

(1)

Statistical Language Model (SLM)

A SLM is a model that assigns probabilities of sequences of words in a language.

Instead of relying on predefined grammar rules,

→ SLM learns these probabilities from large dataset (corpora) & use statistical method to estimate.

→ These models are based on the idea that the probability of a word depends on the words that came before it, & they calculate this probability by analyzing patterns in a text.

e.g:- "The cat chased,"

the model should predict the most probable next word (e.g., "the", "a", "mouse")

Types of SLM :-

* n-grams

* Exponential models

* skip-gram models

(2)

1. N-grams:-

An n-gram is a sequence of n items from a sample of text or speech, such as syllables, letters, words or base pair.

N-gram ^{model} use the frequency of these sequences in a training corpus to predict the likelihood of word sequences.

Ex:- types:-

Uni-gram

bi-gram

tri-gram

& so on.

The n-grams typically are collected from a text or speech corpus.

conditional probability

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

$$P(A, B) = P(B|A) \cdot P(A)$$

more variables :

$$P(A, B, C, D) = P(A) P(B|A) P(C|A, B)$$

$$P(D|A, B, C)$$

(3)

Ex:-

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

$$P(\text{about}) * P(\text{five} | \text{about}) * \\ P(\text{minutes} | \text{about five}) * \\ P(\text{from} | \text{about five minutes})$$

∴ probability of words in sentences:

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

Types :-

Unigram (1-gram): No history is used

bigram (2-gram): ^{previous} one word history

Trigram (3-gram): ^{previous} two words history

fourgram (4-gram): ^{four} words history
^{Previous Three}

Fivegram (5-gram): ^{four} words history
^{previous}

Ex:- I

Unigram (1-gram)

"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

so unigram ^{will} predict dinner

"about five minutes from dinner"

[Note:-

Corpus : is a large collection of text data.
consists of texts, spoken language transcription,
written or even structured document].

Ex :- Estimating Bi-gram probabilities

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

Ex :- 1) <s> Do ?

Soln :-

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

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

(5)

do appear 4 times.

Next Word	Prob. Next word = $\frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$
$P(</s> \text{do})$ (it means how many time end sentence </s> appear after do)	0/4
$P(I \text{do})$	<u>2/4</u> (prob of I appear after do)
$P(\text{am} \text{do})$	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 so

Do I (ie the answer)

Ex ② <s> I like Henry ? {use bi gram.}

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

(5)

$P(\text{college} \text{Henry})$	0
$P(\text{do} \text{Henry})$	0

$\langle \text{Is} \rangle$ is more probable

$\therefore \langle \text{Is} \rangle \text{ I like Henry } \underline{\langle \text{Is} \rangle}$

Ex 3) $\langle \text{s} \rangle$ Do I like ? use Tri-gram

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

Tri-gram use 2 prev. words to predict next word.

Next word prediction probability

$$\text{if } w_{i-2} = \text{I} \text{ & } w_{i-1} = \text{like}$$

Next word	Prob. Next word	$= \frac{\text{count}(w_{i-2}, w_{i-1})}{\text{count}(w_{i-2}, w_{i-1})}$
$P(\langle \text{Is} \rangle \text{I like})$	0/3	(prob. of $\underline{\langle \text{Is} \rangle}$ appear after I like)
$P(\langle \text{I} \rangle \text{I like})$	0/3	($\underline{\text{I}}$ appear after I like)
$P(\langle \text{am} \rangle \text{I like})$	0/3	(prob. of <u>am</u> appear after I like)
$P(\langle \text{Henry} \rangle \text{I like})$	1/3	(prob. of <u>Henry</u> appear after I like)
$P(\langle \text{like} \rangle \text{I like})$	0/3	(prob. of <u>like</u> appear after I like)

(7)

$P(\langle \text{college} \rangle \text{I like})$	2/3
$P(\text{do} \text{I like})$	0/3

college is probable

 $\therefore \text{Do I like } \underline{\text{college}}$

Ex :- 4) :- <ss> Do I like college ? use 4-gram.

4-gram uses previous three words

$$P(\text{I like college}) = 2.$$

Next word prediction probability :-

$$w_{i-3} = \text{I}, w_{i-2} = \text{like}, w_{i-1} = \text{college}$$

Next word	Prob. Next word = $\frac{\text{count}(w_{i-3}, w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-3}, w_{i-2}, w_{i-1})}$
$P(\langle \text{ss} \rangle \text{I like college})$	2/2
$P(\text{I} \text{I like college})$	0/2
$P(\text{am} \text{I like college})$	0/2
$P(\langle \text{Henry} \rangle \text{I like college})$	0/2
$P(\langle \text{like} \rangle \text{I like college})$	0/2
$P(\langle \text{college} \rangle \text{I like college})$	0/2
$P(\text{do} \text{I like college})$	0/2

 $\langle \text{ss} \rangle$ is more probable $\therefore \langle \text{ss} \rangle \text{ Do I like college } \underline{\text{ss}}$

(8)

Ex 5 :- Which of the following sentence is better.

i.e Gets a higher probability in 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 >

Soh :-

Let's :- take any two sentence & find prob.

< s > I like college < /s > &
 < s > Do I like Henry < /s >.

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

(9)

1.) $\langle s \rangle$ I like college $\langle /s \rangle$

$$\begin{aligned}
 &= P(I | \langle s \rangle) \times P(\text{like} | I) \times \\
 &\quad P(\text{college} | \text{like}) \times P(\langle /s \rangle | \text{college}) \\
 &= \frac{3}{7} \times \frac{3}{6} \times \frac{3}{5} \times \frac{3}{3} \\
 &= \frac{9}{70} \\
 &= 0.13
 \end{aligned}$$

2.) $\langle s \rangle$ Do I like Henry $\langle /s \rangle$

$$\begin{aligned}
 &= P(\text{do} | \langle s \rangle) \times P(I | \text{do}) \times P(\text{like} | I) \\
 &\quad \times P(\text{Henry} | \text{like}) \times P(\langle /s \rangle | \text{Henry}) \\
 &= \frac{3}{7} * \frac{2}{4} * \frac{3}{6} * \frac{2}{5} * \frac{3}{5} \\
 &= \frac{9}{350} \\
 &= 0.0257
 \end{aligned}$$

First statement is more probable.