

UNIVERSITY of DELAWARE

# Statistical Language Models

(Search and Data Mining)

Hui Fang  
Department of Electrical and Computer Engineering  
University of Delaware

1



UNIVERSITY of DELAWARE

## Statistical Language Models

- A statistical language model is a probability distribution over word sequences.
- Given a sequence of words, it computes the probability of a sentence of words:
 
$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$
- Examples
 
$$P(\text{Today is Thursday}) = 0.01$$

$$P(\text{Today Thursday is}) = 0.000000000000000001$$


2


UNIVERSITY of DELAWARE

## Language models are useful in many applications

- Text categorization
  - $P(\text{topic} = \text{sports} \mid \text{baseball, baseball, baseball, game})$
- Speech recognition
  - $P(\text{happy} \mid \text{John feels}) \gg P(\text{habit} \mid \text{John feels})$
- Information retrieval
  - $P(\text{an article about dog} \mid \text{puppy}) > P(\text{an article about cat} \mid \text{puppy})$


3


UNIVERSITY of DELAWARE

## Go back to the problem of computing joint probability of words in sentence

$$P(W) = P(w_1 w_2 \dots w_n) = \prod_i P(w_i \mid w_1 w_2 \dots w_{i-1})$$


$P(\text{"its water is so transparent"}) =$   
 $P(\text{its}) \times P(\text{water} \mid \text{its}) \times P(\text{is} \mid \text{its water})$   
 $\times P(\text{so} \mid \text{its water is}) \times P(\text{transparent} \mid \text{its water is so})$


UNIVERSITY of DELAWARE

## Simplest case: Unigram model

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

- Assumption:
  - Each word is generated independently.
- Parameters:
  - $\{p(w_i)\}$
  - $p(w_1) + \dots + p(w_N) = 1$  (N is voc. size)
- This is essentially a multinomial distribution over words
- A piece of text can be regarded as a sample drawn according to this word distribution.


UNIVERSITY of DELAWARE

## Text Generation with Unigram LM

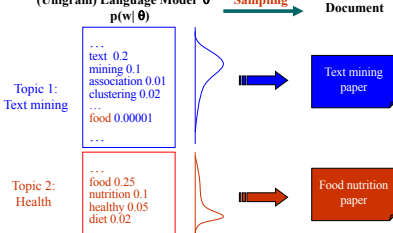
(Unigram) Language Model  $\theta$   $p(w \mid \theta)$  → Sampling Document

Topic 1: Text mining

- ... text 0.2
- ... mining 0.1
- ... association 0.01
- ... clustering 0.02
- ... food 0.00001
- ...

Topic 2: Health

- ... food 0.25
- ... nutrition 0.1
- ... healthy 0.05
- ... diet 0.02
- ...



7

UNIVERSITY of DELAWARE

### Estimation of Unigram LM

(Unigram) Language Model  $\theta$   $p(w_i | \theta) = ?$  ← Estimation Document

10/100 text ?  
5/100 mining ?  
3/100 association ?  
3/100 database ?  
1/100 query ?

text 10 mining 5  
association 3  
database 3  
algorithm 2  
query 1  
efficient 1

A "text mining paper"  
(total #words=100)

8

UNIVERSITY of DELAWARE

### Maximum Likelihood (ML) Estimation

- Finding parameters that yield to maximum likelihood

$$\hat{\theta} = \arg \max_{\theta} P(X | \theta)$$

9

UNIVERSITY of DELAWARE

### Maximum Likelihood Estimation

Data: a document  $d$  with counts  $c(w_1, d), \dots, c(w_N, d)$ , and length  $|d|$   
 Model: multinomial distribution  $\theta$  with parameters  $\{\theta_i = p(w_i)\}$   $\sum_{i=1}^N \theta_i = 1$   
 Likelihood:  $p(d | \theta)$   
 Maximum likelihood estimator:  $\hat{\theta} = \arg \max_{\theta} P(d | \theta)$

$$\hat{\theta} = \arg \max_{\theta} P(d | \theta) = \arg \max_{\theta} \log P(d | \theta)$$

$$\log P(d | \theta) = \log \left( \prod_{i=1}^N P(w_i | \theta) \right) = \sum_{i=1}^N c(w_i, d) \log P(w_i | \theta) = \sum_{i=1}^N c(w_i, d) \log \theta_i$$

Subject to constraints:  $\sum_{i=1}^N \theta_i = 1$

Use Lagrange multiplier approach  $f(\theta | d) = \sum_{i=1}^M c(w_i, d) \log \theta_i + \lambda \left( \sum_{i=1}^M \theta_i - 1 \right)$

11

UNIVERSITY of DELAWARE

### Maximum Likelihood Estimation

$$\hat{\theta} = \arg \max_{\theta} P(d | \theta) = \arg \max_{\theta} \log P(d | \theta)$$

$$\log P(d | \theta) = \log \left( \prod_{i=1}^N P(w_i | \theta) \right) = \sum_{i=1}^N c(w_i, d) \log P(w_i | \theta) = \sum_{i=1}^N c(w_i, d) \log \theta_i$$

Subject to constraints:  $\sum_{i=1}^N \theta_i = 1$

Use Lagrange multiplier approach  $f(\theta | d) = \sum_{i=1}^M c(w_i, d) \log \theta_i + \lambda \left( \sum_{i=1}^M \theta_i - 1 \right)$

Set partial derivatives to zero  $\frac{\partial f(\theta | d)}{\partial \theta_i} = \frac{c(w_i, d)}{\theta_i} + \lambda = 0 \rightarrow \theta_i = -\frac{c(w_i, d)}{\lambda}$

$$\sum_{i=1}^M -\frac{c(w_i, d)}{\lambda} = 1 \rightarrow \lambda = -\sum_{i=1}^M c(w_i, d) \rightarrow \hat{\theta}_i = p(w_i | \hat{\theta}) = \frac{c(w_i, d)}{\sum_{i=1}^M c(w_i, d)} = \frac{c(w_i, d)}{|d|}$$

ML estimator

12

UNIVERSITY of DELAWARE

### Problem with the ML Estimator

- What probability should we give a word that has not been observed?
- If we want to assign non-zero probabilities to unseen words, we need have to discount the probabilities of seen words
- This is what "smoothing" is about ...

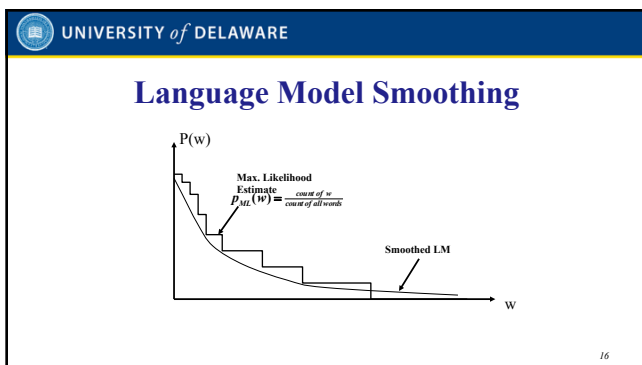
13

UNIVERSITY of DELAWARE

### How to Smooth?

- All smoothing methods try to
  - discount the probability of words seen in a document
  - re-allocate the extra counts so that unseen words will have a non-zero count

14



UNIVERSITY of DELAWARE

## Smoothing Method 1 - Additive Smoothing

- Add a constant  $\delta$  to the counts of each word

$$p(w|d) = \frac{\text{Counts of } w \text{ in } d + 1}{|d| + |V|}$$

↑  
Length of d (total counts)

↑  
Vocabulary size

“Add one”, Laplace smoothing

- Any problems?

17

UNIVERSITY of DELAWARE

## Using Reference Model for Smoothing

- Should all unseen words get equal probabilities?
- We can use a reference model to discriminate unseen words

$$p(w|d) = \begin{cases} p_{\text{seen}}(w|d) & \text{if } w \text{ is seen in } d \\ \alpha_d p(w|REF) & \text{otherwise} \end{cases}$$

Discounted ML estimate

Reference language model

$$\alpha_d = \frac{1 - \sum_{w \in \text{seen}} p_{\text{seen}}(w|d)}{\sum_{w \in \text{unseen}} p(w|REF)}$$

18

UNIVERSITY of DELAWARE

## Smoothing Method 2 - Jelinek Mercer (JM) Smoothing

$$p(w|d) = (1 - \lambda) \frac{c(w,d)}{|d|} + \lambda p(w|REF)$$

ML estimate

parameter

19

UNIVERSITY of DELAWARE

## Smoothing Method 3 - Dirichlet Smoothing

- Assume pseudo counts  $\mu p(w|REF)$

$$p(w|d) = \frac{c(w,d) + \mu p(w|REF)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w,d)}{|d|} + \frac{\mu}{|d| + \mu} p(w|REF)$$

parameter

20