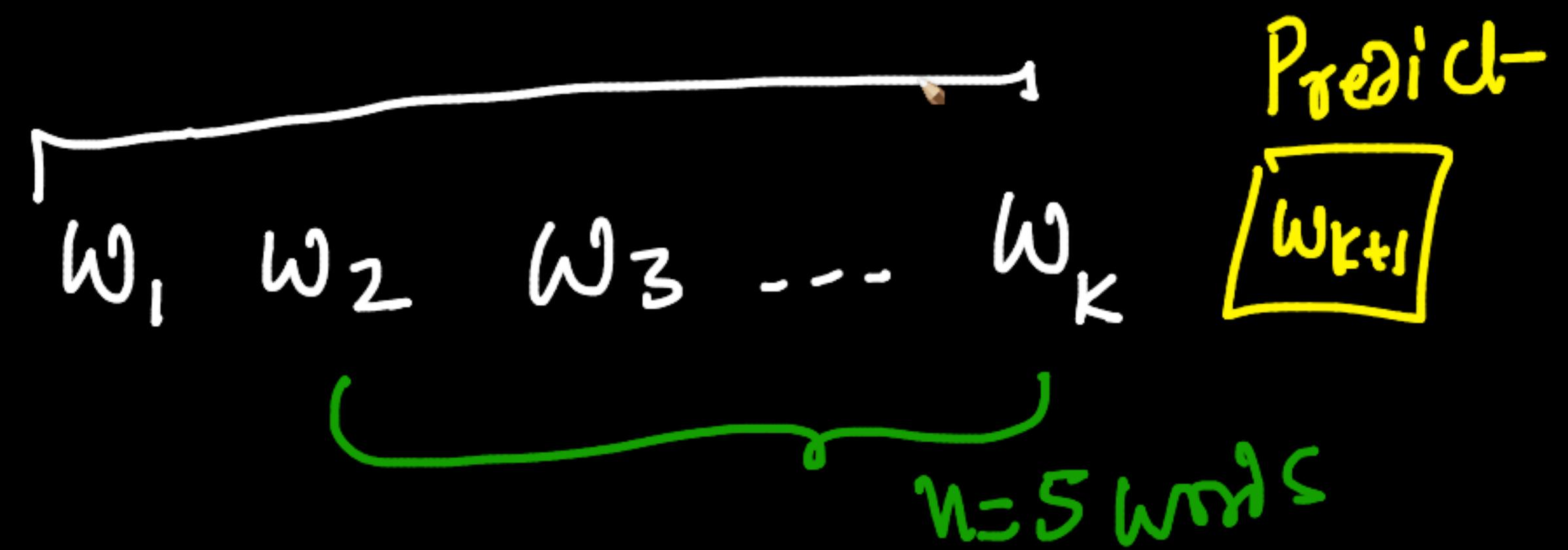


Task:



HINT:

Probability

past  
5 words

①  $P(w_{k+1} | w_k w_{k-1} w_{k-2} w_{k-3} w_{k-4})$

Classical NLP

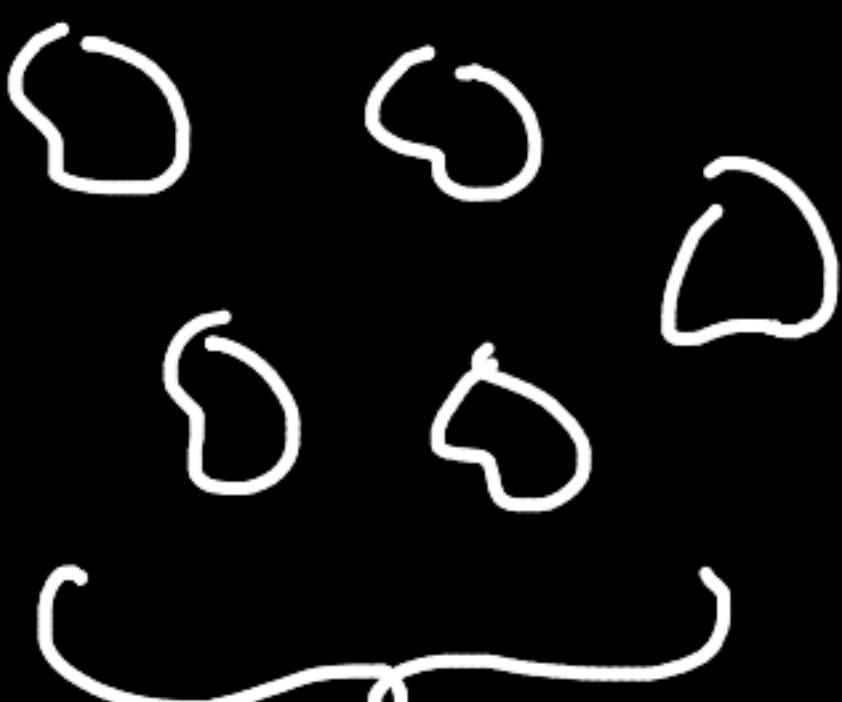
V.V.V.-fast + cheap

Voc: V

② Clustering:

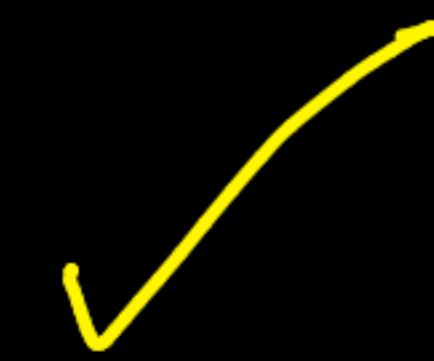
↪ word|sent → vec (Word2Vec + TF-IDF)

(?)  $s_i \rightarrow v_i \rightarrow \text{clustering} \rightarrow$



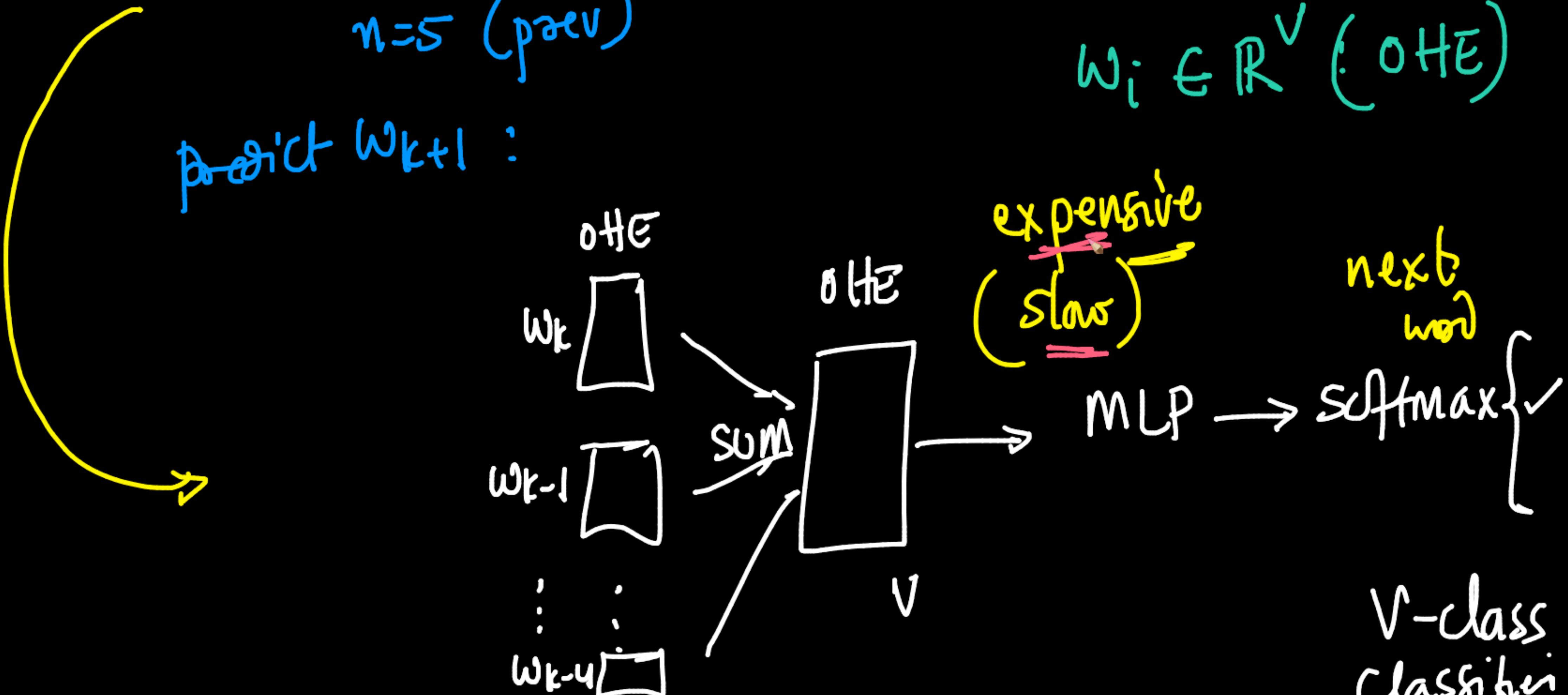
③

multi-class classfn



(MLP)

n=5 (paev)

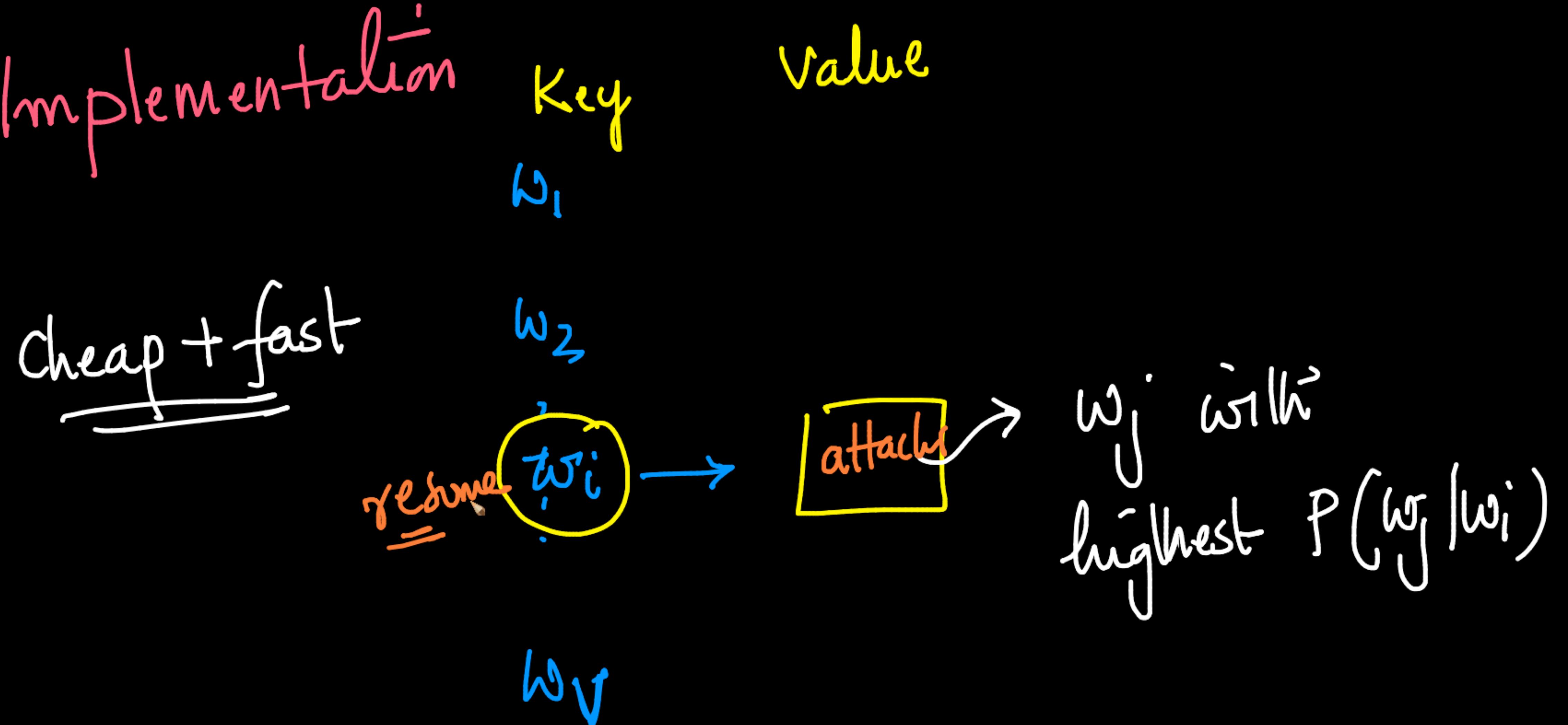
predict  $w_{k+1}$ : $w_i \in \mathbb{R}^V$  (OHE)V-class  
classifier

no past  $\rightarrow P(w_{k+1}) = \frac{n_{k+1}}{\# \text{words}}$  (train data)

n=1  $\rightarrow \underline{\underline{P(w_{k+1} | w_k)}} \leq \frac{P(w_{k+1}, w_k)}{P(w_k)}$

$$= \frac{\text{count}}{\#(w_k)}$$

# Implementation



e.g:

$w_i$   
resume

$w_j$   
attached

$$P(w_1 | w_i)$$

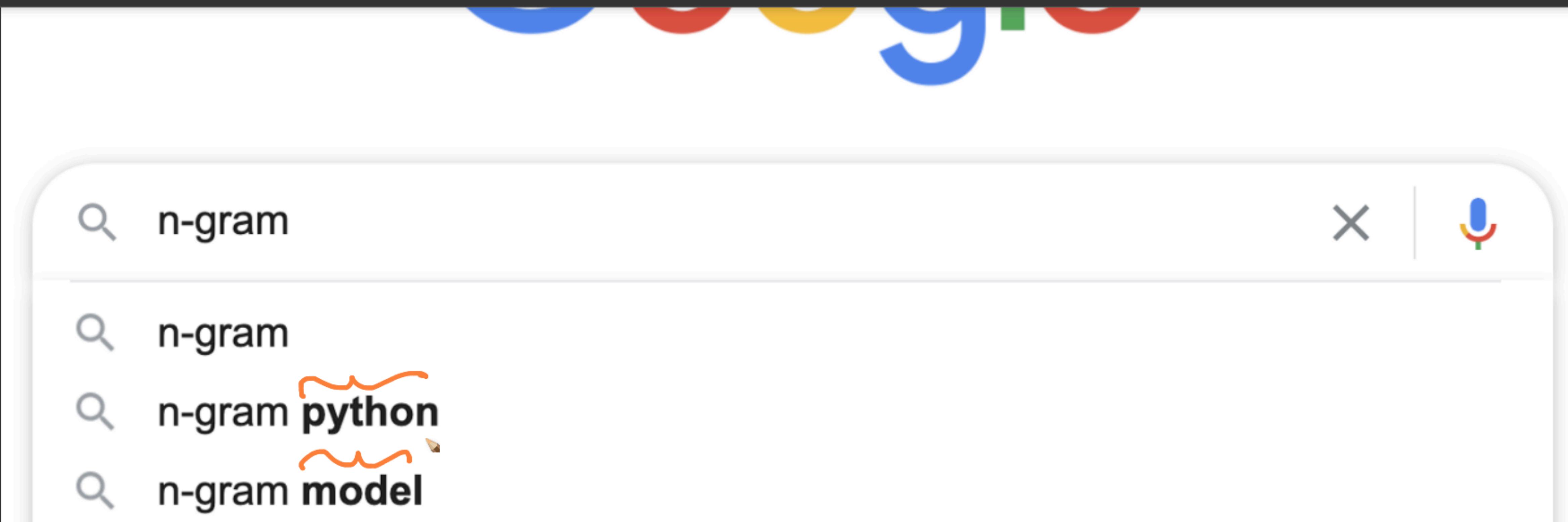
$$P(w_2 | w_i)$$

$$\text{Max} \leftarrow P(w_j = \text{attached} | w_i = \text{resume})$$

$$P(w_v | w_i)$$

+ Code + Text

Connect ▾



How can we implement this feature ?

1. We can get an idea which word sequence came frequently in training data. for example we can see that in most cases "please find my CV" is followed by word "attached".

no Past  $\rightarrow P(w_{k+1}) \xrightarrow{\text{unigram}}$

$n=1 \rightarrow P(w_{k+1} | w_k) \xrightarrow{\text{bigram}}$

$n=2 \rightarrow P(w_{k+1} | w_k w_{k-1}) \xrightarrow{\text{Trigram}}$

$P(w_{k+1} | w_k w_{k-1} \dots w_{k-n+1}) \xrightarrow{n\text{-gram}}$

$$P(w_{k+1} \mid \overbrace{w_k \ w_{k-1}}^{\sim}) = \frac{\#(w_{k+1}, w_k, w_{k-1})}{\#(\underline{w_k \ w_{k-1}})}$$

Problem?

- computationally expensive
- underflow: small probabilities
- Smoothing (Laplace)

language\_modeling\_V1.ipynb - x colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS Connect + Code + Text Update

- We can model this as a conditional probability, wherein, having seen a word we can predict the probability of next word.

## Exploring probabilities of words with an example

- Let us predict the probability of the sentence - "its water is so transparent".
- Using chain rule, the probability can be calculated using:
  - $P(\text{its water is so transparent}) = P(\text{its}) * P(\text{water}|\text{its}) * P(\text{is}|\text{its water}) * P(\text{so}|\text{its water is}) * P(\text{transparent}|\text{its water is so})$

Do we notice pitfalls with the above equation?

- We will have to calculate a lot of probabilities with as the size of sentence scales.

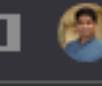
How about we simply estimate by counting and dividing the results?

- $P(\text{transparent} | \text{its water is so}) = \frac{\text{count}(\text{its water is so transparent})}{\text{count}(\text{its water is so})}$

Again, there are few obvious drawbacks with the above:

- There are far too many possible sentences in this method that would need to be calculated.
- This will result in very sparse data making results unreliable.
- Also, language is creative, any particular context might have not occurred before.

language\_modeling\_V1.ipynb - x colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS

+ Code + Text Connect |  

Should we see all the words typed to predict the possible next word ??

1. Find the probability of each word. Do not see any previous word, simply recommend the post probable word in the corpus. (**Unigram**)

$P(I \text{ have a dream}) = P(I) * p(\text{have}) * p(a) * p(\text{dream})$

2. Find the probability of a set of two words (how the given two words come together). - **Bigram**

$P(I \text{ have a dream}) = P(I/) * p(\text{have}/I) * p(a/\text{have}) * p(\text{dream}/a)$

3. Find the probability of a set of three words i.e. next word will be recommended on the basis of previous two words - **Trigram**

$P(I \text{ have a dream}) = P(I/,) * p(\text{have}/, I) * p(a/i,\text{have}) * p(\text{dream}/\text{have},a)$

4. Find the probability of a set of n words i.e. next word will be recommended on the basis of last n-1 words - **n-gram**

Assume Let's assume we have a corpus of only 3 sentences. Let's try to frame the sentence: "I have a dream"

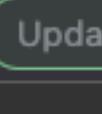
<> text

```
< start > I want a dream job </ start >
< start > I have a dog </ start >
< start > I have a dream company </ start >
```

11 / 11

language\_modeling\_V1.ipynb - x | +

colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS

+ Code + Text Connect |  

Should we see all the words typed to predict the possible next word ??

1. Find the probability of each word. Do not see any previous word, simply recommend the post probable word in the corpus. (**Unigram**)

$P(I \text{ have a dream}) = P(I) * p(\text{have}) * p(a) * p(\text{dream})$

2. Find the probability of a set of two words (how the given two words come together). - **Bigram**

$P(I \text{ have a dream}) = P(I) * p(\text{have}/I) * p(a/\text{have}) * p(\text{dream}/a)$

3. Find the probability of a set of three words i.e. next word will be recommended on the basis of previous two words - **Trigram**

$P(I \text{ have a dream}) = P(I/,) * p(\text{have}/, I) * p(a/i,have) * p(\text{dream}/\text{have},a)$

4. Find the probability of a set of n words i.e. next word will be recommended on the basis of last n-1 words - **n-gram**

Assume Let's assume we have a corpus of only 3 sentences. Let's try to frame the sentence: "I have a dream"

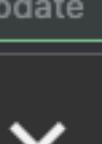
<> text

```
< start > I want a dream job </ start >
< start > I have a dog </ start >
< start > I have a dream company </ start >
```

12 / 12

language\_modeling\_V1.ipynb - x | +

colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS

+ Code + Text Connect |   | 

Should we see all the words typed to predict the possible next word ??

1. Find the probability of each word. Do not see any previous word, simply recommend the most probable word in the corpus. (**Unigram**)

$P(I \text{ have a dream}) = P(I) * p(\text{have}) * p(a) * p(\text{dream})$

2. Find the probability of a set of two words (how the given two words come together). - **Bigram**

$P(I \text{ have a dream}) = P(I/) * p(\text{have}/I) * p(a/\text{have}) * p(\text{dream}/a)$

3. Find the probability of a set of three words i.e. next word will be recommended on the basis of previous two words - **Trigram**

$P(I \text{ have a dream}) = P(I/,) * p(\text{have}/, I) * p(a/i,\text{have}) * p(\text{dream}/\text{have},a)$

4. Find the probability of a set of n words i.e. next word will be recommended on the basis of last n-1 words - **n-gram**

Assume Let's assume we have a corpus of only 3 sentences. Let's try to frame the sentence: "I have a dream"

<> text

```
< start > I want a dream job </ start >
< start > I have a dog </ start >
< start > I have a dream company </ start >
```

13 / 13

language\_modeling\_V1.ipynb - x colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS Connect + Code + Text

Should we see all the words typed to predict the possible next word ??

1. Find the probability of each word. Do not see any previous word, simply recommend the most probable word in the corpus. (**Unigram**)

$P(I \text{ have a dream}) = P(I) * p(\text{have}) * p(a) * p(\text{dream})$

2. Find the probability of a set of two words (how the given two words come together). - **Bigram**

$P(I \text{ have a dream}) = P(I/) * p(\text{have}/I) * p(a/\text{have}) * p(\text{dream}/a)$

3. Find the probability of a set of three words i.e. next word will be recommended on the basis of previous two words - **Trigram**

$P(I \text{ have a dream}) = P(I/,) * p(\text{have}/, I) * p(a/i,\text{have}) * p(\text{dream}/\text{have},a)$

4. Find the probability of a set of n words i.e. next word will be recommended on the basis of last n-1 words - **n-gram**

**n-gram PLMS**

Assume Let's assume we have a corpus of only 3 sentences. Let's try to frame the sentence: "I have a dream"

<> text

```
< start > I want a dream job </ start >
< start > I have a dog </ start >
< start > I have a dream company </ start >
```

14 / 14



P one-time

$$\log P(w_{k+1} |$$

last 5 words)

= bg-likelihood

$$K = V$$

(i; i; i; i)

w<sub>k+1</sub>

last 5 words

Problem:

↓  
size of dict  
↓  
v · v · large

✓ { disb-dict → Memcache  
Redis

language\_modeling\_V1.ipynb - x | +

colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS

+ Code + Text Connect | Update

Assume Let's assume we have a corpus of only 3 sentences. Let's try to frame the sentence: "I have a dream"

{x} **text**

< start > I want a dream job </ start >

< start > I have a dog </ start >

< start > I have a dream company </ start >

< end >

▼ Uni-gram Model

- **Unigram** - For uni-gram models we do not look at the previous words to predict the next word.
- Each word probability is calculated to predict the occurrence

Count table

< start >	I	want	a	dream	job	have	dog	company	< / start >
3/20	3/20	1/20	3/20	2/20	1/20	2/20	1/20	1/20	3/20

- Calculate the probabilities of occurrence of each word in the corpus.

Probability

< start >	I	want	a	dream	job	have	dog	company	< / start >
0.15	0.15	0.05	0.15	0.1	0.05	0.1	0.05	0.05	0.15

16 / 16

Connect



+ Code + Text

## Uni-gram Model

- **Unigram** - For uni-gram models we do not look at the previous words to predict the next word.
- Each word probability is calculated to predict the occurrence

### Count table

< start >	I	want	a	dream	job	have	dog	company	< / start >
3/20	3/20	1/20	3/20	2/20	1/20	2/20	1/20	1/20	3/20

$$P(w_i)$$

- Calculate the probabilities of occurrence of each word in the corpus.

### Probability

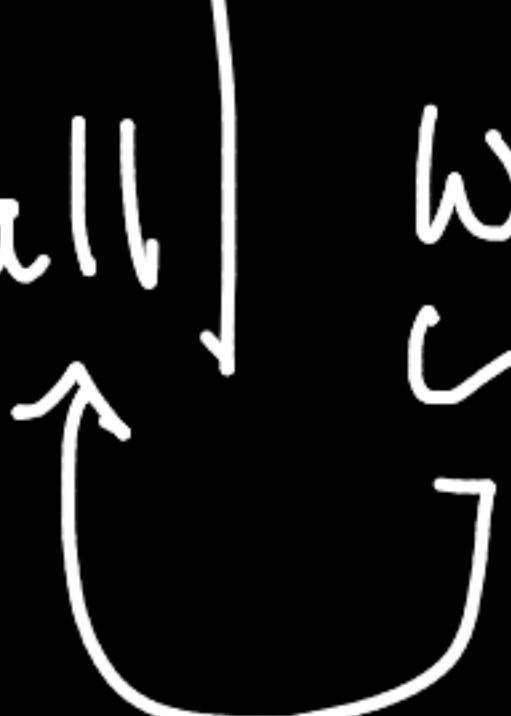
< start >	I	want	a	dream	job	have	dog	company	< / start >
0.15	0.15	0.05	0.15	0.1	0.05	0.1	0.05	0.05	0.15

$$P(I \text{ have a dream}) = P(I) * p(\text{have}) * p(a) * p(\text{dream})$$

- 2. **Bigram** - For Bi-gram models we calculate what is the probability of occurrence of a word given the previous word.

### count table

spell check

$$\hookrightarrow P(wall | wal)$$


①

$$P(l | wa)$$

②

char-level model

+ Code + Text

Connect



## ▼ We will build a N-gram model that will learn word probabilities using MLE

### What is MLE?

- MLE or Maximum Likelihood Estimator is a technique to estimate the parameters of an assumed probability, given some observed data.
- We can compute a particular bi-gram probability of a word  $w_n$ , given the previous word  $w_{n-1}$ , by computing the count of bigram  $C(w_{n-1} w_n)$  and normalize by the sum of all bigrams that share the first word  $w_{n-1}$ .
  - $P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(W_{n-1} w)}$
- Note - Since sum of all bigram counts that start with a given word  $w_{n-1}$  must be equal to the unigram count of that word, we can simplify the above equation to:
  - $P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(W_{n-1})}$
- The above calculation can be extended to other n-grams as well.

## ▼ Unigram Model

+ Code

+ Text

Connect



v

## ▼ Evaluating the model

{x}

**How can we find out the probability of the sentence using word probabilities??**

- **naive assumption** - Each word in the sentence is independent of each other (**markov assumption**)
- each sentence probability is the product of word probabilities.

$$\begin{aligned} P_{\text{eval}}(\text{i have a dream [END]}) &= P_{\text{train}}(\text{i}) P_{\text{train}}(\text{have} \mid \text{i}) P_{\text{train}}(\text{a} \mid \text{have}) P_{\text{train}}(\text{dream} \mid \text{a}) P_{\text{train}}([\text{END}] \mid \text{dream}) \\ &= \underline{\underline{P_{\text{train}}(\text{i})}} \underline{\underline{P_{\text{train}}(\text{have})}} \underline{\underline{P_{\text{train}}(\text{a})}} \underline{\underline{P_{\text{train}}(\text{dream})}} \underline{\underline{P_{\text{train}}([\text{END}])}} \end{aligned}$$

**Now How can we find out the probability of entire email text (having multiple sentences)??**

- **Naive Assumption** - Each sentence is independent of from other sentences (same as words)
- The Probability of entire text is the product of probabilities of sentences

Evaluation text : "I have a dream. We are free at last."

$$P_{\text{eval}}(\text{text}) = P_{\text{eval}}(\text{i have a dream [END]}) P_{\text{eval}}(\text{we are free at last [END]})$$

$$P_{\text{eval}}(\text{i have a dream [END]}) = P_{\text{train}}(\text{i}) P_{\text{train}}(\text{have}) P_{\text{train}}(\text{a}) P_{\text{train}}(\text{dream}) P_{\text{train}}([\text{END}])$$

$$P_{\text{eval}}(\text{we are free at last [END]}) = P_{\text{train}}(\text{we}) P_{\text{train}}(\text{are}) P_{\text{train}}(\text{free}) P_{\text{train}}(\text{at}) P_{\text{train}}(\text{last}) P_{\text{train}}([\text{END}])$$

[+ Code](#) [+ Text](#)

Connect



## Evaluating the model

How can we find out the probability of the sentence using word probabilities??

- **naive assumption** - Each word in the sentence is independent of each other (**markov assumption**)
- each sentence probability is the product of word probabilities.

$$\begin{aligned} P_{\text{eval}}(\text{i have a dream [END]}) &= P_{\text{train}}(\text{i}) P_{\text{train}}(\text{have} \mid \text{i}) P_{\text{train}}(\text{a} \mid \text{have}) P_{\text{train}}(\text{dream} \mid \text{a}) P_{\text{train}}([\text{END}] \mid \text{dream}) \\ &= P_{\text{train}}(\text{i}) P_{\text{train}}(\text{have}) P_{\text{train}}(\text{a}) P_{\text{train}}(\text{dream}) P_{\text{train}}([\text{END}]) \end{aligned}$$

Now How can we find out the probability of entire email text (having multiple sentences)??

- **Naive Assumption** - Each sentence is independent of from other sentences (same as words)
- The Probability of entire text is the product of probabilities of sentences

Evaluation text : "I have a dream. We are free at last."

$$P_{\text{eval}}(\text{text}) = P_{\text{eval}}(\text{i have a dream [END]}) P_{\text{eval}}(\text{we are free at last [END]})$$

$$P_{\text{eval}}(\text{i have a dream [END]}) = P_{\text{train}}(\text{i}) P_{\text{train}}(\text{have}) P_{\text{train}}(\text{a}) P_{\text{train}}(\text{dream}) P_{\text{train}}([\text{END}])$$

$$P_{\text{eval}}(\text{we are free at last [END]}) = P_{\text{train}}(\text{we}) P_{\text{train}}(\text{are}) P_{\text{train}}(\text{free}) P_{\text{train}}(\text{at}) P_{\text{train}}(\text{last}) P_{\text{train}}([\text{END}])$$

$\sum \log P(i)$

Unflatten

+ Code + Text

Connect



## what is the issue with log Probability??

**How to deal with unknown unigrams??** - For a unigram that appears in the evaluation text but not in the training text, its count in the training text – hence its probability – will be zero.

AND

the log of this zero probability is negative infinity, leading to a negative infinity average log likelihood for the entire model

$$n_{\text{train}}(\underline{\text{unigram}}) = \underline{\underline{0}}$$

$$\Rightarrow P_{\text{train}}(\text{unigram}) = \frac{\underline{\underline{0}}}{N_{\text{train}}} = 0$$

$$\Rightarrow \log(P_{\text{train}}(\text{unigram})) = \log(0) = -\infty$$

#Comment: Words not appearing in the trainset but present in the eval set is not a problem only specific to using log. Kindly mention that because of new words in the eval set the probability will become zero so the entire sentence probability becomes zero. Then mention the log(0) result.

language\_modeling\_V1.ipynb - x colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS Connect + Code + Text

1. Add a pseudo-count of k to all the unigrams in our vocabulary.

This will make the count of [UNK] k from 0 (0+K)

The probability of word appearing in eval set but not in training set will not have 0 probability now, why ??

2. The most common value of k is 1, and this goes by the intuitive name of “add-one smoothing”.

$$P_{\text{train}}(\text{unigram}) = \frac{n_{\text{train}}(\text{unigram}) + k}{N_{\text{train}} + kV}$$

unigram vocabulary size

+ Code + Text

Connect



Q

The probability of word appearing in eval set but not in training set will not have 0 probability now, why ??

{x}

2. The most common value of k is 1, and this goes by the intuitive name of “add-one smoothing”.

$$P_{\text{train}}(\text{unigram}) = \frac{n_{\text{train}}(\text{unigram}) + k}{N_{\text{train}} + kV}$$

unigram vocabulary size  
(number of unique unigrams in training text)

The image shows a handwritten-style annotation on a digital whiteboard. A red arrow points from the handwritten 'ntrain' to the term 'n<sub>train</sub>' in the equation. Another red arrow points from the handwritten 'kV' to the term 'kV' in the denominator. Below the equation, the phrase 'unigram vocabulary size' is written in red, followed by the explanatory text '(number of unique unigrams in training text)'.

+ Code + Text

Connect



Q

The probability of word appearing in eval set but not in training set will not have 0 probability now, why ??

{x}

2. The most common value of k is 1, and this goes by the intuitive name of “add-one smoothing”.

*s the stethoscope*

$$P_{\text{train}}(\text{unigram}) = \frac{n_{\text{train}}(\text{unigram}) + k}{N_{\text{train}} + kV} = 1$$

*pencil*

*1*

unigram vocabulary size  
(number of unique unigrams in training text)

+ Code + Text

Connect



Now we can apply Laplace smoothing in Bigram model.

## Effect of Laplace smoothing

Because of the additional pseudo-count  $k$  to each unigram, each time the unigram model encounters an unknown word in the evaluation text, it will convert said unigram to the unigram [UNK].

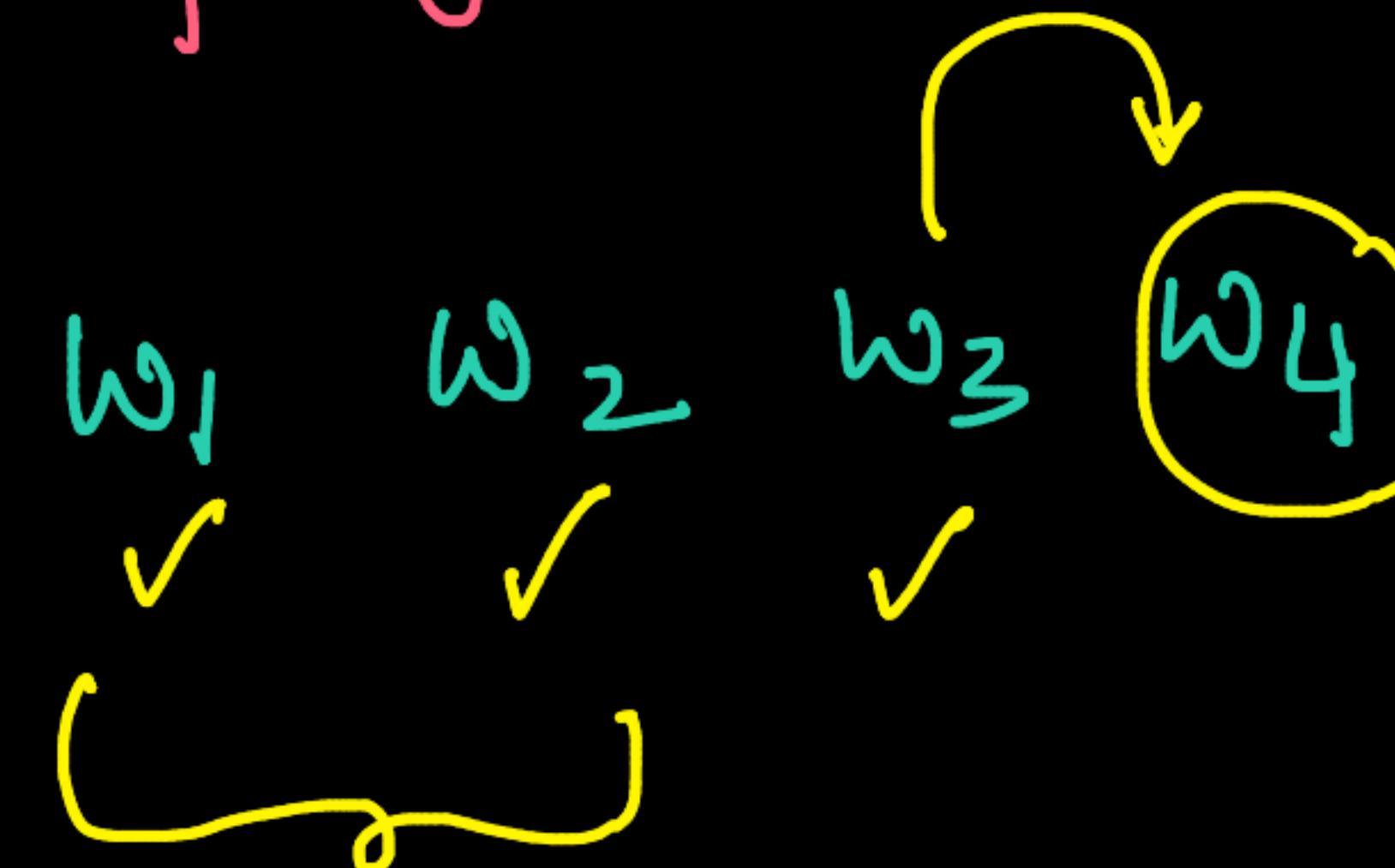
$$P_{\text{train}}(\text{[UNK]}) = \frac{k}{N_{\text{train}} + kV}$$

## N-Gram

- The probability of each word depends on the  $n-1$  words before it. For a trigram model ( $n = 3$ ), for example, each word's probability depends on the 2 words immediately before it.
- This probability is estimated as the fraction of times this  $n$ -gram appears among all the previous  $(n-1)$ -grams in the training set. In other words, training the  $n$ -gram model is nothing but calculating these conditional probabilities from the training text.

HMM

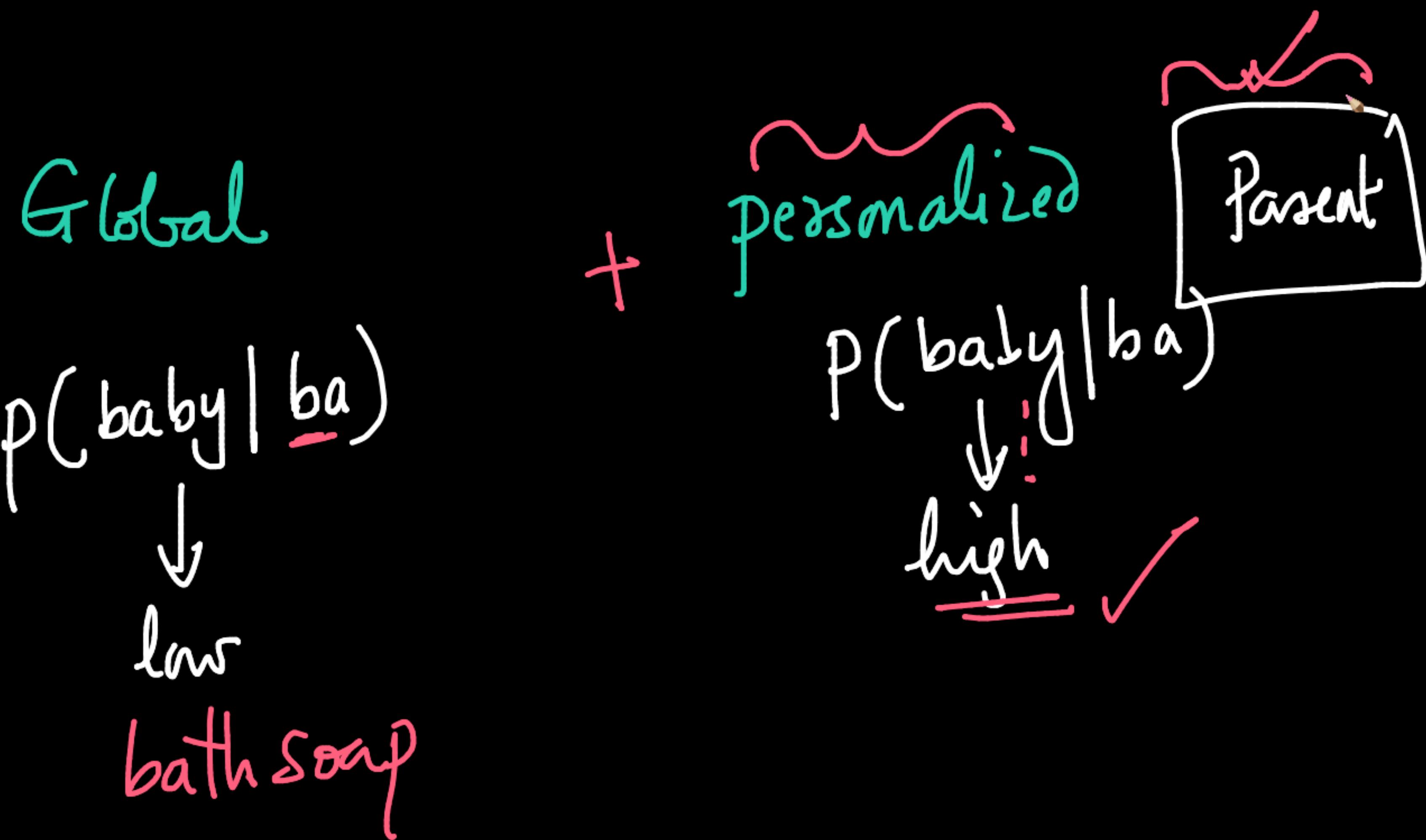
Markov property:  $\rightarrow$  bigram



$$\checkmark \{ P(w_4 | w_2 w_3 w_1) \}$$

$$P(w_4 | w_3)$$

Simple,  
assumption



language\_modeling\_V1.ipynb - x Google Online Shopping site in India: S + colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS#scrollTo=PDZ9tbn36Pw0

+ Code + Text Connect |   

In practice we don't use raw probability as our metric for evaluating language models, but a variant called **perplexity**.

- Perplexity of a language model is the **inverse probability** of the test set, **normalized by the number of words**.



Let us understand Perplexity in detail:

For a test set  $W$ ; containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:
  - $PP(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1/N}$
  - $n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$
- Using chain rule of probability, we can expand:
  - $PP(W) = n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$
- For a bi-gram model, the perplexity of  $W$  can be calculated using:
  - $n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$
- Lower the **perplexity score**, higher the model's **understanding** of the language.

29 / 29

+ Code + Text

Connect



## Let us understand Perplexity in detail:

For a test set  $W$ ; containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:

- $PP(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1/N}$

- $n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$

- Using chain rule of probability, we can expand:

- $PP(W) = n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$

- For a bi-gram model, the perplexity of  $W$  can be calculated using:

- $n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$

- Lower the perplexity score, higher the model's understanding of the language.

```
[ ] def perplexity(n_gram, input_sentence):  
    tokenized_sentence = input_sentence.split()  
    pp_w = []
```

+ Code + Text

Connect



## Let us understand Perplexity in detail:

For a test set  $W$ ; containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:

$$\circ \text{PP}(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1/N}$$

$$\circ n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$$

- Using chain rule of probability, we can expand:

$$\circ \text{PP}(W) = n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

- For a bi-gram model, the perplexity of  $W$  can be calculated using:

$$\circ n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$$

- Lower the perplexity score, higher the model's understanding of the language.

```
[ ] def perplexity(n_gram, input_sentence):  
    tokenized_sentence = input_sentence.split()  
    pp_w = []
```

+ Code + Text

Connect



## Let us understand Perplexity in detail:

For a test set  $\{x\}$  containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:

$$\circ \text{PP}(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1/N}$$

$$\circ n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$$

- Using chain rule of probability, we can expand:

$$\circ \text{PP}(W) = n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

- For a bi-gram model, the perplexity of W can be calculated using:

$$\circ n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$$

- Lower the perplexity score, higher the model's understanding of the language.

```
[ ] def perplexity(n_gram, input_sentence):
    tokenized_sentence = input_sentence.split()
    pp_w = []
```

inverse

$$P(w_1, w_2) = \underbrace{P(w_1)}_{\text{Inverse}} \underbrace{P(w_2)}_{\text{Inverse}}$$

$$P(w_1, w_2) = P(w_1) P(w_1 | w_2)$$

$$P(w_1, w_2, w_3) = P(w_1) \underbrace{P(w_2 | w_1)}_{\text{Inverse}} \underbrace{P(w_3 | w_2)}_{\text{Inverse}}$$

$$P(w_1 | w_2) \\ P(w_2 | w_3)$$

+ Code + Text

Connect



## Let us understand Perplexity in detail:

For a test set  $W$ ; containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:

$$\circ \text{PP}(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1} N$$

$$\circ n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$$

- Using chain rule of probability, we can expand:

$$\circ \text{PP}(W) = n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

- For a bi-gram model, the perplexity of  $W$  can be calculated using:

$$\circ n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$$

- Lower the perplexity score, higher the model's understanding of the language.

```
[ ] def perplexity(n_gram, input_sentence):  
    tokenized_sentence = input_sentence.split()  
    pp_w = [
```

+ Code + Text

Connect



## Let us understand Perplexity in detail:

{x}

For a test set  $W$ ; containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:

- $PP(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1/N}$

- $n \sqrt[N]{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$

- Using chain rule of probability, we can expand:

- $PP(W) = n \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$

- For a bi-gram model, the perplexity of  $W$  can be calculated using:

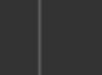
- $n \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$

- Lower the perplexity score, higher the model's understanding of the language.

```
[ ] def perplexity(n_gram, input_sentence):
    tokenized_sentence = input_sentence.split()
    pp_w = [
```

+ Code + Text

Connect



## Let us understand Perplexity in detail:

For a test set  $W$ ; containing words  $w_1, w_2, w_3, \dots, w_n$

- The perplexity can be derived using:

- $PP(W) = P(w_1, w_2, w_3, \dots, w_n)^{-1/N}$

- $n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$

- Using chain rule of probability, we can expand:

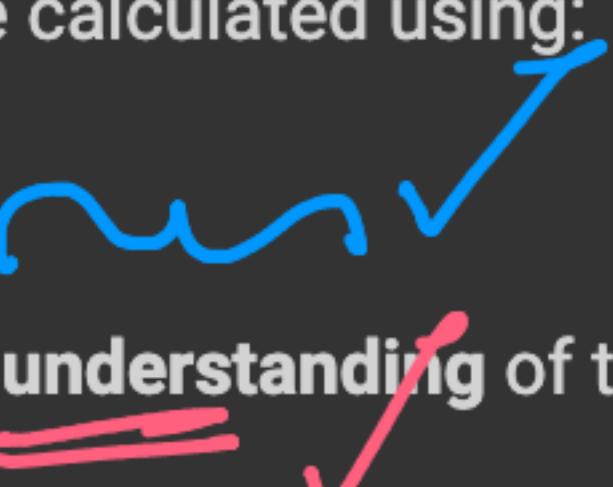
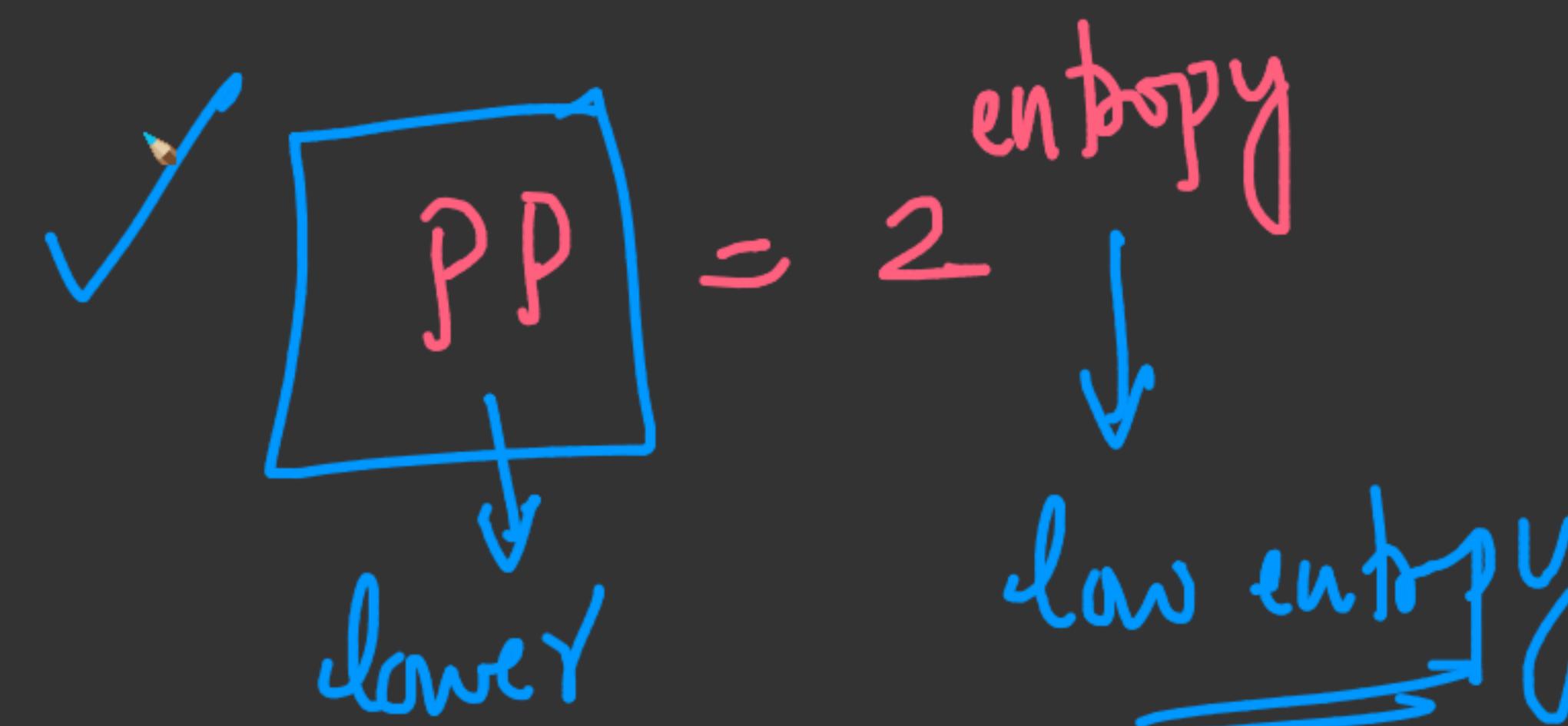
- $PP(W) = n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$

- For a bi-gram model, the perplexity of  $W$  can be calculated using:

- $n \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$

- Lower the perplexity score, higher the model's understanding of the language.

```
[ ] def perplexity(n_gram, input_sentence):  
    tokenized_sentence = input_sentence.split()  
    pp_w = [
```



+ Code + Text

Connect



$$\circ n \sqrt{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}}$$

- Using chain rule of probability, we can expand:

$$\circ PP(W) = n \sqrt{\prod_{i=1}^n \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

- For a bi-gram model, the perplexity of W can be calculated using:

$$\circ n \sqrt{\prod_{i=1}^n \frac{1}{P(w_i | w_{i-1})}}$$

- Lower the perplexity score, higher the model's understanding of the language.

$PP \geq 2^{entropy}$

```
[ ] def perplexity(n_gram, input_sentence):  
    tokenized_sentence = input_sentence.split()  
    pp_w = []  
    N = len(tokenized_sentence)  
  
    if n_gram == 2:  
        tokenized_sentence.insert(0, None)  
        tokenized_sentence.append(None)  
        for i in tokenized_sentence[:-1]:  
            try:  
                pp_w.append(1 / (0.0001+bigram[i][tokenized_sentence[tokenized_sentence.index(i)+1]]))  
            except:  
                continue
```

+ Code + Text

Connect |  

$$\circ P_{Add-k}^*(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + k}{C(w_{n-1}) + kV}$$

## Backoff and Interpolation

- Most applied technique for OOV and produces the best results when compared to smoothing techniques.
- Instead of smoothing and shifting the distribution, we can rely on the knowledge we have already computed.
- In this technique, if we have no examples of a particular tri-gram  $P(w_n|w_{n-2}w_{n-1})$ , we can simply estimate its probability by using the bigram probability  $P(w_n|w_{n-1})$ .
- In case we don't have counts to compute probability of the bi-gram, we can cascade and use the probability of the unigram  $P(w_n)$ .
- Using interpolation helps in generalization across different contexts.
- Techniques to make the interpolation effective:
  - Linear interpolation - Probabilities across each n-gram will have same weightage.
  - " $\lambda$ " interpolation - Weights are assigned to the probabilities based on their computation method.

- Tri-gram probabilities will have the highest weights, followed by bi-gram and uni-gram.

$$P(w_n|w_{n-2}w_{n-1}) = \underbrace{\lambda_1 P(w_n)}_{\text{unigram}} + \underbrace{\lambda_2 P(w_n|w_{n-1})}_{\text{bigram}} + \underbrace{\lambda_3 P(w_n|w_{n-2}w_{n-1})}_{\text{Trigrams}}$$

What are the limitations of N-gram approach to Language Modeling?

(e.g)

$$p(\text{click} \mid \underbrace{\text{HitechCity}}_{\text{loc}}) = \frac{2}{1000}$$

$$p(\text{clk} \mid \text{Hyd}) = \frac{200^4}{100000}$$

↓ backoff

+ Code + Text

Connect



$$\circ P_{Add-k}^*(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + k}{C(w_{n-1}) + kV}$$

## Backoff and Interpolation

- Most applied technique for OOV and produces the best results when compared to smoothing techniques.
- Instead of smoothing and shifting the distribution, we can rely on the knowledge we have already computed.
- In this technique, if we have no examples of a particular tri-gram  $P(w_n|w_{n-2}w_{n-1})$ , we can simply estimate its probability by using the bigram probability  $P(w_n|w_{n-1})$ .
- In case we don't have counts to compute probability of the bi-gram, we can cascade and use the probability of the unigram  $P(w_n)$ .
- Using interpolation helps in generalization across different contexts.
- Techniques to make the interpolation effective:
  - Linear interpolation - Probabilities across each n-gram will have same weightage.
  - " $\lambda$ " interpolation - Weights are assigned to the probabilities based on their computation method.
    - Tri-gram probabilities will have the highest weights, followed by bi-gram and uni-gram.

$$P(w_n|w_{n-2}w_{n-1}) = \lambda_1 P(w_n) + \lambda_2 P(w_n|w_{n-1}) + \lambda_3 P(w_n|w_{n-2}w_{n-1})$$

What are the limitations of N-gram approach to Language Modeling?

+ Code + Text

Connect |   

- n-gram probabilities will have the highest weights, followed by bi-gram and uni-gram.

$$P(w_n|w_{n-2} w_{n-1}) = \lambda_1 P(w_n) + \lambda_2 P(w_n|w_{n-1}) + \lambda_3 P(w_n|w_{n-2} w_{n-1})$$

{x} What are the limitations of N-gram approach to Language Modeling?

□

- Higher the N, the better is the model. But this leads to high computation overhead and would require large computation power, which may not be easily available.
- ✓ {• N-grams are a sparse representation of language. This is because we build the model based on the probability of words co-occurring.  
Zero probabilities are assigned to words that are not present in the training corpus.
- Even with techniques to build non-zero probability models, the results may not be accurate in case of generalization of words over multiple contexts.

We can overcome the limitations of N-gram based Lanugage Models by using Neural Language Models.

- Neural Language Models (NLM) address the n-gram issues through parameterization of words as vectors (word embeddings) and using them as inputs to a neural network.
- We could use robust architectures like RNN and Encoders to understand word contexts better, there-by being more accurate in predicting the next set of words.

40 / 40

+ Code + Text

Connect



multiple contexts.

## We can overcome the limitations of N-gram based Lanugage Models by using Neural Language Models.

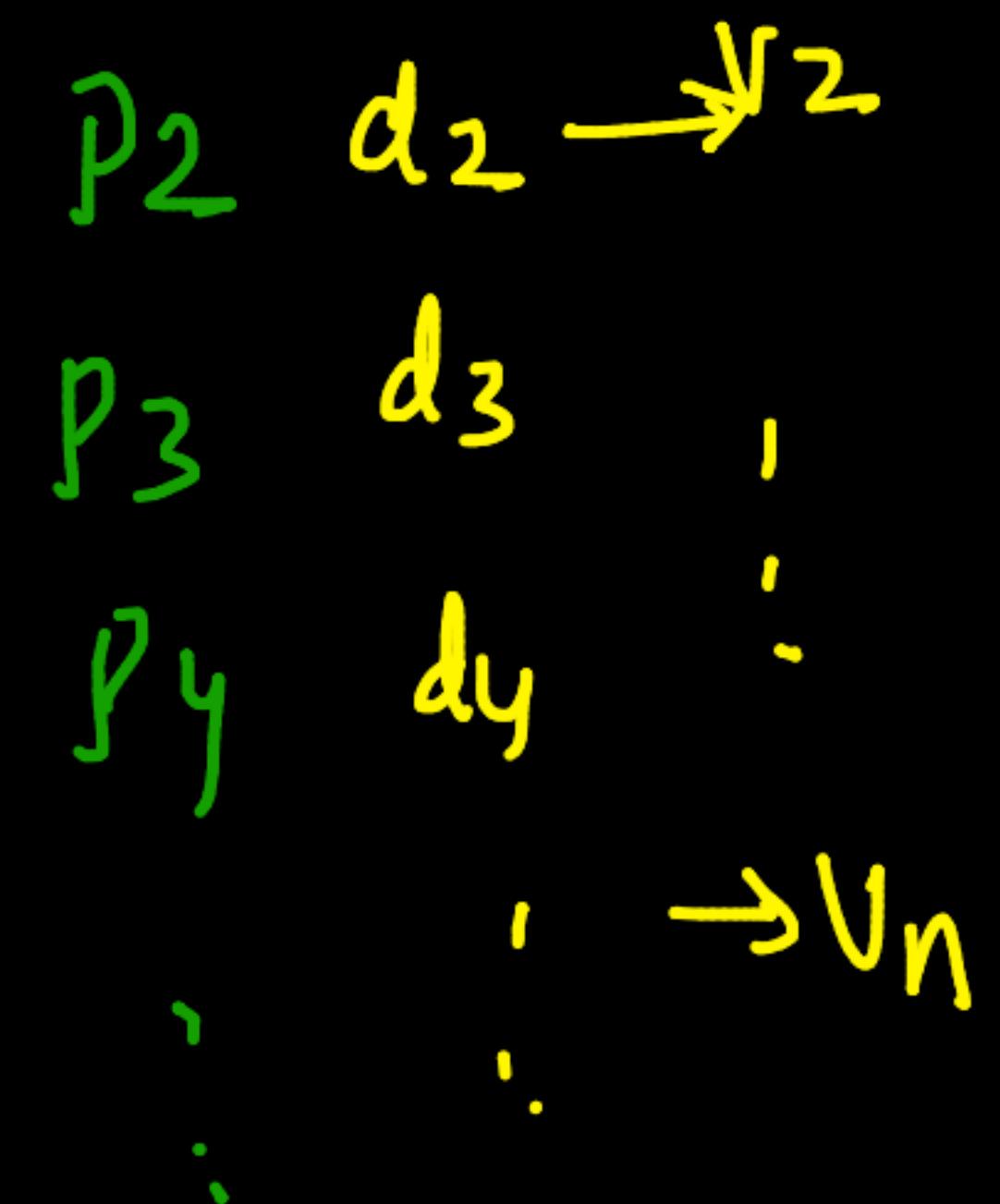
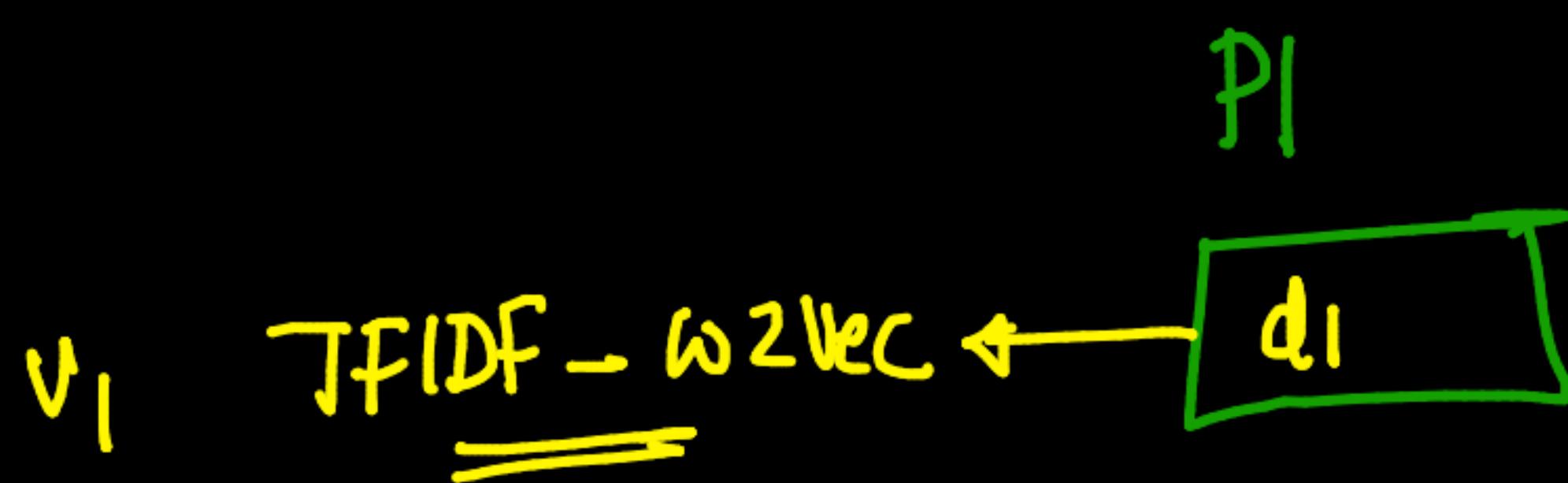
- Neural Language Models (NLM) address the n-gram issues through parameterization of words as vectors (word embeddings) and using them as inputs to a neural network.
- We could use robust architectures like RNN and Encoders to understand word contexts better, there-by being more accurate in predicting the next set of words.

LSM | GRU

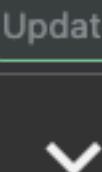
## Conclusion:

- We have modelled the problem of helping people write emails faster as "next-word" prediction problem.
- We have explored Probabilistic Language Models.
- We auto-completed the sentence by building bi-gram and tri-gram language models.
- We explored Perplexity as a technique to ascertain the LM's understanding.
- We saw how Smoothing and interpolation techniques can help over-come the dis-advantages of OOV.

Content-based:

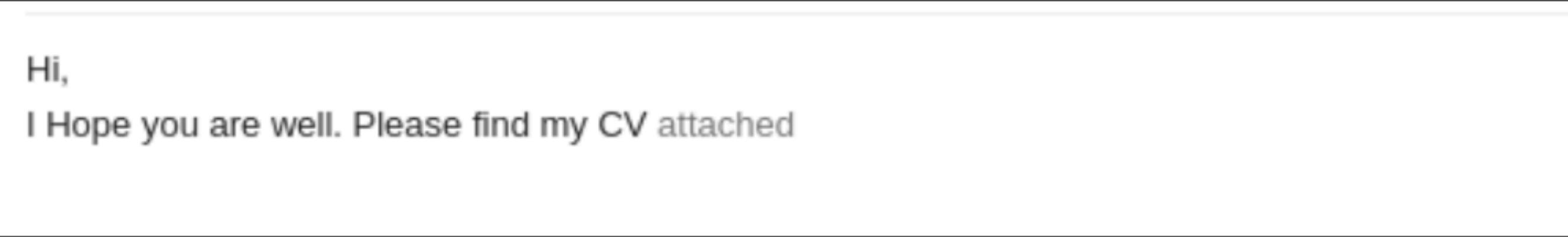


language\_modeling\_V1.ipynb - x Perplexity Intuition (and its deri x Google x | Online Shopping site in India: S x + colab.research.google.com/drive/1FUZIB-plqt1OX6bGqXtvIbpf0lxrb\_BS#scrollTo=tZ9QQZ64hB14

+ Code + Text Connect |   

**Email autofill**

\*You are working as machine learning engineer at Gmail team at Google. Your manager want you to design a solution to suggest the next word in email when they are typing it, to improve their experience and save time. Example as we can see in the below image:-



What are the other real life application of autofill? 1.typing in whatsapp (auto suggest three possible next words) 2.Autofill Google Search



43 / 43