

# Outline

- Natural Language Processing
  - Preprocessing
  - Statistics
  - Language models

# 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*

# 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  $Z_d = \sum_w c(w, d)$

# 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

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

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

- 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

Independence!

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

- 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

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

Markovian assumptions



- $n = 1$ : 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

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

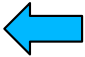
- Sentence: “*the dog ran away*” by trigram ( $n=3$ )

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

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

## Rectify: smoothing

- Sparsity issue:  $\hat{P}[\dots]$  most likely to be 0
- 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  trigram

pseudocount[*dog ran*]  $\approx$  actualcount[*dog ran*] +  $|V|$

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

Since number of bigrams "dog ran"  $\approx$  Number of trigrams containing "dog ran"