

# NLP Basic

## 02 – LANGUAGE MODEL

AI VIET NAM  
Nguyen Quoc Thai

# CONTENT

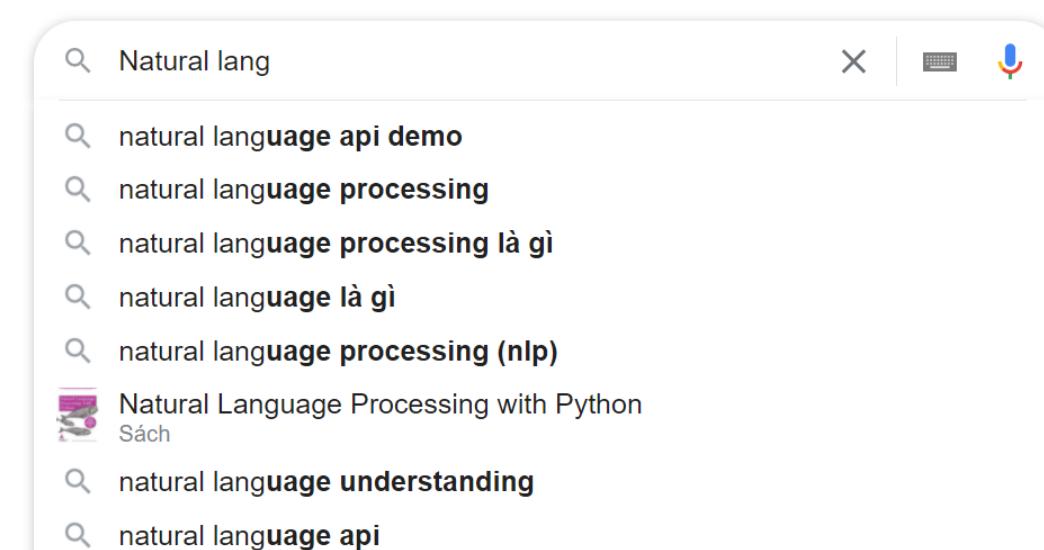
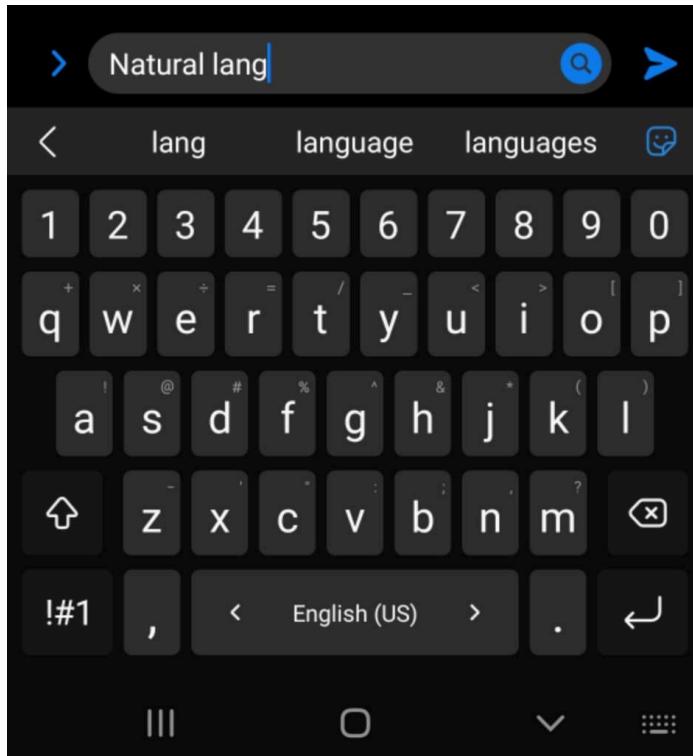
- |   |  |
|---|--|
| 1 | <b>N-grams Language Model</b>              |
| 2 | <b>Evaluating Language Model</b>           |
| 3 | <b>Generalization and Zeros</b>            |
| 4 | <b>Neural Probabilistic Language Model</b> |

# 1 – N-grams Language Model

!

## Language Model

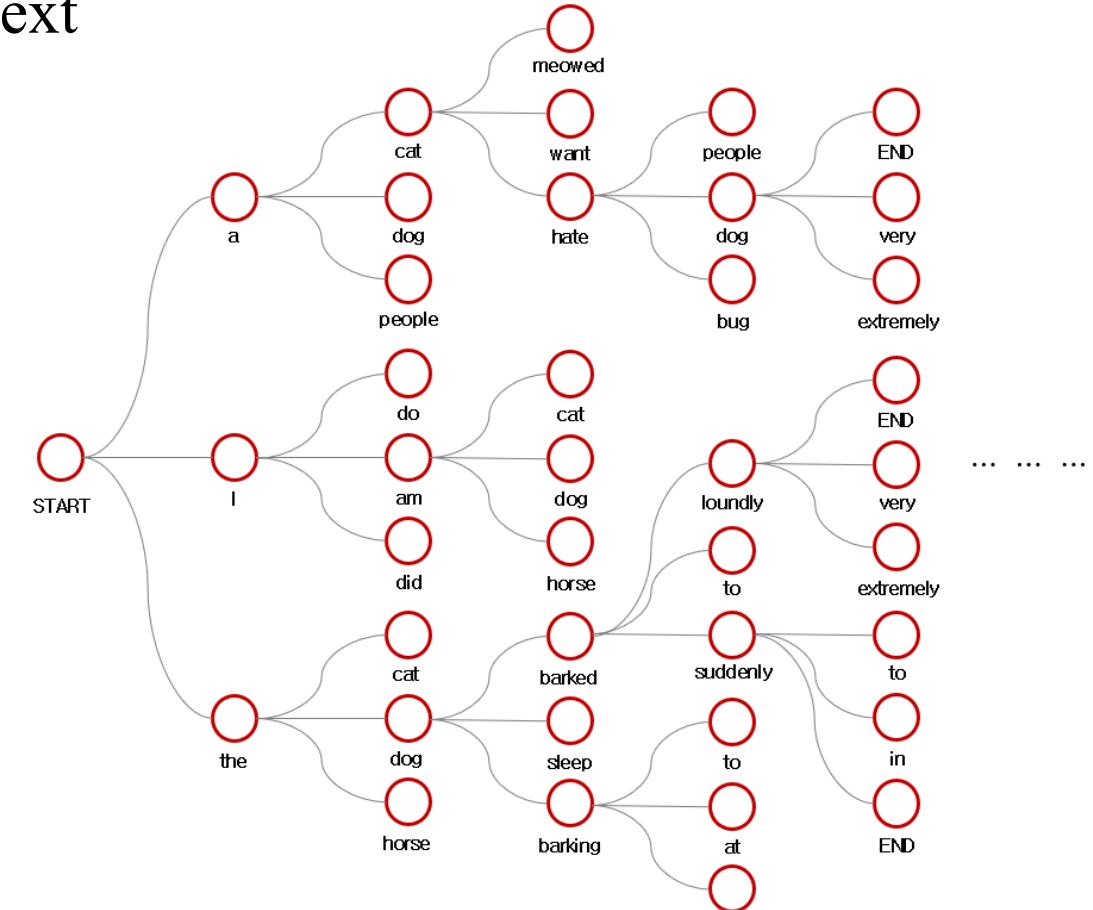
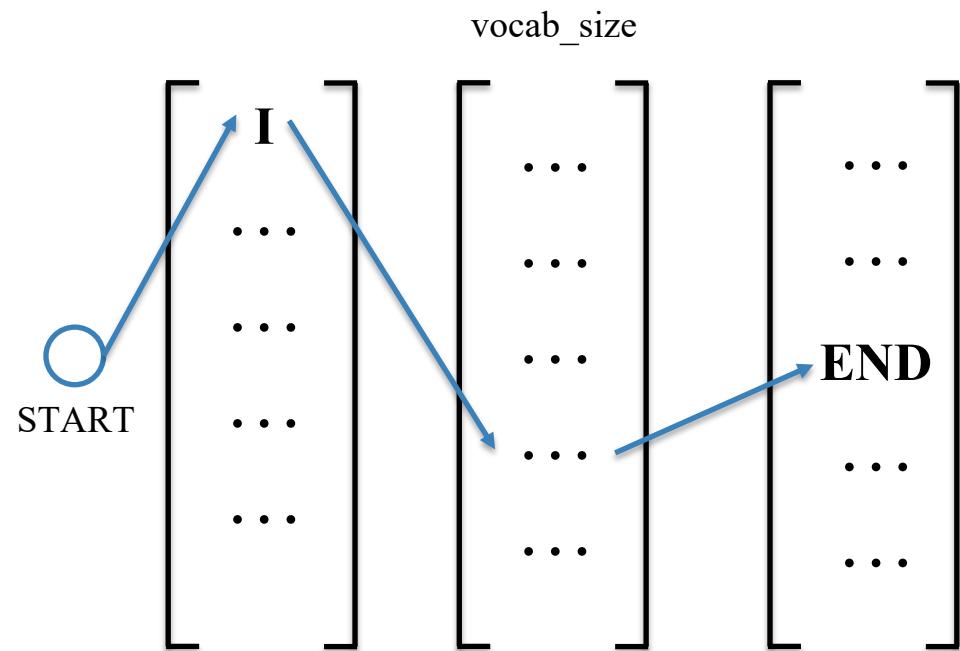
- The task of predicting what word comes next



# 1 – N-grams Language Model



- The task of predicting what word comes next

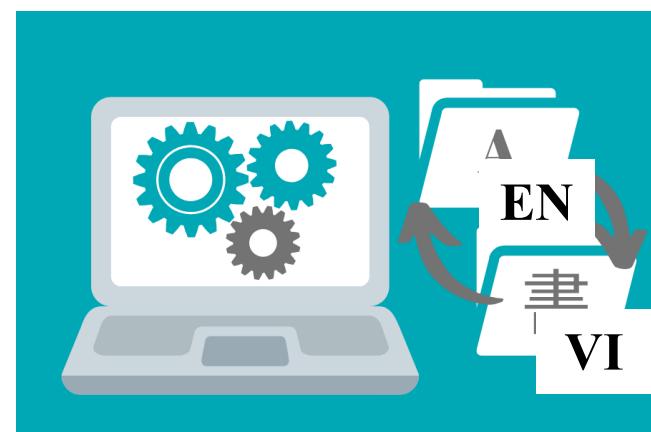


# 1 – N-grams Language Model

!

## Language Model

- Goal: assign a probability to a sentence



Machine Translation

Input Sentence

“Tôi đi học.”

Candidate Sentence

“I go to school.”

$P(\text{"I go to school."})$

>

“I go to work.”

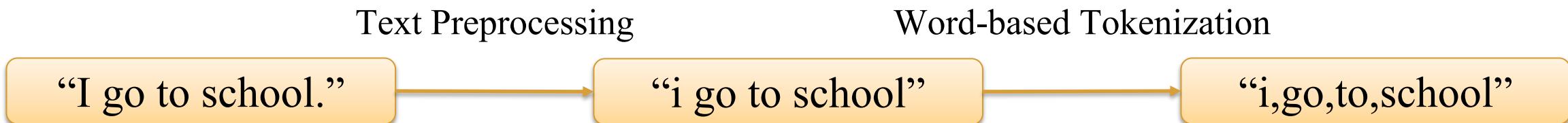
$P(\text{"I go to work."})$

# 1 – N-grams Language Model

!

## Language Model

- The task of predicting what word comes next



- Compute the probability of occurrence of the number of words (tokens) in a sentence (sequence)

$$P(\text{"i,go,to,school"})$$

$$P(W) = P(w_1, w_2, \dots, w_n)$$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- Compute  $P(W)$

$$P(W) = P(w_1, w_2, \dots, w_n)$$

- Conditional Probability

$$P(B|A) = P(AB)/P(A) \Rightarrow P(AB) = P(A).P(B|A)$$

- The chain rule of probability

$$P(A,B,C,D) = P(A).P(B|A).P(C|AB).P(D|ABC)$$

$$\begin{aligned} P(W) &= P(w_1, w_2, \dots, w_n) = P(w_1).P(w_2|w_1).P(w_3|w_{1:2}) \dots P(w_n|w_{1:n-1}) \\ &= \prod_{k=1}^n P(w_k|w_{1:k}) \end{aligned}$$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- Compute  $P(W)$

$$\begin{aligned} P(W) &= P(w_1, w_2, \dots, w_n) = P(w_1).P(w_2|w_1).P(w_3|w_{1:2}) \dots P(w_n|w_{1:n-1}) \\ &= \prod_{k=1}^n P(w_k|w_{1:k}) \end{aligned}$$

- Compute the probability of an upcoming words

“I go to school.”

$P(\text{"i,go,to,school"})$

$= P(i).P(go|i).P(to|i,go).P(school|i,go,to)$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- Compute the probability of an upcoming words

“I go to school.”

$P(\text{"i,go,to,school"})$

$= P(i).P(\text{go}|i).P(\text{to}|i,\text{go}).P(\text{school}|i,\text{go},\text{to})$

- Compute  $P(w|h)$

$P(w|h) = P(\text{school}|i,\text{go},\text{to})$

w: token as word “school”

h: history tokens as “i,go,to”

# 1 – N-grams Language Model

!

## Probabilistic Language Model

➤ Compute  $P(w|h)$

$$P(w|h) = P(\text{school}|i, \text{go}, \text{to})$$

w: token as word “school”

h: history tokens as “i,go,to”

$$P(w|h) = \frac{\text{count}(hw)}{\text{count}(h)}$$

$$P(\text{school}|i, \text{go}, \text{to}) = \frac{\text{count}(i, \text{go}, \text{to}, \text{school})}{\text{count}(i, \text{go}, \text{to})}$$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- The probability of a word depends only on some previous words

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

$$P(w_i | w_{1:i-1}) = P(w_i | w_{i-N+1:i-1})$$

- N-gram model with  $N = \{1, 2, \dots\}$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- N = 1
- Unigram Model (1-gram)

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

$$P(\text{"i,go,to,school"}) \rightarrow = P(i).P(go|i).P(to|i,go).P(school|i,go,to)$$

$$= P(i).P(go).P(to).P(school)$$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- N = 2
- Bigram Model (2-gram)

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

P("i,go,to,school")

$$= \boxed{P(i)} P(go|i).P(to|i,go).P(school|i,go,to)$$

How to compute P(i) ?

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- $N = 2$
- Bigram Model (2-gram)

 $P("i,go,to,school")$ 

$$= P(i).P(go|i).P(to|i,go).P(school|i,go,to)$$

How to compute  $P(i)$  ?

(1) - Remove

$$= P(go|i).P(to|i,go).P(school|i,go,to)$$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- $N = 2$
- Bigram Model (2-gram)

 $P("i,go,to,school")$ 

$$= P(i).P(go|i).P(to|i,go).P(school|i,go,to)$$

How to compute  $P(i)$  ?

(1) - Remove

(2) - Padding

Pad = 1

 $P("< s >, i, go, to, school, </ s >")$ 

$$\begin{aligned} &= P(i|<s>).P(go|<s>,i).P(to|<s>,i,go) \\ &.P(school|<s>,i,go,to).P(</s>|<s>,i,go,to,school) \end{aligned}$$

$$\begin{aligned} &= P(i|<s>).P(go|i).P(to|go) \\ &.P(school|to).P(</s>|school) \end{aligned}$$

# 1 – N-grams Language Model

!

## Probabilistic Language Model

- N = 3
- Trigram Model (3-gram)

$$P(w_{1:n}) = \prod_{i=1}^n P(w_i | w_{i-N+1:i-1}) = \prod_{i=1}^n P(w_i | w_{i-2:i-1})$$



# 1 – N-grams Language Model

!

## Probabilistic Language Model

- $N = 3$
- Trigram Model (3-gram)



$$\begin{aligned} &= P(i|<S>,<S>).P(go|<S>,<S>,i).P(to|<S>,<S>,i,go).P(school|<S>,<S>,i,go,to) \\ &\quad .P(</S>|<S>,<S>,i,go,to,school). P(</S>|<S>,<S>,i,go,to,school,</S>) \end{aligned}$$

$$\begin{aligned} &= P(i|<S>,<S>).P(go|<S>,i).P(to|i,go).P(school|go,to) \\ &\quad .P(</S>|to,school). P(</S>|school,</S>) \end{aligned}$$

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ N-gram Probability:

$$P(w_i | w_{i-N+1:i-1}) = \frac{\text{count}(w_{i-N+1:i-1}, w_i)}{\text{count}(w_{i-N+1:i-1})}$$

➤ Bigram Probability

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

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ Example:

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

< s > tôi đang học < /s >  
< s > tôi đang học lớp nlp < /s >  
< s > lớp nlp có vẻ hơi vui < /s >

$$P(\text{tôi}|<\text{s}>) =$$

$$P(<\text{/s}>|\text{học}) =$$

$$P(\text{đang}|\text{tôi}) =$$

$$P(\text{lớp}|\text{học}) =$$

$$P(\text{học}|\text{đang}) =$$

$$P(\text{nlp}|\text{lớp}) =$$

➤  $P(\text{tôi, đang, học}) =$

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ Example:

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

< s > tôi đang học < /s >  
< s > tôi đang học lớp nlp < /s >  
< s > lớp nlp có vẻ hơi vui < /s >

$$P(\text{tôi}|<\text{s}>) = 2/3$$

$$P(\text{đang}|\text{tôi}) = 2/2$$

$$P(\text{học}|\text{đang}) = 2/2$$

$$P(<\text{s}>|\text{học}) = 1/2$$

$$P(\text{lớp}|\text{học}) = 1/2$$

$$P(\text{nlp}|\text{lớp}) = 2/2$$

❖  $P(\text{tôi, đang, học}) = P(\text{tôi}|<\text{s}>).P(\text{đang}|\text{tôi}).P(\text{học}|\text{đang}).(<\text{s}>|\text{học})$   
 $= 2/3.2/2.2/2.1/2 = 1/3$

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ “Truyện Kiều – Nguyễn Du”:

Trăm năm trong cõi người ta,

Chữ tài chữ mệnh khéo là ghét nhau.

Trải qua một cuộc bể dâu,

Những điều trông thấy mà đau đớn lòng.

Lạ gì bỉ sắc tư phong,

Trời xanh quen thói má hồng đánh ghen.

Cảo thơm lần giờ trước đèn,

Phong tình cổ lục còn truyền sử xanh.

Răng năm Gia Tĩnh triều Minh,

Bốn phương phảng lặng, hai kinh vũng vàng.

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ “Truyện Kiều – Nguyễn Du”:

	<s>	trăm	năm	trong	cõi	người	ta
<s>	0	13	5	23	1	22	1
trăm	0	0	9	0	0	0	0
năm	0	0	4	1	0	1	0
trong	0	0	1	0	3	1	0
cõi	0	0	0	0	0	1	0
người	0	0	0	2	0	2	7
ta	0	0	0	0	0	0	2

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ “Truyện Kiều – Nguyễn Du”:

	<s>	trăm	năm	trong	cõi	người	ta
<s>	0	0.004	0.0015	0.007	0.0003	0.0068	0.0003
trăm	0	0	0.29	0	0	0	0
năm	0	0	0.077	0.019	0	0.019	0
trong	0	0	0.0095	0	0.0285	0.0095	0
cõi	0	0	0	0	0	0.1	0
người	0	0	0	0.009	0	0.009	0.031
ta	0	0	0	0	0	0	0.035

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ “Truyện Kiều – Nguyễn Du”:

Sentence: “Trăm năm trong cõi người ta”

$$\begin{aligned} & P(\text{trăm}, \text{năm}, \text{trong}, \text{cõi}, \text{người}, \text{ta}) \\ &= P(\text{trăm}|\langle s \rangle) \cdot P(\text{năm}|\text{trăm}) \cdot P(\text{trong}|\text{năm}) \cdot P(\text{cõi}|\text{trong}) \\ &\quad \cdot P(\text{người}|\text{cõi}) \cdot P(\text{ta}|\text{người}) \cdot P(\langle /s \rangle|\text{ta}) \\ &= 0.004 * 0.29 * 0.019 * 0.0285 * 0.1 * 0.031 * 0.439 = 8.5e-10 \end{aligned}$$

# 1 – N-grams Language Model

!

## Estimating N-gram Probability

➤ “Truyện Kiều – Nguyễn Du”:

Sentence: “Trăm năm trong cõi người ta”

$$\begin{aligned} & P(\text{trăm}, \text{năm}, \text{trong}, \text{cõi}, \text{người}, \text{ta}) \\ &= P(\text{trăm}|\langle s \rangle) \cdot P(\text{năm}|\text{trăm}) \cdot P(\text{trong}|\text{năm}) \cdot P(\text{cõi}|\text{trong}) \\ &\quad \cdot P(\text{người}|\text{cõi}) \cdot P(\text{ta}|\text{người}) \cdot P(\langle /s \rangle|\text{ta}) \\ &= 0.004 * 0.29 * 0.019 * 0.0285 * 0.1 * 0.031 * 0.439 = 8.5e-10 \end{aligned}$$

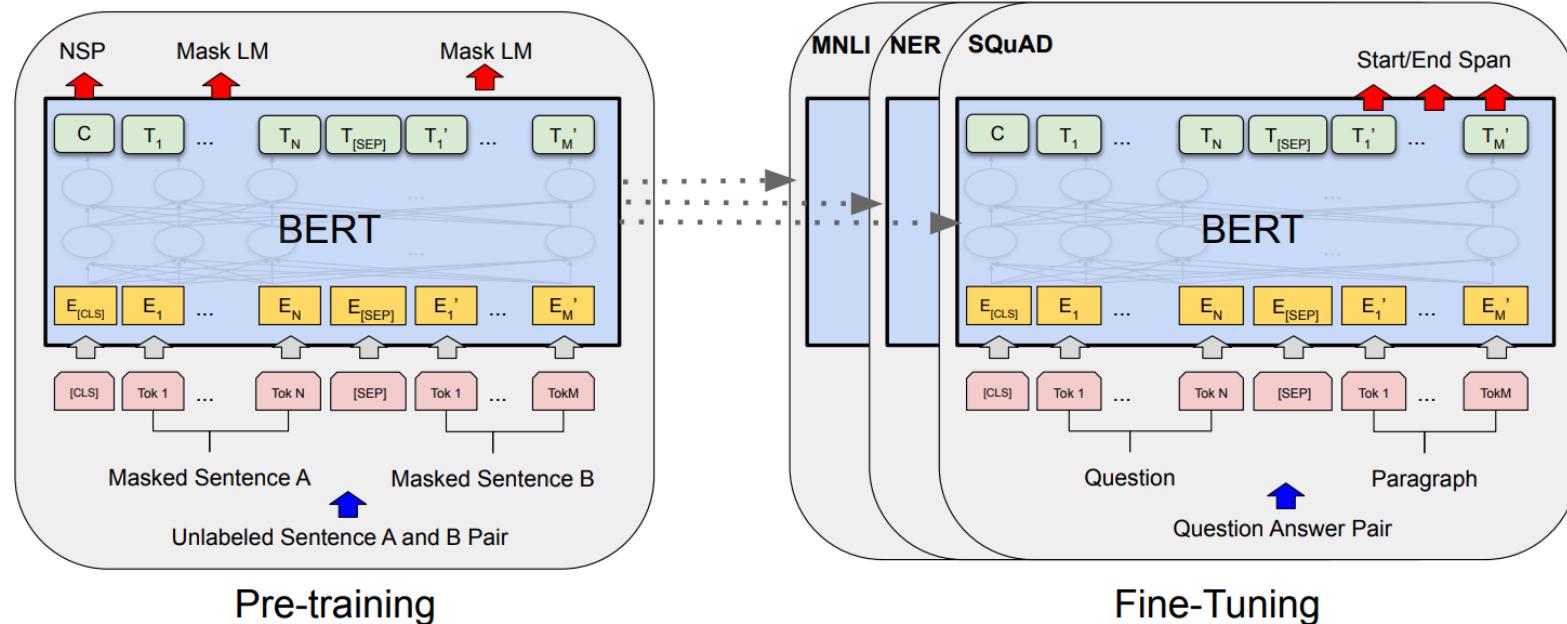
Numerical Underflow => Convert to log space

# 2 – Evaluating Language Model



## Extrinsic Evaluation

- ❖ Best way to evaluate the performance
- ❖ Embedded models into an application (downs task)
- ❖ NLP systems end-to-end: very expensive



# 2 – Evaluating Language Model



## Extrinsic Evaluation

Model	Vocabulary Size(K)	Backbone #Params(M)	MNLI-m/mm ACC	SQuAD v2.0 F1/EM
Base models:12 layers,768 hidden size,12 heads				
BERT <sub>base</sub>	30	86	84.3/84.7	76.3/73.7
RoBERTa <sub>base</sub>	50	86	87.6/-	83.7/80.5
XLNet <sub>base</sub>	32	92	86.8/-	-/80.2
ELECTRA <sub>base</sub>	30	86	88.8/-	-/80.5
DeBERTa <sub>base</sub>	50	100	88.8/88.5	86.2/83.1
DeBERTaV3 <sub>base</sub>	128	86	<b>90.6/90.7</b>	<b>88.4/85.4</b>
Small models:6 layers,768 hidden size,12 heads				
TinyBERT <sub>small</sub>	30	44	84.5/-	77.7/-
MiniLMv2 <sub>small</sub>	30	44	87.0/-	81.6/-
BERT <sub>small</sub>	30	44	81.8/-	73.2/-
DeBERTaV3 <sub>small</sub>	128	44	<b>88.2/87.9</b>	<b>82.9/80.4</b>
XSmall models:12 layers,384 hidden size,6 heads				
MiniLMv2 <sub>xsmall</sub>	30	22	86.9/-	82.3/-
DeBERTaV3 <sub>xsmall</sub>	128	22	<b>88.1/88.3</b>	<b>84.8/82.0</b>

# 2 – Evaluating Language Model

!

## Intrinsic Evaluation: Perplexity

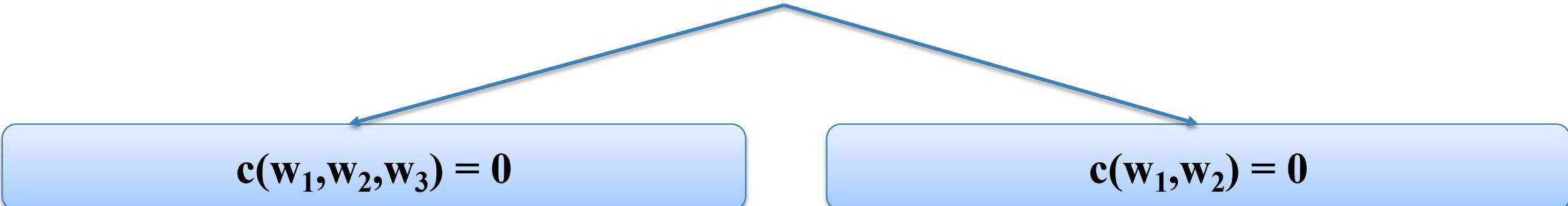
- ❖ Perplexity (PP) is the inverse probability of the test set, normalized by the number of words
- ❖ Lower perplexity => better model
- ❖ Test sample:  $W = \{w_1, w_2, \dots, w_n\}$

$$PP(W) = P(w_1 w_2 \dots w_n)^{-\frac{1}{n}}$$
$$PP(W) = \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_{1:i-1})}}$$

- ❖ Bigram Models:  $PP(W) = \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_{i-1})}}$

# 3 – Generalization and Zeros

$$P(w_3|w_1 w_2) = \frac{c(w_1, w_2, w_3)}{c(w_1, w_2)}$$

 $c(w_1, w_2, w_3) = 0$  $c(w_1, w_2) = 0$

# 3 – Generalization and Zeros

$$P(w_3|w_1 w_2) = \frac{c(w_1, w_2, w_3)}{c(w_1, w_2)}$$

**c(w<sub>1</sub>,w<sub>2</sub>,w<sub>3</sub>) = 0**

**c(w<sub>1</sub>,w<sub>2</sub>) = 0**

**Smoothing**

$$P_{\text{Add}-1}(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

**Backoff**

$$P(w_i|w_{i-2:i-1}) \Rightarrow P(w_i|w_{i-1}) \Rightarrow P(w_i)$$

**Interpolation**

$$\begin{aligned}\widehat{P}(w_n|w_{n-2}w_{n-1}) &= \lambda_1 P(w_n|w_{n-2}w_{n-1}) \\ &\quad + \lambda_2 P(w_n|w_{n-1}) + \lambda_3 P(w_n) \\ \sum_i \lambda_i &= 1\end{aligned}$$

# 3 – Generalization and Zeros

!

## Add-1 Smoothing – “Truyện Kiều”

	<s>	trăm	năm	trong	cõi	người	ta
<s>	0	13	5	23	1	22	1
trăm	0	0	9	0	0	0	0
năm	0	0	4	1	0	1	0
trong	0	0	1	0	3	1	0
cõi	0	0	0	0	0	1	0
người	0	0	0	2	0	2	7
ta	0	0	0	0	0	0	2

# 3 – Generalization and Zeros

!

## Add-1 Smoothing – “Truyện Kiều”

	<s>	trăm	năm	trong	cõi	người	ta
<s>	1	14	6	24	2	23	2
trăm	1	1	10	1	1	1	1
năm	1	1	5	2	1	2	1
trong	1	1	2	1	4	2	1
cõi	1	1	1	1	1	2	1
người	1	1	1	3	1	3	8
ta	1	1	1	1	1	1	3

# 3 – Generalization and Zeros

!

## Add-1 Smoothing – “Truyện Kiều”

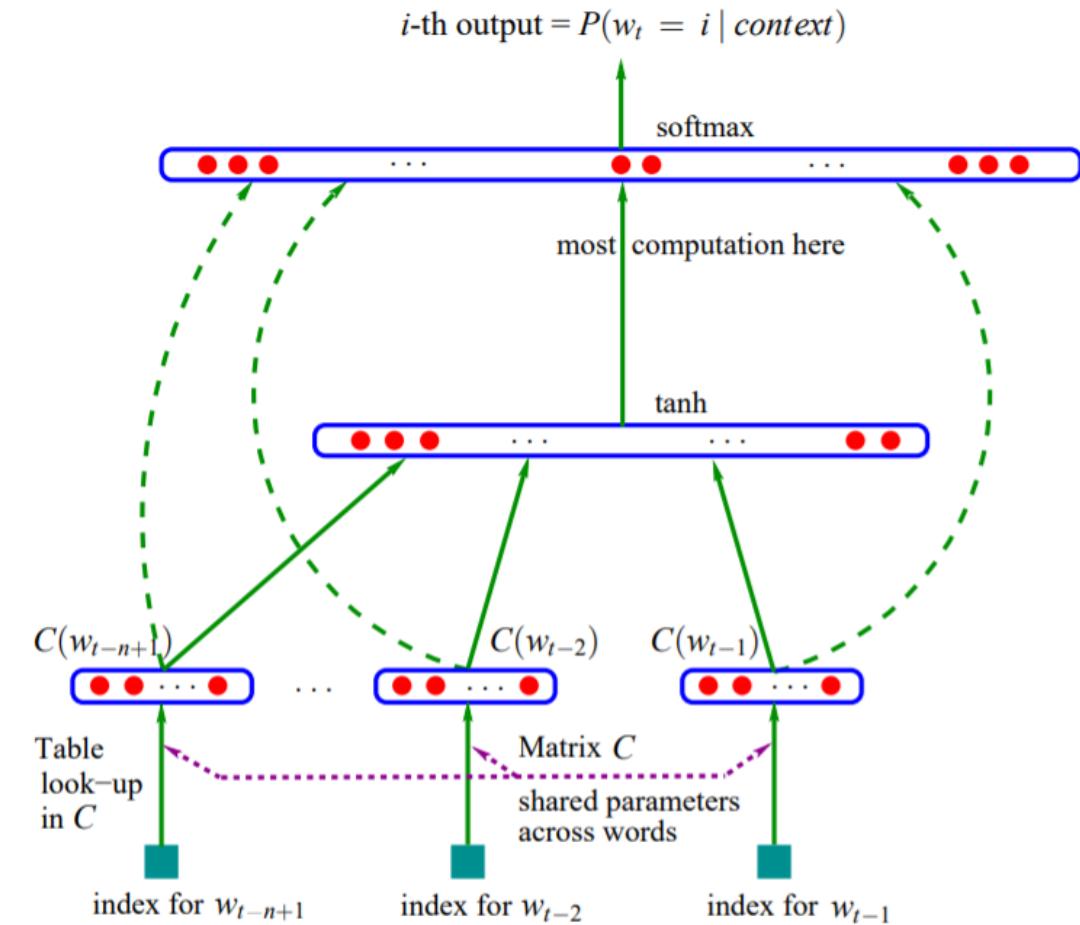
	<s>	trăm	năm	trong	cõi	người	ta
<s>	0.00015	0.0021	0.0009	0.0036	0.0003	0.0035	0.003
trăm	0.0004	0.0004	0.0041	0.0004	0.0004	0.0004	0.0004
năm	0.0004	0.0004	0.002	0.0008	0.0004	0.0008	0.0004
trong	0.0004	0.0004	0.0008	0.0004	0.0015	0.0008	0.0004
cõi	0.0004	0.0004	0.0004	0.0004	0.0004	0.0008	0.0004
người	0.0004	0.0004	0.0004	0.0011	0.0004	0.0011	0.003
ta	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0012

# 4 – Neural Language Model

Source: [Paper]- A Neural Probabilistic Language Model

“trăm năm trong cõi người ta”

Source	Target
trăm	năm
...	...
trăm năm	trong
...	...
trăm năm trong	cõi
...	...
trăm năm trong cõi người	ta



# Thanks!

Any questions?