

# CS 6501 Natural Language Processing

## Statistical Language Modeling

---

Yangfeng Ji

Information and Language Processing Lab  
Department of Computer Science  
University of Virginia



1. Introduction
2.  $N$ -gram Language Models
3. Generation with Bi-gram Models
4. Smoothing Techniques
5. Language Model Evaluation

# Introduction

---

# Word Prediction with Context

Consider the example, what words are likely to follow

Please turn your homework ...

[Jurafsky and Martin, 2019]

# Word Prediction with Context

Consider the example, what words are likely to follow

Please turn your homework ...

Although we do not know the actual word in the original text, we have a good sense about what of these following words are likely to follow

- ▶ in
- ▶ over
- ▶ refrigerator
- ▶ the

[Jurafsky and Martin, 2019]

# Mathematical Formulation

- ▶ Let  $X_1, X_2, \dots, X_{t-1}$  be the random variables representing the words in the context, and  $X_t$  be the **next** word that we would like the model to predict.

# Mathematical Formulation

- ▶ Let  $X_1, X_2, \dots, X_{t-1}$  be the random variables representing the words in the context, and  $X_t$  be the **next** word that we would like the model to predict.
- ▶ Similarly, we can formulate this prediction as a **classification** problem, where  $\mathbf{X}_{1:t-1} = X_1, X_2, \dots, X_{t-1}$  are the input words and  $X_t$  is the output, we can write the classifier in a probabilistic form

$$P(X_t | X_1, \dots, X_{t-1}) \quad \text{or} \quad P(X_t | \mathbf{X}_{1:t-1}) \quad (1)$$

# Mathematical Formulation

- ▶ Let  $X_1, X_2, \dots, X_{t-1}$  be the random variables representing the words in the context, and  $X_t$  be the **next** word that we would like the model to predict.
- ▶ Similarly, we can formulate this prediction as a **classification** problem, where  $\mathbf{X}_{1:t-1} = X_1, X_2, \dots, X_{t-1}$  are the input words and  $X_t$  is the output, we can write the classifier in a probabilistic form

$$P(X_t | X_1, \dots, X_{t-1}) \quad \text{or} \quad P(X_t | \mathbf{X}_{1:t-1}) \quad (1)$$

- ▶ The topics in this and the next lectures will offer **two** modeling methods in equation 1.



# Mathematical Formulation

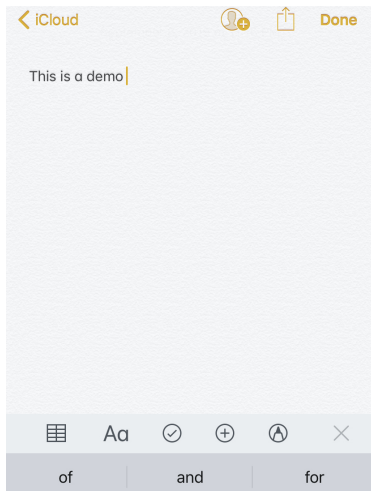
- ▶ Let  $X_1, X_2, \dots, X_{t-1}$  be the random variables representing the words in the context, and  $X_t$  be the **next** word that we would like the model to predict.
- ▶ Similarly, we can formulate this prediction as a **classification** problem, where  $\mathbf{X}_{1:t-1} = X_1, X_2, \dots, X_{t-1}$  are the input words and  $X_t$  is the output, we can write the classifier in a probabilistic form

$$P(X_t | X_1, \dots, X_{t-1}) \quad \text{or} \quad P(X_t | \mathbf{X}_{1:t-1}) \quad (1)$$

- ▶ The topics in this and the next lectures will offer **two** modeling methods in equation 1.
- ▶ Difference with the word embedding methods discussed in the previous lecture
  - ▶ Skip-gram model: predicting the **surrounding** words
  - ▶ Language models: predicting the **next** word

# Word Prediction in Input Methods

Input methods use language models to predict the next likely words, to speed up the typing



# Writing a Poem?

Trevor Noah and Amanda Gorman writing poems with the input methods on their phones



Figure: The Daily Social Distancing Show: Bonus Track feat. Amanda Gorman

Link

# Joint Probability and Chain Rule

Given the conditional probability  $P(X_t | \mathbf{X}_{1:t-1})$ , to evaluate the quality of a text, we need the chain rule in probability to factorize the joint probability  $P(\mathbf{X}_{1:t})$  into a series of conditional probabilities.

$$P(X_1, X_2, \dots, X_t) = P(X_1)P(X_2, \dots, X_k | X_1)$$

(2)

# Joint Probability and Chain Rule

Given the conditional probability  $P(X_t | \mathbf{X}_{1:t-1})$ , to evaluate the quality of a text, we need the chain rule in probability to factorize the joint probability  $P(\mathbf{X}_{1:t})$  into a series of conditional probabilities.

$$\begin{aligned} P(X_1, X_2, \dots, X_t) &= P(X_1)P(X_2, \dots, X_k | X_1) \\ &= P(X_1)P(X_2 | X_1)P(X_3, \dots, X_t | X_1, X_2) \end{aligned}$$

(2)

# Joint Probability and Chain Rule

Given the conditional probability  $P(X_t | \mathbf{X}_{1:t-1})$ , to evaluate the quality of a text, we need the chain rule in probability to factorize the joint probability  $P(\mathbf{X}_{1:t})$  into a series of conditional probabilities.

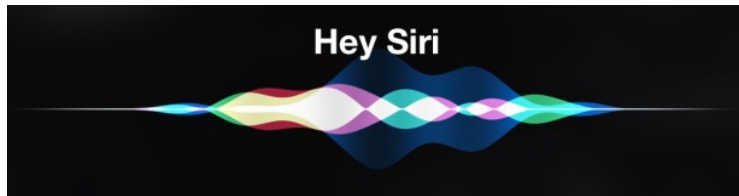
$$\begin{aligned}P(X_1, X_2, \dots, X_t) &= P(X_1)P(X_2, \dots, X_t | X_1) \\ &= P(X_1)P(X_2 | X_1)P(X_3, \dots, X_t | X_1, X_2) \\ &= P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2) \cdots \\ &\quad P(X_t | X_1, \dots, X_{t-1})\end{aligned}$$

(2)

# Joint Probability and Chain Rule

Given the conditional probability  $P(X_t | X_{1:t-1})$ , to evaluate the quality of a text, we need the chain rule in probability to factorize the joint probability  $P(X_{1:t})$  into a series of conditional probabilities.

$$\begin{aligned} P(X_1, X_2, \dots, X_t) &= P(X_1)P(X_2, \dots, X_t | X_1) \\ &= P(X_1)P(X_2 | X_1)P(X_3, \dots, X_t | X_1, X_2) \\ &= P(X_1)P(X_2 | X_1)P(X_3 | X_1, X_2) \cdots \\ &\quad P(X_t | X_1, \dots, X_{t-1}) \\ &= \prod_{i=1}^t P(X_i | X_1, \dots, X_{i-1}) \end{aligned} \tag{2}$$



Given a voice signal, a language model in speech recognition will evaluate the likelihood of decoded texts

$$P(\text{I saw a van}) \gg P(\text{eyes awe of an}) \quad (3)$$

[Jurafsky and Martin, 2019]



## Grammarly:



Rooms that are tiny can be tricky to decorate but they can also be a lot of fun. So when a client challenged us to give her pocket size space a summer makeover for under \$500 dollars, we just couldn't say no. Transforming a very small space doesn't have to blow your budget. Small things like finding a vintage piece of furniture from a relative or adding a fresh coat of paint to your own dated items can add a stylish splash to any abode.

### Correctness

2 alerts

A horizontal progress bar for the 'Correctness' metric, showing a red segment followed by a light pink segment.

### Clarity

A bit unclear

A horizontal progress bar for the 'Clarity' metric, showing a blue segment followed by a light blue segment.

### Engagement

A bit bland

A horizontal progress bar for the 'Engagement' metric, showing a green segment followed by a light green segment.

### Delivery

Slightly off

A horizontal progress bar for the 'Delivery' metric, showing a purple segment followed by a light purple segment.

## Grammarly:



Rooms that are tiny can be tricky to decorate but they can also be a lot of fun. So when a client challenged us to give her pocket size space a summer makeover for under \$500 dollars, we just couldn't say no. Transforming a very small space doesn't have to blow your budget. Small things like finding a vintage piece of furniture from a relative or adding a fresh coat of paint to your own dated items can add a stylish splash to any abode.

### Correctness

2 alerts

### Clarity

A bit unclear

### Engagement

A bit bland

### Delivery

Slightly off

A good writing assistant system involves two tasks

- ▶ evaluate the quality of a text
- ▶ generate revision suggestions

A language model cannot provide support to all functions directly, but is a critical component in the backend system

- ▶ Generative tasks: predicting the next word given a context
  - ▶ Word prediction
  - ▶ Text generation
  - ▶ ...

- ▶ Generative tasks: predicting the next word given a context
  - ▶ Word prediction
  - ▶ Text generation
  - ▶ ...
- ▶ Discriminative tasks: evaluating the quality of texts
  - ▶ Speech recognition
  - ▶ Machine translation
  - ▶ Document summarization
  - ▶ ...

# *N*-gram Language Models

---

## Problem Definition

Given a vocab  $\mathcal{V}$  that contains all the possible word types, then the prediction of  $X_t$  can be formulated as

$$P(X_t \mid X_1, \dots, X_{t-1}) =? \quad (4)$$

# Problem Definition

Given a vocab  $\mathcal{V}$  that contains all the possible word types, then the prediction of  $X_t$  can be formulated as

$$P(X_t \mid X_1, \dots, X_{t-1}) =? \quad (4)$$

The challenges of modeling  $P(X_t \mid X_1, \dots, X_{t-1})$

- ▶ it is a categorical distribution defined on the vocab  $\mathcal{V}$
  
- ▶ it consider the entire context from the very first word  $X_1$  to the previous word  $X_{t-1}$

# Problem Definition

Given a vocab  $\mathcal{V}$  that contains all the possible word types, then the prediction of  $X_t$  can be formulated as

$$P(X_t | X_1, \dots, X_{t-1}) =? \quad (4)$$

The challenges of modeling  $P(X_t | X_1, \dots, X_{t-1})$

- ▶ it is a categorical distribution defined on the vocab  $\mathcal{V}$ 
  - ▶ Unfortunately, we cannot do much on this problem, other than using hierarchical structures in softmax function (e.g., hierarchical softmax and class-factored softmax)
- ▶ it consider the entire context from the very first word  $X_1$  to the previous word  $X_{t-1}$



# Problem Definition

Given a vocab  $\mathcal{V}$  that contains all the possible word types, then the prediction of  $X_t$  can be formulated as

$$P(X_t | X_1, \dots, X_{t-1}) =? \quad (4)$$

The challenges of modeling  $P(X_t | X_1, \dots, X_{t-1})$

- ▶ it is a categorical distribution defined on the vocab  $\mathcal{V}$ 
  - ▶ Unfortunately, we cannot do much on this problem, other than using hierarchical structures in softmax function (e.g., hierarchical softmax and class-factored softmax)
- ▶ it consider the entire context from the very first word  $X_1$  to the previous word  $X_{t-1}$ 
  - ▶ The main topic of this section

# Parameter Estimation

With a collection of texts as training examples, the simple method of estimating the probabilities is using maximum likelihood estimation.

- ▶ In the first lecture, we discussed the MLE of a Bernoulli distribution

$$\hat{P}(X = 1) = \frac{\sum_{i=1}^N \delta(x_i, 1)}{N} = \frac{c(X = 1)}{N} \quad (5)$$

where  $c(X = 1)$  is the number of observations with value 1

# Parameter Estimation

With a collection of texts as training examples, the simple method of estimating the probabilities is using maximum likelihood estimation.

- ▶ In the first lecture, we discussed the MLE of a Bernoulli distribution

$$\hat{P}(X = 1) = \frac{\sum_{i=1}^N \delta(x_i, 1)}{N} = \frac{c(X = 1)}{N} \quad (5)$$

where  $c(X = 1)$  is the number of observations with value 1

- ▶ Similarly, to estimate the conditional probability  $P(X_t | \mathbf{X}_{1:t-1})$ , we have

$$\hat{P}(X_t = x_t | \mathbf{X}_{1:t-1} = \mathbf{x}_{1:t-1}) = \frac{c(\mathbf{x}_{1:t})}{c(\mathbf{x}_{1:t-1})} \quad (6)$$

where  $c(\mathbf{x}_{1:t-1})$  is the number that text  $\mathbf{x}_{1:t-1}$  appears in the training examples, and  $c(\mathbf{x}_{1:t})$  is the number that text  $\mathbf{x}_{1:t}$  appears in the training examples.

[Collins, 2017]

## Parameter Estimation (Cont.)

Imagine we have a **huge** collection of texts for parameter estimation

- ▶ With the sentence “the dog barks”

$$\hat{P}(X_3 = \text{barks} \mid \mathbf{X}_{1:2} = \text{the dog}) = \frac{c(\mathbf{X}_{1:3} = \text{the dog barks})}{c(\mathbf{X}_{1:2} = \text{the dog})} \quad (7)$$

# Parameter Estimation (Cont.)

Imagine we have a **huge** collection of texts for parameter estimation

- ▶ With the sentence “the dog barks”

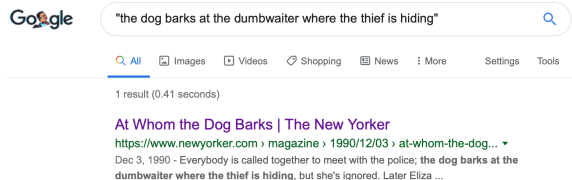
$$\hat{P}(X_3 = \text{barks} \mid \mathbf{X}_{1:2} = \text{the dog}) = \frac{c(\mathbf{X}_{1:3} = \text{the dog barks})}{c(\mathbf{X}_{1:2} = \text{the dog})} \quad (7)$$

- ▶ With the sentence “the dog barks at the dumbwaiter where the thief is hiding”

$$\begin{aligned} \hat{P}(X_{11} = \text{hiding} \mid \mathbf{X}_{1:10} = \text{the dog} \dots \text{is}) \\ = \frac{c(\mathbf{X}_{1:11} = \text{the dog} \dots \text{is hiding})}{c(\mathbf{X}_{1:10} = \text{the dog} \dots \text{is})} \end{aligned} \quad (8)$$

# Parameter Estimation (Cont.)

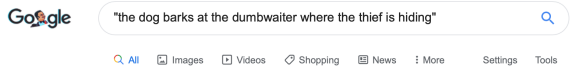
For this specific sentence, we only one training example even if we collect all the texts from the Internet



The image shows a Google search interface. The search bar contains the text "the dog barks at the dumbwaiter where the thief is hiding". Below the search bar, there are navigation options: All, Images, Videos, Shopping, News, More, Settings, and Tools. The search results show 1 result in 0.41 seconds. The result is a link to "At Whom the Dog Barks | The New Yorker" with a URL: <https://www.newyorker.com/magazine/1990/12/03/at-whom-the-dog...>. The snippet below the link reads: "Dec 3, 1990 - Everybody is called together to meet with the police; the dog barks at the dumbwaiter where the thief is hiding, but she's ignored. Later Eliza ..."

# Parameter Estimation (Cont.)

For this specific sentence, we only one training example even if we collect all the texts from the Internet



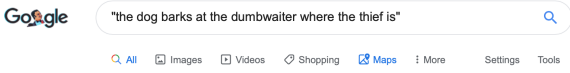
Google search results for the query "the dog barks at the dumbwaiter where the thief is hiding". The search bar contains the full sentence. Below the search bar are navigation links: All, Images, Videos, Shopping, News, More, Settings, and Tools. The search results show 1 result in 0.41 seconds.

1 result (0.41 seconds)

[At Whom the Dog Barks | The New Yorker](#)

<https://www.newyorker.com> › [magazine](#) › 1990/12/03 › [at-whom-the-dog...](#) ▾

Dec 3, 1990 - Everybody is called together to meet with the police; **the dog barks at the dumbwaiter where the thief is hiding**, but she's ignored. Later Eliza ...



Google search results for the query "the dog barks at the dumbwaiter where the thief is". The search bar contains the partial sentence. Below the search bar are navigation links: All, Images, Videos, Shopping, Maps, More, Settings, and Tools. The search results show 1 result in 0.40 seconds.

1 result (0.40 seconds)

[At Whom the Dog Barks | The New Yorker](#)

<https://www.newyorker.com> › [magazine](#) › 1990/12/03 › [at-whom-the-dog...](#) ▾

Dec 3, 1990 - Everybody is called together to meet with the police; **the dog barks at the dumbwaiter where the thief is hiding**, but she's ignored. Later Eliza ...

# Parameter Estimation (Cont.)

For this specific sentence, we only one training example even if we collect all the texts from the Internet

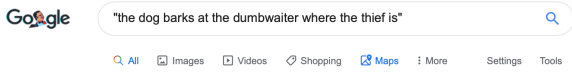


1 result (0.41 seconds)

[At Whom the Dog Barks | The New Yorker](#)

<https://www.newyorker.com> › [magazine](#) › 1990/12/03 › [at-whom-the-dog...](#) ▾

Dec 3, 1990 - Everybody is called together to meet with the police; the dog barks at the dumbwaiter where the thief is hiding, but she's ignored. Later Eliza ...



1 result (0.40 seconds)

[At Whom the Dog Barks | The New Yorker](#)

<https://www.newyorker.com> › [magazine](#) › 1990/12/03 › [at-whom-the-dog...](#) ▾

Dec 3, 1990 - Everybody is called together to meet with the police; the dog barks at the dumbwaiter where the thief is hiding, but she's ignored. Later Eliza ...

$$\hat{P}(X_{11} = \text{hiding} \mid \mathbf{X}_{1:10} = \text{the dog} \cdots \text{is}) = 1.0$$



## Simplification: Uni-gram

- ▶ The main challenge of parameter estimation is the long-term dependence between  $X_t$  and  $X_{1:t-1}$

## Simplification: Uni-gram

- ▶ The main challenge of parameter estimation is the long-term dependence between  $X_t$  and  $X_{1:t-1}$
- ▶ Uni-gram: assume all words are independent with each other. With this assumption, we only need to estimate the probability of each individual word (no conditional probability involved)

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t) \quad (9)$$

# Simplification: Uni-gram

- ▶ The main challenge of parameter estimation is the long-term dependence between  $X_t$  and  $X_{1:t-1}$
- ▶ Uni-gram: assume all words are independent with each other. With this assumption, we only need to estimate the probability of each individual word (no conditional probability involved)

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t) \quad (9)$$

- ▶ For example

$$P(\text{barks} | \text{the dog}) \approx P(\text{barks}) \quad (10)$$

# Simplification: Uni-gram

- ▶ The main challenge of parameter estimation is the long-term dependence between  $X_t$  and  $X_{1:t-1}$
- ▶ Uni-gram: assume all words are independent with each other. With this assumption, we only need to estimate the probability of each individual word (no conditional probability involved)

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t) \quad (9)$$

- ▶ For example

$$P(\text{barks} | \text{the dog}) \approx P(\text{barks}) \quad (10)$$

- ▶ Comments: the tradeoff between prediction power and number of parameters
  - ▶ It has extremely limited prediction power
  - ▶ Number of parameters:  $V = |\mathcal{V}|$

# Bi-gram Models

- ▶ To find a good balance between the prediction power and parameter estimation challenge, we can **limit the contextual information used** in a language modeling.
- ▶ Bi-gram model: uses only one word  $X_{t-1}$  from the previous context to predict the current word  $X_t$

$$P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t | X_{t-1}) \quad (11)$$

# Bi-gram Models

- ▶ To find a good balance between the prediction power and parameter estimation challenge, we can **limit the contextual information used** in a language modeling.
- ▶ Bi-gram model: uses only one word  $X_{t-1}$  from the previous context to predict the current word  $X_t$

$$P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t | X_{t-1}) \quad (11)$$

- ▶ For example, given the text “the dog barks”, the prediction of the last word barks in a bi-gram model is formulated as

$$P(\text{barks} | \text{the dog}) \approx P(\text{barks} | \text{dog}) \quad (12)$$

# Bi-gram Models

- ▶ To find a good balance between the prediction power and parameter estimation challenge, we can **limit the contextual information used** in a language modeling.
- ▶ Bi-gram model: uses only one word  $X_{t-1}$  from the previous context to predict the current word  $X_t$

$$P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t | X_{t-1}) \quad (11)$$

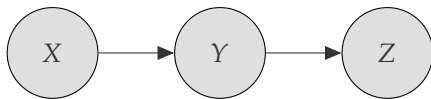
- ▶ For example, given the text “the dog barks”, the prediction of the last word barks in a bi-gram model is formulated as

$$P(\text{barks} | \text{the dog}) \approx P(\text{barks} | \text{dog}) \quad (12)$$

- ▶ In probabilistic modeling, a bi-gram model is an application of the first-order Markov model

# Markov Property

First-order Markov property: given

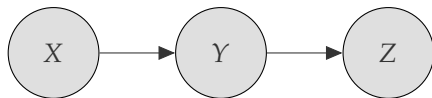


$$P(Z | X, Y) = P(Z | Y) \quad (13)$$



# Markov Property

First-order Markov property: given



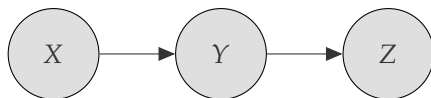
$$P(Z | X, Y) = P(Z | Y) \quad (13)$$

It simplifies the conditional probability

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t | X_{t-1}) \quad (14)$$

# Markov Property

First-order Markov property: given



$$P(Z | X, Y) = P(Z | Y) \quad (13)$$

It simplifies the conditional probability

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t | X_{t-1}) \quad (14)$$

and also the joint probability

$$P(X_1, \dots, X_t) \approx P(X_t | X_{t-1}) \cdot P(X_{t-1} | X_{t-2}) \cdots P(X_2 | X_1) \cdot P(X_1) \quad (15)$$

# Special Tokens

Consider the application of using a bi-gram model

$$P(\text{the dog barks}) = P(\text{the}) \cdot P(\text{dog} \mid \text{the}) \\ P(\text{barks} \mid \text{dog})$$

# Special Tokens

Consider the application of using a bi-gram model

$$P(\text{the dog barks}) = P(\text{the}) \cdot P(\text{dog} \mid \text{the}) \\ P(\text{barks} \mid \text{dog})$$

The model needs

- ▶ a special token ( $\square$ ) to distinguish  $P(\text{the})$  from the marginal distribution of word “the”

# Special Tokens

Consider the application of using a bi-gram model

$$P(\text{the dog barks}) = P(\text{the}) \cdot P(\text{dog} \mid \text{the}) \\ P(\text{barks} \mid \text{dog})$$

The model needs

- ▶ a special token ( $\square$ ) to distinguish  $P(\text{the})$  from the marginal distribution of word “the”
- ▶ another special token ( $\blacksquare$ ) to indicate the end of a sentence

# Special Tokens

Consider the application of using a bi-gram model

$$P(\text{the dog barks}) = P(\text{the}) \cdot P(\text{dog} \mid \text{the}) \\ P(\text{barks} \mid \text{dog})$$

The model needs

- ▶ a special token ( $\square$ ) to distinguish  $P(\text{the})$  from the marginal distribution of word “the”
- ▶ another special token ( $\blacksquare$ ) to indicate the end of a sentence

Factorization with special tokens:

$$P(\square \text{ the dog barks } \blacksquare) = P(\text{the} \mid \square) \cdot P(\text{dog} \mid \text{the}) \\ P(\text{barks} \mid \text{dog}) \cdot P(\blacksquare \mid \text{barks})$$

# Example: Parameter Estimation

## Example sentences

- ▶ □ I am Sam ■
- ▶ □ Sam I am ■
- ▶ □ I do not like green eggs and ham ■

[Jurafsky and Martin, 2019]

# Example: Parameter Estimation

Example sentences

- ▶ □ I am Sam ■
- ▶ □ Sam I am ■
- ▶ □ I do not like green eggs and ham ■

Some of the probabilities:

$$\hat{P}(I \mid \square) = \frac{2}{3} \quad \hat{P}(\blacksquare \mid \text{Sam}) = \frac{1}{2} \quad \hat{P}(\text{do} \mid I) = \frac{1}{3} \quad (16)$$

[Jurafsky and Martin, 2019]



- ▶  $P(X_t | X_{t-1})$  is defined a fixed vocabulary, for normalization purpose

$$P(X_t | X_{t-1}) = \frac{c(X_{t-1}, X_t)}{\sum_{X' \in \mathcal{V}} c(X_{t-1}, X')} \quad (17)$$

- ▶  $P(X_t | X_{t-1})$  is defined a fixed vocabulary, for normalization purpose

$$P(X_t | X_{t-1}) = \frac{c(X_{t-1}, X_t)}{\sum_{X' \in \mathcal{V}} c(X_{t-1}, X')} \quad (17)$$

- ▶ Issues with a fixed vocabulary
  - ▶ Unknown words: word  $x$  is not in the vocabulary
  - ▶ Zero probability: word combination  $(x, x')$  never appears in the training set

# Unknown Words

Replace all words that are not in the vocab with a special token UNK.

For example

- ▶ Original text: “the dog barks at the dumbwaiter where the thief is hiding”
- ▶ After preprocessing: “the dog barks at the UNK where the thief is hiding”

# Unknown Words

Replace all words that are not in the vocab with a special token UNK.

For example

- ▶ Original text: “the dog barks at the dumbwaiter where the thief is hiding”
- ▶ After preprocessing: “the dog barks at the UNK where the thief is hiding”

## Quiz

Can we simply ignore the unknown words? For example, what if the preprocessed text is

“the dog barks at the where the thief is hiding”

We can extend the conditional probability to depend on previous two tokens

$$P(X_t | X_1, \dots, X_{t-1}) \approx P(X_t | X_{t-2}, X_{t-1}) \quad (18)$$

Comments

- ▶ More dependency leads to more accurate predictions
- ▶ Parameter estimation

$$\hat{P}(X_t | X_{t-2}, X_{t-1}) = \frac{c(\mathbf{X}_{t-2:t})}{c(\mathbf{X}_{t-2:t-1})} \quad (19)$$

# Number of Parameters

- ▶ Uni-gram model

- ▶ Ignore context words completely  $P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t)$
- ▶ Number of parameters  $\mathcal{O}(|\mathcal{V}|)$

	$X_1$	$\dots$	$X_{ \mathcal{V} }$
$P(X_t)$			

# Number of Parameters

## ▶ Uni-gram model

- ▶ Ignore context words completely  $P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t)$
- ▶ Number of parameters  $\mathcal{O}(|\mathcal{V}|)$

	$X_1$	$\dots$	$X_{ \mathcal{V} }$
$P(X_t)$			

## ▶ Bi-gram model

- ▶ Use only the adjacent word  $P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t | X_{t-1})$
- ▶ Number of parameters  $\mathcal{O}(|\mathcal{V}|^2)$

$P(X_t   X_{t-1})$	$X_1$	$\dots$	$X_{ \mathcal{V} }$
$X_1$			
$\vdots$			
$X_{ \mathcal{V} }$			

# Number of Parameters

## ▶ Uni-gram model

- ▶ Ignore context words completely  $P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t)$
- ▶ Number of parameters  $\mathcal{O}(|\mathcal{V}|)$

	$X_1$	$\dots$	$X_{ \mathcal{V} }$
$P(X_t)$			

## ▶ Bi-gram model

- ▶ Use only the adjacent word  $P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t | X_{t-1})$
- ▶ Number of parameters  $\mathcal{O}(|\mathcal{V}|^2)$

$P(X_t   X_{t-1})$	$X_1$	$\dots$	$X_{ \mathcal{V} }$
$X_1$			
$\vdots$			
$X_{ \mathcal{V} }$			

## ▶ Tri-gram model

- ▶ Use two preceding words  $P(X_t | \mathbf{X}_{1:t-1}) \approx P(X_t | X_{t-2}, X_{t-1})$
- ▶ Number of parameters  $\mathcal{O}(|\mathcal{V}|^3)$



## Generation with Bi-gram Models

---

- ▶ A bi-gram model with no smoothing
- ▶ Training with the dataset from the arXiv paper abstracts
- ▶ Generating by *randomly* sampling from this bi-gram model

- ▶ A bi-gram model with no smoothing
- ▶ Training with the dataset from the arXiv paper abstracts
- ▶ Generating by *randomly* sampling from this bi-gram model

Checkout the demo code for some examples

# Smoothing Techniques

---

A motivating example:

The `printer` on the 5th floor of Rice hall `crashed`

A motivating example:

The **printer** on the 5th floor of Rice hall **crashed**

$n$ -gram Language Models

- ▶ Uni-gram:  $P(X_t)$
- ▶ Bi-gram:  $P(X_t | X_{t-1})$
- ▶ Tri-gram:  $P(X_t | X_{t-2}, X_{t-1})$
- ▶ 4-gram:  $P(X_t | X_{t-3}, X_{t-2}, X_{t-1})$
- ▶ 5-gram:  $P(X_t | X_{t-4}, X_{t-3}, X_{t-2}, X_{t-1})$

It is the same method used in parameter estimation of naive Bayes classifiers

$$P(X_t | X_{t-1}) = \frac{c(X_{t-1}, X_t) + \alpha}{c(X_{t-1}) + \alpha V} \quad (20)$$

where  $\alpha > 0$  is a hyper-parameter.

Estimate the following three models with MLE:

- ▶ Uni-gram:  $P(X_t)$
- ▶ Bi-gram:  $P(X_t | X_{t-1})$
- ▶ Tri-gram:  $P(X_t | X_{t-2}, X_{t-1})$



Estimate the following three models with MLE:

- ▶ Uni-gram:  $P(X_t)$
- ▶ Bi-gram:  $P(X_t | X_{t-1})$
- ▶ Tri-gram:  $P(X_t | X_{t-2}, X_{t-1})$

Then, the new probability of  $X_t$  given  $X_{t-2}$  and  $X_{t-1}$  is

$$\begin{aligned} P_{LI}(X_t | X_{t-2}, X_{t-1}) &= \lambda_1 \cdot P(X_t) + \lambda_2 \cdot P(X_t | X_{t-1}) \\ &\quad + \lambda_3 \cdot P(X_t | X_{t-2}, X_{t-1}) \end{aligned} \quad (21)$$

$\{\lambda_i\}$  are learned with a held-out corpus (a development set).

# Language Model Evaluation

---

Evaluation with joint probabilities

$$P(\text{I love black coffee}) \text{ vs. } P(\text{black coffee pleases me}) \quad (22)$$

Direct comparison between the probabilities will tell us which sentence is more *fluent*.

Limitation of comparing joint probabilities directly

$$P(\text{I love black coffee}) \text{ vs. } P(\text{I like black coffee very much}) \quad (23)$$

Due to the *length difference*, the second probability may always be smaller than the first.

- ▶ Test data: including the special tokens

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$$

- ▶ Likelihood

$$\text{Log-lik}(\{\mathbf{x}_{m=1}^M\}) = \log_2 \prod_{m=1}^M \prod_{t=1} P(x_{m,t} | \mathbf{x}_{m,1:t-1}) \quad (24)$$

$$= \sum_{m=1}^M \sum_{t=1} \log_2 P(x_{m,t} | \mathbf{x}_{m,1:t-1}) \quad (25)$$

- ▶ Factors
  - ▶ Number of the tokens
  - ▶ No intuitive explanation

The definition of perplexity is

$$\text{Perplexity} = 2^{-\frac{1}{T} \text{Log-lik}(\{x_{m=1}^M\})} \quad (26)$$

where  $T$  is the total number of the log probabilities in  $\text{Log-lik}(\{x_{m=1}^M\})$ .

- ▶ An impossible case

$$P(x_t | \mathbf{x}_{1:t-1}) = 1 \quad (27)$$

- ▶ An impossible case

$$P(x_t | \mathbf{x}_{1:t-1}) = 1 \quad (27)$$

- ▶ Perplexity

$$\begin{aligned} \text{Perplexity} &= 2^{-\frac{1}{T} \sum_{k=1}^M \sum_{m=1} \log_2 1} \\ &= 2^0 \\ &= 1 \end{aligned} \quad (28)$$



- ▶ A trivial case

$$P(x_t \mid x_{1:t-1}) = \frac{1}{|\mathcal{V}|} \quad (29)$$

- ▶ A trivial case

$$P(x_t \mid x_{1:t-1}) = \frac{1}{|\mathcal{V}|} \quad (29)$$

- ▶ Perplexity

$$\begin{aligned} \text{Perplexity} &= 2^{-\frac{1}{T} \sum_{k=1}^M \sum_{m=1} \log_2 \frac{1}{|\mathcal{V}|}} \\ &= 2^{-\frac{1}{T} (T \cdot \log_2 \frac{1}{|\mathcal{V}|})} \\ &= 2^{-\log_2 \frac{1}{|\mathcal{V}|}} \\ &= |\mathcal{V}| \end{aligned} \quad (30)$$

# Typical Values of Perplexity

- ▶  $|\mathcal{V}| = 50K$
- ▶ A uni-gram model: Perplexity = 955
- ▶ A bi-gram model: Perplexity = 137
- ▶ A tri-gram model: Perplexity = 74

Lower is better

[Collins, 2017]

# A Few Comments on Perplexity

## Perplexity

- ▶ is an intrinsic evaluation measurement

## Perplexity

- ▶ is an intrinsic evaluation measurement
- ▶ is not necessarily correlated with the performance of
  - ▶ e.g., lower perplexity does not mean better translation (wrt BLEU score)
- ▶ is not directly comparable even on the same test data
  - ▶ you need the **exactly same** input for comparison



Collins, M. (2017).

Natural language processing: Lecture notes.



Jurafsky, D. and Martin, J. (2019).

Speech and language processing.