# Statistical Language Modelling

D. Thenmozhi

Associate Professor

SSNCE

# Statistical Language Modelling

- Statistical Language Modelling
  - Is a probability distribution P(s) over all possible word sequences

- Dominant approach
  - N-gram model

# Probabilistic Language Models

- Machine Translation:
  - P(**high** winds tonite) > P(**large** winds tonite)
- Spell Correction
  - The office is about fifteen **minuets** from my house
    - P(about fifteen **minutes** from) > P(about fifteen **minuets** from)
- Speech Recognition
  - P(I saw a van) >> P(eyes awe of an)

# Probabilistic Language Models

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

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \ldots w_n)$$

- Related task: probability of an upcoming word:

$$P(w_5 \mid w_1, w_2, w_3, w_4)$$

- A model that computes either of these:

$$P(W) \quad \text{or} \quad P(w_n \mid w_1, w_2 \ldots w_{n-1})$$ is called a **language model**.

# N-gram model

- Decompose sentence probability into a product of conditional probabilities using the chain rule

$$P(x_1,x_2,x_3,...,x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1,x_2)...P(x_n|x_1,...,x_{n-1})$$

$$P(w_1w_2...w_n) = \prod_i P(w_i | w_1w_2...w_{i-1})$$

P("its water is so transparent") =

P(its) × P(water|its) × P(is|its water)

× P(so|its water is) × P(transparent|its water is so)

# Probability Estimation

$$P(\text{the} \mid \text{its water is so transparent that}) =$$

$$\frac{Count(\text{its water is so transparent that the})}{Count(\text{its water is so transparent that})}$$

??

**Markov Assumption**    (Simplifying Assumption)

$$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{that})$$

OR

$$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{transparent that})$$

# Markov Assumption

$$P(w_1 w_2 \ldots w_n) \approx \prod_i P(w_i \mid w_{i-k} \ldots w_{i-1})$$

For each word

$$P(w_i \mid w_1 w_2 \ldots w_{i-1}) \approx P(w_i \mid w_{i-k} \ldots w_{i-1})$$

Uni-gram

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

Bi-gram

$$P(w_i \mid w_1 w_2 \ldots w_{i-1}) \approx P(w_i \mid w_{i-1})$$

# Estimating bi-gram probabilities

- The Maximum Likelihood Estimate

$$P(w_i \mid w_{i-1}) = \frac{count(w_{i-1}, w_i)}{count(w_{i-1})}$$

$$P(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

# Example

$$P(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$P(\texttt{I} \mid \texttt{<s>}) = \frac{2}{3} = .67$     $P(\texttt{Sam} \mid \texttt{<s>}) = \frac{1}{3} = .33$     $P(\texttt{am} \mid \texttt{I}) = \frac{2}{3} = .67$

$P(\texttt{</s>} \mid \texttt{Sam}) = \frac{1}{2} = 0.5$     $P(\texttt{Sam} \mid \texttt{am}) = \frac{1}{2} = .5$     $P(\texttt{do} \mid \texttt{I}) = \frac{1}{3} = .33$

# Example

- Training set
  - The Arabian Knights
  - These are the fairy tales of the east
  - The stories of the Arabian Knights are translated in many languages

- Test sentence
  - The Arabian Knights are the fairy tales of the east

# Bi-gram model

- P(the/<s>)=0.67
- P(are / these)=1.0
- P(tales / fairy)= 1.0
- P(east/the)=0.2
- P(are / Knights)= 1.0
- P(many/in)=1.0

P(Arabian /the)=0.4

P(the/ are)=0.5

P(of/ tales)= 1.0

P(stories/the)=0.2

P(translated / are)=0.5

P(languages /many)= 1.0

P(Knights / Arabian)=1.0

P(fairy /the)=0.2

P(the/of)= 1.0

P(of/stories)= 1.0

P(in/ translated)= 1.0

# Bi-gram model

- The Arabian Knights are the fairy tales of the east

- P(the/<s>) x P(Arabian /the) x P(Knights / Arabian) x P(are / these) x P(the/ are) x  P(fairy /the) x P(tales / fairy) x P(of/ tales) x P(the/of) P(east/the)

- 0.67 x 0.4 x 1.0 x 1.0 x 0.5 x 0.2 x 1.0 x 1.0 x 1.0 x 0.2
        = 0.0067

# Issues in bi-gram model

- Multiplying the probability might cause a numerical underflow
- To avoid this, sum their logs

$$\log(p_1 \times p_2 \times p_3 \times p_4) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

- Suffers from data sparseness
- Zero probability entries in bi-gram matrix
- Soln: Smoothing techniques
  - Refers to the task of re-evaluating zero-probability or low-probability n-grams and assigning them non-zero values
  - Add-one smoothing, good-turing smoothing and catching techniques

# Add-one Smoothing

- Also called Laplace smoothing
- Pretend we saw each word one more time than we did
- Just add one to all the counts!

- MLE estimate:

$$P_{MLE}(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

- Add-1 estimate:

$$P_{Add-1}(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

# Add-one Smoothing

Berkeley Restaurant project sentences

|         | i  | want | to  | eat | chinese | food | lunch | spend |
|---------|----|------|-----|-----|---------|------|-------|-------|
| i       | 5  | 827  | 0   | 9   | 0       | 0    | 0     | 2     |
| want    | 2  | 0    | 608 | 1   | 6       | 6    | 5     | 1     |
| to      | 2  | 0    | 4   | 686 | 2       | 0    | 6     | 211   |
| eat     | 0  | 0    | 2   | 0   | 16      | 2    | 42    | 0     |
| chinese | 1  | 0    | 0   | 0   | 0       | 82   | 1     | 0     |
| food    | 15 | 0    | 15  | 0   | 1       | 4    | 0     | 0     |
| lunch   | 2  | 0    | 0   | 0   | 0       | 1    | 0     | 0     |
| spend   | 1  | 0    | 1   | 0   | 0       | 0    | 0     | 0     |

# Add-one Smoothing

| | i | want | to | eat | chinese | food | lunch | spend |
|---|---|---|---|---|---|---|---|---|
| i | 0.002 | 0.33 | 0 | 0.0036 | 0 | 0 | 0 | 0.00079 |
| want | 0.0022 | 0 | 0.66 | 0.0011 | 0.0065 | 0.0065 | 0.0054 | 0.0011 |
| to | 0.00083 | 0 | 0.0017 | 0.28 | 0.00083 | 0 | 0.0025 | 0.087 |
| eat | 0 | 0 | 0.0027 | 0 | 0.021 | 0.0027 | 0.056 | 0 |
| chinese | 0.0063 | 0 | 0 | 0 | 0 | 0.52 | 0.0063 | 0 |
| food | 0.014 | 0 | 0.014 | 0 | 0.00092 | 0.0037 | 0 | 0 |
| lunch | 0.0059 | 0 | 0 | 0 | 0 | 0.0029 | 0 | 0 |
| spend | 0.0036 | 0 | 0.0036 | 0 | 0 | 0 | 0 | 0 |

# Add-one Smoothing

|         | i  | want | to  | eat | chinese | food | lunch | spend |
|---------|----|------|-----|-----|---------|------|-------|-------|
| i       | 6  | 828  | 1   | 10  | 1       | 1    | 1     | 3     |
| want    | 3  | 1    | 609 | 2   | 7       | 7    | 6     | 2     |
| to      | 3  | 1    | 5   | 687 | 3       | 1    | 7     | 212   |
| eat     | 1  | 1    | 3   | 1   | 17      | 3    | 43    | 1     |
| chinese | 2  | 1    | 1   | 1   | 1       | 83   | 2     | 1     |
| food    | 16 | 1    | 16  | 1   | 2       | 5    | 1     | 1     |
| lunch   | 3  | 1    | 1   | 1   | 1       | 2    | 1     | 1     |
| spend   | 2  | 1    | 2   | 1   | 1       | 1    | 1     | 1     |

# Add-one Smoothing

|         | i        | want     | to       | eat      | chinese  | food     | lunch    | spend    |
|---------|----------|----------|----------|----------|----------|----------|----------|----------|
| i       | 0.0015   | 0.21     | 0.00025  | 0.0025   | 0.00025  | 0.00025  | 0.00025  | 0.00075  |
| want    | 0.0013   | 0.00042  | 0.26     | 0.00084  | 0.0029   | 0.0029   | 0.0025   | 0.00084  |
| to      | 0.00078  | 0.00026  | 0.0013   | 0.18     | 0.00078  | 0.00026  | 0.0018   | 0.055    |
| eat     | 0.00046  | 0.00046  | 0.0014   | 0.00046  | 0.0078   | 0.0014   | 0.02     | 0.00046  |
| chinese | 0.0012   | 0.00062  | 0.00062  | 0.00062  | 0.00062  | 0.052    | 0.0012   | 0.00062  |
| food    | 0.0063   | 0.00039  | 0.0063   | 0.00039  | 0.00079  | 0.002    | 0.00039  | 0.00039  |
| lunch   | 0.0017   | 0.00056  | 0.00056  | 0.00056  | 0.00056  | 0.0011   | 0.00056  | 0.00056  |
| spend   | 0.0012   | 0.00058  | 0.0012   | 0.00058  | 0.00058  | 0.00058  | 0.00058  | 0.00058  |

# Advanced smoothing algorithms

- Intuition used by many smoothing algorithms
  - Good-Turing
  - Kneser-Ney
  - Witten-Bell
- Use the count of things we've **seen once**
  - to help estimate the count of things we've **never seen**

# Good-Turing smoothing

- $N_c$ = the count of things we've seen c times
- Sam I am I am Sam I do not eat

```
I    3
sam  2
am   2
do   1
not  1
eat  1
```

$N_1 = 3$

$N_2 = 2$

$N_3 = 1$

# Good-Turing smoothing

- You are fishing (a scenario from Josh Goodman), and caught:
  - 10 carp, 3 perch, 2 whitefish, 1 trout, 1 salmon, 1 eel = 18 fish
- How likely is it that next species is trout?
  - 1/18
- How likely is it that next species is new (i.e. catfish or bass)
  - Let's use our estimate of things-we-saw-once to estimate the new things.
  - 3/18 (because $N_1=3$)
- Assuming so, how likely is it that next species is trout?
  - Must be less than 1/18
  - How to estimate?

# Good-Turing smoothing

$$P^*_{GT}(\text{things with zero frequency}) = \frac{N_1}{N} \qquad c^* = \frac{(c+1)N_{c+1}}{N_c}$$

- Unseen (bass or catfish)
  - c = 0:
  - MLE p = 0/18 = 0

  - $P^*_{GT}$ (unseen) = $N_1$/N = 3/18

- Seen once (trout)
  - c = 1
  - MLE p = 1/18

  - C*(trout) = 2 * N2/N1
  
    = 2 * 1/3
  
    = 2/3
  - $P^*_{GT}$(trout) = 2/3 / 18 = 1/27

# Summary

- SLM
- N-gram models
- Bi-gram model
- Smoothing techniques

# Questions

- Where will you used statistical language modelling?

- What are the applications of SLM?

- What is the dominant SLM?

- Why do you need smoothing techniques in bi-gram LM?