



# Chinese Restaurants and Backoff

Natural Language Processing: Jordan  
Boyd-Graber

University of Colorado Boulder

SEPTEMBER 10, 2014

# Roadmap

---

After this class, 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

## Outline

---

How does a CRP encode a probability distribution?

How do many CRPs encode backoff?

Language Model Probabilities

## 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

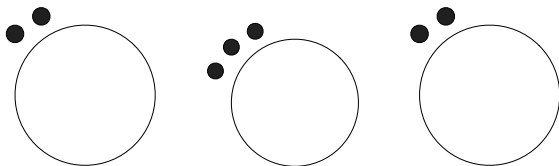
---

- 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.

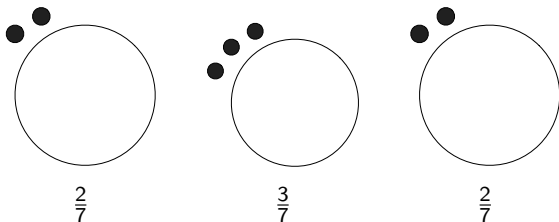




## The Chinese Restaurant as a Distribution

---

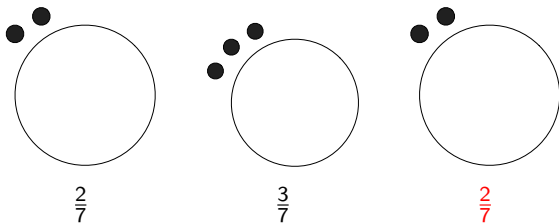
To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



## The Chinese Restaurant as a Distribution

---

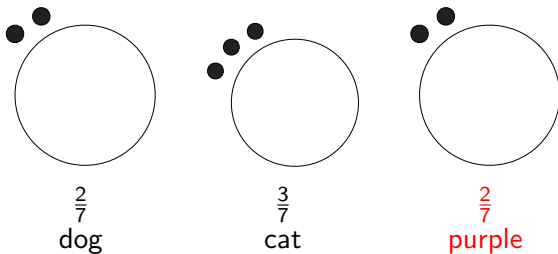
To generate a word, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



## The Chinese Restaurant as a Distribution

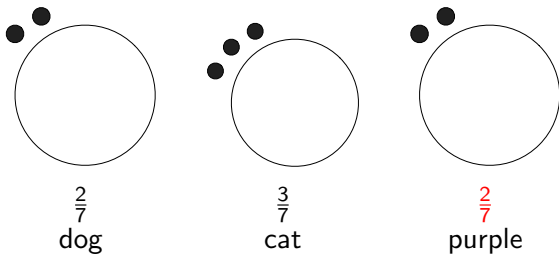
---

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.



## 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.

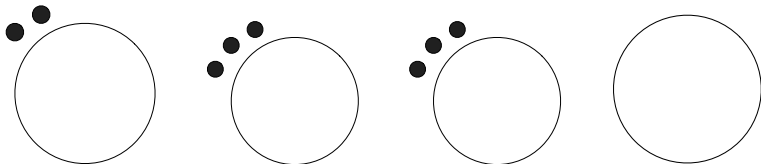


But this is just Maximum Likelihood

Why are we talking about Chinese Restaurants?

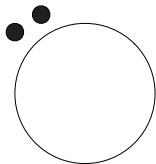
## Always one more table ...

---

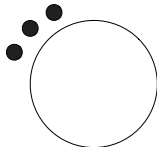


## Always one more table ...

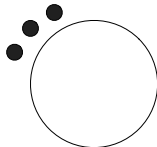
---



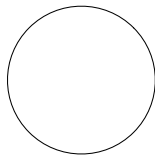
$$\frac{2}{7+\alpha}$$



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



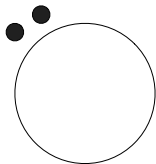
$$\frac{2}{7+\alpha}$$



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

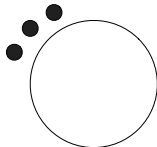
## Always one more table ...

---



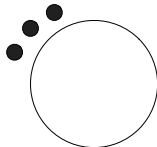
$$\frac{2}{7+\alpha}$$

dog



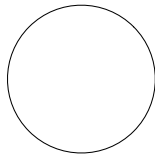
$$\frac{3}{7+\alpha}$$

cat



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

purple

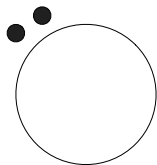


$$\frac{\alpha}{7+\alpha}$$

???

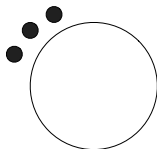
## Always one more table ...

---



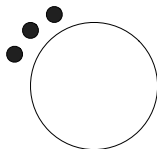
$$\frac{2}{7+\alpha}$$

dog



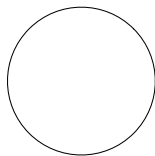
$$\frac{3}{7+\alpha}$$

cat



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

purple



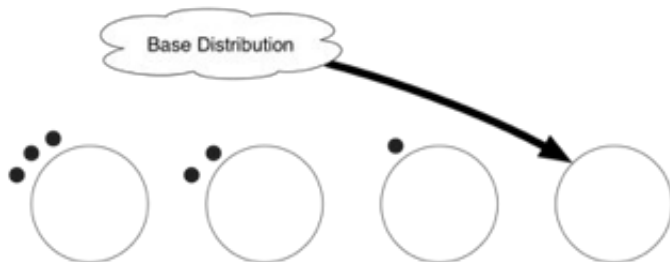
$$\frac{\alpha}{7+\alpha}$$

???



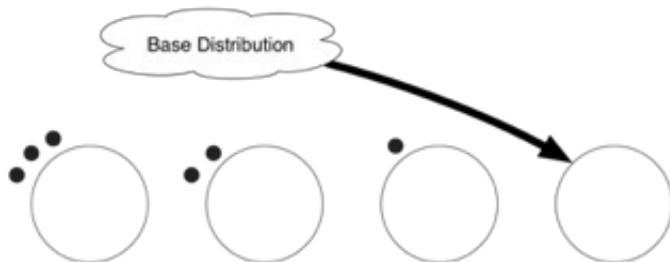
## What to do with a new table?

---



## What to do with a new table?

---

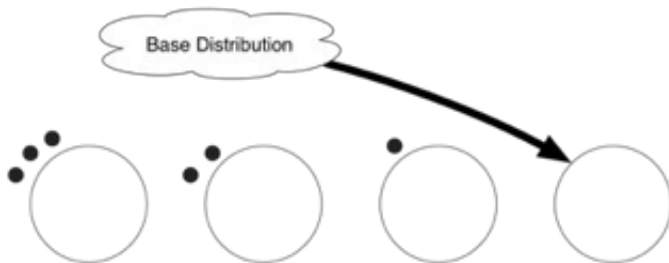


What can be a base distribution?

- 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)

## Outline

---

How does a CRP encode a probability distribution?

How do many CRPs encode backoff?

Language Model Probabilities

## A hierarchy of Chinese Restaurants

---



## Seating Assignments

Dataset:

---

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

## Seating Assignments

Dataset:

---

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

Unigram Restaurant

<s> Restaurant

a Restaurant

b Restaurant

c Restaurant

## Seating Assignments

Dataset:

---

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

Unigram Restaurant

<s> Restaurant

$*$ <sup>1</sup>

b Restaurant

a Restaurant

c Restaurant



## Seating Assignments

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

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

Unigram Restaurant

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

<s> Restaurant

a<sup>1</sup>

a Restaurant

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

b Restaurant

c Restaurant

## Seating Assignments

Dataset:

---

<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

Dataset:

---

<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> b<sup>1</sup>

b Restaurant

c Restaurant

## Seating Assignments

Dataset:

---

<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> b<sup>1</sup>

b Restaurant

c Restaurant

## Seating Assignments

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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

Dataset:

---

<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>

## Outline

---

How does a CRP encode a probability distribution?

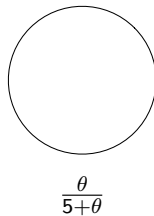
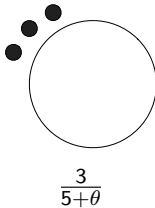
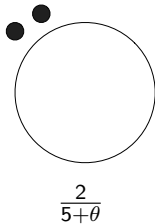
How do many CRPs encode backoff?

Language Model Probabilities



## The rich get richer

---



## Computing the Probability of an Observation

---

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

- Word type  $\textcolor{red}{x}$
- 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

---

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

- Word type  $x$
- Seating assignments  $\vec{s}$
- Concentration  $\theta$
- Context  $\mathbf{u}$
- 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  $\pi(\mathbf{u})$

## 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>

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

## 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)

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

## 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 (GS, later, but not for language models)

## Exercise

---

- Start with restaurant we had before
- Assume you see `<s> b b a c </s>`; 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 `</s>` following `<s>`