

# CS540 Introduction to Artificial Intelligence

## Lecture 9

Young Wu

Based on lecture slides by Jerry Zhu and Yingyu Liang

June 20, 2019

# Feedback, Lectures

## Admin

- Lectures:
  - 1 Concepts: more + RN textbook
  - 2 Examples: more + quiz questions + YouTube videos
  - 3 Applications: more + website and paper links
  - 4 Math: less + YouTube videos
  - 5 Implementation: less + hint file







# Feedback, Assignments

## Admin

- Short instruction: flexibility and creativity
- Long instruction: details, example workflow
- Solutions: details, help with coding
- Each assignment will be re-graded 3 times: all auto-graded, deal with individual submissions after the final no penalty due date.

# Discriminative Model vs Generative Model

## Review

- Week 1 to Week 4 focus on discriminative models.
- Given a training set  $(x_i, y_i)_{i=1}^n$ , the task is classification (machine learning) or regression (statistics), i.e. finding a function  $\hat{f}$  such that given new instances  $x'_i$ ,  $y$  can be predicted as  $\hat{y}_i = \hat{f}(x'_i)$ .
- The function  $\hat{f}$  is usually represented by parameters  $w$  and  $b$ . These parameters can be learned by methods such as gradient descent by minimizing some cost objective function.

# Perceptron

## Review

- Model: LTU Perceptron.

- Objective: minimize mistakes =  $\sum_{i=1}^n \mathbb{1}_{\{y_i \neq a_i\}}$  or maximize accuracy. It is equivalent to minimizing squared error cost, absolute value cost, log cost (cross entropy loss).

- Training: Perceptron algorithm.

- Prediction:  $\hat{y}_i = a'_i = \mathbb{1}_{\{w^T x'_i + b \geq 0\}}$ .





# Neural Network

## Review

- Model: Fully Connected Neural Network
- Objective: minimize squared error cost  $= \sum_{i=1}^n (y_i - a_i^{(L)})^2$ .
- Training: Backpropagation: gradient descent algorithm using chain rule.
- Prediction:  $\hat{y}_i = \mathbb{1}_{\{a_i^{(L)} \geq 0.5\}}$ ,  $a_i^{(l)} = g \left( (w^{(l)})^T a^{(l-1)} + b^{(l)} \right)$   
with  $a^{(0)} = x_i'$ .

# Support Vector Machine

## Review

- Model: Support Vector Machine
- Objective: minimize regularized hinge cost  
$$= \sum_{i=1}^n \max \left\{ 0, 1 - (2y_i - 1) \left( w^T x_i + b \right) \right\} + \lambda \|w\|_2^2$$
 or  
maximize margin.
- Training: Pegasos algorithm: Primal Estimated sub-GrAdient Solver for SVM.
- Prediction:  $\hat{y}_i = a'_i = \mathbb{1}_{\{w^T x'_i + b \geq 0\}}$ .

# Decision Tree

## Review

- Model: Decision Tree
- Objective: recursively minimize negative information gain,  
 $H(Y) - H(Y|X_j)$ .
- Training: ID3: Iterative Dichotomiser 3.
- Prediction:  $\hat{y}_i = \text{label of leaf}$ .

# Nearest Neighbor

## Review

- Model: Nearest Neighbor
- Objective: none.
- Training: memorize the data.
- Prediction:  $\hat{y}_i = \text{mode} \{y_{(1)}, y_{(2)}, \dots, y_{(k)}\}$ .

# Feature Construction

## Review

- Each dimension of  $x_i$  is a feature,  $x_{ij}$ .
- Feature selection is choosing important features to use in predictions: logistic regression regularization, decision tree.
- Feature engineering is creating new features for training: kernelized SVM, convolutional network, traditional computer vision SIFT, HOG, Haar features.

# Applications

## Review

- All classification tasks.
- Homework 1: Handwritten character recognition.
- Homework 2: Facial expression classification.
- Homework 3: Movie box office prediction.
- Homework 4: Face detection in images.
- All recommendation systems: Amazon, Facebook, Google, Netflix, YouTube ...
- Face recognition, object detection, self-driving cars, speech recognition, spam filtering, fraud detection, weather forecast, sports team selection, algorithmic trading, market analysis, gene sequence classification, medical diagnosis ...

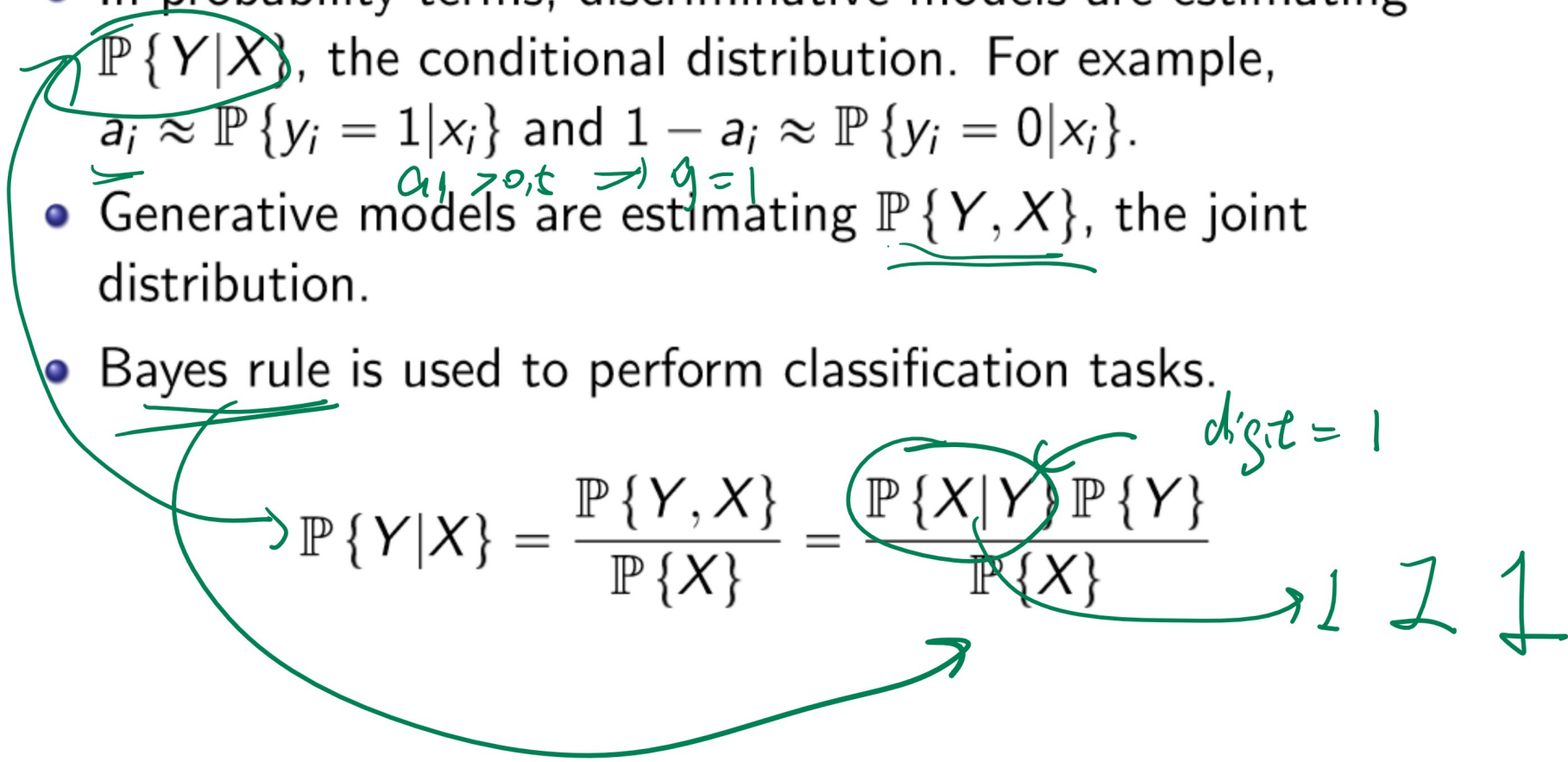
# Generative Models

## Motivation

- In probability terms, discriminative models are estimating  $\mathbb{P}\{Y|X\}$ , the conditional distribution. For example,  $a_i \approx \mathbb{P}\{y_i = 1|x_i\}$  and  $1 - a_i \approx \mathbb{P}\{y_i = 0|x_i\}$ .
- Generative models are estimating  $\mathbb{P}\{Y, X\}$ , the joint distribution.
- Bayes rule is used to perform classification tasks.

$$\mathbb{P}\{Y|X\} = \frac{\mathbb{P}\{Y, X\}}{\mathbb{P}\{X\}} = \frac{\mathbb{P}\{X|Y\}\mathbb{P}\{Y\}}{\mathbb{P}\{X\}}$$

*digit = 1*  
↓ ↓ ↓







# Tokenization

## Motivation

- When processing language, documents (called corpus) need to be turned into a sequence of tokens.

① Split the string by space and punctuations.

② Remove stopwords such as "the", "of", "a", "with" ...

③ Lower case all characters.

④ Stemming or lemmatization words: make "looks", "looked", "looking" to "look".

# Vocabulary

## Motivation

- Word token is an occurrence of a word.
- Word type is a unique token as a dictionary entry.
- Vocabulary is the set of word types.
- Characters can be used in place of words as tokens. In this case, the types are "a", "b", ..., "z", " ", and vocabulary is the alphabet.

HW5

# Bag of Words Features

## Motivation

- Given a document  $i$  and vocabulary with size  $m$ , let  $c_{ij}$  be the count of the word  $j$  in the document  $i$  for  $j = 1, 2, \dots, m$ .
- Bag of words representation of a document has features that are the count of each word divided by the total number of words in the document.

$$x_{ij} = \frac{c_{ij}}{\sum_{j'=1}^m c_{ij'}}$$

fraction of time a word occurs in a doc.

# TF IDF Features

## Motivation

- Another feature representation is called tf-idf, which stands for normalized term frequency, inverse document frequency.

$$tf_{ij} = \frac{c_{ij}}{\max_{j'} c_{ij'}}, \quad idf_j = \log \frac{n}{\sum_{i=1}^n \mathbb{1}_{\{c_{ij} > 0\}}}$$

*Handwritten notes:*  
-  $n$  ← # docs  
-  $\sum_{i=1}^n \mathbb{1}_{\{c_{ij} > 0\}}$  ← # of docs in which  $j$  word appeared in  
-  $x_{ij} = tf_{ij} idf_j$  (underlined)

- $n$  is the total number of documents and  $\sum_{i=1}^n \mathbb{1}_{\{c_{ij} > 0\}}$  is the number of documents containing word  $j$ .

# Bag of Words Features Example

## Motivation

- Given training set, the set of documents is called a corpus. Suppose the set is "I am Groot", "I am Groot", ... (10 times), "We are Groot". The vocabulary is "I" "am" "Groot" "we" "are", then the bag of words features will have the following training set.

I am Groot we are

I am Groot  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 We are Groot

$x_1 =$

$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	0
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	0
...	...	...	...	...
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	0
0	0	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

$x_2 =$

$x_3 =$

# Token Notations

## Definition

- A word (or character) at position  $t$  of a sentence (or string) is denoted as  $z_t$ .
- A sentence (or string) with length  $d$  is  $(z_1, z_2, \dots, z_d)$ .
- $\mathbb{P}\{Z_t = z_t\}$  is the probability of observing  $z_t \in \{1, 2, \dots, j\}$  at position  $t$  of the sentence, usually shortened to  $\mathbb{P}\{z_t\}$ .

$$\begin{aligned} & \text{"I"} \longrightarrow \mathbb{P}_r \{ \text{"I"} \} = 0.3 \\ \mathbb{P}_r \{ z_0 = & \begin{cases} \text{"You"} \longrightarrow \mathbb{P}_r \{ \text{"You"} \} = 0.2 \\ \text{"Great"} \longrightarrow \mathbb{P}_r \{ \text{"Great"} \} = 0 \end{cases} \end{aligned}$$

# Unigram Model

## Definition

- Unigram models assume independence.

$$\mathbb{P}\{z_1, z_2, \dots, z_d\} = \prod_{t=1}^d \mathbb{P}\{z_t\}$$

$P_1(z_1) \cdot P_1(z_2) \cdot \dots$   
 $\cdot P_1(z_d)$

- In general, two events  $A$  and  $B$  are independent if:

$$\mathbb{P}\{A|B\} = \mathbb{P}\{A\} \text{ or } \mathbb{P}\{A, B\} = \mathbb{P}\{A\} \mathbb{P}\{B\}$$

$P_1(A|B) = \frac{P_1(A, B)}{P_1(B)} = P_1(A)$

- For sequence of words, independence means:

$$\mathbb{P}\{z_t | z_{t-1}, z_{t-2}, \dots, z_1\} = \mathbb{P}\{z_t\}$$





# Bigram Model

## Definition

- Bigram models assume Markov property.

$$\mathbb{P}\{z_1, z_2, \dots, z_d\} = \mathbb{P}\{z_1\} \prod_{t=2}^d \mathbb{P}\{z_t | z_{t-1}\}$$

*= P<sub>r</sub>(z<sub>1</sub>) · P<sub>r</sub>(z<sub>2</sub>|z<sub>1</sub>) · P<sub>r</sub>(z<sub>3</sub>|z<sub>2</sub>) ... P<sub>r</sub>(z<sub>d</sub>|z<sub>d-1</sub>)*

- Markov property means the distribution of an element in the sequence only depends on the previous element.

$$\mathbb{P}\{z_t | z_{t-1}, \boxed{z_{t-2}, \dots, z_1}\} = \mathbb{P}\{z_t | z_{t-1}\}$$

*(Note: A green arrow points from the box to the z<sub>t-1</sub> term, and the box is underlined.)*

# Conditional Probability

## Definition

- In general, the conditional probability of an event  $A$  given another event  $B$  is the probability of  $A$  and  $B$  occurring at the same time divided by the probability of event  $B$ .

$$\mathbb{P}\{A|B\} = \frac{\mathbb{P}\{AB\}}{\mathbb{P}\{B\}}$$

- For a sequence of words, the conditional probability of observing  $z_t$  given  $z_{t-1}$  is observed is the probability of observing both divided by the probability of observing  $z_{t-1}$  first.

$$\mathbb{P}\{z_t|z_{t-1}\} = \frac{\mathbb{P}\{z_{t-1}, z_t\}}{\mathbb{P}\{z_{t-1}\}}$$

# Bigram Model Estimation

## Definition

- Using the conditional probability formula,  $\mathbb{P}\{z_t|z_{t-1}\}$ , called transition probabilities, can be estimated by counting all bigrams and unigrams.

$$\hat{\mathbb{P}}\{z_t|z_{t-1}\} = \frac{c_{z_{t-1},z_t}}{c_{z_{t-1}}}$$

*count of #  $z_{t-1}, z_t$*   
*count of #  $z_{t-1}$*

# Unigram MLE Probability

## Quiz (Graded)

Q2

- Given the training data "I am Iron Man", "I love you 3000", "I love you mom", "Tell my family I love them", 18 words in total. With the unigram model, what is the probability of observing a new sentence "I love"?

- A: 0
- B:  $\frac{3}{18}$
- C:  $\frac{3}{4}$
- D:  $\frac{4 \cdot 3}{18 \cdot 4}$
- E:  $\frac{4 \cdot 3}{18 \cdot 18}$

$$P_{\text{r}} \{ \text{"I love"} \} = P_{\text{r}} \{ \text{"I"} \} \cdot P_{\text{r}} \{ \text{"love"} \}$$

MLE

$$= \frac{4}{18} \cdot \frac{3}{18}$$

# Bigram MLE Probability, Part I

## Quiz (Graded)

- Given the training data "I am Iron Man", "I love you 3000", "I love you mom", "Tell my family I love them", 18 words in total. With the bigram model, what is the probability of observing  $Z_2 = \text{"love"}$  given the sentence starts with  $Z_1 = \text{"I"}$ ?

- A: 0
- B:  $\frac{3}{18}$
- C:  $\frac{3}{4}$
- D:  $\frac{4 \cdot 3}{18 \cdot 4}$
- E:  $\frac{4 \cdot 3}{18 \cdot 18}$

$$P(\text{"love"} | \text{"I"}) = \frac{C_{\text{"I love"}}}{C_{\text{"I"}}} = \frac{3}{4}$$

# Bigram MLE Probability, Part II

## Quiz (Graded)

(Q3)

- Given the training data "I am Iron Man", "I love you 3000", "I love you mom", "Tell my family I love them", 18 words in total. With the bigram model, what is the probability of observing a new sentence "I love"?

- A: 0
- B:  $\frac{3}{18}$
- C:  $\frac{3}{4}$
- D:  $\frac{4 \cdot 3}{18 \cdot 4}$
- E:  $\frac{4 \cdot 3}{18 \cdot 18}$

$$\begin{aligned}
 \Pr\{\text{"I love"}\} &= \frac{\Pr\{\text{"I"}\}}{4} \cdot \frac{\Pr\{\text{"love"}|\text{"I"}\}}{3} \\
 &= \frac{4}{18} \cdot \frac{3}{4}
 \end{aligned}$$

# Transition Matrix

## Definition

- These probabilities can be stored in a matrix called transition matrix of a Markov Chain. The number on row  $j$  column  $j'$  is the estimated probability  $\hat{\mathbb{P}}\{j'|j\}$ . If there are 3 tokens  $\{1, 2, 3\}$ , the transition matrix is the following.

$$\begin{bmatrix} \hat{\mathbb{P}}\{1|1\} & \hat{\mathbb{P}}\{2|1\} & \hat{\mathbb{P}}\{3|1\} \\ \hat{\mathbb{P}}\{1|2\} & \hat{\mathbb{P}}\{2|2\} & \hat{\mathbb{P}}\{3|2\} \\ \hat{\mathbb{P}}\{1|3\} & \hat{\mathbb{P}}\{2|3\} & \hat{\mathbb{P}}\{3|3\} \end{bmatrix}$$

- Given the initial distribution of tokens, the distribution of the next token can be found by multiplying it by the transition probabilities.



# Estimating Transition Matrix

## Definition

Suppose the vocabulary is "I", "am", "Groot", "we", "are", and the training set contains 10 "I am Groot" then 1 "We are Groot".

Then the transition matrix is:  $\rightarrow$  I am Groot | I am Groot ... We are Groot

—	I	am	Groot	we	are
I	$\frac{0+1}{10+5}$ 0	$\frac{10+1}{10+5}$ 1	0	$\frac{1}{15}$	0
am	0	0	1	0	0
Groot	0.9	0	0	0.1	0
we	0	0	0	0	1
are	0	0	1	0	0

$C_I \rightarrow$

$\hat{P}_c \{ \text{"am"} | \text{"I"} \}$

Laplace smooth

# Trigram Model

## Definition

- The same formula can be applied to trigram: sequences of three tokens.

$$\hat{\mathbb{P}} \{z_t | z_{t-1}, z_{t-2}\} = \frac{C_{z_{t-2}, z_{t-1}, z_t}}{C_{z_{t-2}, z_{t-1}}}$$

- In a document, it is likely that these longer sequences of tokens never appear. In those cases, the probabilities are  $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$ . Because of this, Laplace smoothing adds 1 to all counts.

$$\hat{\mathbb{P}} \{z_t | z_{t-1}, z_{t-2}\} = \frac{C_{z_{t-2}, z_{t-1}, z_t} + 1}{C_{z_{t-2}, z_{t-1}} + m}$$

← add 1 to each count.  
 → size of vocabulary

# Laplace Smoothing

## Definition

- Laplace smoothing should be used for bigram and unigram models too.

$$\hat{\mathbb{P}} \{z_t | z_{t-1}\} = \frac{c_{z_{t-1}, z_t} + 1}{c_{z_{t-1}} + m} \quad \leftarrow$$

$$\hat{\mathbb{P}} \{z_t\} = \frac{c_{z_t} + 1}{\sum_{z=1}^m c_z + m} \quad \leftarrow$$

- Aside: Laplace smoothing can also be used in decision tree training to compute entropy.

# Smoothing

## Quiz (Graded)

- Fall 2018 Midterm Q12.
- Given a vocabulary of  $10^6$ , a document with  $10^{12}$  tokens with  $C_{\text{zoodles}} = 3$ . What is the MLE estimation of  $P\{\text{zoodles}\}$  with and without Laplace smoothing? (choose 2)

Q5

- ~~A:  $\frac{3}{10^{12}}$~~
- B:  $\frac{3}{10^6}$
- C:  $\frac{3+1}{10^{12}+3}$
- ~~D:  $\frac{3+1}{10^{12}+10^6}$~~
- E:  $\frac{3+1}{10^{12}+10^6-1}$

$P_i\{\text{zoodles}\}$

$$\frac{C_{\text{zoodle}} (+1)}{\sum_{\text{words}} C_{\text{words}} (+m)}$$

↓


$$\frac{3+1}{10^{12}+10^6}$$

$\frac{3}{10^2}$

# N Gram Model

## Algorithm

- Input: series  $\{z_1, z_2, \dots, z_{d_i}\}_{i=1}^n$ .
- Output: transition probabilities  $\hat{\mathbb{P}}\{z_t | z_{t-1}, z_{t-2}, \dots, z_{t-N+1}\}$  for all  $z_t = 1, 2, \dots, m$ .
- Compute the transition probabilities using counts and Laplace smoothing.

$$\hat{\mathbb{P}}\{z_t | z_{t-1}, z_{t-2}, \dots, z_{t-N+1}\} = \frac{C_{z_{t-N+1}, z_{t-N+2}, \dots, z_t} + 1}{C_{z_{t-N+1}, z_{t-N+2}, \dots, z_{t-1}} + m}$$


# Sampling from Discrete Distribution

## Discussion

- In order to generate new sentences given an  $N$  gram model, random realizations need to be generated given the conditional probability distribution.
- Given the first  $N - 1$  words,  $z_1, z_2, \dots, z_{N-1}$ , the distribution of next word is approximated by  $p_x = \hat{\mathbb{P}} \{z_N = x | z_{N-1}, z_{N-2}, \dots, z_1\}$ . This process then can be repeated for on  $z_2, z_3, \dots, z_{N-1}, z_N$  and so on.

# Cumulative Distribution Inversion Method, Part I

## Discussion

- Most programming languages have a function to generate a random number  $u \sim \text{Unif}[0, 1]$ .
- If there are  $m = 2$  tokens in total and the conditional probabilities are  $p$  and  $1 - p$ . Then the following distributions are the same.

$$z_N = \begin{cases} 0 & \text{with probability } p \\ 1 & \text{with probability } 1 - p \end{cases} \Leftrightarrow z_N = \begin{cases} 0 & \text{if } 0 \leq u \leq p \\ 1 & \text{if } p < u \leq 1 \end{cases}$$

# Cumulative Distribution Inversion Method, Part II

## Discussion

- In the general case with  $m$  tokens with conditional probabilities  $p_1, p_2, \dots, p_m$  with  $\sum_{j=1}^m p_j = 1$ . Then the following distributions are the same.

$$z_N = j \text{ with probability } p_j \Leftrightarrow z_N = j \text{ if } \sum_{j'=1}^{j-1} p_{j'} < u \leq \sum_{j'=1}^j p_{j'}$$

- This can be used to generate a random token from the conditional distribution.





