

Estimating n -gram probabilities

- The maximum likelihood estimate of an n -gram model:

$$P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-n)}, \dots, \mathbf{x}^{(t-1)}) \stackrel{MLE}{=} \frac{\text{count}(\mathbf{x}^{(t-n)}, \dots, \mathbf{x}^{(t)})}{\sum_{j=1}^N \text{count}(\mathbf{x}^{(t-n)}, \dots, \mathbf{x}_j^{(t)})}$$

Example from Jurafsky & Martin

<s> I am Sam </s>

<s> Sam I am </s>

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

Here are the calculations for some of the bigram probabilities from this corpus

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

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

Choosing the right n

- Choosing n is a trade off between:
 - Modeling capacity (increases with n)

Examples from Eisenstein (2018)

Gorillas always like to groom **their** friends.

The **computer** that's on the 3rd floor of our office building **crashed**.

- Tractable estimation (harder with larger n ... sparsity)

Example from Jurafsky & Martin

	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

Figure 3.1 Bigram counts for eight of the words (out of $V = 1446$) in the Berkeley Restaurant Project corpus of 9332 sentences. Zero counts are in gray.