

# Language Models

Jordan Boyd-Graber

University of Maryland

Introduction

Slides adapted from Philip Koehn

# Roadmap

After this lecture, you'll be able to:

- Understand probability distributions through the metaphor of the Chinese Restaurant Process
- Be able to calculate Kneser-Ney smoothing
- Understand the role of contexts in language models

# Intuition

- Some words are “sticky”
- “San Francisco” is very common (high ungram)
- But Francisco only appears after one word

# Intuition

- Some words are “sticky”
- “San Francisco” is very common (high ungram)
- But Francisco only appears after one word
- Our goal: to tell a statistical story of bay area restaurants to account for this phenomenon
- How to model this phenomena

# Interpolation

- Higher and lower order  $n$ -gram models have different strengths and weaknesses
  - ▶ high-order  $n$ -grams are sensitive to more context, but have sparse counts
  - ▶ low-order  $n$ -grams consider only very limited context, but have robust counts
- Combine them

$$\begin{aligned} p_I(w_3 | w_1, w_2) = & \lambda_1 p_1(w_3) \\ & + \lambda_2 p_2(w_3 | w_2) \\ & + \lambda_3 p_3(w_3 | w_1, w_2) \end{aligned}$$

## Back-Off

- Trust the highest order language model that contains n-gram

$$p_n^{BO}(w_i | w_{i-n+1}, \dots, w_{i-1}) = \begin{cases} \alpha_n(w_i | w_{i-n+1}, \dots, w_{i-1}) & \text{if } \text{count}_n(w_{i-n+1}, \dots, w_i) > 0 \\ d_n(w_{i-n+1}, \dots, w_{i-1}) p_{n-1}^{BO}(w_i | w_{i-n+2}, \dots, w_{i-1}) & \text{else} \end{cases}$$

- Requires
  - ▶ adjusted prediction model  $\alpha_n(w_i | w_{i-n+1}, \dots, w_{i-1})$
  - ▶ discounting function  $d_n(w_1, \dots, w_{n-1})$

## Let's remember what a language model is

- It is a distribution over the next word in a sentence
- Given the previous  $n - 1$  words

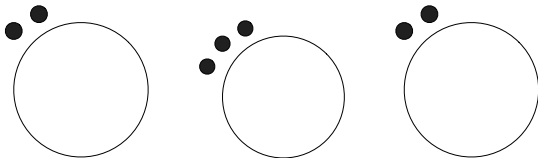
## Let's remember what a language model is

- It is a distribution over the next word in a sentence
- Given the previous  $n - 1$  words
- The challenge: backoff and sparsity



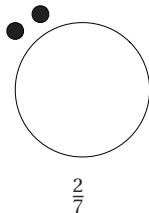
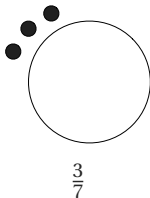
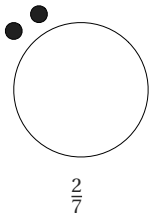
## The Chinese Restaurant as a Distribution

To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



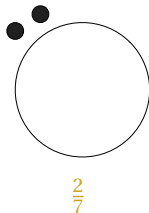
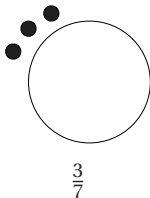
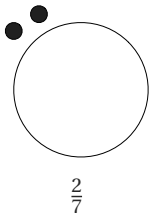
## The Chinese Restaurant as a Distribution

To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



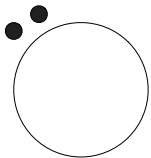
## The Chinese Restaurant as a Distribution

To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.

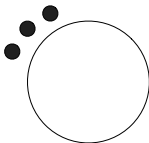


## The Chinese Restaurant as a Distribution

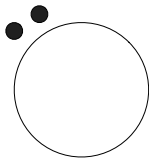
To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



$\frac{2}{7}$   
dog



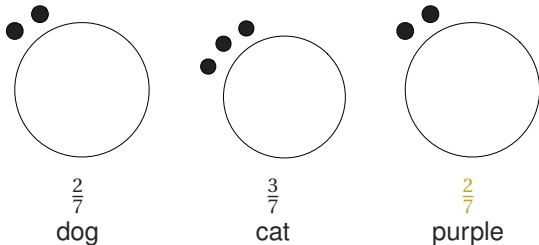
$\frac{3}{7}$   
cat



$\frac{2}{7}$   
purple

## The Chinese Restaurant as a Distribution

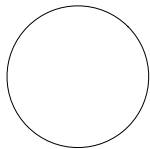
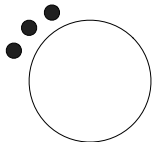
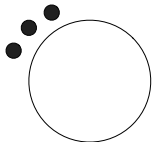
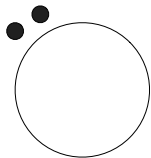
To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



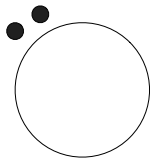
But this is just Maximum Likelihood

Why are we talking about Chinese Restaurants?

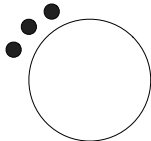
Always one more table ...



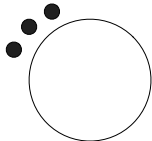
Always one more table ...



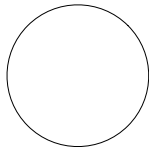
$$\frac{2}{7+\alpha}$$



$$\frac{3}{7+\alpha}$$

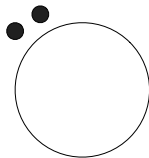


$$\frac{2}{7+\alpha}$$

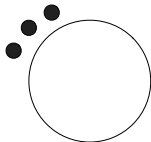


$$\frac{\alpha}{7+\alpha}$$

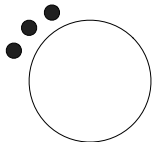
Always one more table ...



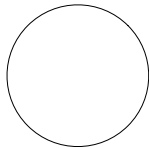
$\frac{2}{7+\alpha}$   
dog



$\frac{3}{7+\alpha}$   
cat



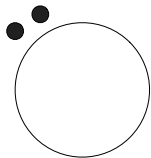
$\frac{2}{7+\alpha}$   
purple



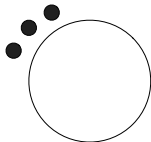
$\frac{\alpha}{7+\alpha}$   
???



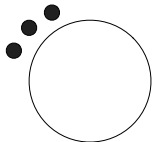
Always one more table ...



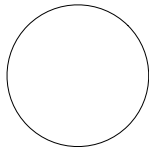
$\frac{2}{7+\alpha}$   
dog



$\frac{3}{7+\alpha}$   
cat

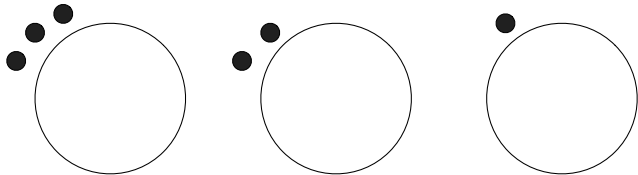


$\frac{2}{7+\alpha}$   
purple

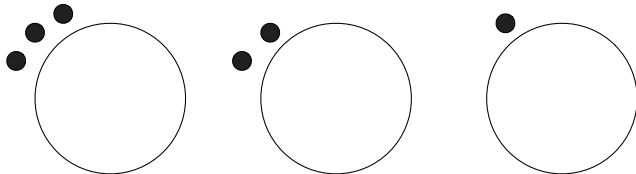


$\frac{\alpha}{7+\alpha}$   
???

## What to do with a new table?



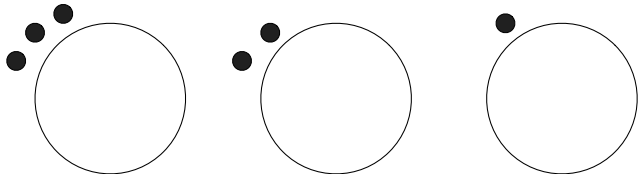
## What to do with a new table?



What can be a base distribution?

- Uniform (Dirichlet smoothing)

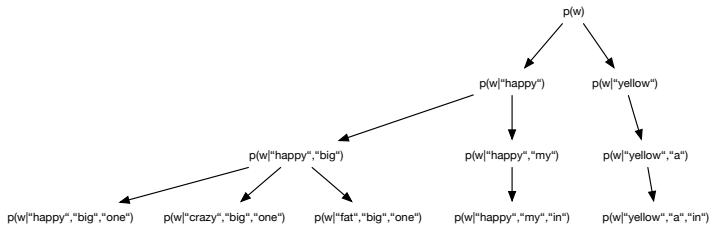
## What to do with a new table?



### What can be a base distribution?

- Uniform (Dirichlet smoothing)
- Specific contexts  $\rightarrow$  less-specific contexts (backoff)

# A hierarchy of Chinese Restaurants



## Seating Assignments

<s> a a a b a c </s>

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

<s> Restaurant

a Restaurant

b Restaurant

c Restaurant

# Seating Assignments

`<s> a a a b a c </s>`

Unigram Restaurant

`<s> Restaurant`

`*`<sup>1</sup>

`b Restaurant`

`a Restaurant`

`c Restaurant`



# Seating Assignments

`<s> a a a b a c </s>`

Unigram Restaurant

`*`<sup>1</sup>

`<s> Restaurant`

`*`<sup>1</sup>

`b Restaurant`

`a Restaurant`

`c Restaurant`

# Seating Assignments

`<s> a a a b a c </s>`

Unigram Restaurant

a<sup>1</sup>

`<s> Restaurant`

a<sup>1</sup>

b Restaurant

a Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

\*<sup>1</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

\*<sup>1</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>2</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>1</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>2</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>1</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>2</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>2</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> \*<sup>1</sup>

b Restaurant

c Restaurant



# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

$a^2$   $*$ <sup>1</sup>

<s> Restaurant

$a^1$

a Restaurant

$a^2$   $*$ <sup>1</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>2</sup> b<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> \*<sup>1</sup>

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

$a^2$   $b^1$

<s> Restaurant

$a^1$

a Restaurant

$a^2$   $b^1$

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

$a^2$   $b^1$

<s> Restaurant

$a^1$

a Restaurant

$a^2$   $b^1$

b Restaurant

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>2</sup> b<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

b Restaurant

\*<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup>

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

$a^2$   $b^1$

<s> Restaurant

$a^1$

b Restaurant

$*$ <sup>1</sup>

a Restaurant

$a^2$   $b^1$

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

b Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup>

c Restaurant



# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

b Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> \*<sup>1</sup>

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

## Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> \*<sup>1</sup>

## <s> Restaurant

a<sup>1</sup>

## b Restaurant

a<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> \*<sup>1</sup>

## c Restaurant

# Seating Assignments

`<s> a a a b a c </s>`

Unigram Restaurant

$a^3$   $b^1$   $c^1$

`<s> Restaurant`

$a^1$

b Restaurant

$a^1$

a Restaurant

$a^2$   $b^1$   $c^1$

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

b Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

c Restaurant

# Seating Assignments

<s> a a a b a c </s>

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

\*<sup>1</sup>

# Seating Assignments

<s> a a a b a c </s>

## Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> \*<sup>1</sup>

## <s> Restaurant

a<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

## b Restaurant

a<sup>1</sup>

## c Restaurant

\*<sup>1</sup>

# Seating Assignments

`<s> a a a b a c </s>`

## Unigram Restaurant

`a`<sup>3</sup> `b`<sup>1</sup> `c`<sup>1</sup> `</s>`<sup>1</sup>

## `<s>` Restaurant

`a`<sup>1</sup>

## a Restaurant

`a`<sup>2</sup> `b`<sup>1</sup> `c`<sup>1</sup>

## b Restaurant

`a`<sup>1</sup>

## c Restaurant

`</s>`<sup>1</sup>

## Real examples

- San Francisco



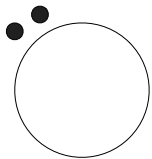
## Real examples

- San Francisco
- Star Spangled Banner

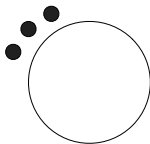
## Real examples

- San Francisco
- Star Spangled Banner
- Bottom Line: Counts go to the context that explains it best

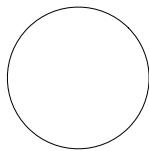
## The rich get richer



$$\frac{2}{5+\theta}$$



$$\frac{3}{5+\theta}$$



$$\frac{\theta}{5+\theta}$$

## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$ .
- The backoff context  $\pi(u)$

## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$ .
- The backoff context  $\pi(u)$

## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$ .
- The backoff context  $\pi(u)$

## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$
- The backoff context  $\pi(u)$

## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$ .
- The backoff context  $\pi(u)$



## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$
- The backoff context  $\pi(u)$

## Computing the Probability of an Observation

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (1)$$

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $u$
- Number seated at table serving  $x$  in restaurant  $u$ ,  $c_{u,x}$
- Number seated at all tables in restaurant  $u$ ,  $c_{u,\cdot}$ .
- The backoff context  $\pi(u)$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{c_{a,b}}{\theta + c_{u,\cdot}} + \frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{c_{a,b}}{\theta + c_{u,\cdot}} + \frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{\theta + c_{u,\cdot}} + \frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{1.0 + c_{u,\cdot}} + \frac{1.0}{1.0 + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{1.0 + 4} + \frac{1.0}{1.0 + 4} p(w = x | \vec{s}, \theta, \pi(u)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{1.0 + 4} + \frac{1.0}{1.0 + 4} p(w = x | \vec{s}, \theta, \pi(u)) \quad (2)$$



Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{1.0 + 4} + \frac{1.0}{1.0 + 4} p(w = x | \vec{s}, \theta, \pi(\emptyset)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{1.0 + 4} + \frac{1.0}{1.0 + 4} p(w = x | \vec{s}, \theta, \pi(\emptyset)) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{5} + \frac{1}{5} \left( \frac{c_{\emptyset, b}}{c_{\emptyset, \cdot} + \theta} + \frac{\theta}{c_{\emptyset, \cdot} + \theta} \frac{1}{V} \right) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{5} + \frac{1}{5} \left( \frac{c_{\emptyset, b}}{c_{\emptyset, \cdot} + \theta} + \frac{\theta}{c_{\emptyset, \cdot} + \theta} \frac{1}{5} \right) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{5} + \frac{1}{5} \left( \frac{c_{\emptyset, b}}{c_{\emptyset, \cdot} + 1.0} + \frac{1.0}{c_{\emptyset, \cdot} + 1.0} \frac{1}{5} \right) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{5} + \frac{1}{5} \left( \frac{1}{c_{\emptyset, \cdot} + 1.0} + \frac{1.0}{c_{\emptyset, \cdot} + 1.0} \frac{1}{5} \right) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{5} + \frac{1}{5} \left( \frac{1}{6 + 1.0} + \frac{1.0}{6 + 1.0} \frac{1}{5} \right) \quad (2)$$

Example:  $p(w = b | \vec{s}, \theta = 1.0, u = a)$

Unigram Restaurant

a<sup>3</sup> b<sup>1</sup> c<sup>1</sup> </s><sup>1</sup>

<s> Restaurant

a<sup>1</sup>

a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

b Restaurant

a<sup>1</sup>

c Restaurant

</s><sup>1</sup>

$$p(w = b | \dots) = \frac{1}{5} + \frac{1}{5} \left( \frac{1}{7} + \frac{1 \cdot 1}{7 \cdot 5} \right) = 0.24 \quad (2)$$



## Discounting

- Empirically, it helps favor the backoff if you have more tables
- Otherwise, it gets too close to maximum likelihood
- Idea is called discounting
- Steal a little bit of probability mass  $\delta$  from every table and give it to the new table (backoff)

## Discounting

- Empirically, it helps favor the backoff if you have more tables
- Otherwise, it gets too close to maximum likelihood
- Idea is called discounting
- Steal a little bit of probability mass  $\delta$  from every table and give it to the new table (backoff)

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x}}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (3)$$

## Discounting

- Empirically, it helps favor the backoff if you have more tables
- Otherwise, it gets too close to maximum likelihood
- Idea is called discounting
- Steal a little bit of probability mass  $\delta$  from every table and give it to the new table (backoff)

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x} - \delta}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta + T\delta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (3)$$

## Discounting

- Empirically, it helps favor the backoff if you have more tables
- Otherwise, it gets too close to maximum likelihood
- Idea is called discounting
- Steal a little bit of probability mass  $\delta$  from every table and give it to the new table (backoff)

$$p(w = x | \vec{s}, \theta, u) = \underbrace{\frac{c_{u,x} - \delta}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta + T\delta}{\theta + c_{u,\cdot}} p(w = x | \vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (3)$$

## Discounting

- Empirically, it helps favor the backoff if you have more tables
- Otherwise, it gets too close to maximum likelihood
- Idea is called discounting
- Steal a little bit of probability mass  $\delta$  from every table and give it to the new table (backoff)

$$p(w = x|\vec{s}, \theta, u) = \underbrace{\frac{c_{u,x} - \delta}{\theta + c_{u,\cdot}}}_{\text{existing table}} + \underbrace{\frac{\theta + T\delta}{\theta + c_{u,\cdot}} p(w = x|\vec{s}, \theta, \pi(u))}_{\text{new table}} \quad (3)$$

### Interpolated Kneser-Ney!

## More advanced models

- Interpolated Kneser-Ney assumes **one table with a dish (word)** per restaurant
- Can get slightly better performance by assuming you can have duplicated tables: **Pitman-Yor** language model
- Requires Gibbs Sampling of the seating assignments
  - ▶ If you walk into a restaurant with your dish on a table, you sample whether to create a new dish or not
  - ▶ Requires more computation, not deterministic
- Bottom line: discrete representation of the context

## More advanced models

- Interpolated Kneser-Ney assumes **one table with a dish (word)** per restaurant
- Can get slightly better performance by assuming you can have duplicated tables: **Pitman-Yor** language model
- Requires Gibbs Sampling of the seating assignments
  - ▶ If you walk into a restaurant with your dish on a table, you sample whether to create a new dish or not
  - ▶ Requires more computation, not deterministic
- Bottom line: discrete representation of the context
- Neural language models use continuous representations to store context, which works better

# Exercises



## Exercise

- Start with restaurant we had before
- Assume you see  $\langle s \rangle$  b b a c  $\langle /s \rangle$ ; add those counts to tables
- Compute probability of b following a ( $\theta = 1.0, \delta = 0.5$ )
- Compute the probability of a following b
- Compute probability of  $\langle /s \rangle$  following  $\langle s \rangle$

# A busy night at the restaurant

## Unigram Restaurant

$a^3$   $b^1$   $c^1$   $\langle /s \rangle^1$

## $\langle s \rangle$ Restaurant

$a^1$

## a Restaurant

$a^2$   $b^1$   $c^1$

## b Restaurant

$a^1$

## c Restaurant

$\langle /s \rangle^1$

# A busy night at the restaurant

## Unigram Restaurant

$a^3$   $b^1$   $c^1$   $\langle /s \rangle^1$

## $\langle s \rangle$ Restaurant

$a^1$   $b^1$

## a Restaurant

$a^2$   $b^1$   $c^1$

## b Restaurant

$a^1$

## c Restaurant

$\langle /s \rangle^1$

# A busy night at the restaurant

## Unigram Restaurant

a<sup>3</sup> b<sup>2</sup> c<sup>1</sup> </s><sup>1</sup>

## <s> Restaurant

a<sup>1</sup> b<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

## b Restaurant

a<sup>1</sup>

## c Restaurant

</s><sup>1</sup>

# A busy night at the restaurant

## Unigram Restaurant

a<sup>3</sup> b<sup>2</sup> c<sup>1</sup> </s><sup>1</sup>

## <s> Restaurant

a<sup>1</sup> b<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

## b Restaurant

a<sup>1</sup> b<sup>1</sup>

## c Restaurant

</s><sup>1</sup>

# A busy night at the restaurant

## Unigram Restaurant

a<sup>3</sup> b<sup>3</sup> c<sup>1</sup> </s><sup>1</sup>

## <s> Restaurant

a<sup>1</sup> b<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>1</sup>

## b Restaurant

a<sup>1</sup> b<sup>1</sup>

## c Restaurant

</s><sup>1</sup>

# A busy night at the restaurant

## Unigram Restaurant

$a^3$   $b^3$   $c^1$   $\langle /s \rangle^1$

## $\langle s \rangle$ Restaurant

$a^1$   $b^1$

## a Restaurant

$a^2$   $b^1$   $c^1$

## b Restaurant

$a^2$   $b^1$

## c Restaurant

$\langle /s \rangle^1$

# A busy night at the restaurant

## Unigram Restaurant

a<sup>3</sup> b<sup>3</sup> c<sup>1</sup> </s><sup>1</sup>

## <s> Restaurant

a<sup>1</sup> b<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>2</sup>

## b Restaurant

a<sup>2</sup> b<sup>1</sup>

## c Restaurant

</s><sup>1</sup>



# A busy night at the restaurant

## Unigram Restaurant

a<sup>3</sup> b<sup>3</sup> c<sup>1</sup> </s><sup>1</sup>

## <s> Restaurant

a<sup>1</sup> b<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>2</sup>

## b Restaurant

a<sup>2</sup> b<sup>1</sup>

## c Restaurant

</s><sup>2</sup>

# A busy night at the restaurant

## Unigram Restaurant

a<sup>3</sup> b<sup>3</sup> c<sup>1</sup> </s><sup>1</sup>

## <s> Restaurant

a<sup>1</sup> b<sup>1</sup>

## a Restaurant

a<sup>2</sup> b<sup>1</sup> c<sup>2</sup>

## b Restaurant

a<sup>2</sup> b<sup>1</sup>

## c Restaurant

</s><sup>2</sup>

As you see more data, bottom restaurants do more work.

b following a

$$= \frac{1-\delta}{\theta+5} + \frac{\theta+3\delta}{\theta+5} p^{(b)} \quad (4)$$

$$= \frac{1-\delta}{\theta+5} + \frac{\theta+3\delta}{\theta+5} \left( \frac{3-\delta}{\theta+8} + \frac{\theta+4\delta}{\theta+8} \frac{1}{V} \right) \quad (5)$$

(6)

b following a

$$= \frac{1-\delta}{\theta+5} + \frac{\theta+3\delta}{\theta+5} p^{(b)} \quad (4)$$

$$= \frac{1-\delta}{\theta+5} + \frac{\theta+3\delta}{\theta+5} \left( \frac{3-\delta}{\theta+8} + \frac{\theta+4\delta}{\theta+8} \frac{1}{V} \right) \quad (5)$$

(6)

b following a

$$= \frac{1-\delta}{\theta+5} + \frac{\theta+3\delta}{\theta+5} p^{(b)} \quad (4)$$

$$= \frac{1-\delta}{\theta+5} + \frac{\theta+3\delta}{\theta+5} \left( \frac{3-\delta}{\theta+8} + \frac{\theta+4\delta}{\theta+8} \frac{1}{V} \right) \quad (5)$$

(6)

0.23

a following b

$$= \frac{2-\delta}{\theta+3} + \frac{\theta+2\delta}{\theta+3} p(a) \quad (7)$$

$$= \frac{2-\delta}{\theta+3} + \frac{\theta+2\delta}{\theta+3} \left( \frac{3-\delta}{\theta+8} + \frac{\theta+4\delta}{\theta+8} \frac{1}{V} \right) \quad (8)$$

(9)

a following b

$$= \frac{2-\delta}{\theta+3} + \frac{\theta+2\delta}{\theta+3} p(a) \quad (7)$$

$$= \frac{2-\delta}{\theta+3} + \frac{\theta+2\delta}{\theta+3} \left( \frac{3-\delta}{\theta+8} + \frac{\theta+4\delta}{\theta+8} \frac{1}{V} \right) \quad (8)$$

$$(9)$$

a following b

$$= \frac{2-\delta}{\theta+3} + \frac{\theta+2\delta}{\theta+3} p(a) \quad (7)$$

$$= \frac{2-\delta}{\theta+3} + \frac{\theta+2\delta}{\theta+3} \left( \frac{3-\delta}{\theta+8} + \frac{\theta+4\delta}{\theta+8} \frac{1}{V} \right) \quad (8)$$

(9)

0.55



$\langle /s \rangle$  following  $\langle s \rangle$

$$= \frac{\theta + 2\delta}{\theta + 2} p(\langle /s \rangle) \quad (10)$$

$$= \frac{\theta + 2\delta}{\theta + 2} \left( \frac{1 - \delta}{\theta + 8} + \frac{\theta + 4\delta}{\theta + 8} \frac{1}{V} \right) \quad (11)$$

(12)

$\langle /s \rangle$  following  $\langle s \rangle$

$$= \frac{\theta + 2\delta}{\theta + 2} p(\langle /s \rangle) \quad (10)$$

$$= \frac{\theta + 2\delta}{\theta + 2} \left( \frac{1 - \delta}{\theta + 8} + \frac{\theta + 4\delta}{\theta + 8} \frac{1}{V} \right) \quad (11)$$

$$(12)$$

$\langle /s \rangle$  following  $\langle s \rangle$

$$= \frac{\theta + 2\delta}{\theta + 2} p(\langle /s \rangle) \quad (10)$$

$$= \frac{\theta + 2\delta}{\theta + 2} \left( \frac{1 - \delta}{\theta + 8} + \frac{\theta + 4\delta}{\theta + 8} \frac{1}{V} \right) \quad (11)$$

(12)

0.08