

< Notes



Language models and n-grams

Today's lecture is based on Chapter 4 of the textbook
 Wednesday's quiz will cover chapter 4+8.

Language models assign a probability to a sentence or a sequence of words.

Applications

① Auto complete for texting

② Spelling correction

③ Speech recognition

$p(\text{it's hard to wreck a nice beach})$

④ Machine Translation

$f_1 \nearrow$ high \swarrow large $p(\text{high winds...})$

⑤ Other NLG

$> p(\text{large winds})$

- summarization

$p(\text{high gusts})$

- question answering

$\leftarrow p(\text{large gusts})$

- dialog system

Goal: compute the probability of a sentence or sequence of words

$$P(w) = P(w_1, w_2, w_3, w_4, w_5 \dots w_N)$$

Also want to calculate prob of an upcoming word

$$P(w_5 | w_1, w_2, w_3, w_4)$$

"Language model" LM

Compute the joint prob:

$$P(\text{the underdog Philadelphia Eagles won})$$

Manipulate probabilities

< Notes



Manipulate probabilities

conditional prob

$$p(B|A) = \frac{p(A, B)}{p(A)}$$

↑ given ↑ joint prob ↑ prior prob

$$p(A, B) = p(B|A) \cdot p(A)$$

$$p(A, B, C, D) = \frac{p(A) \cdot p(B|A) \cdot p(C|A, B)}{p(D|A, B, C)}$$

↑ Chain Rule

$$p(x_1, x_2, x_3 \dots x_N) = \underbrace{p(x_1)}_{\text{---}} \circ \underbrace{p(x_2|x_1)}_{\text{---}} \circ \\ \underbrace{p(x_3|x_1, x_2)}_{\text{---}} \circ \dots \circ \underbrace{p(x_N|x_1, \dots x_{N-1})}_{\text{---}}$$

$$\Rightarrow p(w_1, w_2, w_3, w_4 \dots w_n) = \prod_i p(w_i | w_1, w_2, w_3, \dots, w_{i-1})$$

How do we estimate these probabilities for a sentence?

$$p(\text{won} | \text{the underdog team})$$

Maximum likelihood estimation

MLE

$$\frac{\text{count}(\text{the underdog team won})}{\text{count}(\text{the underdog team})}$$

< Notes

Problems

- sentence may not appear
- need a large data set
- many ways combining words into sentences
- Never going to have enough data to estimate the prob. of whole sentences

Simplifying assumption = The Markov assumption

$p(\text{won} \mid \text{the underdog team})$

$\approx p(\text{won} \mid \text{team})$

$\Rightarrow p(w_i)$ only depends on prev k words
NOT the whole context

N-gram
 $\uparrow \uparrow$ $n=2$ bigram
tri gram

$\approx p(\text{won} \mid \text{underdog team})$

$$p(w_1, w_2, w_3, w_4, \dots, w_n) \approx \prod_i p(w_i \mid w_{i-k}, \dots, w_{i-1})$$

$$p(w_i \mid w_1, w_2, w_3) \approx p(w_i \mid w_{i-k}, \dots, w_{i-1})$$

num context words that we take into account

How much history should we use in our Markov assumption?

n-gram models
 ↑
 sequence size
 (= history plus current word)

unigram = 1 word, no history

$$p(w_1, w_2, w_3, \dots, w_n)$$

$$\prod_i p(w_i) = \frac{\text{count}(w_i)}{\text{total words in our corpus}}$$

< Notes

sequence size
history plus current word)

unigram = 1 word, no history

$$\prod p(w_i) = \frac{\text{count}(w_i)}{\text{total words in our corpus}}$$

bigram = 1 word, 1 as history

$$\prod p(w_i | w_{i-1})$$

trigram = 1 word, 2 ..

$$\prod p(w_i | w_{i-2}, w_{i-1})$$

4 gram = 1 word, 3 ..

$$\prod p(w_i | w_{i-3}, w_{i-2}, w_{i-1})$$

MLE

$$\text{unigram } p(w_i) = \frac{\text{count}(w_i)}{\text{len(corpus)}} = \frac{100}{1 \text{ billion}}$$

$$\text{bigram } p(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1} w_i)}{\text{count}(w_{i-1} *)}$$

$$\text{trigram } p(w_i | w_{i-2} w_{i-1}) = \frac{\text{count}(w_{i-2} w_{i-1} w_i)}{\text{count}(w_{i-2} w_{i-1} *)}$$

Is MLE effective here?

$$\text{count} = 0 \leftarrow$$

Historical Notes

Markov chains/assumptions

1910s = Andrei Markov 1913

20k letters Eugene Onegin

1948 Claude Shannon uses n-gram to approximate English

1950s/60s Chomsky decries finite-state Markov models

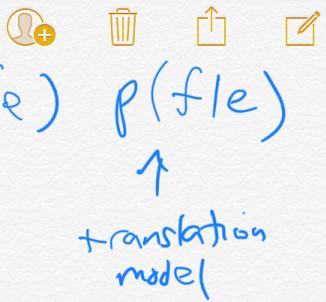
1980s Fred Jelinek at IBM Watson uses n-grams for ASR

① Machine Translation ≈

$$\arg \max_e p(e|f) = \arg \max_e p(e) p(f|e)$$

↑ ↑

< Notes



② Stock market prediction LM

> Fred left IBM to JHU

- Renaissance Tech

- Peter Brown

- Robert Mercer

Trump's book / Brexit
Cambridge Analytica
Breitbart News
ACL Lifetime achievement award

Problems for MLE

zeros

Train
denied the allegations
.. .. report
.. .. claims
.. .. requests

Test
denied the memo

$p(\text{memo} | \text{denied the}) = 0$ and we assign
0 probability to sentences containing such words

Out of Vocab (OOV)

<UNK> to deal with OOVs

fixed lexicon L of size V

normalize training data by replacing any word
not in L w/ <UNK>

Avoid zeros with smoothing

~~Simplest~~ Simplest smoothing ← Kneser-Ney smoothing, Stupid



< Notes



~~Simplest~~ Simplest smoothing \leftarrow Kneser-Ney smoothing, stupid backoff
 add one aka Laplace smoothing

$$P_{add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1} w_i) + 1}{c(w_{i-1}) + V}$$

Backoff

Interpolation

$$\begin{aligned} P(w_i | w_{i-2} w_{i-1}) &= \lambda_1 P(w_i | w_{i-2} w_{i-1}) \\ &\quad + \lambda_2 P(w_i | w_{i-1}) \\ &\quad + \lambda_3 P(w_i) \\ \lambda_1 + \lambda_2 + \lambda_3 &= 1 \end{aligned}$$