

Special Topics in Natural Language Processing

CS6980

Ashutosh Modi
CSE Department, IIT Kanpur



Lecture 4: Language Models 1
Jan 10, 2020

PROBABILITY REVIEW



Probability Review

Bayes Rule:

$$\begin{aligned} P(X, Y) &= P(Y \mid X)P(X) \\ &= P(X \mid Y)P(Y) \end{aligned}$$

$$P(Y \mid X) = \frac{P(X, Y)}{P(X)}$$

$$P(X \mid Y) = \frac{P(X, Y)}{P(Y)}$$



Probability Review

Bayes Rule:

$$\begin{aligned} P(X, Y) &= P(Y \mid X)P(X) \\ &= P(X \mid Y)P(Y) \end{aligned}$$

$$P(Y \mid X) = \frac{P(X, Y)}{P(X)}$$

$$P(X \mid Y) = \frac{P(X, Y)}{P(Y)}$$

Marginalization:

$$P(X) = \sum_Y P(X, Y)$$

$$P(Y) = \sum_X P(X, Y)$$



Probability Review

Bayes Rule:

$$\begin{aligned} P(X, Y) &= P(Y | X)P(X) \\ &= P(X | Y)P(Y) \end{aligned}$$

$$P(Y | X) = \frac{P(X, Y)}{P(X)}$$

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Marginalization:

$$P(X) = \sum_Y P(X, Y)$$

$$P(Y) = \sum_X P(X, Y)$$

Independence:

Given Joint Distribution:

$$P(Y_1, Y_2, \dots, Y_N)$$

Suppose, Y_1 is independent of $\{Y_3, \dots, Y_N\}$

$$\implies P(Y_1 | Y_2, Y_3, \dots, Y_N) = P(Y_1 | Y_2)$$



Probability Review

Chain Rule:

$$p(X_1, X_2, \dots, X_n) = p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1)$$



Probability Review

Chain + Independence Rule:

Suppose, R.V.s occur in a sequence: X_1, X_2, \dots, X_n ,

Assume that every R.V. only depends on the previous one in the sequence then,

$$p(X_1, X_2, \dots, X_n) = ?$$



Probability Review

Chain + Independence Rule:

Suppose, R.V.s occur in a sequence: X_1, X_2, \dots, X_n ,

Assume that every R.V. only depends on the previous one in the sequence then,

$$p(X_1, X_2, \dots, X_n) = p(X_n | X_{n-1}) \times p(X_{n-1} | X_{n-2}) \times \dots \\ p(X_3 | X_2) \times p(X_2 | X_1) \times p(X_1)$$



LANGUAGE MODELS

Based on Michael Collin's NLP Lecture Notes: <http://www.cs.columbia.edu/~mcollins/lm-spring2013.pdf>



Guess the next word

Lake Zurich is well known, its water is so transparent that _____



Guess the next word

Lake Zurich is well known, its water is so transparent that _____

The dog ran after the _____



Guess the next word

Lake Zurich is well known, its water is so transparent that _____

The dog ran after the _____

The book is in the _____



Guess the next word

Lake Zurich is well known, its water is so transparent that you

The dog ran after the mouse

The book is in the car



Guess the next word

Lake Zurich is well known, its water is so transparent that you

The dog ran after the mouse

The book is in the car

Language Model (LM) at work

$$p(w_k \mid w_{k-1}, w_{k-2}, \dots, w_1)$$



Terminology

Vocabulary/Lexicon: Set of all tokens/words in the language

$$\mathcal{V} = \{\text{apple, banana, the, . . . , zebra}\}$$

Sentence: Sequence of words

$$\mathcal{S} = w_1 w_2 w_3 \dots w_n$$

where $w_i \in \mathcal{V} \quad \forall i = 1, \dots, n - 1; \quad n \geq 1$

and $w_n = STOP$



Terminology

Sentence: Sequence of words

$$\mathcal{S} = w_1 w_2 w_3 \dots w_n$$

where $w_i \in \mathcal{V}$ $\forall i = 1, \dots, n - 1;$ $n \geq 1$

and $w_n = STOP$

the dog ran STOP

the cat is brown STOP

the class is tomorrow STOP

the the the STOP

apple STOP

STOP

.....



Terminology

Sentence: Sequence of words

$$\mathcal{S} = w_1 w_2 w_3 \dots w_n$$

where $w_i \in \mathcal{V}$ $\forall i = 1, \dots, n - 1$; $n \geq 1$

and $w_n = STOP$

the dog ran STOP

the cat is brown STOP

the class is tomorrow STOP

the the the STOP

apple STOP

STOP

.....

Sentence Set:

$$\mathcal{V}^\dagger = \{\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_3, \dots\}$$



Language Model (LM)

LM Definition: Probability Distribution over sentences in \mathcal{V}^\dagger



Language Model (LM)

LM Definition: Probability Distribution over sentences in \mathcal{V}^\dagger

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$



Language Model (LM)

LM Definition: Probability Distribution over sentences in \mathcal{V}^\dagger

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n) \geq 0 \quad \forall w_i \in \mathcal{V}$$

$$\sum_{\langle w_1, w_2, \dots, w_n \rangle \in \mathcal{V}^\dagger} p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n) = 1$$



Language Model (LM)

LM Definition: Probability Distribution over sentences in \mathcal{V}^\dagger

$$p(w_1, w_2, \dots, w_n)$$

$$p(w_1, w_2, \dots, w_n) \geq 0$$

$$\sum_{\langle w_1, w_2, \dots, w_n \rangle \in \mathcal{V}^\dagger} p(w_1, w_2, \dots, w_n) = 1$$



Why Language Models (LMs)?

- LMs used in several applications like speech recognition
- Neural LMs have shown SOTA results on many NLP tasks
- The LM techniques can be generalized to many other scenarios.



Language Model (LM)

LM Problem: How do we learn probability distribution over sentences?

$$p(w_1, w_2, \dots, w_n) = ?$$



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there?



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there? $|\mathcal{V}|^n$



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there? $|\mathcal{V}|^n$

Assume number of words in the vocabulary is 40000 (<< english vocabulary),
Suppose, each sentence is of length 10

Number of possible sentences = $1.048576e + 46$



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there? $|\mathcal{V}|^n$

Assume number of words in the vocabulary is 40000 (<< english vocabulary),
Suppose, each sentence is of length 10

Exponential Growth

Number of possible sentences = $1.048576e + 46$



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there? $|\mathcal{V}|^n$

Assume number of words in the vocabulary is 40000 (<< english vocabulary),

Suppose, each sentence is of length 10

Exponential Growth

Number of possible sentences = $1.048576e + 46$

2. What about sentences never seen?



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there? $|\mathcal{V}|^n$

Assume number of words in the vocabulary is 40000 (<< english vocabulary),

Suppose, each sentence is of length 10

Exponential Growth

Number of possible sentences = $1.048576e + 46$

2. What about sentences never seen?

$$p(w_1, w_2, \dots, w_n) = \frac{C(w_1, w_2, \dots, w_n)}{N}$$



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

Naïve Approach: Count all sentences in a corpus.

1. How many sentences of length n are there? $|\mathcal{V}|^n$

Assume number of words in the vocabulary is 40000 (<< english vocabulary),

Suppose, each sentence is of length 10

Exponential Growth

Number of possible sentences = $1.048576e + 46$

2. What about sentences never seen?

Poor Generalization

$$p(w_1, w_2, \dots, w_n) = \frac{C(w_1, w_2, \dots, w_n)}{N}$$



Probability Distribution

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

$$\begin{aligned} p(X_1, X_2, \dots, X_n) &= p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ &\quad p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ &\quad p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1) \end{aligned}$$



Probability Distribution: Independence Assumption

$$p(X_1 = w_1, X_2 = w_2, \dots, X_n = w_n)$$

$$\begin{aligned} p(X_1, X_2, \dots, X_n) &= p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ &\quad p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ &\quad p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1) \end{aligned}$$

Assume, each word depends on only the previous word.

→ Each word is independent of the entire history given the previous word



Probability Distribution: Independence Assumption

$$p(X_1, X_2, \dots, X_n) = p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1)$$

Assume, each word depends on only the previous word.

→ Each word is independent of the entire history given the previous word

$$p(X_k \mid X_{k-1}, \dots, X_2, X_1) = p(X_k \mid X_{k-1})$$

First Order Markov Assumption



Markov Assumption

First Order Markov Assumption

$$p(X_k \mid X_{k-1}, \dots, X_2, X_1) = p(X_k \mid X_{k-1})$$



Markov Assumption

First Order Markov Assumption

$$p(X_k \mid X_{k-1}, \dots, X_2, X_1) = p(X_k \mid X_{k-1})$$

Second Order Markov Assumption

$$p(X_k \mid X_{k-1}, \dots, X_2, X_1) = p(X_k \mid X_{k-1}, X_{k-2})$$



Probability Distribution: Independence Assumption

$$\begin{aligned} p(X_n, X_{n-1}, \dots, X_1) &= p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ &\quad p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ &\quad p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1) \end{aligned}$$



Probability Distribution: Independence Assumption

$$p(X_n, X_{n-1}, \dots, X_1) = p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1)$$

$$p(X_n, X_{n-1}, \dots, X_1) = p(X_n \mid X_{n-1}, X_{n-2}) \times p(X_{n-1} \mid X_{n-2}, X_{n-3}) \times \\ \dots \times p(X_k \mid X_{k-1}, X_{k-2}) \times \dots \dots \\ p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1)$$



Probability Distribution: Independence Assumption

$$p(X_n, X_{n-1}, \dots, X_1) = p(X_n \mid X_{n-1}, \dots, X_2, X_1) \times p(X_{n-1} \mid X_{n-2}, \dots, X_2, X_1) \times \\ p(X_{n-2} \mid X_{n-3}, \dots, X_2, X_1) \times \dots \dots \\ p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1)$$

$$p(X_n, X_{n-1}, \dots, X_1) = p(X_n \mid X_{n-1}, X_{n-2}) \times p(X_{n-1} \mid X_{n-2}, X_{n-3}) \times \\ \dots \times p(X_k \mid X_{k-1}, X_{k-2}) \times \dots \dots \\ p(X_3 \mid X_2, X_1) \times p(X_2 \mid X_1) \times p(X_1)$$

$$p(X_n, X_{n-1}, \dots, X_1) = p(X_1) \times p(X_2 \mid X_1) \times \prod_{i=3}^n p(X_i \mid X_{i-1}, X_{i-2})$$



Language Model

$$p(X_n, X_{n-1}, \dots, X_1) = p(X_1) \times p(X_2 \mid X_1) \times \prod_{i=3}^n p(X_i \mid X_{i-1}, X_{i-2})$$

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$



N-Gram Model



N-Gram Model

Unigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i)$$



N-Gram Model

Unigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i)$$

Bigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1})$$



N-Gram Model

Unigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i)$$

Bigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1})$$

Trigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$



N-Gram Model

Unigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i)$$

Bigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1})$$

Trigram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

N-gram Model

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, \dots, X_{i-(n-1)} = w_{i-(n-1)})$$



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$

STEP 1: Initialize $w_0 = w_{-1} = \text{START}$, $i = 1$



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$

STEP 1: Initialize $w_0 = w_{-1} = \text{START}$, $i = 1$

STEP 2: Sample (Generate) a word: $w_i \sim p(X_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$

STEP 1: Initialize $w_0 = w_{-1} = \text{START}$, $i = 1$

STEP 2: Sample (Generate) a word: $w_i \sim p(X_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$

STEP 3: If $w_i = \text{STOP}$



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$

STEP 1: Initialize $w_0 = w_{-1} = \text{START}$, $i = 1$

STEP 2: Sample (Generate) a word: $w_i \sim p(X_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$

STEP 3: If $w_i = \text{STOP}$

 Terminate and Return sequence: $w_1, w_2, \dots, w_{i-1}, \text{STOP}$



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$

STEP 1: Initialize $w_0 = w_{-1} = \text{START}$, $i = 1$

STEP 2: Sample (Generate) a word: $w_i \sim p(X_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$

STEP 3: If $w_i = \text{STOP}$

 Terminate and Return sequence: $w_1, w_2, \dots, w_{i-1}, \text{STOP}$

else:

 Increment i and Go to STEP 2



Language Generation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

where,

$$w_0 = w_{-1} = \text{START}; \quad w_n = \text{STOP}; \quad w_i \in \mathcal{V} \quad \forall 1 \leq i \leq n - 1$$

STEP 1: Initialize $w_0 = w_{-1} = \text{START}$, $i = 1$

STEP 2: Sample (Generate) a word: $w_i \sim p(X_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$

STEP 3: If $w_i = \text{STOP}$

 Terminate and Return sequence: $w_1, w_2, \dots, w_{i-1}, \text{STOP}$

else:

 Increment i and Go to STEP 2



Language Model Example

I took the book when the book was available



Language Model Example

I took the book when the book was available

$$\begin{aligned} p(\text{I took the book when the book was available}) &= p(STOP \mid \text{available, was}) \\ &\quad \times p(\text{available} \mid \text{was, book}) \\ &\quad \times p(\text{was} \mid \text{book, the}) \\ &\quad \times p(\text{book} \mid \text{the, when}) \\ &\quad \times p(\text{the} \mid \text{when, book}) \\ &\quad \times p(\text{when} \mid \text{book, the}) \\ &\quad \times p(\text{book} \mid \text{the, took}) \\ &\quad \times p(\text{the} \mid \text{took, I}) \\ &\quad \times p(\text{took} \mid \text{I, } START) \\ &\quad \times p(\text{I} \mid START, START) \end{aligned}$$



Language Model Estimation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$



Language Model Estimation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

Let, $p(w_i \mid w_{i-1}, w_{i-2}) = \theta(w_i \mid w_{i-1}, w_{i-2})$

$\theta(q \mid r, s)$ is the model parameter corresponding to each trigram $\{s, r, q\}$



Language Model Estimation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

Trigram Language Model

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n \theta(w_i \mid w_{i-1}, w_{i-2})$$



Language Model Estimation

$$p(X_n = w_n, X_{n-1} = w_{n-1}, \dots, X_1 = w_1) = \prod_{i=1}^n p(X_i = w_i \mid X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2})$$

Trigram Language Model

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n \theta(w_i \mid w_{i-1}, w_{i-2})$$

We would like to estimate $\theta(q \mid r, s)$



Language Model Estimation

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n \theta(w_i \mid w_{i-1}, w_{i-2})$$

We would like to estimate $\theta(q \mid r, s)$

Constraints:

$$\theta(q \mid r, s) \geq 0$$

$$\sum_{q \in \mathcal{V} \cup \{\text{STOP}\}} \theta(q \mid r, s) = 1$$



Language Model Estimation

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n \theta(w_i \mid w_{i-1}, w_{i-2})$$

We would like to estimate $\theta(q \mid r, s)$

Constraints:

$$\theta(q \mid r, s) \geq 0$$

$$\sum_{q \in \mathcal{V} \cup \{\text{STOP}\}} \theta(q \mid r, s) = 1$$

How many parameters in the model?



Language Model Estimation

How many parameters in the model?

$$\theta(q \mid r, s)$$

$$q \in \mathcal{V} \cup STOP$$

$$r, s \in \mathcal{V} \cup START$$



Language Model Estimation

How many parameters in the model?

$$\theta(q \mid r, s)$$

$$q \in \mathcal{V} \cup STOP$$

$$r, s \in \mathcal{V} \cup START$$

About $|\mathcal{V}|^3$ parameters

If $|\mathcal{V}| \sim 40000 \implies 6.4 \times 10^{13}$ parameters



Language Model Estimation

How many parameters in the model?

$$\theta(q \mid r, s)$$

$$q \in \mathcal{V} \cup STOP$$

$$r, s \in \mathcal{V} \cup START$$

About $|\mathcal{V}|^3$ parameters

If $|\mathcal{V}| \sim 40000 \implies 6.4 \times 10^{13}$ parameters

Number of parameters for n-gram mode $\sim |\mathcal{V}|^n !!$



Maximum Likelihood Estimate (MLE)

Given a corpus,

$$\mathcal{S}_1 = \{w_1^{(1)}, w_2^{(1)}, \dots, w_n^{(1)}\}$$

$$\mathcal{S}_2 = \{w_1^{(2)}, w_2^{(2)}, \dots, w_n^{(2)}\}$$

⋮

$$\mathcal{S}_N = \{w_1^{(N)}, w_2^{(N)}, \dots, w_n^{(N)}\}$$



Maximum Likelihood Estimate (MLE)

Given a corpus,

$$\mathcal{S}_1 = \{w_1^{(1)}, w_2^{(1)}, \dots, w_n^{(1)}\}$$

$$\mathcal{S}_2 = \{w_1^{(2)}, w_2^{(2)}, \dots, w_n^{(2)}\}$$

⋮

$$\mathcal{S}_N = \{w_1^{(N)}, w_2^{(N)}, \dots, w_n^{(N)}\}$$

$$p(\mathbf{S} \mid \Theta) = \prod_{k=1}^N \prod_{i=1}^n \theta(w_i^{(k)} \mid w_{i-1}^{(k)}, w_{i-2}^{(k)})$$

$$\Theta = \{\theta(q \mid r, s), \theta(r \mid s, t), \dots\}$$



Maximum Likelihood Estimate (MLE)

$$\mathcal{L}(\Theta) = \prod_{k=1}^N \prod_{i=1}^n \theta(w_i^{(k)} \mid w_{i-1}^{(k)}, w_{i-2}^{(k)})$$



Maximum Likelihood Estimate (MLE)

$$\mathcal{L}(\Theta) = \prod_{k=1}^N \prod_{i=1}^n \theta(w_i^{(k)} \mid w_{i-1}^{(k)}, w_{i-2}^{(k)})$$

Optimization Problem

$$\operatorname{argmax}_{\Theta} \prod_{k=1}^N \prod_{i=1}^n \theta(w_i^{(k)} \mid w_{i-1}^{(k)}, w_{i-2}^{(k)})$$

Constraints

$$\theta(q \mid r, s) \geq 0 \quad \sum_{q \in \mathcal{V} \cup \{\text{STOP}\}} \theta(q \mid r, s) = 1$$



Maximum Likelihood Estimate (MLE)

Optimization Problem

$$\operatorname{argmax}_{\Theta} \prod_{k=1}^N \prod_{i=1}^n \theta(w_i^{(k)} | w_{i-1}^{(k)}, w_{i-2}^{(k)})$$

Constraints

$$\theta(q | r, s) \geq 0 \quad \sum_{q \in \mathcal{V} \cup \{\text{STOP}\}} \theta(q | r, s) = 1$$

$$\theta(q | r, s) = \frac{C(q, r, s)}{\sum_{q_i} C(q_i, r, s)} = \frac{C(q, r, s)}{C(r, s)}$$

$C(q, r, s)$ = Number of times trigram {s,r,q} occurs in the corpus

$C(r, s)$ = Number of times bigram {s,r} occurs in the corpus



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$\theta(STOP \mid ran, cat) = ?$$

$$\theta(STOP \mid am, I) = ?$$

$$\theta(ran \mid cat, The) = ?$$

$$\theta(The \mid START, START) = ?$$



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$\theta(STOP \mid ran, cat) = \frac{2}{3}$$

$$\theta(STOP \mid am, I) = \frac{2}{4}$$

$$\theta(ran \mid cat, The) = \frac{2}{2}$$

$$\theta(The \mid START, START) = \frac{2}{7}$$



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$\theta(out \mid ran, cat) = ?$$



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$\theta(out \mid ran, cat) = \frac{0}{3}$$



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$p(S = \{\text{John ran}\}) = ?$$



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$p(S = \{John\ ran\}) = \theta(STOP \mid ran, John) \\ \times \theta(ran \mid John, START) \\ \times \theta(John \mid START, START)$$



Example

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

I am John STOP

I am out today STOP

John I am STOP

Mary I am STOP

The cat ran STOP

John and cat ran STOP

The cat ran after the mouse STOP

$$p(S = \{\text{John ran}\}) = \underbrace{\theta(\text{STOP} \mid \text{ran, John})}_{= \text{indeterminate}} \\ \times \underbrace{\theta(\text{ran} \mid \text{John, START})}_{= 0} \\ \times \underbrace{\theta(\text{John} \mid \text{START, START})}_{= \frac{2}{7}}$$



Problems with MLE

$$\theta(q \mid r, s) = \frac{C(q, r, s)}{C(r, s)}$$

Overfitting (Sparsity):

Many trigram estimates ($\theta(q \mid r, s)$) are 0

Hapax Legomenon

Indeterminate Estimates:

$C(r, s)$ can be 0



Zipf's Law

- Frequency of usage of a word is inversely proportional to its rank in the frequency table.



Zipf's Law

- Frequency of usage of a word is inversely proportional to its rank in the frequency table.
- Most frequent word occurs twice more often than second most frequent word and thrice than the third most frequent word.



Zipf's Law

- Frequency of usage of a word is inversely proportional to its rank in the frequency table.
- Most frequent word occurs twice more often than second most frequent word and thrice than the third most frequent word.

$$f_i \propto \frac{1}{i}$$

$$f_i = ci^k; \quad k = -1$$

$$\log(f_i) = k \log(i) + C$$



Zipf's Law

- Frequency of usage of a word is inversely proportional to its rank in the frequency table.
- Most frequent word occurs twice more often than second most frequent word and thrice than the third most frequent word.

$$f_i \propto \frac{1}{i}$$

$$f_i = ci^k; \quad k = -1$$

$$\log(f_i) = k \log(i) + C$$

- Word usage in languages follow *power law*
- Plot of *log frequency vs log rank* is linear



Summary

- Language Models provide a way to estimate probability of a sentence
- Markov assumptions make it possible to estimate the sentence probability
- Trigram models are trivial but work well in practice
- MLE for LM gives simple way to estimate probabilities in terms of word counts.
- MLE has overfitting problems.



References

1. Michael Collin's NLP Lecture Notes:
<http://www.cs.columbia.edu/~mcollins/lm-spring2013.pdf>
2. Chapter 4, Speech and Language Processing, Dan Jurafsky and James Martin



- Next week
 - Language Models Smoothing
 - Language Model Evaluation
 - Neural Language Models

