

N-gram Model [underflow problem, (smoothing) zero-probability problem] N-gram Evaluation

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.

Condition 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 :-

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

Types :-

Unigram (1-gram): No history is used

Bigram (2-gram): One word history is used

Trigram (3-gram): Two words " " "

Four-gram (4-gram): Three " " " "

Five-gram (5-gram): Four " " " "

Adv :-

- Easy to understand, implement
- can be easily convert to any gram

Disadv :-

→ Underflow due to multiplication of probabilities

Solution: use log. Add prob.

→ Zero probability problem

Solution: use Laplace smoothing

Ex :-

Consider the corpus and find the most probable statement using 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>

Ans :-  
First statement  
is more  
probable.

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

$$= P(I | \langle S \rangle) \times P(\text{like} | I) \times P(\text{college} | \text{like}) \times P(\langle S \rangle | \text{college})$$

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

$$= \log(\frac{3}{7}) + \log(\frac{3}{6}) + \log(\frac{3}{5}) + \log(\frac{3}{3}) = -2.0513.$$

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

$$= P(\text{do} | \langle S \rangle) \times P(I | \text{do}) \times P(\text{like} | I) \times P(\text{Henry} | \text{like}) \times P(\langle S \rangle | \text{Henry})$$

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

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

Find the probability for the following statements

1.  $\langle s \rangle$  like college  $\langle l s \rangle$
2.  $\langle s \rangle$  do I like Henry  $\langle l s \rangle$

1.  $\langle s \rangle$  like college  $\langle l s \rangle$

$$= P(\text{like} | \langle s \rangle) \times P(\text{college} | \text{like}) \times P(\langle l s \rangle | \text{college})$$

$$= 0.7 \times 3/5 \times 3/3$$

$$= 0.$$

2.  $\langle s \rangle$  do I like Henry  $\langle l s \rangle$

$$= P(\text{do} | \langle s \rangle) \times P(\text{I} | \text{do}) \times P(\text{like} | \text{I}) \times$$

$$P(\text{Henry} | \text{like}) \times P(\langle l s \rangle | \text{Henry})$$

$$= 3/7 \times 2/4 \times 3/6 \times 2/5 \times 3/5$$

$$= 0.0257. \quad (\text{Second statement is probable})$$

In the first stmt, the other words have prob. have more frequency but it doesn't appear after start so it becomes zero.

To overcome zero probability problem, we can apply Laplace smoothing

(4)

Laplace Smoothing:-

Find the unique words in the corpus.

<S></S> & I, Henry, do, like, am, college.

= 8 unique words

But we exclude <S> as it never comes <sup>in bigram</sup> here.

Total unique words = 7.

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

add 1 at numerator & add 7 in denominator.

① <S> like college </S>

$$= p(\text{like} | \langle S \rangle) \times p(\text{college} | \text{like}) \times p(\langle /S \rangle | \text{college})$$

$$= \frac{(0+1)}{(7+7)} * \frac{(3+1)}{5+7} * \frac{3+1}{3+7}$$

$$= 1/14 * 4/12 * 4/10 = 0.0095$$

② <S> do + like Henry </S>

$$= p(\text{do} | \langle S \rangle) * p(\text{I} | \text{do}) * p(\text{like} | \text{I}) * p(\text{Henry} | \text{like}) \\ * p(\langle /S \rangle | \text{Henry})$$

$$= \frac{(8+1)}{(7+7)} * \frac{(2+1)}{(4+7)} * \frac{(3+1)}{(6+7)} * \frac{(2+1)}{(5+7)} * \frac{(3+1)}{(5+7)}$$

$$= 4/14 * 3/11 * 4/13 * 3/12 * 4/12 = 0.0020$$

First statement is more probable.

## Language Model Evaluation:- N-gram Evaluation

Language model 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:

- 1) Extrinsic
- 2) Intrinsic.

### Extrinsic Evaluation :-

It evaluates the language model when solving a specific task.

For e.g. speech recognition accuracy, Machine translation accuracy, spelling correction accuracy - compare 2 (or more) models, and check which works best.

#### Disadvantage:-

- ↳ Expensive
- ↳ Time consuming.

### Intrinsic Evaluation :-

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

It uses perplexity:

#### Definition of perplexity:-

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

(6)

- Lower the value of perplexity : Better model  
 → More-value of perplexity : confused for prediction

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

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

WSJ Corpus :-

contains 38 million words : Training  
 1.5 million words : Test

N-gram order	unigram	Bigram	Trigram
Perplexity Measure	962	170	109.

Perplexity of tri-gram is very less than this other it is the best model for predicting the next word.

Ex:-

Perplexity for Bi-gram < s > I like college < /s >

$$= P(I | s) \times P(\text{like} | I) \times P(\text{college} | \text{like}) \times$$

$$\text{so } N=4. = \frac{①}{3} \times \frac{②}{1} \times \frac{③}{3} \times \frac{④}{5} = 0.13.$$

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

Perplexity of Trigram < s > I like college < /s >

$$P(W) = P(\text{like} | s) \times P(\text{college} | \text{like}) + P(s | \text{like})$$

$$\text{so } N=3. = \frac{①}{3} \times \frac{②}{2} \times \frac{③}{3} = 0.22$$

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

Note :- This is a small corpus, so we can't take decision. Since both values are same. for large corpus - the values vary.