

NGram Language Model Log of Probabilities Laplace Smoothing Perplexity

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:



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

>

P("about five minutes from")=P(about) × P(five |about) × P(minutes | about five) × P(from | about five minutes)

Probability of words in sentences:

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

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



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

"about five minutes from....."

Assumption: Next word may be college, class

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

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

"about five minutes from...."



Count(about five minutes from)= P(about | <S>) × P(five | about) × P(minutes | five) × P(from | minutes)

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

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

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

=P(college | from)

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

=P(class | from)

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



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

Count(about five minutes from)= P(five | <S>, about) × P(minutes | about, five) × P(from | five, minutes)

Count(about five minutes from college) = P(five | <S>, about) × P(minutes | about, five) × P(from | five , minutes) × P(college | minutes from)

Count(about five minutes from class) = P(five | <S>, about) × P(minutes | about, five) × P(from | five , minutes) × P(class | minutes from)

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

=P(college | minutes from)

P(class | about five minutes from) = count(about five minutes from class) count(about five minutes from) =P(class | minutes from)

19-02-2024

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



=P(college | five minutes from)

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

=P(college | 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>I like college

<S> Do Henry like college

<S> Henry I am

<S> Do I like Henry

<S> Do I like college

<S>I do like Henry

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

1)	<\$>	Do	?
----	------	----	---

-	100					Stock
<\$>	2 34	-	-	-		c 196
		-	-		ν.	~

<\$>1 like college </\$>

<S> Do Henry like college

<S> Henry I am

<S> Do I like Henry

<S> Do I like college

<\$>1 do like Henry </\$>

Word	Frequency
<\$>	7
5	7
1	6
am	2
Henry	5
like	5
college	3
do	4

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



	count(w _{i-1} , w _i)
Next word	Probability Next Word = 1
P(do)	0/4
P(<i> do)</i>	2/4
P(<am> do)</am>	0/4
P(<henry> do)</henry>	1/4
P(<like do)<="" td="" =""><td>1/4</td></like>	1/4
P(<college do)<="" td="" =""><td>0/4</td></college>	0/4
P(do do)	0/4

I is more probable



2) <S> I like Henry ?

<S>I am Henry

<\$>1 like college </\$>

<S> Do Henry like college

<S> Henry I am

<S> Do I like Henry

<S> Do I like college

<S>1 do like Henry

Word	Frequency
<s></s>	7
\$	7
- 1	6
am	2
Henry	5
like	5
college	3
do	4

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

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

Use Tri-gram

P<I like>=3



<\$>1 like college </\$>

<S> Do Henry like college

<S> Henry I am

<S> Do I like Henry

<S> Do I like college

<S>I do like Henry

Next word prediction probability

W_{i-2}=I and W_{i-1}=like

Next word	Probability Next Word= count(w _{i-2} , w _{i-1} , w _i
ivextword	count(w_{i-2}, w_{i-1})
P(I like)	0/3
P(<i> I like)</i>	0/3
P(<am> I like)</am>	0/3
P(<henry> I like)</henry>	1/3
P(<like i="" like)<="" td="" =""><td>0/3</td></like>	0/3
P(<college i="" like)<="" td="" =""><td>2/3</td></college>	2/3
P(do I like)	0/3

College is probable



4) <S> Do I like college ?

Use Four-gram





<\$> Henry I am </\$>

<S> Do I like Henry

<S> Do I like college

<\$> I do like Henry </\$>

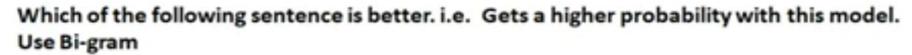
Next word prediction probability

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

	$count(w_{i-3}, w_{i-2}, w_{i-1}, w_i)$
Next word	Probability Next Word= $count(w_{i-3}, w_{i-2}, w_{i-1})$
P(I like college)	2/2
P(<i> I like college)</i>	0/2
P(<am> I like college)</am>	0/2
P(<henry> I like college)</henry>	0/2
P(<like college)<="" i="" like="" td="" =""><td>0/2</td></like>	0/2
P(<college college)<="" i="" like="" td="" =""><td>0/2</td></college>	0/2
P(do I like college)	0/2









D

>11				

<\$>1 like college </\$>

<S> Do Henry like college

<\$> Henry I am </\$>

<S> Do I like Henry

<S> Do I like college

<\$>1 do like Henry </\$>

Word	Frequency
<\$>	7
5	7.
1	6
am	2
Henry	5
like	5
college	3
do	- 4

<S> I like college

<S> like college =?

$$=P(1|~~) \times P(like | 1) \times P(college | like) \times P(| college)~~$$

=3/7 × 3/6 × 3/5 ×3/3 = 9/70=0.13

2. <S> Do I like Henry

=P(do |
$$<$$
S>) × P(I | do) × P(like | I) × P(Henry | like) × P($<$ /S> | Henry)
=3/7 × 2/4 × 3/6 ×2/5 ×3/5 = 9/350=0.0257

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



Word	Frequency
<\$>	7
	7
1	6
am	2
Henry	5
like	5
college	3
do	4

First statement is more probable

1. <S> I like college

$$=P(1| ~~) \times P(like | 1) \times P(college | like) \times P(| college)~~$$

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

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

2. <S> Do I like Henry

$$=P(do | ~~) \times P(I | do) \times P(like | I) \times P(Henry | like) \times P(| Henry)~~$$

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

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



<s>1</s>	am	Henry	\$
-		T. T	-9

<S> Henry I am

<S> Do I like Henry

<S> Do I like college

<S>I do like Henry

Word	Frequency		
<\$>	7		
	7		
- 1	₽ 6		
am	2		
Henry	5		
like	5		
college	3		
do	4		

Second statement is more probable

1. <S> like college

=P(like |
$$\langle S \rangle$$
) × P(college | like) × P($\langle S \rangle$ | college)

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

2. <S> Do I like Henry

=P(do |
$$\langle S \rangle$$
) × P(I | do) × P(like | I) × P(Henry | like) × P($\langle S \rangle$ | Henry)

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

Laplace Smoothing



Corpus

<s></s>	lam	Н	enrv	/ </th <th>/S></th>	/ S>

<S>I like college

<S>Do Henry like college

<S>Henry I am

<S>Do I like Henry

<S>Do I like college

<S>I do like Henry

Word	Frequency
<\$> \$	7
	7
I	6
am	2
Henry	5
like	5
college	3
Do	4

Features: <S>, , I, am, Henry, like, college, Do

Total number of unique words: 8

But we exclude <S> as it is not used in bigram.

So, total unique word is 7.

Add one Smoothing

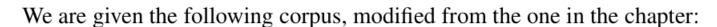


<S> like college

2. <S> Do I like Henry

=0.0020

=P(do |
$$<$$
S>) × P(I | do) × P(like | I) × P(Henry | like) × P($<$ /S> | Henry)
=(3+1)/(7+7) × (2+1)/(4+7) × (3+1)/(6+7) × (2+1)/(5+7) × (3+1)/(5+7)
= $4/14 \times 3/11 \times 4/13 \times 3/12 \times 4/12$





$$\langle s \rangle$$
 I am Sam $\langle s \rangle$

$$<$$
s $>$ Sam I am $s $>$$

$$\langle s \rangle$$
 I am Sam $\langle s \rangle$

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

Using a bigram language model with add-one smoothing, what is P(Sam am)? Include <s> and </s> in your counts just like any other token.

Word	Frequency
I	4
am	3
do	1
not	1
like	1
green	1
eggs	1
Sam	4
and -2024	1

P(Sam|am) =
$$\frac{2}{3}$$
 (Bigram Model)

P(Sam|am) =
$$\frac{3}{14}$$
 (Bigram Model with add-one smoothing)

Language Model Evaluation



- LM is better if it is assigning a high probability to the real, frequently observed and grammatical sentence over false, rarely observed and ungrammatical sentences.
- Two different criteria for evaluation
 - ✓ Extrinsic
 - ✓ Intrinsic
- Extrinsic Evaluation
- It evaluates the language model when solving a specific task.
- Ex: Speech recognition accuracy, Machine translation accuracy, spelling correction accuracy
- Compare 2 or more models and check which one is better.
- Disadvantages



Language Model Evaluation

- Intrinsic Evaluation:
- 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.
- Lower the value of perplexity : Better model
- More value of perplexity : Confused for prediction

Perplexity (Intrinsic Evaluation model)



• 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.

$$PP(w) = P(w_1, w_2, w_3, w_4, ..., wN)^{\frac{1}{N}}$$

$$PP(w) = \left(\prod_{i} \frac{1}{P(w_{i} \mid w_{1}, w_{2}, \dots, w_{i-1})}\right)^{\frac{1}{N}} \qquad PP(w) = \left(\prod_{i} \frac{1}{P(w_{i} \mid w_{i-1})}\right)^{\frac{1}{N}}$$

- Lower the value of perplexity : **Better model**
- More value of perplexity : Confused for prediction



Perplexity for Bigram <S> I like college

=P(I| ~~) × P(like | I) × P(college | like) × P(| college)
=
$$3/7 \times 3/6 \times 3/5 \times 3/3 = 9/70 = 0.13$$~~

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

Perplexity for Trigram <S> I like college

$$P(w)=P(like | ~~1) \times P(college | 1 like) \times P(| like college)~~$$

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

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



Advantages:

- Easy to understands, implement
- Can be easily convert to any gram

Þ

Disadvantages:

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



References:

Daniel Jurafsky, James H. Martin —Speech and Language Processing Second Edition, Prentice Hall, 2008.

Christopher D.Manning and Hinrich Schutze, — Foundations of Statistical Natural Language Processing, MIT Press, 1999.



THANK YOU