

# Natural Language Processing Basics

Yingyu Liang

University of Wisconsin-Madison

# Natural language Processing (NLP)

- The processing of the **human** languages by **computers**
- One of the oldest AI tasks
- One of the most important AI tasks
- One of the hottest AI tasks nowadays

# Difficulty


- Difficulty 1: ambiguous, typically no formal description
- Example: “*We saw her duck.*”

How many different meanings?

# Difficulty

- Difficulty 1: ambiguous, typically no formal description
- Example: “*We saw her duck.*”
- 1. We looked at a duck that belonged to her.
- 2. We looked at her quickly squat down to avoid something.
- 3. We use a saw to cut her duck.

# Difficulty

- Difficulty 2: computers do not have human concepts
- Example: “She like *little animals*. For example, yesterday we saw her *duck*.”
- 1. We looked at a duck that belonged to her.
- 2. We looked at her quickly squat down to avoid something.
- 3. We use a saw to cut her duck.

# Words

Preprocess

Zipf's Law

# Preprocess

- Corpus: often a set of text documents
  - Tokenization or text normalization: turn corpus into sequence(s) of tokens
1. Remove unwanted stuff: HTML tags, encoding tags
  2. Determine word boundaries: usually white space and punctuations
    - Sometimes can be tricky, like Ph.D.

# Preprocess

- Tokenization or text normalization: turn data into sequence(s) of tokens
  1. Remove unwanted stuff: HTML tags, encoding tags
  2. Determine word boundaries: usually white space and punctuations
    - Sometimes can be tricky, like Ph.D.
  3. Remove stopwords: the, of, a, with, ...



# Preprocess

- Tokenization or text normalization: turn data into sequence(s) of tokens
  1. Remove unwanted stuff: HTML tags, encoding tags
  2. Determine word boundaries: usually white space and punctuations
    - Sometimes can be tricky, like Ph.D.
  3. Remove stopwords: the, of, a, with, ...
  4. Case folding: lower-case all characters.
    - Sometimes can be tricky, like US and us
  5. Stemming/Lemmatization (optional): looks, looked, looking → look

# Vocabulary

Given the preprocessed text

- Word token: occurrences of a word
- Word type: unique word as a dictionary entry (i.e., unique tokens)
- Vocabulary: the set of word types
  - Often 10k to 1 million on different corpora
  - Often remove too rare words

# Zipf's Law

- Word count  $f$ , word rank  $r$
- Zipf's law:  $f * r \approx \text{constant}$

| Word       | Count $f$ | rank $r$ | $fr$  |
|------------|-----------|----------|-------|
| the        | 3332      | 1        | 3332  |
| and        | 2972      | 2        | 5944  |
| a          | 1775      | 3        | 5235  |
| he         | 877       | 10       | 8770  |
| but        | 410       | 20       | 8400  |
| be         | 294       | 30       | 8820  |
| there      | 222       | 40       | 8880  |
| one        | 172       | 50       | 8600  |
| two        | 104       | 100      | 10400 |
| turned     | 51        | 200      | 10200 |
| comes      | 16        | 500      | 8000  |
| family     | 8         | 1000     | 8000  |
| brushed    | 4         | 2000     | 8000  |
| Could      | 2         | 4000     | 8000  |
| Applausive | 1         | 8000     | 8000  |

Zipf's law on the corpus *Tom Sawyer*

# Text: Bag-of-Words Representation

Bag-of-Words

tf-idf

# Bag-of-Words

How to represent a piece of text (sentence/document) as numbers?

- Let  $m$  denote the size of the vocabulary
- Given a document  $d$ , let  $c(w, d)$  denote the #occurrence of  $w$  in  $d$
- Bag-of-Words representation of the document

$$v_d = [c(w_1, d), c(w_2, d), \dots, c(w_m, d)]/Z_d$$

- Often  $Z_d = \sum_w c(w, d)$

# Example

- Preprocessed text: *this is a good sentence this is another good sentence*
- BoW representation:  
 $[c('a', d)/Z_d, c('is', d)/Z_d, \dots, c('example', d)/Z_d]$
- What is  $Z_d$ ?
- What is  $c('a', d)/Z_d$ ?
- What is  $c('example', d)/Z_d$ ?

# tf-idf

- tf: normalized term frequency

$$tf_w = \frac{c(w, d)}{\max_v c(v, d)}$$

- idf: inverse document frequency

$$idf_w = \log \frac{\text{total \#documents}}{\text{\#documents containing } w}$$

- tf-idf:  $tf-idf_w = tf_w * idf_w$

- Representation of the document

$$v_d = [tf-idf_{w_1}, tf-idf_{w_2}, \dots, tf-idf_{w_m}]$$

# Cosine Similarity

How to measure similarities between pieces of text?

- Given the document vectors, can use any similarity notion on vectors
- Commonly used in NLP: cosine of the angle between the two vectors

$$\text{sim}(x, y) = \frac{x^T y}{\sqrt{x^T x} \sqrt{y^T y}}$$



# Text: statistical Language Model

Statistical language model

N-gram

Smoothing

# Probabilistic view

- Use probabilistic distribution to model the language
- Dates back to Shannon (information theory; bits in the message)

# Statistical language model

- Language model: probability distribution over sequences of tokens
- Typically, tokens are words, and distribution is discrete
- Tokens can also be characters or even bytes
- Sentence: *“the quick brown fox jumps over the lazy dog”*

Tokens:  $x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$   $x_7$   $x_8$   $x_9$

# Statistical language model

- For simplification, consider fixed length sequence of tokens (sentence)

$$(x_1, x_2, x_3, \dots, x_{\tau-1}, x_{\tau})$$

- Probabilistic model:

$$P [x_1, x_2, x_3, \dots, x_{\tau-1}, x_{\tau}]$$

# Unigram model

- Unigram model: define the probability of the sequence as the product of the probabilities of the tokens in the sequence

$$P[x_1, x_2, \dots, x_\tau] = \prod_{t=1}^{\tau} P[x_t]$$

- Independence!

# A simple unigram example

- Sentence: “*the dog ran away*”

$$\hat{P}[\textit{the dog ran away}] = \hat{P}[\textit{the}] \hat{P}[\textit{dog}] \hat{P}[\textit{ran}] \hat{P}[\textit{away}]$$

- How to estimate  $\hat{P}[\textit{the}]$  on the training corpus?

# A simple unigram example

- Sentence: “*the dog ran away*”

$$\hat{P}[\textit{the dog ran away}] = \hat{P}[\textit{the}] \hat{P}[\textit{dog}] \hat{P}[\textit{ran}] \hat{P}[\textit{away}]$$

- How to estimate  $\hat{P}[\textit{the}]$  on the training corpus?

| Word  | Count $f$ |
|-------|-----------|
| the   | 3332      |
| and   | 2972      |
| a     | 1775      |
| he    | 877       |
| but   | 410       |
| be    | 294       |
| there | 222       |
| one   | 172       |

# n-gram model

- $n$ -gram: sequence of  $n$  tokens
- $n$ -gram model: define the conditional probability of the  $n$ -th token given the preceding  $n - 1$  tokens

$$P[x_1, x_2, \dots, x_\tau] = P[x_1, \dots, x_{n-1}] \prod_{t=n}^{\tau} P[x_t | x_{t-n+1}, \dots, x_{t-1}]$$



# n-gram model

- $n$ -gram: sequence of  $n$  tokens
- $n$ -gram model: define the conditional probability of the  $n$ -th token given the preceding  $n - 1$  tokens

$$P[x_1, x_2, \dots, x_\tau] = P[x_1, \dots, x_{n-1}] \prod_{t=n}^{\tau} P[x_t | x_{t-n+1}, \dots, x_{t-1}]$$

Markovian assumptions



# Typical $n$ -gram model

- $n = 1$ : unigram
- $n = 2$ : bigram
- $n = 3$ : trigram

# Training $n$ -gram model

- Straightforward counting: counting the co-occurrence of the grams

For all grams  $(x_{t-n+1}, \dots, x_{t-1}, x_t)$

1. count and estimate  $\hat{P}[x_{t-n+1}, \dots, x_{t-1}, x_t]$
2. count and estimate  $\hat{P}[x_{t-n+1}, \dots, x_{t-1}]$
3. compute

$$\hat{P}[x_t | x_{t-n+1}, \dots, x_{t-1}] = \frac{\hat{P}[x_{t-n+1}, \dots, x_{t-1}, x_t]}{\hat{P}[x_{t-n+1}, \dots, x_{t-1}]}$$

# A simple trigram example

- Sentence: “*the dog ran away*”

$$\hat{P}[\textit{the dog ran away}] = \hat{P}[\textit{the dog ran}] \hat{P}[\textit{away}|\textit{dog ran}]$$

$$\hat{P}[\textit{the dog ran away}] = \hat{P}[\textit{the dog ran}] \frac{\hat{P}[\textit{dog ran away}]}{\hat{P}[\textit{dog ran}]}$$

# Drawback

- Sparsity issue:  $\hat{P}[\dots]$  most likely to be 0
- Bad case: “*dog ran away*” never appear in the training corpus, so  $\hat{P}[\textit{dog ran away}] = 0$
- Even worse: “*dog ran*” never appear in the training corpus, so  $\hat{P}[\textit{dog ran}] = 0$

# Rectify: smoothing

- Basic method: adding non-zero probability mass to zero entries
- Example: **Laplace smoothing** that adds one count to all  $n$ -grams  
 $\text{pseudocount}[\textit{dog}] = \text{actualcount}[\textit{dog}] + 1$

# Rectify: smoothing

- Basic method: adding non-zero probability mass to zero entries
- Example: **Laplace smoothing** that adds one count to all  $n$ -grams  
 $\text{pseudocount}[dog] = \text{actualcount}[dog] + 1$

$$\hat{P}[dog] = \frac{\text{pseudocount}[dog]}{\text{pseudo length of the corpus}} = \frac{\text{pseudocount}[dog]}{\text{actual length of the corpus} + |V|}$$

# Rectify: smoothing

- Basic method: adding non-zero probability mass to zero entries
- Example: **Laplace smoothing** that adds one count to all  $n$ -grams  
pseudocount[*dog ran away*] = actualcount[*dog ran away*] + 1  
pseudocount[*dog ran*] = ?



# Rectify: smoothing

- Basic method: adding non-zero probability mass to zero entries
- Example: **Laplace smoothing** that adds one count to all  $n$ -grams  
pseudocount[*dog ran away*] = actualcount[*dog ran away*] + 1  
pseudocount[*dog ran*] = actualcount[*dog ran*] +  $|V|$

$$\hat{P}[\textit{away}|\textit{dog ran}] \approx \frac{\text{pseudocount}[\textit{dog ran away}]}{\text{pseudocount}[\textit{dog ran}]}$$

since #bigrams  $\approx$  #trigrams on the corpus

# Example

- Preprocessed text: *this is a good sentence this is another good sentence*
- How many unigrams?
- How many bigrams?
- Estimate  $\hat{P}[is|this]$  without using Laplace smoothing
- Estimate  $\hat{P}[is|this]$  using Laplace smoothing ( $|V| = 10000$ )

# Rectify: smoothing

- Basic method: adding non-zero probability mass to zero entries
  - Example: Laplace smoothing
- Back-off methods: restore to lower order statistics
  - Example: if  $\hat{P}[\textit{away}|\textit{dog ran}]$  does not work, use  $\hat{P}[\textit{away}|\textit{ran}]$  as replacement
- Mixture methods: use a linear combination of  $\hat{P}[\textit{away}|\textit{ran}]$  and  $\hat{P}[\textit{away}|\textit{dog ran}]$

# Another drawback

- High dimension: # of grams too large
- Vocabulary size: about  $10^k = 2^{14}$
- #trigram: about  $2^{42}$

# Rectify: clustering

- Class-based language models: cluster tokens into classes; replace each token with its class
- Significantly reduces the vocabulary size; also address sparsity issue
- Combinations of smoothing and clustering are also possible