

The goal in this chapter is to assign a probability to a sentence. \rightarrow machine translation

applications + machine translation

- ↳ speech recognition

- summary generation, etc

goal: compute the probability of a sentence (or seq. of words)

$$P(w) = P(w_1, w_2, \dots, w_n)$$

↓
related task: compute the probability of the next word.

→ A model that computer either of these is called a language model

$P(\omega)$? \downarrow

$$p(x_1) \times p(x_2, \dots, x_n | x_1)$$

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

I [chain rule $P(A, B) = P(A) \times P(B|A)$]

applying chain rule recursively on $P(w)$:

$$P(w) = P(w_1) P(w_2|w_1) P(w_3|w_1, w_2) \dots P(w_n|w_1, \dots, w_{n-1})$$

$$P(w_1, w_2, \dots, w_n) = \prod P(w_i | w_1, \dots, w_{i-1})$$

Let's say we have M words in a .

$$P(w_1, w_2, \dots, w_M) = p(w_1) p(w_2 | w_1) \dots p(w_M | w_1, \dots, w_{M-1})$$

to compute the probability

$$p(w_M | w_1, \dots, w_{M-1})$$

I
we need to model V^M contents.

this is not possible with finite sample data.

To solve this problem, n -gram models, which make a simplifying approximation to condition on only the past $(n-1)$ words
i.e.,

$$p(w_M | w_1, \dots, w_{M-1}) \approx p(w_M | w_{M-1}, \dots, w_{M-n+1})$$

simplest case: unigram model

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

bigram model:

$$p(w_1, \dots, w_n) \approx \prod_i P(w_i | w_{i-1})$$

→ can extend to trigrams, 4-grams, 5-grams, ...

in general, this is an insufficient model of language.

because language has long distance dependencies.

estimating bi-gram probabilities

$$p(w_i | w_{i-1}) : \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

by for n-grams

$$p(w_m | w_{m-1} \dots w_{m-n+1}) : \frac{\text{count}(w_m, w_{m-1} \dots w_{m-n+1})}{\text{count}(w_{m-1} \dots w_{m-n+1})}$$

Evaluating language models:

intrinsic evaluation:

intrinsic evaluation: but whether our model prefers
"good sentences" to "bad ones".

Perplexity:

w_1, \dots, w_M + held-out data for validating

$$l(w) = \sum_{m=1}^M \log P(w_m | w_{m-1}, \dots, w_1)$$

$$\text{ppl}(w) = 2^{-\frac{l(w)}{M}}$$

→ minimizing perplexity is the same as maximizing probability.

There might be a chance that probability $P(w_m | w_{m-1}, \dots, w_1)$ will be zero, so to avoid this we could do smoothing

Smoothing

! intuition: steal some probability mass & distribute among the zero-probability ones

Laplace smoothing:

$$P(w_i | w_{1:i-1}) = \frac{C(w_{1:i-1}, w_i) + 1}{C(w_{1:i-1}) + V}$$

$$C^V(w_{1:i-1}, w_i) = \frac{(C(w_{1:i-1}, w_i) + 1) \cdot C(w_{1:i-1})}{C(w_{1:i-1}) + V}$$

amount
 $d_C = \frac{C^V}{C}$

Interpolation: mix unigram, bigram, trigram statistics

$$P(w_i | w_{i-1}, w_{i-2}) : \lambda_1 P(w_i | w_{i-1}, w_{i-2}) +$$

$$\lambda_2 P(w_i | w_{i-1}) + \lambda_3 P(w_i)$$

Dealing with unseen words:

↳ create a token $\langle \text{UNK} \rangle$