

# Lecture 21

Advanced smoothing techniques (N-gram models)

# N-gram probability calculation

eg. Of next word prediction (bigram model) : this is ...?....

- Probability of the Nth word given the previous N-1 words

(bigram : N=2) 
$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

(N-gram) 
$$P(w_n|w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}w_n)}{C(w_{n-N+1}^{n-1})}$$

max probability → winner among all candidate words = next word

- Probability of the whole sentence containing n words (eg. for bigrams)

$$P(w_1^n) \approx \prod_{k=1}^n P(w_k|w_{k-1})$$

# OOV (Out of Vocabulary) words and why do you require smoothing?

- Formula 1 (prob of next word)

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

- Formula 2 (prob of whole sentence)

$$P(w_1^n) \approx \prod_{k=1}^n P(w_k|w_{k-1})$$

- $w_{k-2}w_{k-1}w_k$  (trigram)
- This is... $w_k$ ...

- Count of denom=0
- 0/0 (avoid divide by zero)
- Count of denom in non zero
- Count of num=0 (prob in Formula 2)
- Even if 1 prob in chain is 0 the whole sent. prob=0

# Laplace Smoothing (recap)

- **Laplace smoothing** (of last slide probability formula)
- Also called add-one smoothing
- To prevent 0/0 situation
- Add 1 to numerator and V to the denominator

$$p_i^* = \frac{c_i + 1}{N + V}$$

Let  $c^*$  be the modified count in the numerator

V=size of the vocabulary = number of unique unigram (denominator) in the training corpus

# Class Assignment (deadline 5 pm on 21/10/20)

Training corpus → build a probabilistic model (N-gram) → Test input → prediction of the next word → prob(sentence) [LM: Language modelling]

*Training corpus:*

*N= 2 (bigram) V=13*

<s>the start of the day was good</s><s> if the start is good the whole day is good</s><s> the goodness of the day matters</s>

Text prediction (LM):

<s> the </s><s>the.....

# Advanced smoothing techniques

# 1. Good-Turing discounting

- Proposed by Good (1953)
- Ques: How to handle OOV (Out of Vocabulary) words whose count in the training corpus is zero?
- Let  $N_c$  be the number of N-grams that occur  $c$  times (in the training corpus)
- (EXCEPT N0)
- $$c^* = (c + 1) \frac{N_{c+1}}{N_c}$$
- $N_2=6$   $N_1=14$   $N_0=0$   $N=20$
- For N0 case (OOV),  $c^*=N_1/N=14/20$

## 2. Knesser-Ney smoothing

$$P_{\text{KN}}(w_i|w_{i-1}) = \begin{cases} \frac{C(w_{i-1}w_i) - \mathbf{D}}{C(w_{i-1})}, & \text{if } C(w_{i-1}w_i) > 0 \\ \alpha(w_i) \frac{|\{w_{i-1}: C(w_{i-1}w_i) > 0\}|}{\sum_{w_i} |\{w_{i-1}: C(w_{i-1}w_i) > 0\}|} & \text{otherwise.} \end{cases}$$

## 3. Interpolation

$$P_{\text{base}}(w_i|w_{i-n+1}^{i-1}) = \lambda_n P_{\text{ML}}(w_i|w_{i-n+1}^{i-1}) + (1 - \lambda_n) P_{\text{base}}(w_i|w_{i-n+2}^{i-1})$$



## 4. Katz backoff

- Similar to interpolation
- Backoff to lower N-grams in case higher N-grams not present in training corpus

$$P_{\text{katz}}(z|x,y) = \begin{cases} P^*(z|x,y), & \text{if } C(x,y,z) > 0 \\ \alpha(x,y)P_{\text{katz}}(z|y), & \text{else if } C(x,y) > 0 \\ P^*(z), & \text{otherwise.} \end{cases}$$
$$P_{\text{katz}}(z|y) = \begin{cases} P^*(z|y), & \text{if } C(y,z) > 0 \\ \alpha(y)P^*(z), & \text{otherwise.} \end{cases}$$