$=3/7 \times 3/6 \times 3/5 \times 3/3 = 9/70 = \mathbf{0.13}$$

$$= \log(3/7) + \log(3/6) + \log(3/5) + \log(3/3) = \mathbf{-2.0513}$$

2. <S> Do I like Henry </S>

$$=P(\text{do} | <S>) \times P(I | \text{do}) \times P(\text{like} | I) \times P(\text{Henry} | \text{like}) \times P(</S> | \text{Henry})$$

$$=3/7 \times 2/4 \times 3/6 \times 2/5 \times 3/5 = 9/350 = \mathbf{0.0257}$$

$$= \log(3/7) + \log(2/4) + \log(3/6) + \log(2/5) + \log(3/5) = \mathbf{-3.6607}$$

<S> I am Henry </S>

<S> I like college </S>

<S> Do Henry like college </S>

<S> Henry I am </S>

<S> Do I like Henry </S>

<S> Do I like college </S>

<S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

Second statement is more probable

1. <S> like college </S>

$$=P(\text{like} \mid \text{<S>}) \times P(\text{college} \mid \text{like}) \times P(\text{</S>} \mid \text{college})$$

$$=0/7 \times 3/5 \times 3/3 = 0$$

2. <S> Do I like Henry </S>

$$=P(\text{do} \mid \text{<S>}) \times P(\text{I} \mid \text{do}) \times P(\text{like} \mid \text{I}) \times P(\text{Henry} \mid \text{like}) \times P(\text{</S>} \mid \text{Henry})$$

$$=3/7 \times 2/4 \times 3/6 \times 2/5 \times 3/5 = 9/350 = \mathbf{0.0257}$$

Laplace Smoothing

<S> I am Henry </S>
<S> I like college </S>
<S> Do Henry like college </S>
<S> Henry I am </S>
<S> Do I like Henry </S>
<S> Do I like college </S>
<S> I do like Henry </S>

Word	Frequency
<S>	7
</S>	7
I	6
am	2
Henry	5
like	5
college	3
do	4

Unique words are : <S>, </S>, I, Henry do, like, am, college

Total unique words: 8

But we exclude <S> as it never comes in bi-gram calculations

Total unique words: 7

Give the following bi-gram probabilities estimated by Laplace model.

1. <S> like college </S>

$$=P(\text{like} \mid \text{<S>}) \times P(\text{college} \mid \text{like}) \times P(\text{</S>} \mid \text{college})$$

$$=(0+1)/(7+7) \times (3+1)/(5+7) \times (3+1)/(3+7)$$

$$=1/14 \times 4/12 \times 4/10$$

$$=0.0095$$

2. <S> Do I like Henry </S>

$$=P(\text{do} \mid \text{<S>}) \times P(\text{I} \mid \text{do}) \times P(\text{like} \mid \text{I}) \times P(\text{Henry} \mid \text{like}) \times P(\text{</S>} \mid \text{Henry})$$

$$=(3+1)/(7+7) \times (2+1)/(4+7) \times (3+1)/(6+7) \times (2+1)/(5+7) \times (3+1)/(5+7)$$

$$=4/14 \times 3/11 \times 4/13 \times 3/12 \times 4/12$$

$$=0.0020$$

First statement is more probable

Perplexity

The language model is best when it predicts an unseen test set.

Definition of Perplexity:

It is the inverse probability of the test data which is normalized by the number of words.

$$PP(w) = P(w_1, w_2, w_3, \dots, w_N)^{-\frac{1}{N}}$$

$$PP(w) = \left(\prod_i \frac{1}{P(w_i | w_1, w_2, \dots, w_{i-1})} \right)^{\frac{1}{N}} \quad PP(w) = \left(\prod_i \frac{1}{P(w_i | w_{i-1})} \right)^{\frac{1}{N}}$$

Lower the value of perplexity: **Better Model**

More value of perplexity: **Confused for prediction**

WSJ Corpus

Training: 38 million words

Test: 1.5 million words

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

Perplexity for Bigram <S> I like college </S>

$$=P(I | <S>) \times P(\text{like} | I) \times P(\text{college} | \text{like}) \times P(</S> | \text{college})$$

$$=3/7 \times 3/6 \times 3/5 \times 3/3 = 9/70 = \mathbf{0.13}$$

$$\mathbf{PP(w) = (1/0.13)^{1/4} = 1.67}$$

Perplexity for Trigram <S> I like college </S>

$$P(w) = P(\text{like} | <S> I) \times P(\text{college} | I \text{ like}) \times P(</S> | \text{like college})$$

$$P(w) = 1/3 \times 2/3 \times 3/3 = 2/9 = \mathbf{0.22}$$

$$\mathbf{PP(w) = (1/0.22)^{1/3} = 1.66}$$