Nearest Neighbor
○○○○○○○○○○○○○

Natural Language Processing
○○○○○○○○○○○○○○○○○○○○○

Sampling
○○○○

# CS540 Introduction to Artificial Intelligence
## Lecture 7

### Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

June 20, 2022

# Lecture Feedback, Additional Examples, Solutions

Admin

# Midterm Details

## Admin

# Midterm Coverage

## Admin

# $K$ Nearest Neighbor

### Description

- Given a new instance, find the $K$ instances in the training set that are the closest.
- Predict the label of the new instance by the majority of the labels of the $K$ instances.

Nearest Neighbor
0000●00000000

Natural Language Processing
0000000000000000000000

Sampling
0000

# Distance Function
### Definition

- Many distance functions can be used in place of the Euclidean distance.

$$\rho\left(x, x'\right) = \left\|x - x'\right\|_2 = \sqrt{\sum_{j=1}^{m}\left(x_j - x_j'\right)^2}$$

- An example is Manhattan distance.

$$\rho\left(x, x'\right) = \sum_{j=1}^{m}\left|x_j - x_j'\right|$$

# Manhattan Distance Diagram

### Definition

# 1 Nearest Neighbor

Quiz

Nearest Neighbor

Natural Language Processing

Sampling

○○○○○○○●○○○○○

○○○○○○○○○○○○○○○○○○○○○○

○○○○

# 3 Nearest Neighbor

Quiz

# K Fold Cross Validation

Discussion

Nearest Neighbor

○○○○○○○○○●○○○

Natural Language Processing

○○○○○○○○○○○○○○○○○○○○○

Sampling

○○○○

# 5 Fold Cross Validation Example

Discussion

# Leave One Out Cross Validation

Discussion

# Cross Validation

Quiz

# Cross Validation 2

## Quiz

Nearest Neighbor
○○○○○○○○○○○○○

Natural Language Processing
●○○○○○○○○○○○○○○○○○○○○○

Sampling
○○○○

# 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".

Nearest Neighbor
ooooooooooooo

Natural Language Processing
oooooooooooooooooooo

Sampling
oooo

# 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.

# Bag of Words Features

### Definition

- Given a document $i$ and vocabulary with size $m$, let $c_{ij}$ be the count of the word $j$ in the document $i$ for $j = 1, 2, ..., m$.

- 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.

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

# Bag of Words Features Example

## Motivation

# TF IDF Features

Definition

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

$$\text{tf}_{ij} = \frac{c_{ij}}{\max_{j'} c_{ij'}}, \ \text{idf}_j = \log \frac{n}{\sum_{i=1}^{n} \mathbb{1}_{\{c_{ij} > 0\}}}$$

$$x_{ij} = \text{tf}_{ij} \, \text{idf}_j$$

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

# Unigram Model
Definition

- Unigram models assume independence.
$$\mathbb{P}\left\{z_1, z_2, ..., z_d\right\} = \prod_{t=1}^{d} \mathbb{P}\left\{z_t\right\}$$

- In general, two events $A$ and $B$ are independent if:
$$\mathbb{P}\left\{A|B\right\} = \mathbb{P}\left\{A\right\} \text{ or } \mathbb{P}\left\{A, B\right\} = \mathbb{P}\left\{A\right\}\mathbb{P}\left\{B\right\}$$

- For a sequence of words, independence means:
$$\mathbb{P}\left\{z_t|z_{t-1}, z_{t-2}, ..., z_1\right\} = \mathbb{P}\left\{z_t\right\}$$

# Maximum Likelihood Estimation
### Definition

- $\mathbb{P}\{z_t\}$ can be estimated by the count of the word $z_t$.
$$\hat{\mathbb{P}}\{z_t\} = \frac{c_{z_t}}{\displaystyle\sum_{z=1}^{m} c_z}$$

- This is called the maximum likelihood estimator because it maximizes the probability of observing the sentences in the training set.

Nearest Neighbor
○○○○○○○○○○○○○

Natural Language Processing
○○○○○○○●○○○○○○○○○○○○○○

Sampling
○○○○

# Bigram Model

### Definition

- Bigram models assume Markov property.

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

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

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

# Markov Chain Demo

## Motivation

# 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.

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

# Unigram MLE Probability

Quiz

# Bigram MLE Probability

Quiz

# Unigram MLE Probability
## Quiz

# Bigram MLE Probability

Quiz

Nearest Neighbor
000000000000

Natural Language Processing
00000000000000000000000

Sampling
0000

# 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

Nearest Neighbor
○○○○○○○○○○○○○

Natural Language Processing
○○○○○○○○○○○○○○○○○●○○○○

Sampling
○○○○

# Trigram Model
### Definition

- The same formula can be applied to trigram: sequences of three tokens.
$$\hat{\mathbb{P}}\left\{z_t|z_{t-1}, z_{t-2}\right\} = \frac{c_{z_{t-2}, z_{t-1}, z_t}}{c_{z_{t-2}, z_{t-1}}}$$

- In a document, likely, these longer sequences of tokens never appear. In those cases, the probabilities are $\frac{0}{0}$. Because of this, Laplace smoothing adds 1 to all counts.
$$\hat{\mathbb{P}}\left\{z_t|z_{t-1}, z_{t-2}\right\} = \frac{c_{z_{t-2}, z_{t-1}, z_t} + 1}{c_{z_{t-2}, z_{t-1}} + m}$$

Nearest Neighbor
○○○○○○○○○○○○○

Natural Language Processing
○○○○○○○○○○○○○○○○○○●○○○

Sampling
○○○○

# Laplace Smoothing

Definition

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

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

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

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

# Smoothing Example

## Quiz

# Smoothing Example 2

## Quiz

# Smoothing Example 3

Quiz

Nearest Neighbor
0000000000000

Natural Language Processing
0000000000000000000000

Sampling
●000

# Sampling from Discrete Distribution
### Discussion

- 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, $z_1, z_2, ..., z_{N-1}$, the distribution of next word is approximated by $p_x = \hat{\mathbb{P}} \{z_N = x | z_{N-1}, z_{N-2}, ..., z_1\}$. This process then can be repeated for on $z_2, z_3, ..., z_{N-1}, z_N$ and so on.

# CDF Inversion Method Diagram

Discussion

Nearest Neighbor

○○○○○○○○○○○○○

Natural Language Processing

○○○○○○○○○○○○○○○○○○○○○○

Sampling

○○●○

# Generating New Words 1

Quiz

# Generating New Words 2

## Quiz