

CS183 Review Linguistic

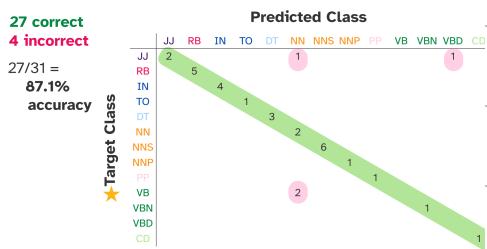
L2. Tasks, Data, Evaluation

- Tasks and Evaluation Metric

- Task: Classification : input text data, output label from a closed set
E.g. Part of Speech Tagging (词性标注)

For this task, metric can be ? Prediction class.

Def: Confusion Matrix : Target class | O O O-
Diagonal : Corrects!
⇒ Accuracy:



⇒ Precision:

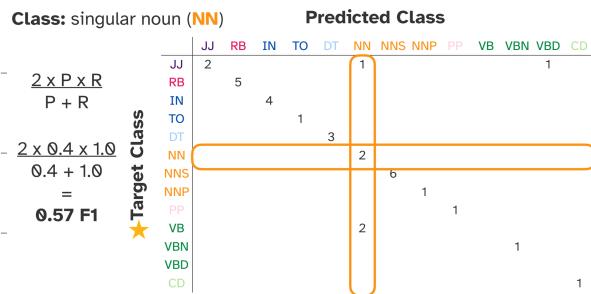
在 class 上 predict 正确率
For a particular class, how many of the model's predictions of that class are accurate? (how precise is the model?)

Class: singular noun (NN) Predicted Class

	JJ	RB	IN	TO	DT	NN	NNS	NNP	PP	VB	VBN	VBD	CD
Target Class													
JJ	2	5	4	1	1	1	3	2	6	1	1	1	1
RB													
IN													
TO													
DT													
NN													
NNS													
NNP													
PP													
VB													
VBN													
VBD													
CD													

⇒ Recall:

With Precision & Recall ⇒ F1 Score:



$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

For examples with a particular correct class label, how often does the model accurately predict that label? (how well does the model recall the items for a particular class?)

Class: singular noun (NN) Predicted Class

	JJ	RB	IN	TO	DT	NN	NNS	NNP	PP	VB	VBN	VBD	CD
Target Class													
JJ	2	5	4	1	1	1	3	2	6	1	1	1	1
RB													
IN													
TO													
DT													
NN													
NNS													
NNP													
PP													
VB													
VBN													
VBD													
CD													

- Task: Automatic Speech Recognition: Utterance ⇒ text

Metric: Word Error Rate (WER)

Possible Error: Substitution (S), Deletion (D) & Insertion (I)

$$WER = (S+D+I)/N. \quad N \text{ is the \# of words in GT script}$$

Key: Calculate Corrects 'S+D+I' ⇒ Levenshtein Distance

Levenshtein Distance :

Hypothesis (i)

	0.	1. Take	2. me	3. to	4. you	5. see	6. campus
Reference (j)	0.	0	1	2			
1. Take	1	0	$\min(2,1,3)$				
2. me							
3. to							
4. UC							
5. campus							

$$d[i, j] = \min(d[i-1, j] + 1,$$

$d[i, j-1] + 1,$

`d [i -1 , j -1] + local_substitution (i , j))`

“Grow” the edit matrix iteratively, by accumulating the cost for each element $d[i,j]$

D P !

local(i,j) =

} 1, if $H[i] \neq RE[j]$

{ 0, else

Insertion

Deletion

Substitution

Hypothesis (i)

Reference (j)	0.	1. Take	2. me	3. to	4. you	5. see	6. campus
0.	0	1	2	3	4	5	6
1. Take	1	0	1	2	3	4	5
2. me	2	1	0	1	2	3	4
3. to	3	2	1	0	1	2	3
4. UC	4	3	1	1	1	2	2
5. campus	5	4	2	2	2	2	2

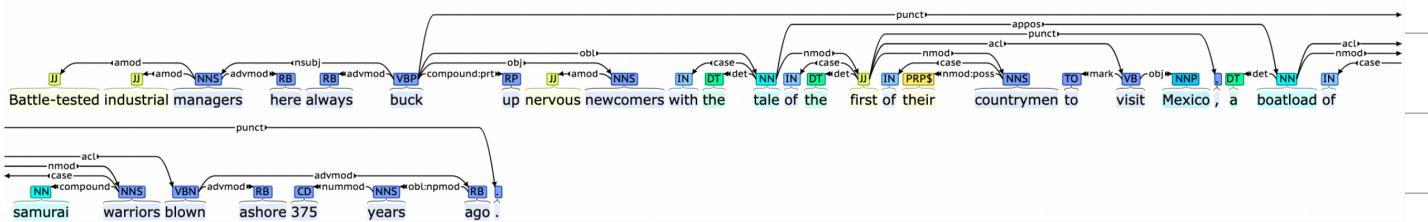
Optimal path:
tells us what
the errors were

you ≠ UC

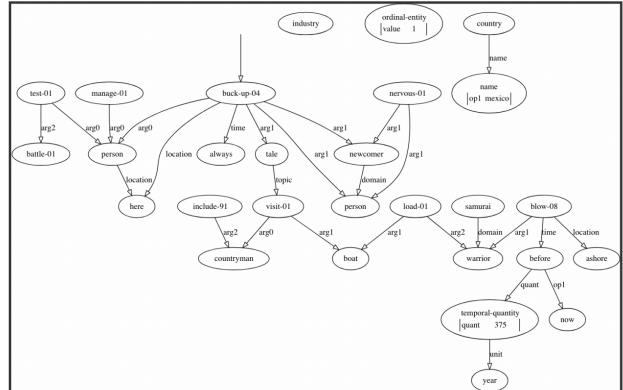
$$\therefore \text{local}(i,j) = 1$$

四

- Task: Structured Prediction: text data → ^{structured combination}_{of labels from a closed set}



E.g. Semantic Parsing (语义分析)
Abstract Meaning Representation)



Metric: For them: Task-specific

E.g. Semantic Parsing : match accuracy
Text SQL : execution accuracy

• Task: Text Comprehension : Text → Short piece of text

Information Extraction , Question Answering - - -

• Task: Reference Based Text Generation : like machine translation

Metric: ① BLEU !

$G_n(x)$, n-gram

$C(s, x)$

- System outputs / candidates \hat{y}
 \hat{y}_1 The most common form of language use is conversation.
 \hat{y}_2 The most common form of language is speech.
(里面无重复字)
- $G_n(x)$ is the set of n-grams in sequence x
 $G_1(y) = \{\text{the, most, natural, form, of, language, use, is, dialogue, .}\}$
 $G_2(y) = \{\text{the most, most natural, natural form, form of, ...}\}$
- $C(s, x)$ is the number of occurrences of n-gram s in x
 $C(\text{the}, y) = 1 \quad C(\text{most}, y) = 1 \quad C(\text{natural form}, y) = 1$

n-gram precision

How many of the n-grams actually occur in the reference?

$$P_n(\hat{y}, y) = \frac{\sum_{s \in G_n(\hat{y})} \min(C(s, \hat{y}), C(s, y))}{\sum_{s \in G_n(\hat{y})} C(s, \hat{y})}$$

y The most natural form of language use is dialogue.

\hat{y}_1 The most common form of language use is conversation.

分子：所有 n-gram 元组，有重复，的数量，在 GT 中出现

分子：对于每一个 n 元组， $\min(C(s, \hat{y}), C(s, y))$

n-gram precision

How many of the n-grams actually occur in the reference?

Calculation

Example

注意取出的 $G_n(x)$

来自于 prediction !

$$P_n(\hat{y}, y) = \frac{\sum_{s \in G_n(\hat{y})} \min(C(s, \hat{y}), C(s, y))}{\sum_{s \in G_n(\hat{y})} C(s, \hat{y})}$$

y The most natural form of language use is dialogue.

\hat{y}_1 The most common form of language use is conversation.

$G_1(\hat{y}_1) = \{\text{the, most, common, form, of, language, use, is, conversation, .}\}$

s	$C(s, \hat{y})$	$C(s, y)$	$\min(C(s, \hat{y}), C(s, y))$
the	1	1	1
most	1	1	1
common	1	0	0
form	1	1	1
of	1	1	1
language	1	1	1
use	1	1	1
is	1	1	1
conversation	1	0	0
.	1	1	1
Sum	10		8

$$P_1(\hat{y}_1, y) = 8 / 10 = 0.8$$

To prevent rather short sentences for BLEU score 'stealing'

\Rightarrow Brevity Penalty

$$BP(\hat{y}, y) = \begin{cases} 1, & |\hat{y}| \geq |y| \\ e^{1 - \frac{|y|}{|\hat{y}|}}, & \text{else} \end{cases}$$

n-gram precision

How many of the n -grams actually occur in the reference?

$$P_n(\hat{y}, y) = \frac{\sum_{s \in G_n(\hat{y})} \min(C(s, \hat{y}), C(s, y))}{\sum_{s \in G_n(\hat{y})} C(s, \hat{y})}$$

y The most natural form of language use is dialogue.
 \hat{y}_3 language

$$\underline{BP(\hat{y}, y)} = \begin{cases} 1 & |\hat{y}| \geq |y| \\ e^{(1 - (|y|/|\hat{y}|))} & |\hat{y}| < |y| \end{cases}$$

In practice: Average n -gram precision, for up to $N = 4$

Final BLEU formula:

$$BP(\hat{y}, y) \exp \left(\sum_{n=1}^N \frac{1}{N} \ln P_n(\hat{y}, y) \right)$$

Average over the log results

$$P_n(\hat{y}, y) = \frac{\sum_{s \in G_n(\hat{y})} \min(C(s, \hat{y}), C(s, y))}{\sum_{s \in G_n(\hat{y})} C(s, \hat{y})}$$

$$\underline{BP(\hat{y}, y)} = \begin{cases} 1 & |\hat{y}| \geq |y| \\ e^{(1 - (|y|/|\hat{y}|))} & |\hat{y}| < |y| \end{cases}$$

$$\text{BLEU} = BP(\hat{y}, y) \exp \left(\sum_{n=1}^N \frac{1}{N} \ln P_n(\hat{y}, y) \right)$$

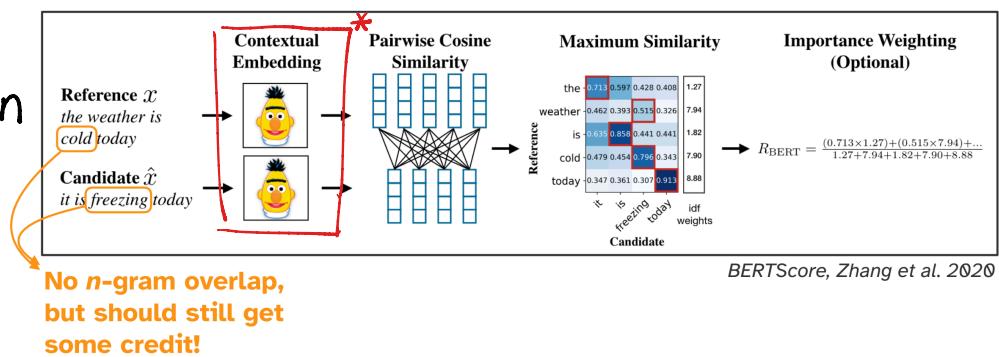
② Representation Similarity: Synonyms should work or suffice, but this can't be reflected on BLEU

Need Semantic or Contextual Representation

\Rightarrow Embedding

Then calculate

similarity -----.



③ Human Eval: Especially for long-text generation.

Δ : There've been study showing that automatic metric can reflect human judgement!

④ LLM-as-a-Judge: Use LLM to score or judge

• Task: Communicative Success: System take action via language

E.g. Input Image & Object, Output referring expression

Easy to eval since they take actions toward given object!

• Task: Dialogue & Interactive Systems

Hard to eval! \Leftrightarrow { Define slot filling task, but reduce expressivity
Use another model as simulator rather than human user

- Building Language Technology — A paradigm:

Task — Implementation (Training included) — Eval (with Data!)

Tasks (in General)



Most tasks can be thought of as a mapping from some input X to some output Y , where one or both of these have to do with language

Task	Input (X)	Output (Y)
Text classification	Text	Label from a fixed class
Automatic speech recognition	Audio signal	Text
Dependency parsing	Text	Dependency tree
Code generation	Text	Executable code
Question answering	Document and question (both text)	Answer (text)
Translation	Text (in source language)	Text (in target language)
Open-ended NLG	Optional prompt (text)	Text
Referring expression generation	Image and target object	Referring expression (text)
Dialogue	Conversation history (text)	Next utterance (text)

How can we implement a classification model?

- Manual rules

```
def classify(x: str) -> str:
    sports_keywords = ["baseball", "soccer", "football", "tennis"]
    if any(keyword in x for keyword in sports_keywords):
        return "sports"
    else:
        return "other"
```

- Prompting a model without training

If the following text is about sports, reply "sports". Otherwise, reply "other".

"Cal football is set to lose its entire starting ..."



- Machine learning

I love to play baseball.
The stock price is going up.
He got a hat-trick yesterday.
He is wearing tennis shoes.

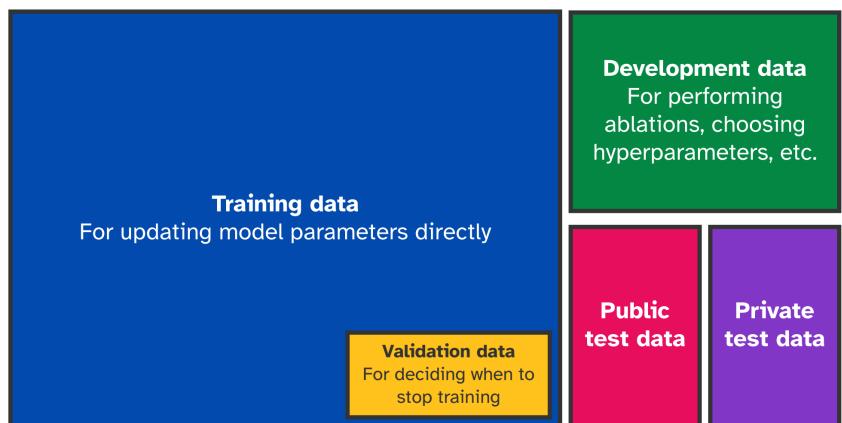
sports
other
sports
other



Data



Before pretrained models, nearly all datasets came with splits, assumed to be IID:



- Benchmark and Model Eval

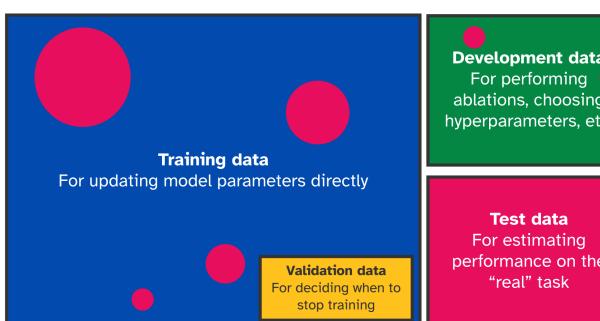


Really important:

- Estimating how well our models will work on real-world data
- Shared understanding of model performance with standardized evaluations
- Building trust within a community in proving how well a new model does
- Driving progress towards specific tasks and capabilities

There can be pitfall as well !

Does the benchmark actually evaluate what we want it to?



Correct 80% of the time without looking at an image!

傍相关

- Learned a spurious correlation that gives the model high accuracy: bananas are yellow
- Doesn't actually test visual understanding

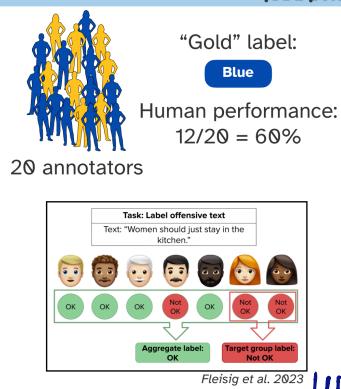
'learn'

Model 'memorize' Test Data!

Defining "Human Performance"

Is the task actually evaluable?

- Is there an objectively correct label?
- What color should we make our Berkeley t-shirts?
- Do annotators disagree? If so, how are we processing their disagreements?



The Long-Tail Paradox

- Which of the following words is most rare? myriad, solipsist, anachronistic, apricate
- Can you think of a rarer word than these? (Check in the [Google Books Ngram Viewer](#))
- LLMs can look really impressive when we are trying to challenge them with tricks, but maybe this is because we're bad at coming up with long-tail (rare) challenging tasks out of the blue!
- Can seem even more impressive when the task is unverifiable, and any model response could satisfy you!

tasks

We focus on few popular tasks, but ignore long-tail

And for long-tail task, we are easy to be content

L3: Speech & Lexical Semantics (词汇语义)

Linguistic Units: Smallest: Phonemes (音位) & Graphemes (字位)

语音中区分意义最小声音单位 文字系统中代表一个音位最小单位

Speech can be represented as acoustic waveform (声波)

compute

⇒ Signal Processing can be helpful

① Analyze short-time window

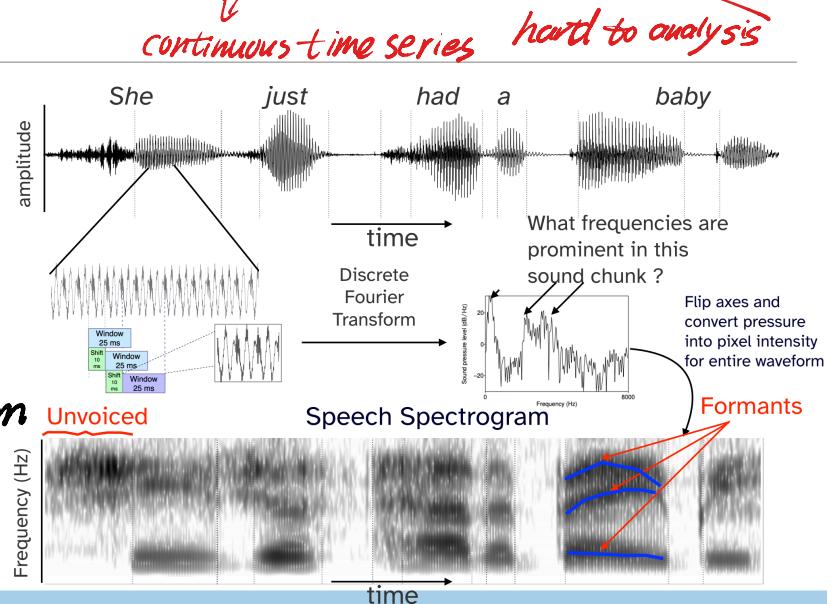
② Run Fourier Transform



Spectrogram

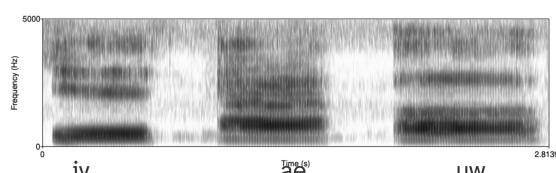
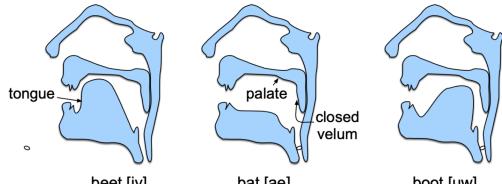
而又有发音对应的 Spectrogram

即可进行 'decode':



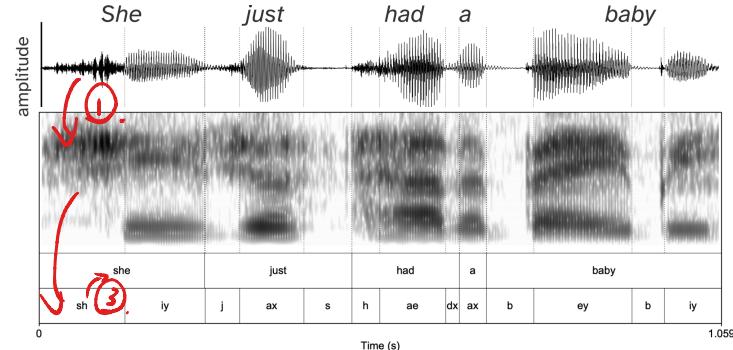
语音学: phonetics

Study of speech sounds — their physical production, spectral and perceptual properties



Articulatory Phonetics

Acoustic Phonetics



- Spectrogram reveals some segmental structure with distinct properties

- These are phonemes — perceptually distinct speech sounds

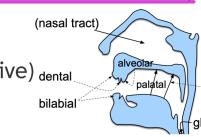
发音

International Phonetic Alphabet (IPA)

- Phoneticians compiled a common set of sounds used to codify different speech sounds (across languages)

IPA: 音 - 行. - 行 - 音

THE INTERNATIONAL PHONETIC ALPHABET (n)										
CONSONANTS (PULMONIC)										
	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Vocalic	Palatal	V-		
Plosive	p b	t d	t d	t d	k					
Nasal	m n	n	n	n	n					
Troll	r	r	r	r	r					
Tap or Flap	t	t	t	t	t					
Fricative	f v	θ ð	s z	ʃ ʒ	ʂ ʐ					
	ɸ β	θ	θ	θ	θ					



- Consonants are characterized by place and manner of articulation
- /p/ is caused by constriction at lips (labial)
- /p/ is caused by sudden release of air (plosive)
- Vowels are characterized by jaw position and tongue shape
- Some vowels also use lips (eg. sound uw in cool)

音位学

Phonology is the study of rules that govern the organization of sounds in a language (Phonemes → Syllables → Words)

- Pronunciation dictionaries (often made by linguists) give the syllables and phonemes within each word in vocabulary

音位 → 音节 → 单词

发音字典 (CMU)

Conclusion: →

- Words are composed of atomic units based on sound (for spoken languages)
- Sounds are a function of how we move our vocal tracts and mouth anatomy
- Languages have distinct sets of possible sounds (phonemic inventory)
- And “rules” governing which sound sequences are likely (syllable structure)

接下来介绍 Morphemes / Lexemes (语素 / 词位 (字典词)) .

- Text data can be viewed as a sequence of words
- First step in building a language technology: building a function that maps from arbitrary text data to that sequence

['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills', '']
tokenized text

- Type-token distinction: (word type).
- Type: a unique word in a text corpus
- Token: an instance of a word type, appearing in a particular context

Example . ↓ .

['., ',', :, :, 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic', 'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form', 'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening', 'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches', 'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young']
wordtypes (vocabulary)

instances (tokens) of wordtype ' , '

['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills', '']
tokenized text

- Simplest tokenizer (for English): splitting on spaces

```
tokenized = s.split(' ')
```

```
[ 'The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language',
  'use', 'is', 'dialogue:', 'Every', 'language', 'user', 'including',
  'young', 'children', 'and', 'illiterate', 'adults', 'hold', 'a',
  'conversation', 'whereas', 'reading', 'writing', 'preparing',
  'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far',
  'from', 'universal', 'skills.]
```

- But this gets us some weird wordtypes:

```
'dialogue:'
'user'
'skills.'
```

Not really words different from
dialogue, user, skills

这又了 token, how to tokenize?

Simplest: whitespace!

→ punctuations mixed,
maybe need rules for that

- nltk tokenizers, with special rules for punctuation

```
import nltk
tokenized = nltk.word_tokenize(s)
```

```
[ 'The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language',
  'use', 'is', 'dialogue:', 'Every', 'language', 'user', '',
  'including', 'young', 'children', 'and', 'illiterate', 'adults', '',
  'can', 'hold', 'a', 'conversation', '', 'whereas', 'reading', '',
  'writing', '', 'preparing', 'speeches', 'and', 'even', 'listening',
  'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills', '.']
```

- But this still loses similarity between wordtypes

```
'skills'
'reading'
'speeches'
```

Lexically similar to, but
morphologically distinct from
skill, read, speech

- Once you've "trained" your tokenizer, you're stuck with it

```
vocab = ['.', ',', ':', 'Every', ..., 'writing', 'young']
vocab.index('ChatGPT') → not found!
```

- Once you've "trained" your tokenizer, you're stuck with it

```
vocab = ['.', ',', ':', 'Every', ..., 'writing', 'young', '<UNK>']
tokenized_indices =
[vocabulary.index(token) for token in tokenized_text]
if token in vocabulary
else vocabulary.index('<UNK>')
```

Not Robust enough!

项目 2: Character Unit → Encode, using Unicode

- Strings are sequences of characters (bytes)!

```
tokenized = s.encode()
```

```
54 68 65 20 6d 6f 73 74 20 6e 61 74 75 72 61 6c 20 61 6e 64 20 62 61 73 69 63 20 66 6f 72 6d 20 6f 66 20
6e 61 6e 67 75 61 69 65 20 75 73 65 20 69 73 20 64 69 61 6e 6f 67 75 65 3a 20 45 76 65 72 79 20 6c 61 65
67 61 6e 67 75 61 69 65 20 75 73 65 20 69 73 20 64 69 61 6e 6f 67 75 65 3a 20 45 76 65 72 79 20 6c 61 65
66 20 6e 64 20 69 62 6c 69 74 65 72 61 74 75 72 60 61 64 66 74 72 20 63 67 6e 20 66 6f 66 64 20 61 6c 61 65
20 63 66 67 76 65 72 73 74 69 66 6c 69 72 63 73 20 77 68 65 72 65 61 73 20 72 65 61 64 69 66 67 72 20 77 72 69
74 69 66 67 2c 20 70 72 65 70 61 72 69 66 20 73 70 65 63 68 65 73 20 61 66 64 20 65 72 65 66 20 6c 61 65
69 73 74 65 66 69 66 20 73 60 70 73 70 65 63 68 65 73 20 61 72 65 20 66 61 72 20 66 72 6f 6d 20 75
6e 69 76 65 72 73 61 6c 20 73 6b 69 6e 6c 73 2e
```

- Now our vocabulary is a fixed size (all possible Unicode characters)



- But: individual characters are not meaningful

```
[., ',', ':', 'Every', ..., 'user', 'whereas', 'writing', 'young']
```

VS.

```
[1, 'a', 'n', 'g', 'u', 'b', 'c', ..., 'a', 'b', 'c', ..., '7', '8', '9']
```

- But: input sequences are much longer

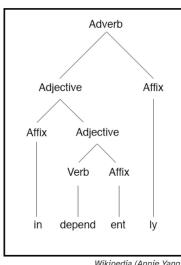
language

VS.

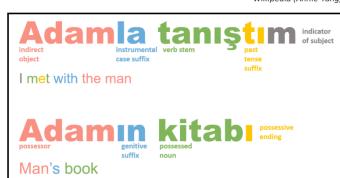
```
'1', 'a', 'n', 'g', 'u', 'b', 'c', ..., 'a', 'g', 'e'
```

Compromise: Character ← word → subparts

How to set the subparts vocab? ⇒ BPE



Wikipedia (Annie Yang)



- Main principle: words are (often) composed of subparts (morphemes)
- Our vocabulary should have entries for frequent words kept whole, because we have a lot of data about those words
- But it should also have entries for parts of less-frequent words, so our ML models can learn how to compose parts of words into whole words (especially unfamiliar words!)

- Modern standard for building a tokenizer

- Inputs: collection of texts and target vocabulary size

- Initial vocabulary is the set of all bytes (characters) across the texts

- Until the target vocabulary size is reached, repeat the following:

- Tokenize all of the texts using the current vocabulary
- Find the most common bigram in the tokenized texts, then add it to the vocabulary as a new wordtype *

每个单词 letter level 看 freq, 看 2-gram freq, 找 freq max, 进 vocab, 之后 updated

vocab tokenize, 找 2-gram, 看 freq -

直至 vocab 中 token 数量 达到 目标 n 为止

* Begin: Every letter is a token △: 1 => 2 tokens

- Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary

tokenized texts

'h' 'u'	15	most frequent bigram; add to vocabulary
'u' 'g'	20	
'p' 'u'	17	
'u' 'n'	16	
'b' 'u'	4	
'g' 's'	5	

- Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary

tokenized texts

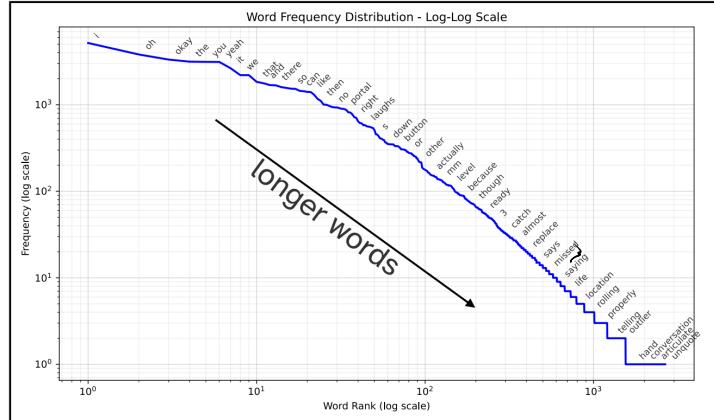
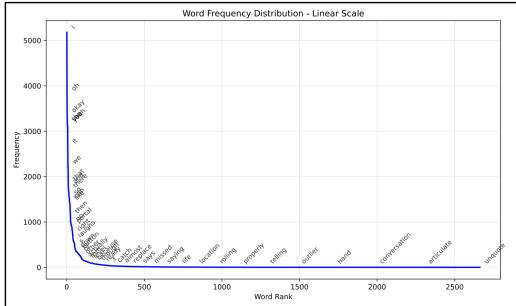
```
('h', 'u', 'g', 'p', 'n', 'b', 's', 't') → ('h' 'u', 'g', 'p', 10), ('p' 'u', 'g', 5), ('p' 'u', 'n', 12), ('b' 'u', 'n', 4), ('h' 'u', 'g', 's', 5)
```

```
('h', 'u', 'g', 'p', 'n', 'b', 's', 't', 'ug') → ('h' 'u', 'g', 'p', 10), ('p' 'u', 'g', 5), ('p' 'u', 'n', 12), ('b' 'u', 'n', 4), ('h' 'u', 'g', 's', 5)
```

```
('h', 'u', 'g', 'p', 'n', 'b', 's', 't', 'ug') → ('h' 'u', 'g', 'p', 10), ('p' 'u', 'g', 5), ('p' 'u', 'n', 12), ('b' 'u', 'n', 4), ('h' 'u', 'g', 's', 5)
```

```
('h', 'u', 'g', 'p', 'n', 'b', 's', 't', 'ug') → ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hug', 's', 5)
```

Long-tail ("Zipfian") distribution of wordtypes (from Portal 2 dialogues)



词汇 freq 基本规律: Zipfian Dist:

只有少数词汇会被 frequently 使用，绝大多数用得非常少

Word Meaning:

Core question: if a machine learning system is processing text, how should words be represented as input to and within the ML system?

一个词或符号的意义，就是它在现实世界或特定语境中所指代的那个具体事物、概念或实体

Hypernym: 上义词

更广泛

Hyponym: 下义词

更狭窄

Synset: 同义词集

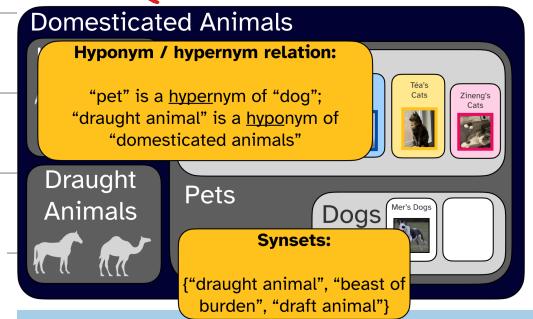


- Missing context-dependence (polysemy)
- Missing meaning of new words
- Requires human labor
- Subjective and culturally dependent

(指称语义学)

Denotational semantics: the symbol *refers to* something in the context in which the language is used

Word Meaning:
本体论



Machine learning models expect numerical inputs
独热编码

- Represent each word as a unique **one-hot vector** (vector with values 0, except for one value of 1)

$$\begin{aligned}\phi(\text{cat}) &= [0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0] \\ \phi(\text{kitten}) &= [0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]\end{aligned}$$

- Vector dimension: the same size as the vocabulary

Problem: no notion of similarity

How can I make my house safe for a new kitten? | How can I make my house safe for a new cat?

- (Is there some way we could use WordNet to get more informative numerical representations of words?)

- Instead: represent words as continuous vectors.

$$\begin{aligned}\phi(\text{cat}) &= [-0.6 \ 1.3 \ -0.1 \ 0.3] \\ \phi(\text{kitten}) &= [-0.7 \ -1.5 \ 0.0 \ 0.3]\end{aligned}$$

Implicitly provides notions of similarity:

$$\|\phi(\text{cat}), \phi(\text{kitten})\| < \|\phi(\text{cat}), \phi(\text{dog})\|$$

- Similarity is learnable from text at scale:

... feeding time in multiple **cat** households can often be ...
... is to relieve the **cat** from the stress of ...
... effect on the feral **cat** population in rural areas ...

... So I wrapped the **kitten** up in a towel ...
... they noticed that the **kitten** population was getting high ...
... were looking for a **kitten**, hoping that someone ...

Meaning \propto Context 中而来



Core principle: **distributional hypothesis**

- Words that are used in similar contexts have similar meanings
- Context: typically, other words in a text, but really anything can be context!

- Input: corpus D including pairs (w, c) where w is a word and c is a context,

e.g.:

- "Everybody likes tesgüino"
 $(w = \text{"tesgüino"}, c = \text{"likes"})$
e.g., context is every word in the sentence
 $(w = \text{"tesgüino"}, c = \text{"everybody"})$

- We want to model: given that we observe word w , what's likely in the context c ?

$$p(c | w; \theta)$$

- Our objective: find parameters that maximize the corpus probability

$$\arg \max_{\theta} \prod_{(w,c) \in D} p(c | w; \theta)$$

- (There are other methods for word2vec (e.g., GloVe, CBOW), but we will only look at Skip-Gram in this class)

Objective: find

$$\arg \max_{\theta} \prod_{(w,c) \in D} p(c | w; \theta)$$

where

$$p(c | w; \theta) = \frac{\exp(\text{score}(c, w))}{\sum_{c' \in C} \exp(\text{score}(c', w))}$$

score: similarity of
we use cosine similarity because, representations of
in contrast to dot product, it is unbiased towards vector

magnitude: $\text{score}(c, w) = \cos(\phi(c), \phi(w)) = \frac{\phi(c) \cdot \phi(w)}{\|\phi(c)\| \|\phi(w)\|}$

Training Paradigm : maximize prob of (c, w) pair, and prob is evaluated via

Objective: find

$$\arg \max_{\theta} \prod_{(w,c) \in D} \frac{\exp(\text{score}(c, w))}{\sum_{c' \in C} \exp(\text{score}(c', w))}$$

Problem: intractable to sum over all contexts

But ~~it~~ is intractable ! Vocab C is huge ! So, instead of 'multiclass classification'
(Don't need prob dist of w over the entire vocab) *Negative Sampling*

we do binary classification \Rightarrow Given (c, w) where they don't form a context and (c, w) from D , tell each sample form context or not,

NEW objective: find

$$\arg \max_{\theta} \prod_{(w,c) \in D} p((w, c) \in D) \prod_{(w,c) \in D'} p((w, c) \notin D)$$

What is D ?

$$p((w, c) \in D) = \frac{1}{1 + \exp(-\text{score}(c, w))}$$

$$p((w, c) \notin D) = 1 - p((w, c) \in D)$$

Binary Classification Prob Representation: Sigmoid

NEW objective: find

$$\arg \max_{\theta} \prod_{(w,c) \in D} p((w, c) \in D) \prod_{(w,c) \in D'} p((w, c) \notin D)$$

\Leftarrow How $c, w, c, w \in D$ is built

- * Negative samples: sample and train on $l * |D|$ pairs (w', c') where

$$w' \sim p(W) \quad c' \sim p(C)$$

unigram prior
over words unigram prior
over contexts

- * How to train?

- * Gradient descent

- * In practice, we work with log probabilities, not direct probabilities, to avoid float underflow

$$a : b :: c : ? \equiv d = \arg \max_{i \in W} \frac{(\phi(b) - \phi(a) + \phi(c)) \cdot \phi(i)}{\|(\phi(b) - \phi(a) + \phi(c)) \cdot \phi(i)\|}$$

cosine the expected and some
similarity of representation of the other word
analogy word i

With embedding, we can even ① vector arithmetic

$D(A \overset{?}{=} B, B \overset{?}{=} C \overset{?}{=})$

② Translation: key Assumption:

Two Language Embedding is the different mapping of the same semantic space

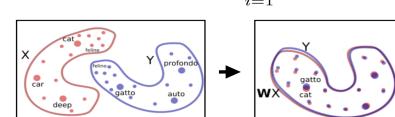
- What if we have embeddings in different languages and want to align them? Given:

- ϕ^A : language A's embedding function

- ϕ^B : language B's embedding function

- $D = \{w_i^A, w_i^B\}_{i=1}^M$: dataset pairing M word translations

- Find a matrix W such that $\min_W \sum_{i=1}^M \|W\phi(w_i^A) - \phi(w_i^B)\|^2$



But still has challenge: like Polysemy (-多义)

Word Embeddings in the Age of LLMs



- If we have representations for parts of sentences (wordtypes), then can we get a representation of the whole sentence?
- One option: bag-of-words representation

$$\phi(\bar{x}) = \frac{1}{|\bar{x}|} \sum_{i=0}^{|\bar{x}|} \phi(x_i)$$

↑

Try to represent a sentence, but miss a lot of things. (Syntax, Semantic)

Word Embedding Today ⇒

- Building a language technology pre-LLMs, for a target task:
 - Download some pre-trained word embeddings from the Internet
 - Initialize your language model with these word embeddings (all other parameters randomized)
 - Use your task-specific data to fine-tune the other parameters (possibly also fine-tuning the word embeddings)
 - Having a good starting point via word embeddings is **really** important, especially with little task-specific data

- But we don't do this anymore
- Language models still learn embeddings specific to each wordtype
 - But we don't download them from the Internet — we download the whole language model Download as a whole

L4. Syntax

Grammaticality:

- Grammaticality judgments are not universal!
- No such thing as "bad grammar" — disagreements over what feels grammatical or not are due to language variation
- Grammatical sentences don't need to be meaningful
- Our main task: we have some language L , and we want to know whether a new sentence $x \in L$
- Should we represent L as a finite-sized set of possible sentences? **Apparently no!**
- What about a regular expression?

Some sentences sound grammatical, some not. 正则表达式

Can regular expression work?

- Can't handle center-embedding

- Lack of Memory

⇒ No!

⇒ Introducing: Context-Free Grammar

The cat that thinks the cow thinks
the rabbits hid are wrong. X

The rabbits hid.

Regex:
DT N VB

The cow thinks the rabbits hid.

DT N VB DT N VB (DT N VB) *

The cat that thinks the cow thinks the rabbits hid is wrong.

DT N IN VB DT N VB DT N VB JJ

DT N IN (DT N VB) * VB JJ?

Is there another formalism that tells us whether
a string is "accepted" in a language or not?

Context-Free Grammars



语法抽象概念或中间变量 ←

词汇，不能再分 ←

* 定义如何用 nonterminal

& terminal to re-write

a non-terminal

△ Can be re-written

as --- !

- Set of nonterminal symbols 非终结符集合
- Set of terminal symbols (wordtypes) 终结符
- Set of production rules defining how nonterminal symbols could be expressed via the composition of other nonterminal and terminal symbols

Nonterminal symbols

DT N VB JJ IN

S NP VP SBAR

Terminal symbols (vocabulary)
are, cat, cow, hid, is,
rabbits, that, the, thinks,
wrong, ...



Production rules

- DT → {the, a, an, ...}
- N → {cat, cow, rabbits, dogs, ...}
- VB → {hid, is, thinks, was, are, ...}
- JJ → {wrong, right, blue, red, ...}
- IN → {that, in, of, because, ...}
- NP → {DT N, NP IN VP}
- VP → {VB, VB JJ, VB S}
- S → {NP VP, VP}

How to use production rule? Arrow:

Right → Left : Parsing / Reduction

Left → Right : Generation

E.g.:

The cat that thinks the cow thinks the rabbits hid is wrong.

The cat that thinks the cow thinks the rabbits hid is wrong.

Nonterminal symbols

DT N VB JJ IN

S NP VP SBAR

Terminal symbols (vocabulary)

are, cat, cow, hid, is,
rabbits, that, the, thinks,
wrong, ...

Production rules

DT → {the, a, an, ...}
N → {cat, cow, rabbits, dogs, ...}
VB → {hid, is, thinks, was, are, ...}
JJ → {wrong, right, blue, red, ...}
IN → {that, in, of, because, ...}
NP → {DT N, NP IN VP}
VP → {VB, VB JJ, VB S}
S → {NP VP, VP}

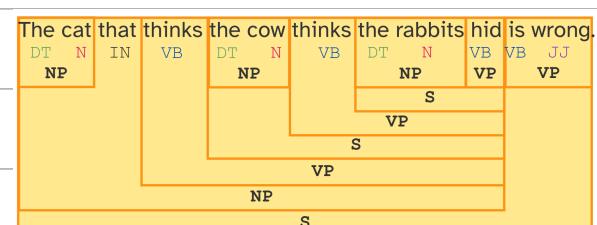
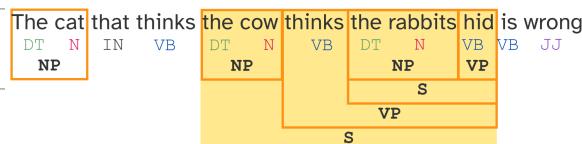
Nonterminal symbols

DT N VB JJ IN

S NP VP

Production rules

NP → {DT N, NP IN VP}
VP → {VB, VB JJ, VB S}
S → {NP VP, VP}



Nonterminal symbols

DT N VB JJ IN

S NP VP

Production rules

NP → {DT N, NP IN VP}
VP → {VB, VB JJ, VB S}
S → {NP VP, VP}

Nonterminal symbols

DT N VB JJ IN

S NP VP

Production rules

NP → {DT N, NP IN VP}
VP → {VB, VB JJ, VB S}
S → {NP VP, VP}

NP is singular, so its corresponding VP should be too

- CFG is often called phrase structure or constituency grammar

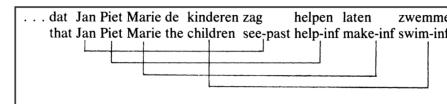
- Each production rule describes a **constituent**

若 CFG 中加入了 'agreement', 如
NP&NPs, VP&VPs, 则不再

Constituent constructions are independent of one another
(this is why the grammar is context-free)

Augmenting a CFG with agreement (e.g., distinguishing plural vs. singular NPs and plural vs. singular VPs plural) means it is no longer context-free

- Also, some languages aren't even context-free, not even considering agreement:



Kaplan et al.

Probabilistic CFG



Production rules

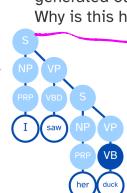
augmented with probabilities

DT → {the, a, an, ...} 0.02 0.02 0.01 0.01
N → {cat, cow, rabbits, dogs, ...} 0.01 0.05 0.05 0.03 0.02
VB → {hid, is, thinks, was, are, ...} 0.02 0.02 0.02 0.02 0.01
JJ → {wrong, right, blue, red, ...} 0.04 0.03 0.02 0.05
IN → {that, in, of, because, ...} 0.01 0.01 0.01 0.01
NP → {DT N, NP IN VP} 0.3 0.2 0.2 0.1
VP → {VB, VB JJ, VB S} 0.3 0.2 0.1 0.1
S → {NP VP, VP} 0.05 0.05 0.05 0.05

p(rule|nonterminal)

learn from data

syntactic parsing: given a PCFG, which rule applications generated our sentence?
Why is this hard?



CFG 中 Rule 规定了叶节点的 rule, 而 PCFG 对句

能 "进行了量化" - 树权概率 = 构成此树的所有

歧义语法分析 产生式规则的概率之乘积

Combinatory Categorial Grammar (CCG)



- Another way of representing a constituency grammar: bottom-up

从它出发

- Elements of a CCG:

• Lexical items (wordtypes) (每个词)

• Each paired with a syntactic type (≈ nonterminal or composition thereof)

the : NP/N dog : N John : NP bit : (S\NP)/NP

If a Noun appears to the right, then it creates a Noun Phrase N → dog Noun Phrase Noun Phrase If a Noun Phrase appears to the right, then it creates an element with the type S\NP

NP → PRO PRO → John If an NP appears to the left of that element, it creates a Sentence

Introduce : CCG :

Syntactic type: 定义词如何与其他词依据

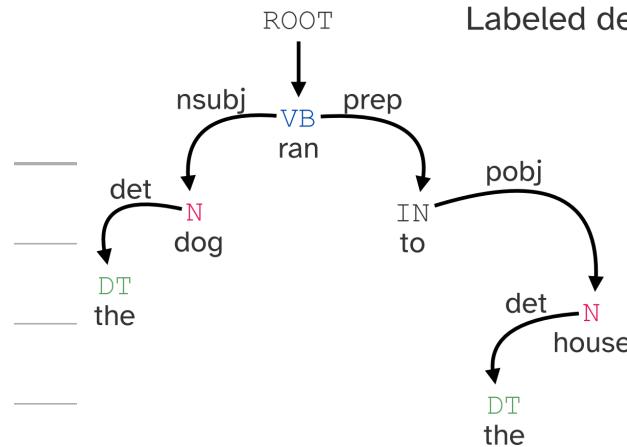
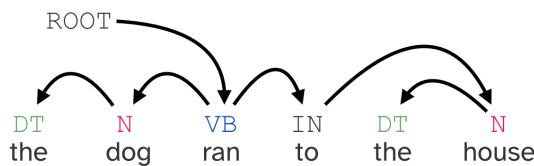
词性进行组合

/ : 需要一个参数 ; (S\NP)/NP : 左右均需NP,

& \" : 最终组成 S

Introduce : Dependency Grammar

- A dependency grammar includes:
 - Terminal symbols (wordtypes) 约束(规则).
 - Parts of speech ← 词性 ↴
 - Some constraints on which parts of speech can be attached to other parts of speech



Labeled dependencies

Why Syntax is Important? ① 模型学

Lec5. Compositional Semantics

Inductive Bias

- Denotational semantics: tokens are references to things in the real world
 - Ontologies: tokens are references to nodes in some knowledge graph
 - Word embeddings: tokens are represented by continuous vectors

Lexical semantics: we can get word meanings

 - Denotational semantics: tokens are references to things in the real world
 - Ontologies: tokens are references to nodes in some knowledge graph
 - Word embeddings: tokens are represented by continuous vectors

Syntax: we can determine what sequences of word types are possible or not possible in a language by modeling latent structure

Main
how do
of the
(a) the
(b) the

ever

Main challenge of semantic parsing:
how do we get a single representation
of the entire sentence's meaning from
(a) the meanings of its words, and
(b) their order and latent structure?

0.5	0.0	0.9
0.3	0.8	0.9
0.2	0.4	0.0
0.6	0.6	0.0
0.4	0.5	0.6
0.7	0.8	0.4
0.3	0.8	0.4

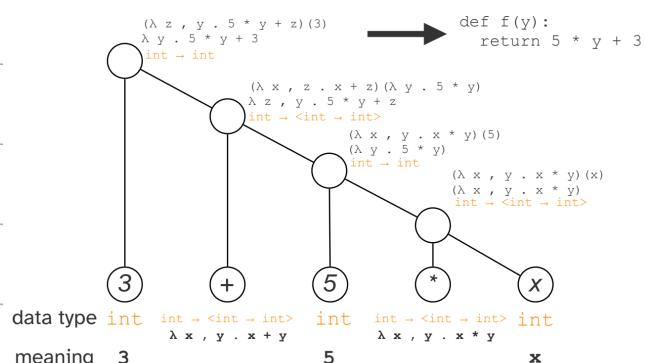
⇒ With word meaning & Syntax,
how to semantic parsing?

Analogy:

CCG and Lambda Calculus

- CFG ⇒ 组合结构

Lambda : "+" $\Rightarrow \lambda x. \lambda y. x + y$



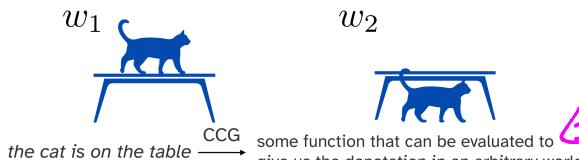
- Truth-Conditional Semantics

 In the context of their use, statements are either true or false, i.e., they have the type `t` (aka, `bool` in python)

- Let's call this context a world w
 - We'd like the outcome of our semantic parsing to be a some that can evaluate to true or false (i.e., $\mathcal{W} \rightarrow [0, 1]$)

→ Goal : Sentence Semantic Parsing → Function. ↗

This J: Input a word w, output whether this sentence is true or not in C.



everyone	likes	Pepper
$s / (s \setminus NP)$ $\langle e \rightarrow t \rangle \rightarrow t$ $\lambda f . \forall x$ (person(x) \rightarrow f(x))	$(s \setminus NP) / NP$ $e \rightarrow \langle e \rightarrow t \rangle$ $\lambda x, y . \text{likes}(y, x)$	NP e Pepper
		$\frac{s \setminus NP}{\lambda y . \text{likes}(y, Pepper)}$

	everyone	likes	Pepper
Syntactic type	$S / (S \setminus NP)$	$(S \setminus NP) / NP$	NP
Semantic type	$\langle e \rightarrow t \rangle \rightarrow t$	$e \rightarrow \langle e \rightarrow t \rangle$	e
λ -expression	$\lambda f . \forall x$ (person(x) \rightarrow f(x))	$\lambda x, y . \text{likes}(y, x)$	Pepper

Final Result \Leftarrow

Example:

CCG

Pepper 语义上为实体

λ -exp \Rightarrow function

everyone	likes	Pepper
$s / (s \setminus NP)$ $\langle e \rightarrow t \rangle \rightarrow t$ $\lambda f . \forall x$ (person(x) \rightarrow f(x))	$(s \setminus NP) / NP$ $e \rightarrow \langle e \rightarrow t \rangle$ $\lambda x, y . \text{likes}(y, x)$	NP e Pepper
		$\frac{s \setminus NP}{\lambda y . \text{likes}(y, Pepper)}$
		$\frac{s}{t}$

$\lambda f . \forall x$ (person(x) \rightarrow f(x)) $\vdash (\lambda y . \text{likes}(y, Pepper))$
 $\forall x$ (person(x) \rightarrow ($\lambda y . \text{likes}(y, Pepper)$))(x)
 $\forall x$ (person(x) \rightarrow likes(x, Pepper))

Modern: //

- In NLP, nobody is really mapping from sentences to lambda calculus representations anymore
- However, many of our problems still take the form of mapping from language to some meaningful structured representation

Formal Semantics

- Logical operators, like \vee , \wedge , and \neg
Pepper is clever and curious
- Quantifiers like \forall and \exists
Some cats like water
- Relationships between functions \Rightarrow and \Leftrightarrow
Squares are rectangles ($\forall x (\text{square}(x) \rightarrow \text{rectangle}(x))$)
- Verbs can have tenses, and can be modified with adverbs
- We can talk about beliefs others have
- Some combinations of meanings are nonsensical (unevaluable) green ideas
- Sentences aren't just statements — sometimes they are commands, questions, etc.
- Sentences exist in the context of previous sentences and their meanings

Lec5: Dialogue, Conversation & Interaction

- Pragmatic: 语用学，研究语言在特定情景(context)中如何被使用。
理解和产生意义

三个核心概念: Speech Act, Presupposition, Implication

语言行为
预设
隐含

Speech Acts

- Our interpretation of utterances used in context often goes beyond the literal (formal) meaning
- By interpreting speech as action, we can ascribe intent to utterances that isn't obvious from their formal representation

Do you mind if I sit next to you?

Yeah (go ahead)

No (I don't mind)

Sorry, someone is coming

What color is Pepper?

\downarrow

$[\lambda f . (\text{Color}(f) \wedge f(\text{Pepper}))]^i$

{black}

Semantics: mapping from surface form (sequence of tokens) to formal executable representation

Pragmatics: executing the logical form against some context to acquire its denotation

当我们将语言诠释为行动时，我们能够推断潜在意图

Presupposition



- Propositions that must be true about a world in order to compute the denotation of a particular sentence
- In other words: implicit assumptions made by utterances

Pepper owns a house Can be true or false

Pepper's house is big Computing the truth value requires that, in our world, there is an entity x such that house (x) and owns (Pepper, x)

Awareness of presupposition in speech is very useful for critical analysis of persuasive speech, e.g. in politics

一句话要成立，必依赖某些未言明的背景假设

Principles behind Implicature: Gricean Maxims



- General principles we believe we mutually hold about how what kinds of utterances we should add to conversation given what's been said so far:
- **Quantity:** utterances should contain just the right amount of information — not too little or too much
- **Truth:** utterances should not contain falsehoods
- **Relation:** utterances should be relevant to what's been said before
- **Manner:** utterance form and meaning should be clear
- What happens when we contribute utterances that break these principles?



Principles behind Implicature: Gricean Maxims



- **Flouting** conversation maxims is "breaking" them under the assumption the listener knows the speaker is intentionally breaking them
 - E.g., flouting relevance:
- **Violating** maxims is breaking them under the assumption the listener won't believe a maxim has been broken

Do you know what they weather will be like today?

You should bring your umbrella.

My dog ate my homework.



上面是四项会话原则 分 打破原则两种方式 . //

下面介绍会话中的共同基础: (Common Grounds)

During interaction, we maintain some representation of what we believe is mutually known by conversation participants

- Mutually known: I know it, I believe you know it, I believe you believe I know it, I believe you believe I believe you know it, ...
- What can be in the common ground?
- Principles guiding how we interact with one another (e.g., Gricean maxims, shared lexicon)
- Propositions about the world, values and beliefs
- Things in our shared environment, including things we are paying attention to

Cooperative Interaction:

Multi-Agent Foundations of Interaction



- You're playing Minecraft, and want to build a house
- Your actions influence the state of the world
- But there's also stochasticity (e.g., animals, enemies, villagers appear and might destroy your house, but their behavior is predictable)
- But what happens when another player comes to the game?
 - They also want to build a house
 - Maybe they will take some of your resources from you
 - Their behavior is not completely stochastic, though!
- Let's assume my partner and I share the same (high-level) goal, and it's in the common ground
 - Goal: build a house
- It's likely optimal for me to choose actions in a way that depends on my partner's actions, to avoid redundancy and execute the goals more efficiently
 - I'll gather wood while my partner places it in the right spots to create a foundation
- Some environments might require that I and my partner take different actions at the same time
 - E.g., pressing paired switches

To best model another player, we may want to keep track of their:

- **Beliefs:** what information are they using to make decisions? How do they perceive the world, and how do they build an internal model of the world as they act in it?
- **Goals:** what do they want to get done? Do they share the same goals, are their goals orthogonal, or are they trying to sabotage mine?
- **Intentions:** how will they attempt to execute their goals? What skills do they have and what strategies are they likely to take?
- **Model of me:** if it's useful for me to reason about them, they are probably reasoning about me, too — how does this influence their actions?
- How can we more successfully coordinate with one another, especially under uncertainty over how the environment works?
- How can I better model my partner?
 - What they observe and know?
 - What they are trying to do?
 - What are their plans to act?
- How can I influence my partner?
 - By sharing information with them?
 - By telling them what to do?
 - By teaching them about how to act?

在 multi-agent 研究中：

- At the beginning of an interaction, we might have significant uncertainty over other agents
- • Over their goals, beliefs, and skills
- • And also over how they use language
- • But over time, we converge to more similar representations
- • By building models of one another from observing their behavior
- • By explicitly resolving uncertainties via language use
- • This refines our expectations of other agents
- • Certainty in our expectations allows us to take communicative shortcuts

• Challenges for learning in multi-agent environments:

- Agents need to learn their first-order policy (goal, observation → action)
- But they also might need to model how their adaptations influence the behavior of other agents, including the introduction of new shared abstractions (i.e., words and conventions)!

• Challenges for evaluating multi-agent systems:

- Dynamics depend heavily on initial conditions: (uncertainty over variation across agent partners' beliefs, goals, and intentions)
- E.g., in teaching contexts, a teacher will adapt their language and pedagogical strategy to the learner's existing knowledge and skills

Simple Multi-Agent Scenario:

- **Environment:** set of candidate referents, available to both players
- **Players**
 - Speaker: knows the identity of a target object
 - Listener: no privileged information
- **Shared goal:** listener picks out the target object
- **Actions:**
 - Speaker: natural language
 - Listener: selection of candidate referent
- **Communication channel:** unidirectional, single utterance (no dialogue)

Observation space now includes utterances made by other agent(s)

Action space now includes the ability to produce an utterance

Dynamic Interaction.

The influence of learning on interaction dynamics:

- Agent policies adapt as they learn about other agents
- This influences the observations other agents make of them
- Which in turn influences how other agents adapt

Element of Scenario

- Interaction dynamics also depend heavily on the properties of the context itself:
 - Incentive structure
 - Environment design — perception and action, novelty
 - Participants — how many, any existing structures among them, roles, a priori asymmetries
 - Communication channel
- Work in computational linguistics, psycholinguistics, and cognitive science aims to characterize the relationship between scenario design and linguistic behavior

Conversation as a multi-agent game.

- No immediate physical environment, but we still exist in a social world, with common (or uncommon) knowledge, goals, and values
- One can work towards social goals by using language as action, e.g. through:
 - Education
 - Persuasion
 - Hate speech and dogwhistles

Lec 6: Multilingual NLP

- For any task we expect out of language technologies, they should work for any language
- • Question answering, information retrieval, summarization
- • Dialogue systems and chatbots
- • Language generation
- • Language technologies can also support cross-language communication
 - Machine translation
 - Language learning

Challenges:

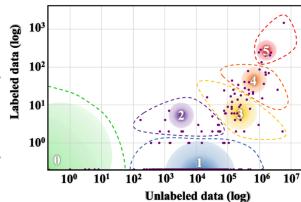
- Not all languages have writing systems, with many languages where:

- Only audio recording is possible or available
- Developing a consistent writing system is difficult, and finding or creating written records is extremely time-consuming relative to how much the language is used

- Some languages are rarely written or have inconsistent uses of writing systems, and only used conversationally

- Some writing systems are not-yet digitized, or all documents are handwritten

② Data Scarcity



- For most languages, very little data is available for training or evaluating language technologies

- There's even less labeled or parallel data!
- 1.2B total speakers, virtually no available data for building language technologies

③ Dialectical Variation

- The same languages can vary significantly depending on who is speaking it where
 - Regional differences
 - Formality differences

- What can vary? Any linguistic features!

- Speakers often mix different dialects with one another in the same conversation (code-mixing)

In some languages, syllables may be distinguished not only by which specific vowel is used, but also by:

- Syllable length
- Pitch contour (spectrogram)
- Pitch height
- Phonation (breathy, creaky, etc.)

Differences between vowels may be more subtle

← ④ Speech System

⑤ Morphology ⇒

⑥ Lexical Semantics

(like how to express

'92'!)

- Synthetic languages denote syntactic relationships between words using inflection (modification of a word, e.g., conjugating a word) or agglutination (adding particles to a word)

- Issues with tokenization

⑦ Syntax : (E.g. English & Chinese. English & Japanese)

⑧ Semantics : Wide variety of possible language feature

⑨ Idioms & Figurative Speech ⑩ Difference in Language use

⑪ Language Change over time.

Lec6: Sequence Modeling

For now: let's assume utterances are sequences of tokens from our vocabulary

Our vocabulary has a fixed size and consists of discrete wordtypes (we'll get to modeling continuous language signals, like speech, in a few weeks!) 词表大小固定, 且离散

A sequence is denoted as:

$$\bar{x} = \langle x_1, \dots, x_n \rangle \quad x \in \mathcal{V} \quad x_n = \text{EOS}$$

We can also consider writing out all possible sequences given our vocabulary (though this set is infinitely large): \mathcal{V}^+

Example: Count-Based Language Model



Documents + frequencies: $\mathcal{D} = \{('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)\}$

$$\mathcal{V} = \{b, g, h, n, p, s, u\}$$

$$\mathcal{V}^+ = \{b, g, h, \dots, bb, bg, \dots, bug, bun, \dots, ssssss, \dots\}$$

Learning problem: we want to estimate the probability distribution $p(\bar{x}) \in \Delta^{\mathcal{V}^+}$ that generated our observations \mathcal{D}

One simple option: just count based on document occurrence

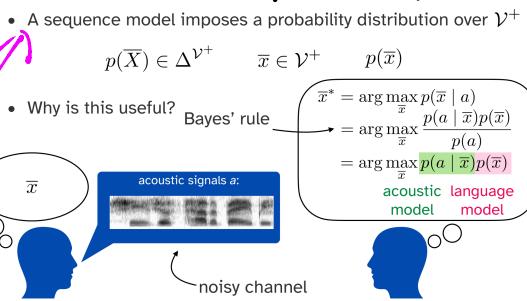
$$p(\bar{x}) = \frac{C(\bar{x})}{|\mathcal{D}|}$$

How to model a sentence, i.e., a token sequence?

← Inductive Bias

建模目标本质:

- 种建模方式,



Count-Based ⇒ frequency

Training Data ⇒ Corpus. But ...

Language Modeling is Hard

$$p(\bar{x}) \in \Delta^{\mathcal{V}^+} \quad \bar{x} \in \mathcal{V}^+ \quad p(\bar{x})$$

- How might we go about assigning a probability to any possible sequence, even ones we've never seen before?
- One intuition: sentences have internal consistencies!

可见概率建模不能以句子为unit.

Autoregressive language modeling:

- The probability of a sequence is a product of local token probabilities
- The probability of a token depends on the ones that came before it

$$p(\bar{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

Introduce: Autoregressive Language Modeling:

Sampling Procedure,

$$p(\bar{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

Sampling from an Autoregressive Language Model



$$p(\bar{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \\ = p(x_1)p(x_2 | x_1) \dots p(x_n | x_1, \dots, x_{n-1})$$

Let's sample a sequence from this approximation:

- Sample the first word $x_1 \sim p(X_1) \in \Delta^{\mathcal{V}}$ $p(X_2 | x_1)$ = the
- Sample the second word $x_2 \sim p(X_2 | x_1) \in \Delta^{\mathcal{V}}$

*途中EOS, 停止!

$$\bar{x} = \langle \text{the, same} \rangle$$

X	p(x)
-ir	0.03
same	0.02
impact	0.02
line	0.02
...	

Shortback:

- ③ Without a big n , cannot handle long-distance dependencies

In sequence modeling, long & short term dependency is significant

N -Gram works, but we dream of:

' θ ' can be MN!

- Let's say we have a language model that can give us a probability of any text $p_\theta(\bar{x})$
- We created this language model using a corpus $\mathcal{D} = \{\bar{x}_i\}_{i=1}^M$
- We care how well this generalizes to some held-out dataset \mathcal{D}'

Lec 7: Generation.

In this section, we explore more possible way to generate sequence in a proper and effective manner.

Sequence sampling method in Lec 6 \Rightarrow

Operation: Temperature.

Operation: modify the logits before computing probabilities



Sneak peek: computing probabilities over wordtypes using pretty much any modern language model

- Score each wordtype independently
 $s(w) = f(w | x_1, \dots, x_{i-1}; \theta)$ ← logits
- Renormalize using softmax
 $p(X_i = w | x_1, \dots, x_{i-1}) = \frac{\exp(s(w))}{\sum_{w' \in \mathcal{V}} \exp(s(w'))}$

\Leftrightarrow Bring in Temperature!

Temperature parameter controls the "smoothness" of this distribution:

$$p(X_i = w | x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

τ Entropy ↑

$\tau \downarrow, \downarrow$; Sentences

are more deterministic

($\tau=0$, argmax)

- Temperature allows us to control the entropy of the output distribution without changing its relative ranking of items

$\tau \uparrow$, entropy ↑

- Higher temperature: closer to a uniform distribution

- Lower temperature: "peakier" distribution (in the limit, gives all probability mass to the most probable item)

$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

No notion of similarity

- What if the count of the target n -gram is 0? \rightarrow smoothing
- What if our $n-1$ -gram prefix has a count of 0? \rightarrow backoff
- Storage is at worst exponential wrt $|\mathcal{V}|$
- Can't learn anything from the counts of n -grams containing similar words

$$p(\text{bike} | \text{I bought a}) \approx p(\text{bicycle} | \text{I purchased a})$$

- Likelihood: probability of the data under our model

$$\prod_{i=1}^M p_\theta(\bar{x}_i)$$

- Negative log likelihood (fixes float underflow)

$$-\sum_{i=1}^M \log p_\theta(\bar{x}_i) = -\sum_{i=1}^M \sum_{j=1}^N \log p_\theta(x_j^i | x_1^i, \dots, x_{j-1}^i)$$

- Perplexity: inverse probability of data, normalized by number of tokens in the dataset

$$= \exp \left(-\frac{1}{\sum_{i=1}^M |\bar{x}_i|} \sum_{i=1}^M \sum_{j=1}^N \log p_\theta(x_j^i | x_1^i, \dots, x_{j-1}^i) \right)$$

$$p(\bar{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

- As we generate, we build our output sequence \bar{x} , which starts as an empty sequence $\bar{x} = \langle \rangle$
- At each step i , we choose an item from the vocabulary \mathcal{V} by performing some operation \downarrow on the local probability distribution $p(X_i | x_1, \dots, x_{i-1}) \in \Delta^{\mathcal{V}}$
- Then, we append this to our running sequence $\bar{x} \leftarrow \bar{x} + \langle x_i \rangle$
- If we ever choose EOS, we stop generation \downarrow 如何利用这条件
- Main point: generation methods differ wrt the operation they perform on $p(X_i | x_1, \dots, x_{i-1})$ \downarrow 哪里决定了生成方法

logits (before softmax) \downarrow
 $s(w) = f(w | x_1, \dots, x_{i-1}; \theta)$ parameter

$$p(X_i = w | x_1, \dots, x_{i-1}) = \frac{\exp(s(w))}{\sum_{w' \in \mathcal{V}} \exp(s(w'))}$$

softmax operation

Softmax-Logit
- Paradigm

- $\tau = 1$: no changes to the probability distribution
- $\tau \rightarrow 0$: relative probability assigned to highest-probability item in distribution increases
- in practice, setting a temperature of 0 recovers "argmax", putting all of the mass on the highest-probability item

$\tau \uparrow \Rightarrow$

uniform distribution

- Temperature allows us to control the entropy of the output distribution without changing its relative ranking of items

$\tau \uparrow$, entropy ↑

- Higher temperature: closer to a uniform distribution

- Lower temperature: "peakier" distribution (in the limit, gives all probability mass to the most probable item)

Suppose we have task:

Finding the Most Probable Sequence

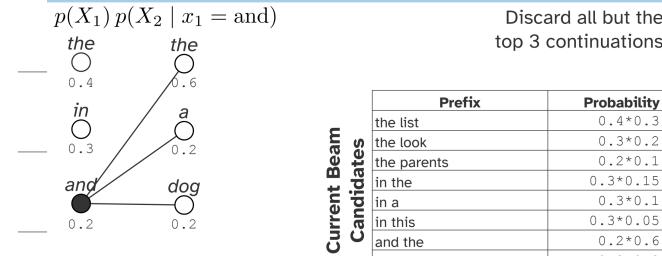
Operation: find $\arg \max_{\bar{x} \in \mathcal{V}^+} p(\bar{x})$

- Why is this hard?

Greedy doesn't guarantee global optima!

- An approximation: greedy "sampling"
 $x_i \leftarrow \arg \max_{x \in \mathcal{V}} p(X_i | x_1, \dots, x_{i-1})$
- Just choose the most probable wordtype at each generation step (no random sampling needed)

⇒ Beam Search to ease this issue! Example :



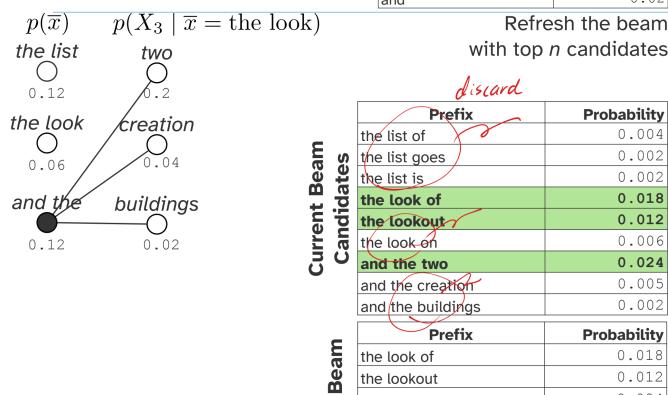
Discard all but the top 3 continuations

Current Beam Candidates

Prefix	Probability
the list	0.4 * 0.3 = 0.12
the look	0.3 * 0.2 = 0.06
the parents	0.2 * 0.1 = 0.02
in the	0.3 * 0.15 = 0.045
in a	0.3 * 0.1 = 0.03
in this	0.3 * 0.05 = 0.015
and the	0.2 * 0.6 = 0.12
and a	0.2 * 0.2 = 0.04
and dog	0.2 * 0.2 = 0.04

Beam

Prefix	Probability
the	0.04
in	0.02
and	0.02



- What if our sequences have different lengths?

Length normalization: $p(\bar{x}) \propto -\frac{1}{|\bar{x}|^\alpha} \sum_{i=1}^{|\bar{x}|} \log p(x_i | x_1, \dots, x_{i-1})$

Argmax can be problematic!

Argmax produces repetitive, less diverse, and overall too-probable output sequences

What's missing?

Beam Search

...to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and...

Human

...which grant increased life span and three years warranty. The Antec HCG series consists of five models with power ratings spanning from 400W to 900W. Now, we should note that we have already tested the HCG-620 in a previous review and were quite satisfied with its performance. In today's review we will rigorously test the Antec HCG-520, which as its model number implies, has 520W capacity and contrary to Antec's strong beliefs in multi-rail PSUs is equipped...

But: length may matter! Item < 1 !

⇒ Length may matter!

words!

* Argmax "only" chose those highly freq

⇒ Change picking strategy!

Operation: top- p (nucleus) sampling

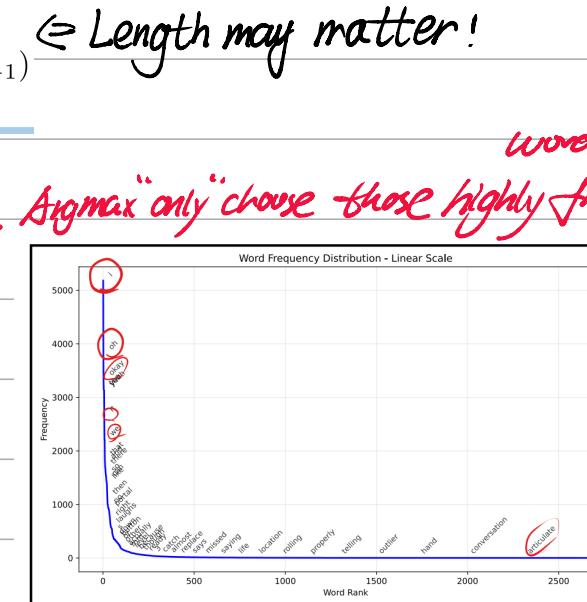
- Identify the set of n tokens \mathcal{P} that have the highest probabilities under $p(X_i | x_1, \dots, x_{i-1})$ and their cumulative probability is p
- Set the probabilities of all but these tokens to 0
- Renormalize by dividing remaining probabilities by $\sum_{x \in \mathcal{P}} p(X_i | x_1, \dots, x_{i-1}) = p$



Operation: choose $x_i \leftarrow \arg \max_{x \in \mathcal{V}} p(X_i | x_1, \dots, x_{i-1})$

- Why isn't this guaranteed to get us the highest-probability sequence?
- A better approximation for global argmax: **beam search**
 - During generation, we maintain a "beam" of n sequences instead of just one
 - At each generation step i ,
 - We select the n most likely next tokens X_i for each prefix, and create n more sequences
 - Then we look at all the n^2 sequences so far, and discard all but the n most likely sequences
 - At the end, we select the sequence that has the highest probability among the set

Prefix	Probability
the list	0.12
the look	0.06
the parents	0.02
in the	0.05
in a	0.03
in this	0.02
and the	0.12
and a	0.04
and dog	0.04



Operation: constrained decoding

Operation: e sampling

- Identify the set of n tokens \mathcal{E} such that $\forall x \in \mathcal{E}, p(x | x_1, \dots, x_{i-1}) \geq \epsilon$ 只考虑概率高于 ϵ (threshold) 的 token.
- Set the probabilities of all but these tokens to 0 and then sample
- Renormalize by dividing remaining probabilities by $\sum_{x \in \mathcal{E}} p(X_i | x_1, \dots, x_{i-1})$



- For some tasks, we have additional information about what wordtypes can or cannot be next

- E.g., in code generation, I can't generate more than I have (对于任意输出空间上的需求，人工取出)

- While modern LLMs can learn these patterns from data at scale, it can sometimes still be useful to constrain our output space $C \subseteq \mathcal{V}$, 并用 $\sum_{x \in C} p(X_i | x_1, \dots, x_{i-1})$ 进行 renorm.

- Similar to before: given a set of possible continuations $C \subseteq \mathcal{V}$ we will set the probabilities of all other tokens to 0, then renormalize using $\sum_{x \in C} p(X_i | x_1, \dots, x_{i-1})$

Lec 8 : Neural Sequence Modeling

Previously, we've mentioned the $p(\bar{x} | x_1, \dots, x_{i-1}) = \prod_{i=1}^{|x|} p(x_i | x_1, \dots, x_{i-1})$ we dream of:

$$p(\bar{x}) = \prod_{i=1}^{|x|} p(x_i | x_1, \dots, x_{i-1})$$

Autoregressive approximation

← Formulation

$$\approx \prod_{i=1}^{|x|} p(x_i | x_{i-n+1}, \dots, x_{i-1})$$

N-gram approximation

← Care about previous $N-1$ words

$$\approx \prod_{i=1}^{|x|} \frac{C(x_{i-n+1}, \dots, x_{i-1}, x_i)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

Count-based approximation

← Use n & $(n-1)$ gram to calculate

$$\approx \prod_{i=1}^{|x|} p(x_i | x_{i-n+1}, \dots, x_i; \theta)$$

Neural n-grams

← what we dream of.

$$p(x_i | x_{i-n+1}, \dots, x_i; \theta)$$

$$= \frac{\exp(s(x_i | x_{i-n+1}, \dots, x_{i-1}; \theta))}{\sum_{x' \in \mathcal{V}} \exp(s(x' | x_{i-n+1}, \dots, x_{i-1}; \theta))}$$

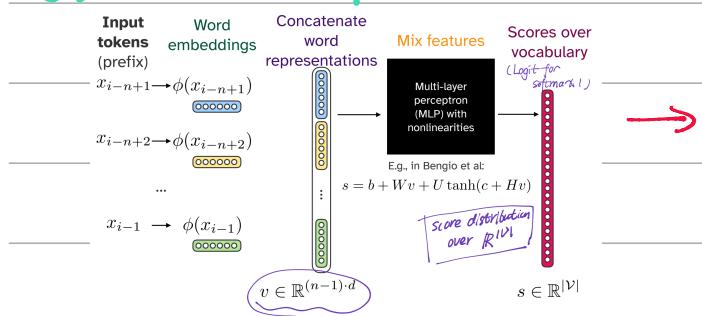
Softmax over wordtype scores

$$s(x | x_{i-n+1}, \dots, x_{i-1}; \theta) = f_\theta(\phi(x_{i-n+1}), \dots, \phi(x_{i-1}))[x]$$

Neural model takes $n-1$ inputs: embeddings of each prefix wordtype

Utilize 'softmax-logit paradigm' & Embedding to bring in semantic information

If can be MLP, then our model can be:



→ softmax → sample.

$$\phi \in \mathbb{R}^{|\mathcal{V}| \times d}$$

$$H \in \mathbb{R}^{d(n-1) \times h}$$

$$c \in \mathbb{R}^h$$

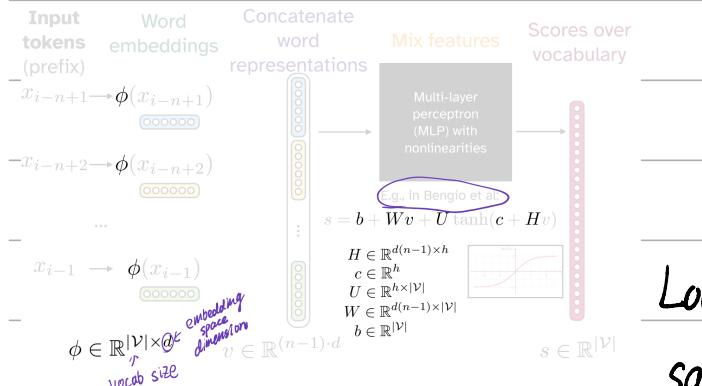
$$U \in \mathbb{R}^{h \times |\mathcal{V}|}$$

$$W \in \mathbb{R}^{d(n-1) \times |\mathcal{V}|}$$

$$b \in \mathbb{R}^{|\mathcal{V}|}$$

$$C_n \in \mathbb{N}^{|\mathcal{V}|^n}$$

$$C_{n-1} \in \mathbb{N}^{|\mathcal{V}|^{n-1}}$$



Neural N-grams

$$n = 6$$

$$|\mathcal{V}| = 20,000$$

$$d = 100$$

$$h = 60$$

→ ~13M parameters

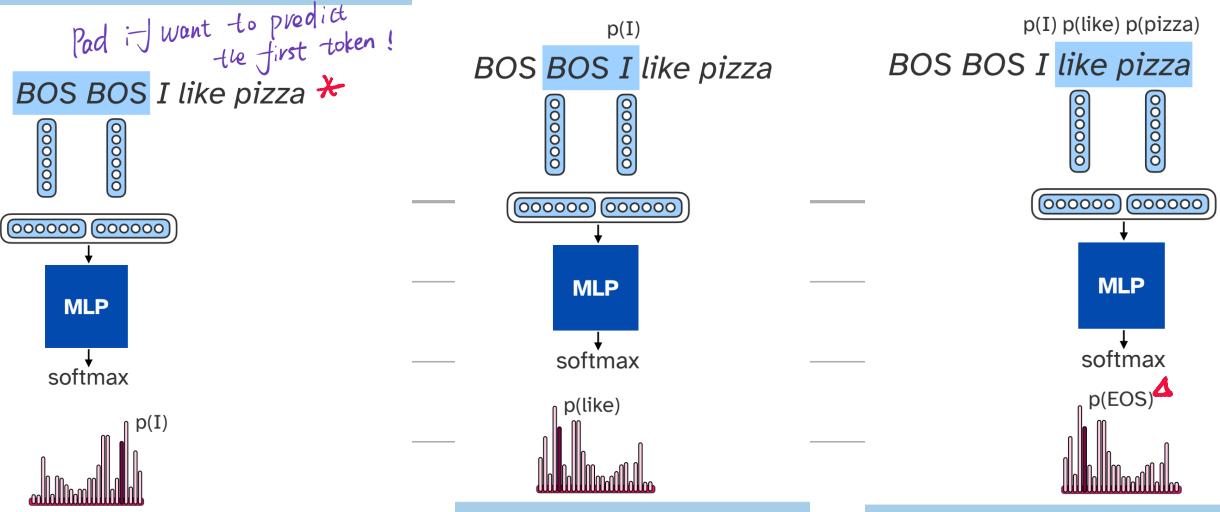
save parameter!

Count-based N-grams

$$6.4 \times 10^{25} \dots$$

Looks nice. Moreover, it even can save storage !

(learnable parameters)



Example : $n=3$, for generating one word , look at previous (3-1) words

* Pad for generating first (3-1) words ! ⚡ STOP !

$$\exp\left(-\frac{1}{3}(\log p(I) + \log p(\text{like}) + \log p(\text{pizza}) + p(\text{EOS}))\right)$$

↑ 4

Generated Sentence's Perplexity

Training Stuff :

- We have some training dataset $\mathcal{D} = \{\bar{x}\}_{i=1}^M$
- We want to find the parameters of our neural model θ that maximize the likelihood of this dataset
- Train using gradient descent in log-probability space

$$\theta^* = \arg \min_{\theta} \left(- \sum_{\bar{x} \in \mathcal{D}} \sum_{i=1}^{|\bar{x}|} \log p(x_i | x_{i-n+1}, \dots, x_{i-1}; \theta) \right)$$

$\phi \in \mathbb{R}^{|\mathcal{V}| \times d}$
 $H \in \mathbb{R}^{d(n-1) \times h}$
 $c \in \mathbb{R}^h$
 $U \in \mathbb{R}^{h \times |\mathcal{V}|}$
 $W \in \mathbb{R}^{d(n-1) \times |\mathcal{V}|}$
 $b \in \mathbb{R}^{|\mathcal{V}|}$

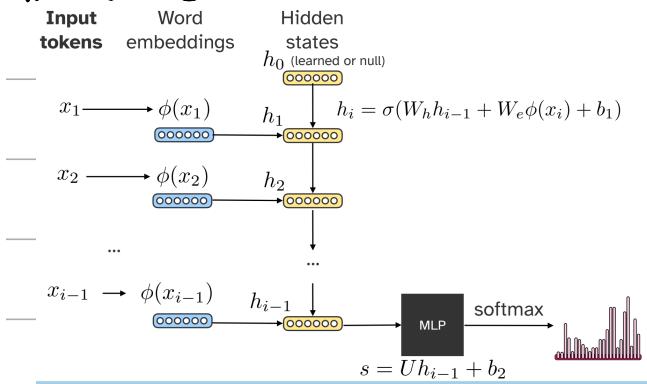
- ✓ What if the count of the target n -gram is 0?
- ✓ What if our $n-1$ -gram prefix has a count of 0?
- ✓ Storage is at worst exponential wrt $|\mathcal{V}|$
- ✓ No notion of similarity
- ✓ Intervening words Can learn about it
- ✗ Without a big n , cannot handle long-distance dependencies still need $n-1$ token prefix !!

- N-gram models (count-based or neural) impose fixed windows of context
- Could we instead truly compute $p(x_i | x_1, \dots, x_{i-1})$ for any i ?

Many previous issues are solved. but n-gram's Markov Property still violates long-dist dependency.

⇒ To tackle this, introduce : Recurrent Neural Network:

- Key idea:** add a latent space whose values are both:
 - determined by computations performed during previous "timesteps", and
 - taken in as input at each timestep
- The value of this latent space is called a "hidden state" ✘



shared.

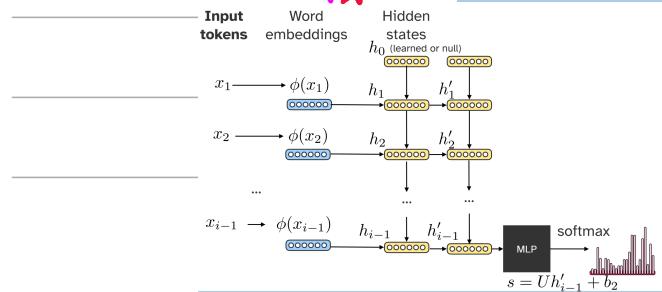
$$h_i = \sigma(W_h h_{i-1} + W_e \phi(x_i) + b_i)$$

↓

6: sigmoid

↖ Single layer

⇐ Multi-Layer



$$h_i = \sigma(W_h h_{i-1} + W_e \phi(x_i) + b_1)$$

$$s = U h'_{i-1} + b_2$$

$$\phi(\cdot) \in \mathbb{R}^{|\mathcal{V}| \times d}$$

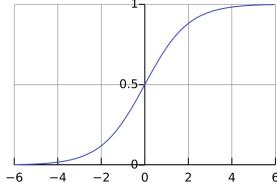
$$W_e \in \mathbb{R}^{d \times d'}$$

$$b_1 \in \mathbb{R}^{d'}$$

$$W_h \in \mathbb{R}^{d' \times d'}$$

$$U \in \mathbb{R}^{d' \times |\mathcal{V}|}$$

$$b_2 \in \mathbb{R}^{|\mathcal{V}|}$$



$$\phi(\cdot) \in \mathbb{R}^{|\mathcal{V}| \times d}$$

$$H \in \mathbb{R}^{d(n-1) \times h}$$

$$c \in \mathbb{R}^h$$

$$U \in \mathbb{R}^{h \times |\mathcal{V}|}$$

$$W \in \mathbb{R}^{d(n-1) \times |\mathcal{V}|}$$

$$b \in \mathbb{R}^{|\mathcal{V}|}$$

$$\phi(\cdot) \in \mathbb{R}^{|\mathcal{V}| \times d}$$

$$W_e \in \mathbb{R}^{d \times d'}$$

$$b_1 \in \mathbb{R}^{d'}$$

$$W_h \in \mathbb{R}^{d' \times d'}$$

$$U \in \mathbb{R}^{d' \times |\mathcal{V}|}$$

$$b_2 \in \mathbb{R}^{|\mathcal{V}|}$$

Good!

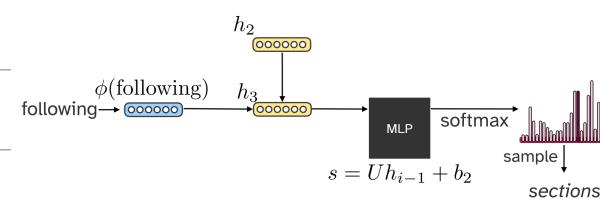
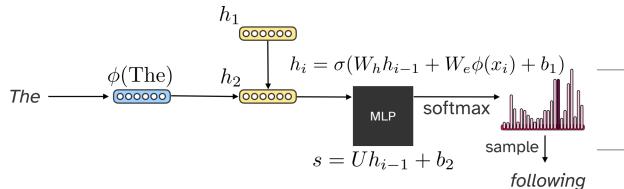
Neural N-grams

$$\begin{aligned} n &= 6 \\ |\mathcal{V}| &= 20,000 \\ d &= 100 \\ h &= 60 \end{aligned}$$

$$\rightarrow \sim 13M \text{ parameters}$$

$$|\mathcal{V}| = 20,000 \rightarrow \sim 2.1M \text{ parameters}$$

Plus, we can (theoretically)
keep information from any
point in the past!



- We have some training dataset $\mathcal{D} = \{\bar{x}\}_{i=1}^M$
- We want to find the parameters of our neural model θ that maximize the likelihood of this dataset
- Train using gradient descent in log-probability space

$$\theta^* = \arg \min_{\theta} \left(- \sum_{\bar{x} \in \mathcal{D}} \sum_{i=1}^{|\bar{x}|} \log p(x_i | x_1, \dots, x_{i-1}; \theta) \right)$$

这就是 loss!

During training, we always compute the probability of each observed token by conditioning on its observed prefix. This is also called "teacher forcing".

Exploding gradients: if the activations are too big, the gradients become too big, and the values of the weights become too big, eventually getting values closer to ∞

- Can solve using gradient clipping (if the norm of the gradient is greater than a threshold, clamp its value or scale it down)

Vanishing Gradients

$$\frac{\delta \mathcal{L}(x_i | x_1, \dots, x_{i-1})}{\delta h_j} = \frac{\delta \mathcal{L}(x_i | x_1, \dots, x_{i-1})}{\delta h_i} \prod_{j < t \leq i} \text{diag}(\sigma'(f(h_{i-1}))) W_h$$

$$\prod_{j < t \leq i} W_h \prod_{j < t \leq i} \text{diag}(\sigma'(f(h_{i-1})))$$

$$W_h^{i-j} \prod_{j < t \leq i} \text{diag}(\sigma'(f(h_{i-1})))$$

$$\begin{aligned} \text{BOS} &\rightarrow \phi(\text{BOS}) \rightarrow h_0 \\ \text{I} &\rightarrow \phi(\text{I}) \rightarrow h_1 \\ \text{like} &\rightarrow \phi(\text{like}) \rightarrow h_2 \\ \text{pizza} &\rightarrow \phi(\text{pizza}) \rightarrow h_3 \\ &\quad \vdots \\ &\rightarrow \text{MLP} \rightarrow \text{softmax} \rightarrow p(\text{EOS}) \end{aligned}$$

$$\mathcal{L}(x_n | x_1, \dots, x_{n-1}) = -\log p(x_n | x_1, \dots, x_{n-1}; \theta)$$

Problem:

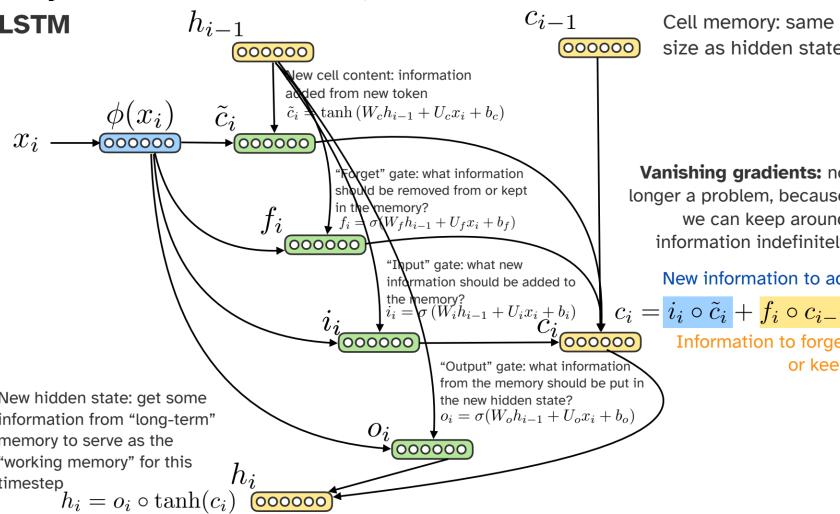
Vanishing
Gradients



- Vanishing gradients:** gradient signal from many timesteps in the future is negligible compared to gradient signal from nearby tokens *

⇒ Introduce Advanced Model:

LSTM



Vanishing gradients: no longer a problem, because we can keep around information indefinitely

New information to add

$$i_i = \sigma(W_i h_{i-1} + U_i x_i + b_i)$$

Information to forget or keep

$$f_i = \sigma(W_f h_{i-1} + U_f x_i + b_f)$$

"Input" gate: what new information should be added to the memory?

$$o_i = \sigma(W_o h_{i-1} + U_o x_i + b_o)$$

"Output" gate: what information from the memory should be put in the new hidden state?

$$c_i = i_i \circ \tilde{a}_i + f_i \circ c_{i-1}$$

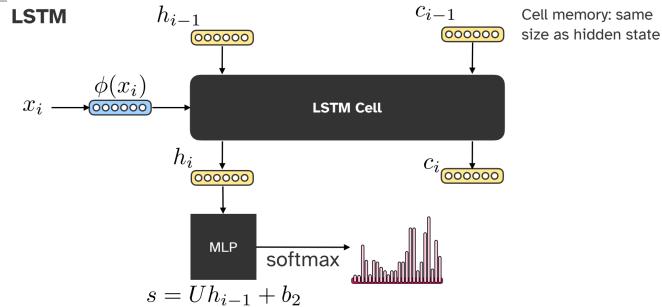
"Forget" gate: what information should be removed from or kept in the memory?

$$\tilde{a}_i = \sigma(W_c h_{i-1} + U_c x_i + b_c)$$

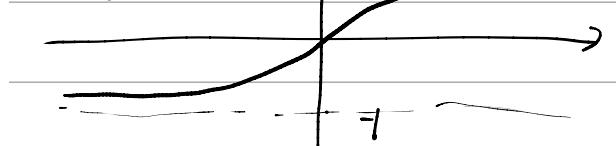
New hidden state: get some information from "long-term" memory to serve as the "working memory" for this timestep

$$h_i = o_i \circ \tanh(c_i)$$

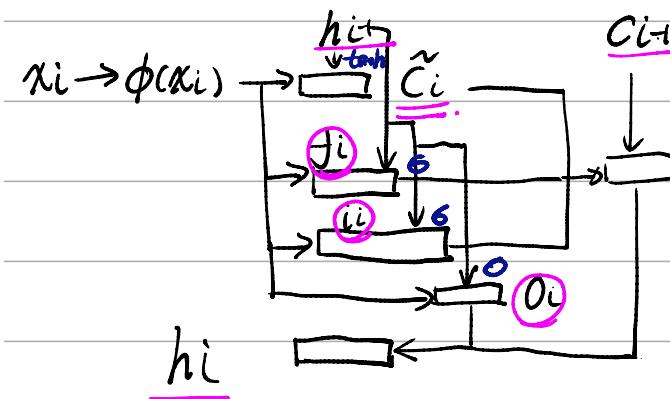
LSTM



tanh



tanh allow negative number



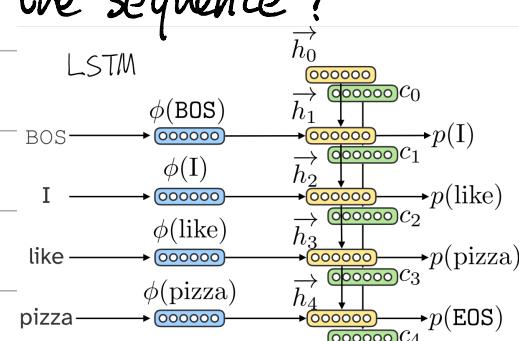
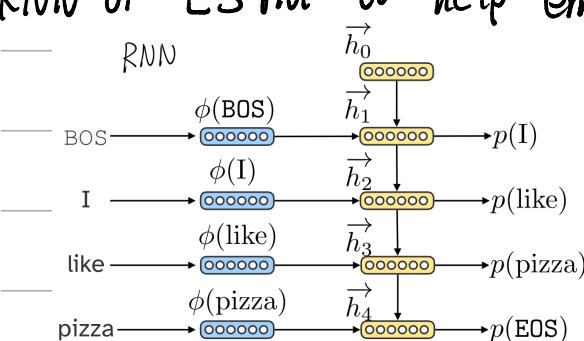
$$c_i = i_i \circ \tilde{a}_i + f_i \circ c_{i-1}$$

i_i f_i O_i are scalars

$$h_i = o_i \circ \tanh(c_i)$$

Lec 8: Sequence Embedding

We have word embedding. What about sentence? Can we use RNN or LSTM to help embed the sequence?



$$\begin{aligned}\phi(\bar{x}) &= \phi(I \text{ like pizza}) \\ &= \frac{1}{|\bar{x}|} \sum_{i=1}^{|\bar{x}|} \phi(x_i)\end{aligned}$$

Option 0: average word embeddings

But we lose all notions of order

$$\begin{aligned}\phi(\bar{x}) &= \phi(I \text{ like pizza}) \\ &= h_{-1}\end{aligned}$$

Option 1: take the last hidden state

- But this is overly weighted by information from later tokens

$$\begin{aligned}\phi(\bar{x}) &= \phi(I \text{ like pizza}) \\ &= \frac{1}{|\bar{x}|} \sum_{i=1}^{|\bar{x}|} h_i\end{aligned}$$

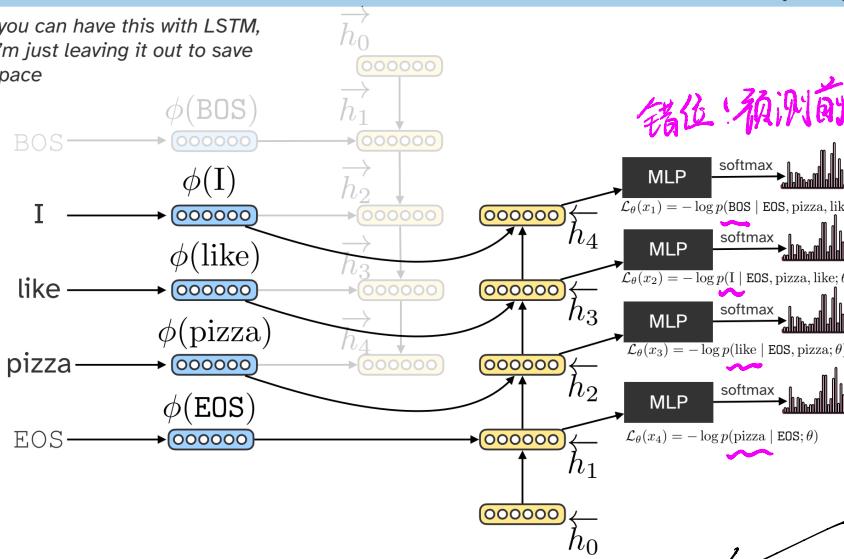
Option 2: average hidden states over time

But earlier states have no information about later ones, while later ones have full information of the sequence

* To tackle this => Bidirectional RNN

Bidirectional RNN

*you can have this with LSTM, I'm just leaving it out to save space



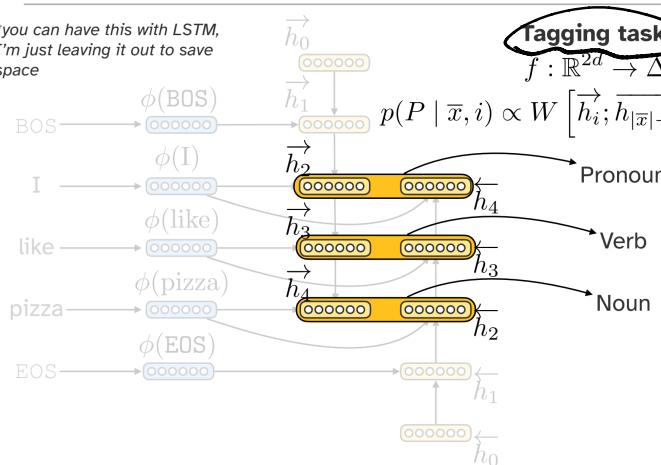
$$\begin{aligned}\phi(\bar{x}) &= \phi(I \text{ like pizza}) \\ &= \frac{1}{|\bar{x}|} \sum_{i=1}^{|\bar{x}|} [\vec{h}_i; \vec{h}_i]\end{aligned}$$

① 加入 BOS & EOS

② 正向做完一遍，反向做一遍

③ \vec{h}_i 拼接，并 normalize.

*you can have this with LSTM, I'm just leaving it out to save space



- Let's say we have some embedding function (aka encoder)

$$Enc : \mathcal{V}^+ \rightarrow \mathbb{R}^n$$

Note: $Enc \equiv \phi$ is the notation we'll use for encoding sequences

- And we want to learn to classify text

E.g.: spam vs. not spam

- We want to learn the optimal parameters of some function

$$f_\theta : \mathbb{R}^n \rightarrow \Delta \{\text{spam, not spam}\}$$

such that for some labeled dataset $\mathcal{D} = \{\bar{x}_i, y_i\}_{i=1}^M$

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^M f(y_i | Enc(\bar{x}_i); \theta)$$

- Our classifier might just look like a linear transformation or MLP, e.g.,

$$f_\theta(\bar{x}) = \sigma(W Enc(\bar{x}) + b)$$

$$\theta = \{W, b\}$$

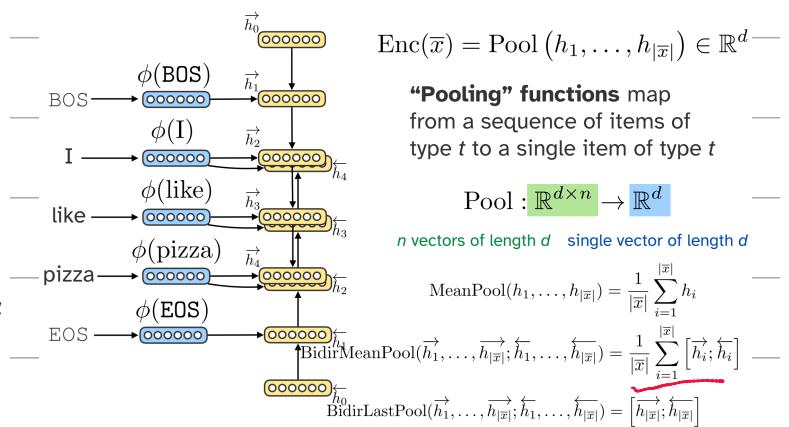
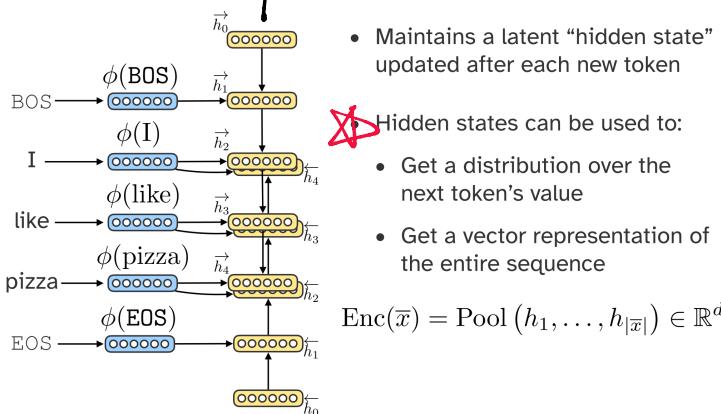
And we can train this classifier by fixing the parameters of the Encoder

- Or, we could also backpropagate into the encoder itself

$$\theta = \{W, b\} \cup \theta_{RNN}$$

Lec 9. Sequence to Sequence Modeling

Recap:



补充一个新 Sequence Embedding 的应用

Natural Language Inference

$$f : \mathcal{V}^+ \times \mathcal{V}^+ \rightarrow \Delta^{\{\text{entailment, contradiction, neutral}\}}$$

Output: 哪种关系

vector concatenation

$$f(\bar{x}_p, \bar{x}_h) = \text{softmax}(W[\text{Enc}(\bar{x}_p); \text{Enc}(\bar{x}_h)] + b)$$

text encoding text encoding
of premise of hypothesis

Train using labeled data

$$\bar{x}_p, \bar{x}_h, y \in \{\text{entailment, contradiction, neutral}\}$$

Can be used to measure quality of sentence representations!

Important!

Want to implement a new task!

- Task: map from text in one language (L1) to a distribution over texts in another language (L2), assigning high probability to texts that preserves the input text meaning

$$f : \mathcal{V}_{\text{L1}}^+ \rightarrow \Delta^{\mathcal{V}_{\text{L2}}^+}$$

$\bar{x}_{\text{English}} = \text{Oula A. Alrifai is a Syrian emigrant to the United States and writer for various Washington-based think tanks.}$

↓ Google Translate

$\bar{x}_{\text{Malay}} = \text{Oula A. Alrifai ialah seorang pendatang Syria ke Amerika Syarikat dan penulis untuk pelbagai badan pemikir yang berpangkalan di Washington.}$

End-to-end, we can think about the translation problem as generating an output sequence token-by-token by

[mapping from some input word to the target vocabulary, using the target language syntax.]

- What can we do with sequence embeddings?

- Text classification

$$f_\theta(\bar{x}) = \sigma(W\text{Enc}(\bar{x}) + b) \quad f : \mathcal{V}^+ \rightarrow \mathcal{C}$$

- Multi-sentence classification, e.g., natural language inference

$$f : \mathcal{V}^+ \times \mathcal{V}^+ \rightarrow \Delta^{\{\text{entailment, contradiction, neutral}\}}$$

$$f(\bar{x}_p, \bar{x}_h) = \text{softmax}(W[\text{Enc}(\bar{x}_p); \text{Enc}(\bar{x}_h)] + b)$$

- Retrieval

$$\arg \max_{\bar{x}' \in \mathcal{D}} \text{sim}(\text{Enc}(\bar{x}), \text{Enc}(\bar{x}'))$$

Conditional sequence generation!

- Aka "sequence transduction"

- E.g., machine translation, response generation in dialogue

- What do we get from sequence models? Recall.

- Autoregressive (token-by-token) generation

- Sentence embeddings $\text{Enc}(\bar{x})$

↓ How to translate

- Key idea: encode some input into a vector, decode it using autoregressive generation

- Aka: sequence-to-sequence, encoder-decoder, sequence transduction, conditional language model...

We can use RNN to encode & decode:

← ↗ [Core idea!]

- Encode input sentence: $\text{Enc}(\bar{x}_{\text{L1}})$
- Generate output sequence token-by-token, conditioning on previously-generated tokens and the input sentence embedding:

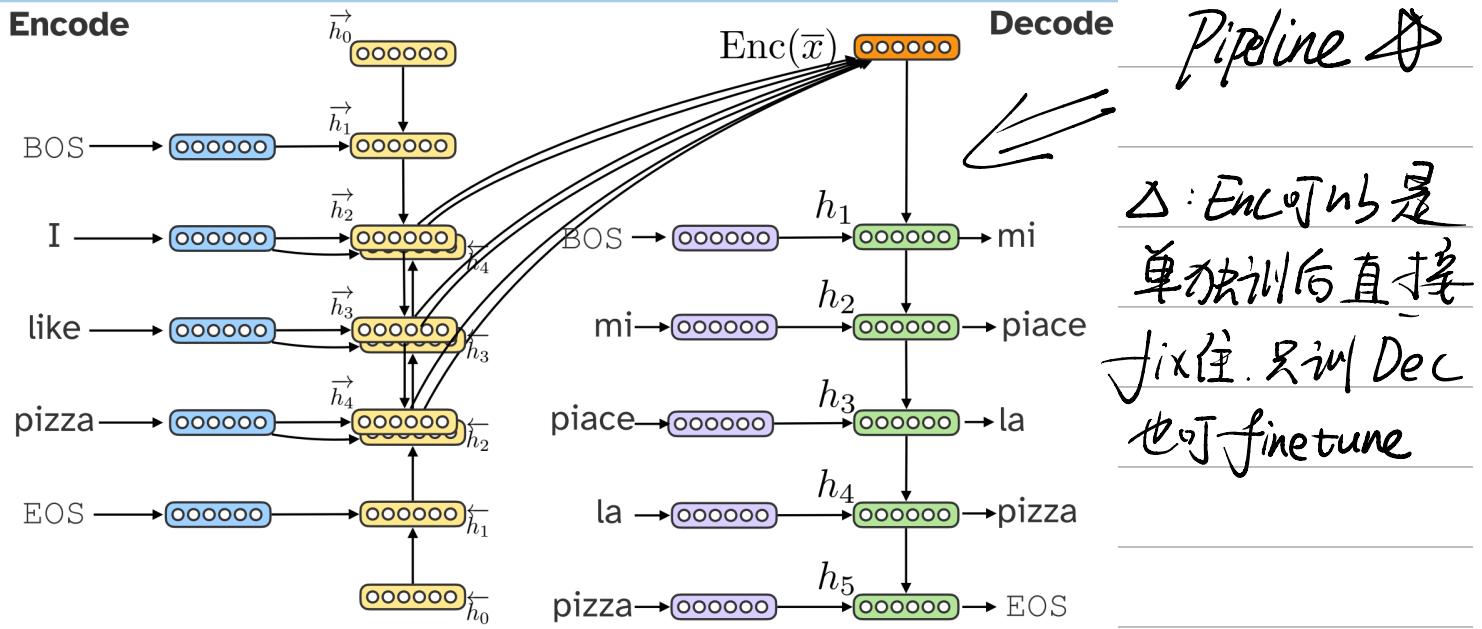
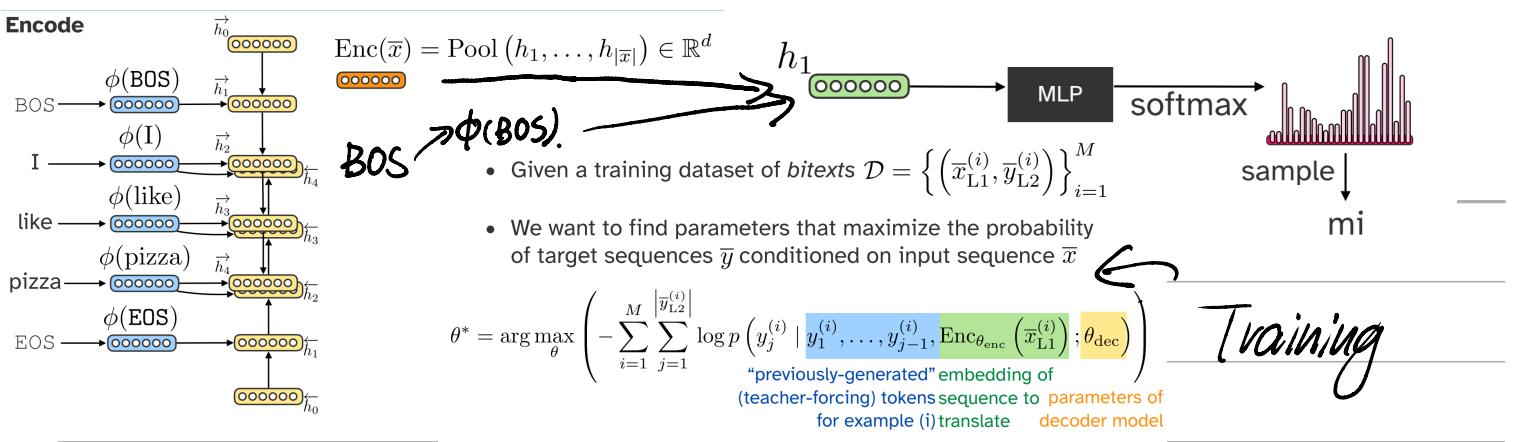
$$\bar{y}_{\text{L2}} = \langle y_1, \dots, y_n \rangle \quad y_n = \text{EOS}$$

parameters of decoder model (e.g., RNN parameters)

$$y_i \sim p(Y_i | y_1, \dots, y_{i-1}, \text{Enc}(\bar{x}_{\text{L1}}); \theta_{\text{dec}})$$

previously-generated tokens embedding of (empty when starting translation) sequence typically, we use y to denote translate

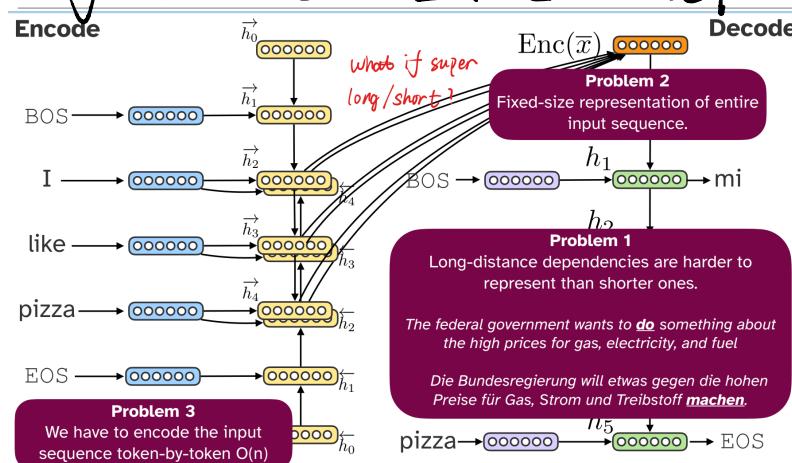
output tokens (you can condition on other things, too! e.g., images)



Options for backprop: use pre-trained θ_{enc} and freeze, or finetune by backprop directly $\theta = \theta_{\text{enc}} \cup \theta_{\text{dec}}$

$$p(y | x)$$

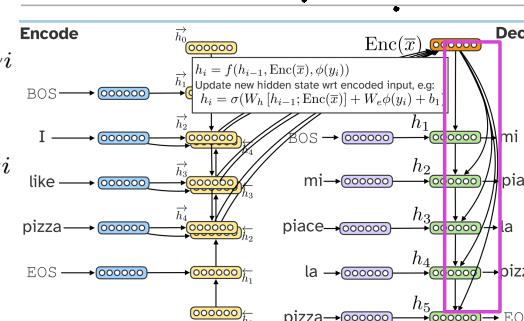
- Machine translation
- Summarization (long text \rightarrow short text)
- Dialogue (previous utterances \rightarrow next utterance)
- Syntactic parsing (input sentence \rightarrow serialized parse)
- Code generation (natural language \rightarrow Python code)



但并非 perfect:
 ⇐ Prob 1: 有些 translation 翻译, 其 dependency 经常很长 (RNN毛病)

Prob 2: Pooling ($h_1, \dots, h_{|\bar{x}|}$)
 是否单一? ; Prob 3: RNN毛病

Prob2 MeanPool($h_1, \dots, h_{|\bar{x}|}$) = $\frac{1}{|\bar{x}|} \sum_{i=1}^{|\bar{x}|} h_i$
 $= \sum_{i=1}^{|\bar{x}|} \frac{1}{|\bar{x}|} h_i$
 equally weight contribution of each input token / index



Prob1 简单解决:
 Enc(\bar{x}) 经每个 decode
 过程都补上,

$$\text{Enc}(\bar{x}) = \sum_{i=1}^{|\bar{x}|} p(i) h_i$$

$p(I) \in \Delta^{\mathbb{N}_{1:|\bar{x}|}}$

- Generalization of pooling: assign a **probability distribution** over input indices, and compute **weighted sum** over hidden states **need: ①**
- What if this probability distribution could change **②** during generation?

With respect to prob₂, can we average with weight? And the weight can be in form of prob dist?

⇒ Attention. Coming Soon

Lec 10 : Attention

$$\text{Enc}(\bar{x}, g) = \sum_{i=1}^{|\bar{x}|} p(i | g) h_i$$

$g \in \mathbb{R}^d \quad p : \mathbb{R}^d \rightarrow \Delta^{\mathbb{N}_{1:|\bar{x}|}}$

Our need: ② ①

- Formulation We want weight: $p(i | g)$.

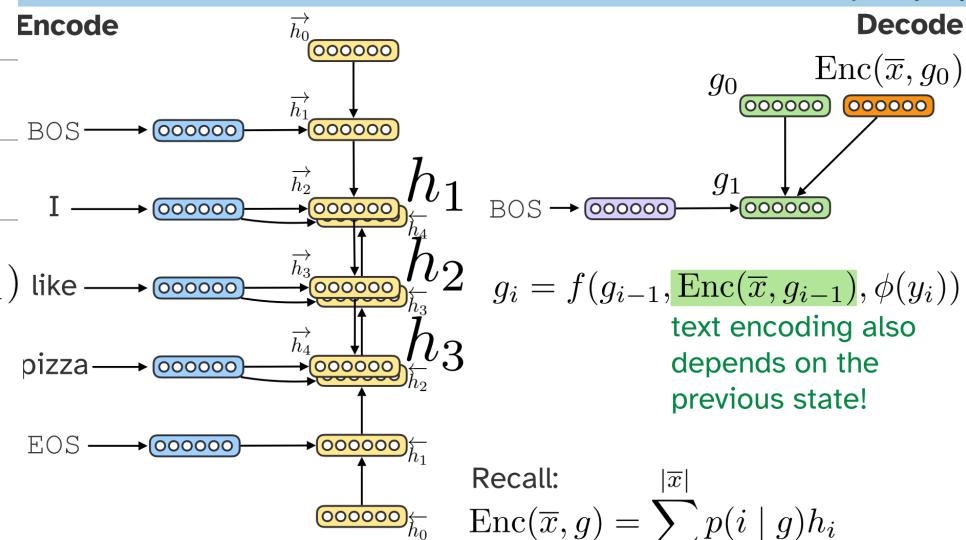
* (of the previously mentioned problems)

$g \in \mathbb{R}^d$

$$p : \mathbb{R}^d \rightarrow \Delta^{\mathbb{N}_{1:|\bar{x}|}}$$

g can be modeled as hidden state:

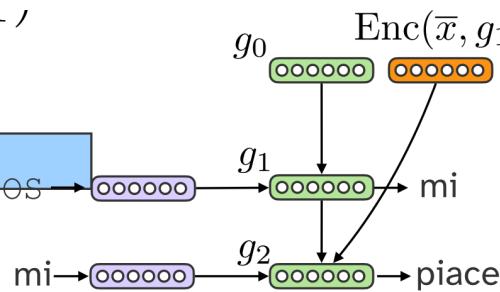
Encode



Two questions:

- Where does g come from?
- How to implement p ?

1



- Update decoding hidden state based on an updated version of the input sequence encoding

$$g_i = f(g_{i-1}, \text{Enc}(\bar{x}, g_{i-1}), \phi(y_i))$$

- The updated sequence is determined via a weighted sum over input hidden states

$$\text{Enc}(\bar{x}, g) = \sum_{i=1}^{|\bar{x}|} p(i | g) h_i \quad p : \mathbb{R}^d \rightarrow \Delta^{\mathbb{N}_{1:|\bar{x}|}}$$

- Next time: how do we compute $p(I | g)$?

- One solution: Attention

Query: what are some words for mammals in Turkish?

i	Key	Value
1	sheep	koyun
2	river	ırımkar
3	author	yazar
4	bull	boğa
5	fly	sinek
6	radio	radyo
7	farmer	çiftçi
8	plant	bitki
9	actress	aktris
10	home	ev
11	snail	salyangoz

$$p(i | q, k_i) \propto \text{sim}(\text{Enc}(q), \phi(k_i))$$

Weight of each item is proportional to the similarity of its key with the query

hidden state

$$- p(i | q, k_i) v_i \text{ value}$$

- Weight each value accordingly

- In our RNN Scenario : $p(i | q, k_i) \propto \text{sim}(\text{Enc}(q), \phi(k_i)) \equiv p(i | g, h_i) \propto \text{sim}(g, h_i)$

- **Query:** current decoder hidden state g

$$\text{Query Key } |\bar{x}|$$

$$\text{Enc}(\bar{x}, g) = \sum_{i=1}^{|\bar{x}|} p(i | g) h_i$$

- **Key:** encoder hidden state (at some index i) h_i

- **Value:** also the encoder hidden state at index i

$$p(i | g, h_i) \propto \text{sim}(g, h_i)$$

Query Key

$$\text{Enc}(\bar{x}, g) = \sum_{i=1}^{|\bar{x}|} p(i | g) h_i$$

Value

There are multiple ways to eval $\text{sim}(q, k_i)$...

$$p(i | q, k_i) \propto \text{sim}(q, k_i) \quad q \in \mathbb{R}^d \quad \text{sim}(q, \mathbf{k}) = q^\top W \mathbf{k} \quad \text{Bilinear attention}$$

$$= \frac{e^{\text{sim}(q, k_i)}}{\sum_{i=1}^{|\bar{x}|} e^{\text{sim}(q, k_i)}} \quad \mathbf{k} \in \mathbb{R}^{d' \times |\bar{x}|} \quad W \in \mathbb{R}^{d \times d'}$$

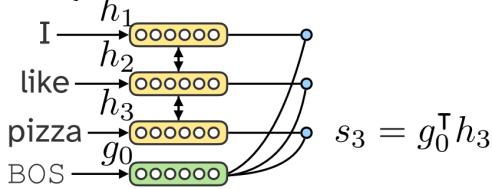
$$\text{sim}(q, \mathbf{k}) = w_2^\top \tanh(W_1[q; \mathbf{k}]) \quad \text{MLP attention}$$

$$W_1 \in \mathbb{R}^{(d+d') \times d'} \quad w_2 \in \mathbb{R}^{d''}$$

$$\text{sim}(q, \mathbf{k}) = q^\top \mathbf{k} \quad \text{Dot product attention}$$

$$d = d' \quad \text{sim}(q, \mathbf{k}) = \frac{q^\top \mathbf{k}}{\sqrt{d'}} \quad \text{Scaled dot product attention}$$

Pipeline !!



α is calculated via $s \Leftarrow$

$$p(Y_1 | \bar{x}) = \text{softmax}(f(c_0, g_0))$$

usage
sample!

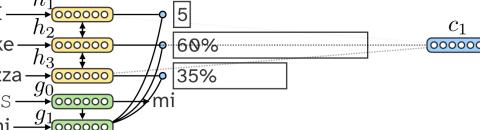
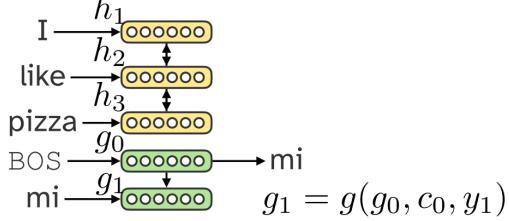
Similarities between query and keys
 $s_0 = g_0^\top \mathbf{h} \in \mathbb{R}^{|\bar{x}|}$

Weights over keys (weights sum to 1)
 $\alpha_0 = p(\cdot | g_0, \mathbf{h}) = \text{softmax}(s_0) \in \Delta^{N_1:|\bar{x}|}$

Encoding is weighted sum over values

$$c_0 = \text{Enc}(\bar{x}, g_0) = \sum_{i=1}^{|\bar{x}|} \alpha_{0,i} h_i \in \mathbb{R}^d$$

I				
like				
pizza				



Similarities between query and keys
 $s_1 = g_1^\top \mathbf{h} \in \mathbb{R}^{|\bar{x}|}$

Weights over keys (weights sum to 1)
 $\alpha_1 = p(\cdot | g_1, \mathbf{h}) = \text{softmax}(s_1) \in \Delta^{N_1:|\bar{x}|}$

Encoding is weighted sum over values

$$c_1 = \text{Enc}(\bar{x}, g_1) = \sum_{i=1}^{|\bar{x}|} \alpha_{1,i} h_i \in \mathbb{R}^d$$

Updating and
repeat

	mi			
I				
like				

Overview:

- We encoded our input sequence into hidden states: $h_1, \dots, h_{|\bar{x}|} = \mathbf{h} \in \mathbb{R}^{d \times |\bar{x}|}$
- Now we want to predict the next word y_{i+1}
- We have access to the previous decoder hidden state g_i
- First, compute attention scores for each input hidden state using similarity between query (g_i) and keys (\mathbf{h}): $s_i = a(g_i, \mathbf{h}) \in \mathbb{R}^{|\bar{x}|}$
- Then, take softmax of attention scores to get a distribution over keys:
 $\alpha_i = \text{softmax}(s_i) \in \Delta^{N_1:|\bar{x}|}$
- Finally, compute a weighted sum of values (\mathbf{h}) using this distribution
 $c_i = \sum_{j=1}^{|\bar{x}|} \alpha_{i,j} h_j \in \mathbb{R}^d$
- Use the weighted sum to predict the next word:
 $p(Y_{i+1} | \bar{x}, y_1, \dots, y_i) = \text{softmax}(f(c_i, g_i))$
- Use the weighted sum to update the decoder hidden state: $g_i = g(g_{i-1}, c_{i-1}, y_i)$

Complexity

$$h_1, \dots, h_{|\bar{x}|} = \mathbf{h} \in \mathbb{R}^{d \times |\bar{x}|}$$

$$O(n) s_i = a(g_i, \mathbf{h}) \in \mathbb{R}^{|\bar{x}|}$$

$$\alpha_i = \text{softmax}(s_i) \in \Delta^{N_1:|\bar{x}|}$$

$$c_i = \sum_{j=1}^{|\bar{x}|} \alpha_{i,j} h_j \in \mathbb{R}^d$$

$$p(Y_{i+1} | \bar{x}, y_1, \dots, y_i) = \text{softmax}(f(c_i, g_i)) \quad O(m)$$

$$g_i = g(g_{i-1}, c_{i-1}, y_i)$$

O(n m)

m: dimension of vocab

n: input length

- Main problem: fixed-size representation ("bottleneck") of input sequence at each stage of decoding**

Attention fixes this by allowing access to different parts of the input sequence based on where we are in decoding

- While decoding token-by-token, we can "focus" on different input parts

- Helps with vanishing gradient problem via shortcut to far-away states

- Provides some kind of interpretability ("soft" alignment) learned without any supervision

Attention was initially developed for sequence-to-sequence problems

→ Self-Attention

But we can also use it during single-sequence generation to solve analogous problems:

- Problem 1: long-distance dependencies
- Problem 2: fixed-size representation of sequence so far

$$h_{j+1} = g(h_j, c_j, x_{j+1}) \quad h_{j+1} = \phi(x_{j+1})$$

How attention solved this problem.

And other benefits.

- Keys, values, and queries are transformations of the same underlying representations

- For each token generated at index i , we have:

- A key: $k_i = Kh_i \in \mathbb{R}^d$ $K \in \mathbb{R}^{d \times d}$

- A value: $v_i = Vh_i \in \mathbb{R}^d$ $V \in \mathbb{R}^{d \times d}$

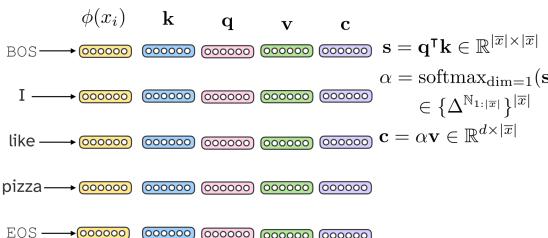
- When trying to generate our next token at index j , we need a query: $q_j = Qh_j \in \mathbb{R}^d$ $Q \in \mathbb{R}^{d \times d}$

- Scores over keys: $s_{ji} = q_j^T k_i$ $\alpha_{ji} = \frac{\exp(s_{ji})}{\sum_{j'=1}^J \exp(s_{ji'})}$

- Attention distribution over keys: $\alpha_{ji} = \frac{\exp(s_{ji})}{\sum_{j'=1}^J \exp(s_{ji'})}$

- Weighted sum of values: $c_j = \sum_{i=1}^I \alpha_{ji} v_i$

↓ Crucial for supporting parallel computation



We can get context-sensitive representations of each input token in parallel!

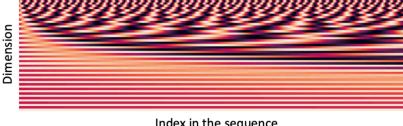
Two problems →

Prob1:

- Currently, our model is operating only on embedding of individual wordtypes $\phi(x_i) = \phi(x)$
- Let's also embed their indices: $\phi(x_i) = \phi(x) + \phi(i)$

- Sinusoidal embeddings

$$\mathbf{p}_i = \begin{pmatrix} \sin(i/1000^{2+1/d}) \\ \cos(i/1000^{2+1/d}) \\ \vdots \\ \sin(i/1000^{2+\frac{d}{2}/d}) \\ \cos(i/1000^{2+\frac{d}{2}/d}) \end{pmatrix}$$



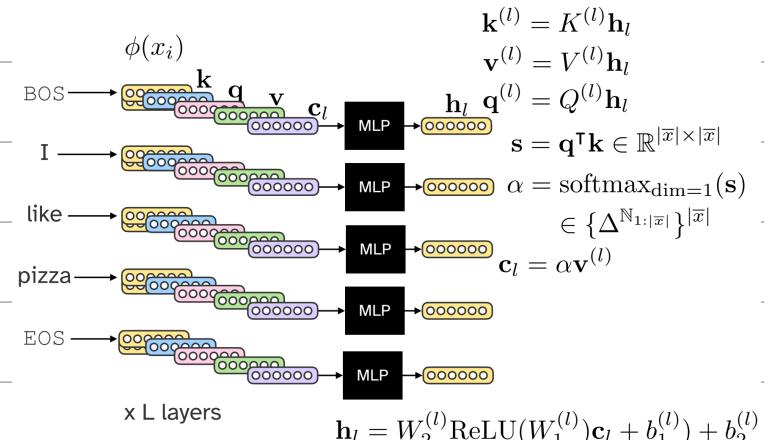
- (in practice, most methods use learned embeddings)

Prob2: $h_l = \text{MLP}(C_l)$

Calculate $q_k v$: First block

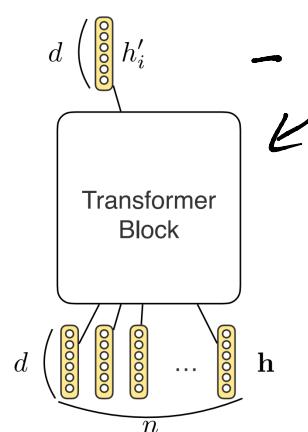
use embedding $\phi(\alpha_i)$; Following

blocks use h_{l-1} from the last block



Lec 11: Transformer

- A transformer block takes as input a sequence of input vectors $\mathbf{h} \in \mathbb{R}^{d \times n}$
- It generates a sequence of output vectors $\mathbf{h}' \in \mathbb{R}^{d \times n}$
- For now, let's focus on how it computes the output vector corresponding to the input at index i , $h'_i \in \mathbb{R}^d$



- Encoder

Compute attention over input vectors

$$\mathbf{k} = K\mathbf{h} \in \mathbb{R}^{d_k \times n}$$

$$q_i = Qh_i \in \mathbb{R}^{d_k}$$

$$s_i = \frac{q_i^T \mathbf{k}}{\sqrt{d_k}} \in \mathbb{R}^n$$

$$\alpha_i = \text{softmax}(s_i) \in \Delta^{N_1:n}$$

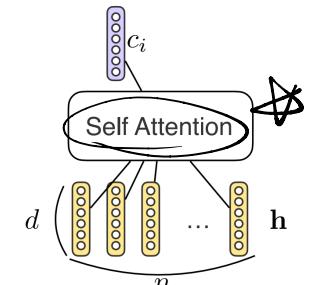
$$\mathbf{v} = V\mathbf{h} \in \mathbb{R}^{d_k \times n}$$

$$c_i = \sum_{i'=1}^n \alpha_i v_{i'}$$

$$c_i = \text{SelfAttention}(\mathbf{h})_i$$

(single head.)

$$\theta_{\text{SelfAttention}} = \{K, Q, V\}$$



MultiHead Attention

Compute attention over input vectors

$$\mathbf{k}^{(j)} = K_j \mathbf{h} \in \mathbb{R}^{d_k \times n}$$

$$q_i^{(j)} = Q_j h_i \in \mathbb{R}^{d_k}$$

$$s_i^{(j)} = \frac{q_i^{(j)T} \mathbf{k}^{(j)}}{\sqrt{d_k}} \in \mathbb{R}^n$$

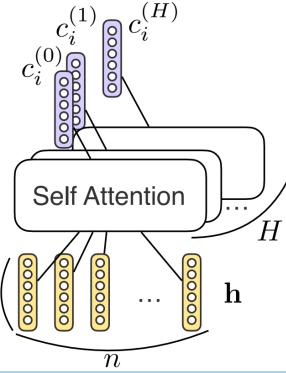
$$\alpha_i^{(j)} = \text{softmax}(s_i^{(j)}) \in \Delta^{N_{1:n}}$$

$$\mathbf{v}^{(j)} = V_j \mathbf{h} \in \mathbb{R}^{d_k \times n}$$

$$c_i^{(j)} = \sum_{i'=1}^n \alpha_i^{(j)} v^{(j)}$$

$$c_i^{(j)} = \text{SelfAttention}_j(\mathbf{h})_i$$

$$\theta_{\text{SelfAttention}_j} = \{K_j, Q_j, V_j\}$$

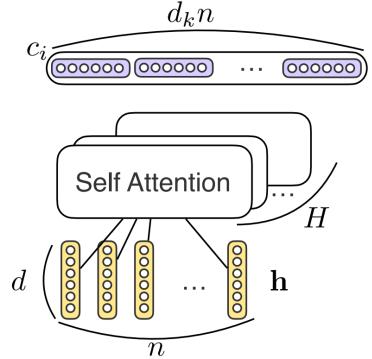


Compute attention over input vectors

$$c_i^{(j)} = \text{SelfAttention}_j(\mathbf{h})_i$$

$$c_i = [c_i^{(1)}; \dots; c_i^{(H)}]$$

concat

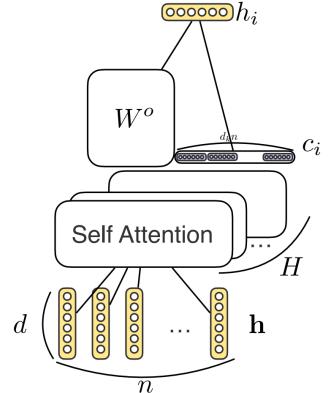


Compute attention over input vectors

$$c_i^{(j)} = \text{SelfAttention}_j(\mathbf{h})_i$$

$$c_i = [c_i^{(1)}; \dots; c_i^{(H)}]$$

$$\hat{h}_i = W^o c_i$$



residual + layernorm

1. Compute attention over input vectors

Layer normalization cuts down on uninformative variation in activations, speeding up training.

$$h_i^{\text{MHA}} = \text{MultiHeadAttention}^H(\mathbf{h})_i$$

2. Add and norm

$$h_i^{\text{AddNorm1}} = \text{LayerNorm}(h_i + h_i^{\text{MHA}})$$

LayerNorm : $\mathbb{R}^d \rightarrow \mathbb{R}^d$

$$\text{LayerNorm}(x) = \frac{g}{\sigma(x)}(x - \mu(x))$$

$$\mu(x) = \frac{1}{d} \sum_{i=1}^d x_i \quad \sigma(x) = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2}$$

$$h_i^{\text{AddNorm1}} = \text{LayerNorm}(h_i + h_i^{\text{MHA}})$$

1. Compute attention over input vectors

$$h_i^{\text{MHA}} = \text{MultiHeadAttention}^H(\mathbf{h})_i$$

2. Add and norm

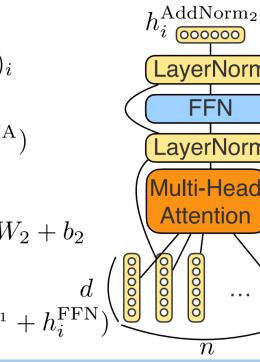
$$h_i^{\text{AddNorm1}} = \text{LayerNorm}(h_i + h_i^{\text{MHA}})$$

3. Feedforward layer

$$h_i^{\text{FFN}} = \text{ReLU}(h_i^{\text{AddNorm1}} W_1 + b_1) W_2 + b_2$$

4. Another add and norm

$$h_i^{\text{AddNorm2}} = \text{LayerNorm}(h_i^{\text{AddNorm1}} + h_i^{\text{FFN}})$$



1. Compute attention over input vectors

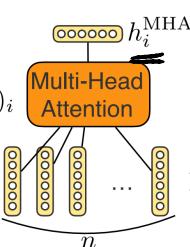
$$\theta_{\text{MultiHeadAttention}^H} = \{W^o\} \cup \{K_j, Q_j, V_j\}_{j=1}^H$$

$$c_i^{(j)} = \text{SelfAttention}_j(\mathbf{h})_i$$

$$c_i = [c_i^{(1)}; \dots; c_i^{(H)}]$$

$$\hat{h}_i = W^o c_i$$

Multi-headed attention allows us to look at multiple places in the input sequence at once.



1. Compute attention over input vectors

$$h_i^{\text{MHA}} = \text{MultiHeadAttention}^H(\mathbf{h})_i$$

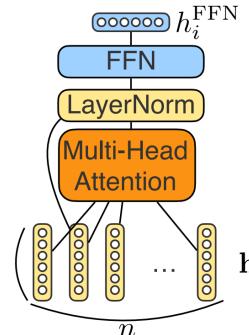
2. Add and norm

$$h_i^{\text{AddNorm1}} = \text{LayerNorm}(h_i + h_i^{\text{MHA}})$$

3. Feedforward layer

$$h_i^{\text{FFN}} = \text{ReLU}(h_i^{\text{AddNorm1}} W_1 + b_1) W_2 + b_2$$

Here is a nonlinearity that makes learning easier!



Residual Connection

$$h^{\text{MHA}} = \text{MultiHeadAttention}^H(\mathbf{h})$$

$$h^{\text{AddNorm}} = \text{LayerNorm}(\mathbf{h} + h^{\text{MHA}})$$

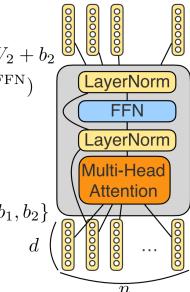
$$h^{\text{FFN}} = \text{ReLU}(\mathbf{h}^{\text{AddNorm}} W_1 + b_1) W_2 + b_2$$

$$\mathbf{h}' = \text{LayerNorm}(\mathbf{h}^{\text{AddNorm}} + h^{\text{FFN}})$$

$$\mathbf{h}' = \text{TransformerBlock}(\mathbf{h})$$

$$\theta_{\text{TransformerBlock}}$$

$$= \theta_{\text{MultiHeadAttention}^H} \cup \{W_1, W_2, b_1, b_2\}$$

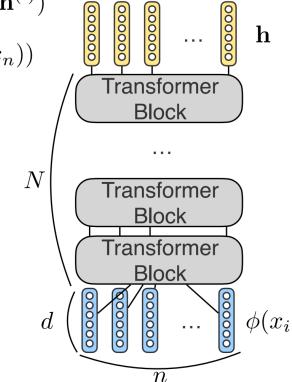


$$\mathbf{h}^{(l+1)} = \text{TransformerBlock}_{l+1}(\mathbf{h}^{(l)})$$

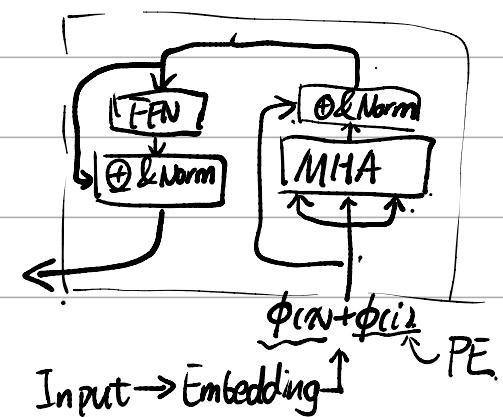
$$\mathbf{h}^N = \text{TransformerEncoder}(\phi(x_1), \dots, \phi(x_n))$$

$$\mathbf{h}^N \in \mathbb{R}^{d \times n}$$

Input: $\phi(x_i) = \phi(x) + \phi(i)$
word + position embeddings
 $\bar{x} = \langle x_1, \dots, x_n \rangle$



XN Encoder.

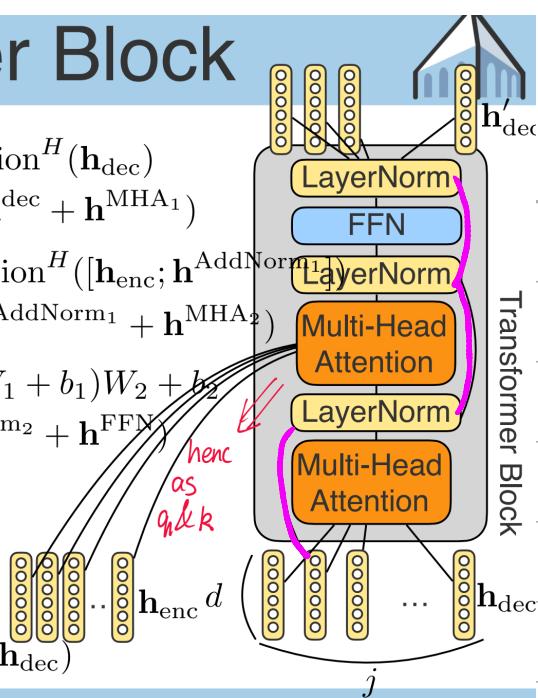


Decoder Block

$$\begin{aligned} \mathbf{h}^{\text{MHA}_1} &= \text{MultiHeadAttention}^H(\mathbf{h}_{\text{dec}}) \\ \mathbf{h}^{\text{AddNorm}_1} &= \text{LayerNorm}(\mathbf{h}_{\text{dec}} + \mathbf{h}^{\text{MHA}_1}) \\ \mathbf{h}^{\text{MHA}_2} &= \text{MultiHeadAttention}^H([\mathbf{h}_{\text{enc}}; \mathbf{h}^{\text{AddNorm}_1}]) \\ \mathbf{h}^{\text{AddNorm}_2} &= \text{LayerNorm}(\mathbf{h}^{\text{AddNorm}_1} + \mathbf{h}^{\text{MHA}_2}) \\ \mathbf{h}^{\text{FFN}} &= \text{ReLU}(\mathbf{h}^{\text{AddNorm}_2} W_1 + b_1) W_2 + b_2 \\ \mathbf{h}'_{\text{dec}} &= \text{LayerNorm}(\mathbf{h}^{\text{AddNorm}_2} + \mathbf{h}^{\text{FFN}}) \end{aligned}$$

One last layer norm to get our final output vectors for each token generated so far!

$$\mathbf{h}'_{\text{dec}} = \text{DecoderBlock}(\mathbf{h}_{\text{enc}}, \mathbf{h}_{\text{dec}})$$



- Decoder
First MHA: self.

Second MHA: $q, k = h_{\text{enc}}$

\triangleright : LN's result

同时, MHA & LN 数量 +1

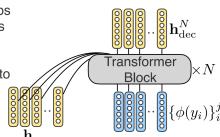
Res 连接都跟上

$$\mathbf{h}_{\text{dec}}^{(l+1)} = \text{DecoderBlock}_{l+1}(\mathbf{h}_{\text{enc}}, \mathbf{h}_{\text{dec}}^{(l)})$$

$$\mathbf{h}_{\text{dec}}^N = \text{Decoder}(\mathbf{h}_{\text{enc}}, \phi(y_1), \dots, \phi(y_j))$$

$$\text{Transformer} : \mathbb{R}^{d \times n} \times \mathcal{V}^j \rightarrow \mathbb{R}^{d \times j}$$

The transformer decoder maps from a set of input encodings and a sequence of tokens generated so far to a set of vectors, each corresponding to a token in the currently-generated sequence



How to generate index $j+1$'s word?

Decoder 前输入前 $j+1$ word,

算出 h_{dec}^N 后, 取 h_j^N :

WGR $^{m \times d}$,

$$\text{prob}(Y_{j+1} | y_1, \dots, y_j) = \text{softmax}(W h_j^N)$$

Given some paired data (\bar{x}, \bar{y}) :

$$\mathbf{h}_{\text{enc}} = \text{Transformer}_{\text{Encoder}}(\bar{x}) \quad \text{We can compute these in parallel on the GPU}$$

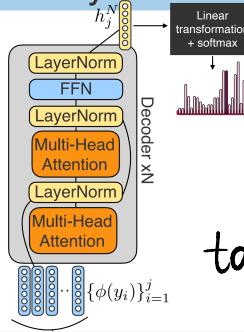
$$p(y_{i+1}) = \text{Decoder}(\mathbf{h}_{\text{enc}}, y_1, \dots, y_i)$$

Probability of a token should only depend on the ones that come before it! But we still want to take advantage of parallelism...

Masking:

- When doing the forward pass on the decoder, self-attention can be parallelized
- But we don't want the model to learn to rely on "future" words in the sequence!
- Solution: during training, set $s_{ij} = -\infty$ if i (query index) $<$ j (key/value index)

Decoder-Only Transformer



- Transformer for sequence generation

Decoder only.

take the last h_j^N !

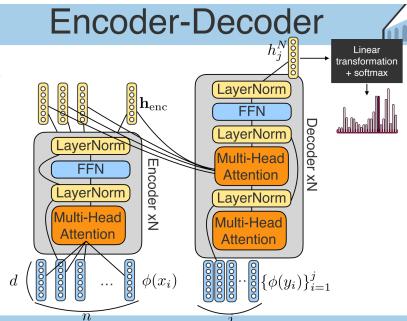
- Transformer for Sequence modeling
Encoder only

- Let's say we want to predict the word at index $j+1$
- All we need to do is transform this into a distribution over the vocabulary

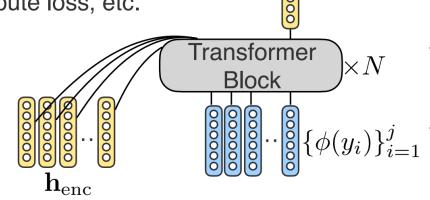
$$p(Y_{j+1} | x_1, \dots, x_n, y_1, \dots, y_j) = \text{softmax}(W h_j^N)$$

$$W \in \mathbb{R}^{|\mathcal{V}| \times d}$$

- Then we can sample, compute loss, etc.



$$h_j^N$$



(Training of EnkDec).

\neg -Training : for parallelism, masking is important !! Mask # 在 Decoder 中第

- / T MHA

Recall the autoregressive approximation of sequence probability:

$$p(\bar{x}) = \prod_{i=1}^{|\bar{x}|} p(x_i | x_1, \dots, x_{i-1})$$

Can we do this with a transformer model?

Yes! We just don't have an encoder (i.e., decoder-only LLM).

We still have to be careful about masking while training.

What if we just want some representation of a sentence, and we don't really care about being able to generate sentences? Enc(\bar{x})

We get this with the transformer encoder:

$$\mathbf{h}_{\text{enc}} = \text{Transformer}_{\text{Encoder}}(\bar{x})$$

$$p(\bar{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}, x_{i+1} \dots x_n)$$

Training can be problematic.

- From our transformer encoder, we get

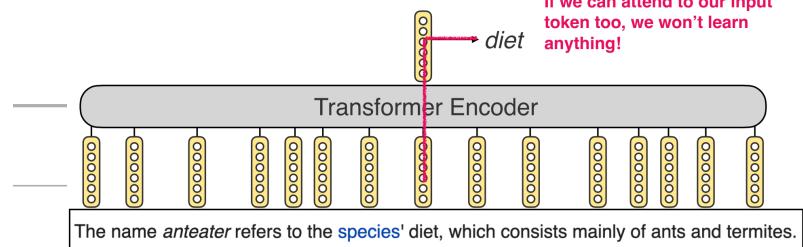
$$h_i = \text{Transformer}_{\text{Encoder}}(\bar{x})_i \in \mathbb{R}^d$$

- We could just transform this to the vocabulary space, and compute the loss from there:

$$p(X_i) = \text{softmax}(W h_i)$$

- Do you anticipate any problems?

If we can attend to our input token too, we won't learn anything!



For word with index j , output h_j can see x_j !!

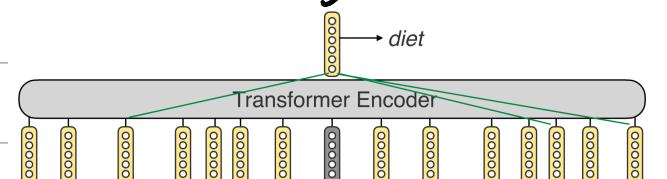
And our training goal is : $p(X_j = x_j) = \text{softmax}(W h_j)$

Information leakage!

\Rightarrow Mask word j ! For encoder

to infer with $x_i (i \neq j)$! ^{*}(mask some tokens).

\Rightarrow BERT



Instead: randomly sample some tokens to replace with a MASK placeholder token in the input, then train the model to predict those words