

---

**11-442 / 11-642:  
Search Engines**

**Best-Match Retrieval:  
Statistical Language Models**

Jamie Callan  
Carnegie Mellon University  
[callan@cs.cmu.edu](mailto:callan@cs.cmu.edu)

---

**Introduction**

**The last lecture introduced two best match retrieval models**

- The vector space model (VSM), Okapi BM25

**This lecture introduces statistical language models**

- A newer form of probabilistic retrieval model
- Inspired by work in speech recognition and machine translation

**Very different underlying theory**

- But, implemented using many of the same values used by the vector space and Okapi BM25 retrieval models
  - Different explanations, but similar effects

## Outline

### Statistical language models

- Introduction to language models
- Query likelihood
- Kullback-Leibler (KL) Divergence
- Indri

3

© 2017 Jamie Callan

## Statistical Language Models

**A statistical language model uses a probability distribution to determine the probability of terms or term sequences**

### Common types of models:

- Unigrams:  $p(t_i | \theta)$        $p(\text{"search"} | \theta)$
- Bigrams:  $p(t_i | t_{i-1}, \theta)$        $p(\text{"engines"} | \text{"search"}, \theta)$
- Trigrams:  $p(t_i | t_{i-2}, t_{i-1}, \theta)$        $p(\text{"class"} | \text{"search"}, \text{"engines"}, \theta)$

**Language model**



**Bigrams, trigrams, etc., haven't helped much for IR (so far)**

- But they are a key idea in speech recognition

4

© 2017 Jamie Callan

## Statistical Language Models

### Word histogram

Term	$tf_d$
camera	17
image	13
picture	11
up	8
movie	8
like	7
mode	7
software	7
red	6
digital	5
eye	5
shutter	5
sony	5

### Unigram language model

Term	$P(t \theta_d)$
camera	0.09551
image	0.07303
picture	0.06180
up	0.04494
movie	0.04494
like	0.03933
mode	0.03933
software	0.03933
red	0.03371
digital	0.02809
eye	0.02809
shutter	0.02809
sony	0.02809

$$P(t|d) = \frac{tf_{t,d}}{length_d}$$

5

© 2017 Jamie Callan

## A Language Model can be Created From Any Language Sample

### Examples

- A document collection
- A document
- Also sentence, paragraph, chapter, ...
- A query

### This is similar to the vector space

- A vector can be created for any sample of text

Term	$P(t \theta_d)$
camera	0.09551
image	0.07303
picture	0.06180
up	0.04494
movie	0.04494
like	0.03933
mode	0.03933
software	0.03933
red	0.03371
digital	0.02809
eye	0.02809
shutter	0.02809
sony	0.02809

6

© 2017 Jamie Callan

## Terminology

**d:** A document

**$\theta_d$ :** The language model for document d

**q:** A query

**$\theta_q$ :** The language model for query q

### **d and $\theta_d$ are not the same thing**

- The document and the model of the document are different
- However, to simplify notation, we will often treat them as the same
  - i.e.,  $p(d | q)$  instead of  $p(d | \theta_q)$  or  $p(\theta_d | \theta_q)$

7

© 2017 Jamie Callan

## Retrieval Model Based Upon Statistical Language Models

### **A document d defines a probability distribution $\theta_d$ over index terms**

- E.g., the probability of generating/observing an index term

### **A query q also defines a probability distribution $\theta_q$**

- A sparse distribution (more on this later)

<u>t</u>	<u>p(t)</u>
Apple	0.0204
Apple	0.0001
Campaign	0.0034
Obama	0.0102
:	:

### **How is a language model used to rank document d?**

- Rank d by  $p(d | \theta_q)$  “query likelihood”
- Rank d by the similarity of  $\theta_d$  and  $\theta_q$  “KL divergence”

8

© 2017 Jamie Callan

## Outline

### Statistical language models

- Introduction to language models
- Query likelihood
- Kullback-Leibler (KL) Divergence
- Indri

9

© 2017 Jamie Callan

## Rank by $P(d|q)$ : The Query-Likelihood Approach

**Task:** Rank documents by  $p(d|q)$

- Given a query  $q$ , what is the probability of document  $d$ ?

- **Problems**

- $q$  is a very sparse language model (few terms)
- $q$  contains very little frequency information

Term	$P(t \theta_Q)$
digital	0.5
camera	0.5

- **Solution**

$$p(d|q) = \frac{p(q|d)p(d)}{p(q)}$$

**Bayes rule**

$$\propto p(q|d)p(d)$$

**Drop document-independent term**

**Key issues**

- How are  $p(q|d)$  and  $p(d)$  estimated?

10

© 2017 Jamie Callan

## Query-Likelihood Approach: Estimating $p(d)$

It is simple and convenient to assume that  $p(d)$  is uniform

- Later we will consider non-uniform  $p(d)$

So...

$$\begin{aligned} p(d | q) &\propto p(q | d)p(d) \\ &\propto p(q | d) \end{aligned} \quad \text{Drop the constant term (for now)}$$

11

© 2017 Jamie Callan

## What is $p(q|d)$ ? Simple Unigram Approach

Assume that a query is composed of independent terms

- Unjustified, but convenient
- Doesn't hurt effectiveness of other models (e.g., vector space)

Then ...

$$p(q | d) = \prod_{q_i \in q} p(q_i | d)$$

This should look a little familiar (although the notation differs)

- The score of  $(q, d)$  is based on the scores of  $(q_i, d)$
- Similar to what we saw with BM25 and the vector space
- How is  $p(q_i|d)$  calculated?

12

© 2017 Jamie Callan

## Estimating $p(q_i|d)$

**Maximum likelihood estimation (MLE) is a simple approach**

$$p_{MLE}(q_i | d) = \frac{tf_{q_i,d}}{length(d)}$$

**Is this a good estimate?**

- Estimates are based on small samples (a single document)
  - So, perhaps not very accurate
- $p_{MLE}(q_i|d) = 0$  if  $q_i$  isn't in document  $d$

$$p(q | d) = \prod_{q_i \in q} p_{MLE}(q_i | d) = 0 \quad \text{Boolean AND}$$

- $q_i$  could be a good description of document  $d$ , even if it does not occur in document  $d$

13

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Smoothing

**Smoothing is used to solve two problems in the language modeling framework**

- Imprecise probability estimates from MLE
- Probability estimates for unobserved terms
  - E.g., query terms that don't occur in the document

**But ... there are many smoothing methods ... which to use?**

14

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Jelinek-Mercer (“Mixture Model”) Smoothing

### Linear interpolation with a reference language model

$$p(q_i | d) = (1 - \lambda)p_{MLE}(q_i | d) + \lambda p_{MLE}(q_i | C)$$

### Smoothing decreases as $\lambda \rightarrow 0$

#### What value of $\lambda$ is best?

- Small  $\lambda$  (little smoothing) is best for short queries
- Larger  $\lambda$  (more smoothing) is better for long queries
- We will see later in the lecture why this is true

15

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Bayesian Smoothing With Dirichlet Priors

### Method #2: Bayesian smoothing using Dirichlet priors

$$p_{MLE}(q_i | d) = \frac{tf_{q_i,d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu}$$
$$p_{MLE}(q_i | C) = \frac{ctf_{q_i}}{|length_{tokens}(C)|}$$

#### What value of $\mu$ is best?

- $\mu$  in [1,000-10,000] seems best
  - Some people think  $\mu$  is approximately related to average document length ... others disagree

16

© 2017 Jamie Callan



## Estimating $p(q_i|d)$ : Smoothing

We know that tf.idf weights are effective  
...how do they relate to statistical language models

tf is easy to see

$$p_{MLE}(q_i | d) = \frac{tf_{q_i, d}}{length(d)}$$

idf is a little harder to see

- Jelenick-Mercer smoothing provides an idf-like effect
- Next slide...

17

© 2017 Jamie Callan

## No Smoothing Example:

$$p(q | d) = \prod_{q_i \in q} p_{MLE}(q_i | d)$$

Two query terms:  $p(\text{apple} | C) = 0.01$ ,  $p(\text{ipod} | C) = 0.001$   
“frequent” “rare”

**Document  $d_1$**

doclen = 50,  $tf_{\text{apple}} = 2$ ,  $tf_{\text{ipod}} = 3$

$(2/50) \times$

$(3/50) =$

$0.04 \times 0.06 =$

$p(q|d_1) = 0.0024$

**Document  $d_2$**

doclen = 50,  $tf_{\text{apple}} = 3$ ,  $tf_{\text{ipod}} = 2$

$(3/50) \times$

$(2/50) =$

$0.06 \times 0.04 =$

$p(q|d_2) = 0.0024$

**The model does not distinguish between frequent and rare terms**

18

© 2017 Jamie Callan

### Jelinek-Mercer Smoothing Example:

$$p(q|d) = \prod_{q_i \in q} (1-\lambda)p_{MLE}(q_i|d) + \lambda p_{MLE}(q_i|C)$$

Two query terms:  $p(\text{apple} | C) = 0.01$ ,  $p(\text{ipod} | C) = 0.001$   
 $\lambda = 0.4$       “frequent”      “rare”

**Document  $d_1$**

**doclen = 50,  $tf_{\text{apple}} = 2$ ,  $tf_{\text{ipod}} = 3$**

$$(0.6 \times (2/50) + 0.4 \times 0.01) \times$$

$$(0.6 \times (3/50) + 0.4 \times 0.001) =$$

$$(0.6 \times 0.04 + 0.004) \times$$

$$(0.6 \times 0.06 + 0.0004) =$$

$$(0.024 + 0.004) \times (0.036 + 0.0004) =$$

$$p(q|d_1) = 0.001019$$

**Document  $d_2$**

**doclen = 50,  $tf_{\text{apple}} = 3$ ,  $tf_{\text{ipod}} = 2$**

$$(0.6 \times (3/50) + 0.4 \times 0.01) \times$$

$$(0.6 \times (2/50) + 0.4 \times 0.001) =$$

$$(0.6 \times 0.06 + 0.004) \times$$

$$(0.6 \times 0.04 + 0.0004) =$$

$$(0.036 + 0.004) \times (0.024 + 0.0004) =$$

$$p(q|d_2) = 0.000976$$

**The model does distinguish between frequent and rare terms**

19

© 2017 Jamie Callan

### Why Does Smoothing Work?

$$p(q|d) = \prod_{q_i \in q} (1-\lambda)p_{MLE}(q_i|d) + \lambda p_{MLE}(q_i|C)$$

**doclen = 50,  $tf_{\text{apple}} = 2$ ,  $tf_{\text{ipod}} = 2$**

$$p(\text{apple} | d) = 0.6 \times (2/50) + 0.4 \times 0.010 = 0.0280$$

$$p(\text{ipod} | d) = 0.6 \times (2/50) + 0.4 \times 0.001 = 0.0244$$

**What is the effect of matching one additional instance of a term?**

$$p_8(\text{apple} | d) = 0.6 \times (1/50) = 0.012$$

$$p_8(\text{ipod} | d) = 0.6 \times (1/50) = 0.012$$

– The unsmoothed effect of each match is the same

**The incremental value of matching a term is multiplied by the  $p(q|d)$  of other query terms**

$$- tf_{\text{apple}} = 3, tf_{\text{ipod}} = 2: p(q|d) = (0.028 + \underline{0.012}) \times 0.0244 = 0.000976$$

$$- tf_{\text{apple}} = 2, tf_{\text{ipod}} = 3: p(q|d) = 0.028 \times (0.0244 + \underline{0.012}) = 0.001019$$

20

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : The idf-like Effect of Jelenick-Mercer Smoothing

$$\begin{aligned}
 p(q|d) &= \prod_{q_i \in q} p(q_i|d) \\
 &= \prod_{q_i \in q} ((1-\lambda)p_{MLE}(q_i|d) + \lambda p_{MLE}(q_i|C)) && \text{J-M smoothing} \\
 &= \prod_{q_i \in q} ((1-\lambda)p_{MLE}(q_i|d) + \lambda p_{MLE}(q_i|C)) \frac{\lambda p_{MLE}(q_i|C)}{\lambda p_{MLE}(q_i|C)} && \text{Multiply by 1} \\
 &= \prod_{q_i \in q} \left( \left( \frac{(1-\lambda)p_{MLE}(q_i|d)}{\lambda p_{MLE}(q_i|C)} + 1 \right) \lambda p_{MLE}(q_i|C) \right) && \text{Recombine} \\
 &= \prod_{q_i \in q} \left( \frac{(1-\lambda)p_{MLE}(q_i|d)}{\lambda p_{MLE}(q_i|C)} + 1 \right) \prod_{q_i \in q} \lambda p_{MLE}(q_i|C) && \text{Recombine} \\
 &\propto \prod_{q_i \in q} \left( \frac{(1-\lambda)p_{MLE}(q_i|d)}{\lambda p_{MLE}(q_i|C)} + 1 \right) && \text{Drop constant}
 \end{aligned}$$

21

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Jelinek-Mercer (“Mixture Model”) Smoothing

**HW2 requires you to think about the effects of parameters**

$$p(q_i|d) = (1-\lambda)p_{MLE}(q_i|d) + \lambda p_{MLE}(q_i|C)$$

**What happens when  $\lambda$  approaches 0?**

- There is no smoothing (no “idf like” effect)

$$p(q_i|d) = (1-\lambda)p_{MLE}(q_i|d)$$

22

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Jelinek-Mercer (“Mixture Model”) Smoothing

HW2 requires you to think about the effects of parameters

$$p(q_i | d) = (1 - \lambda)p_{MLE}(q_i | d) + \lambda p_{MLE}(q_i | C)$$

How important is smoothing to queries of different lengths?

- Short queries
  - Few query terms, so (usually) every query term must match
  - “idf weighting” is less important, so small  $\lambda$  is best
- Long queries
  - Many query terms, so (usually) most query terms must match
  - “idf weighting” is more important
  - Give priority to the less common terms, so larger  $\lambda$  is best

23

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Bayesian Smoothing With Dirichlet Priors

HW2 requires you to think about the effects of parameters

$$p(q_i | d) = \frac{tf_{q_i,d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu}$$

What happens when  $\mu$  approaches 0?

- There is no smoothing (maximum likelihood only)

$$p(q_i | d) = \frac{tf_{q_i,d}}{length(d)}$$

24

© 2017 Jamie Callan

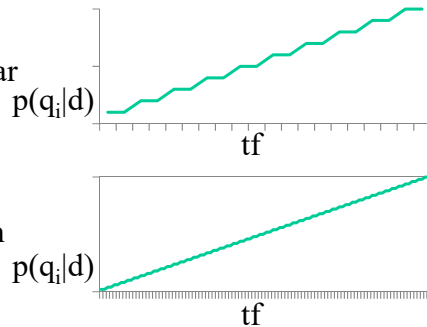
## Estimating $p(q_i|d)$ : Bayesian Smoothing With Dirichlet Priors

HW2 requires you to think about the effects of parameters

$$p(q_i | d) = \frac{tf_{q_i,d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu}$$

How important is smoothing to documents of different lengths?

- Short documents
  - Probabilities are more granular
  - Larger  $\mu$  is more important
- Long documents
  - Probabilities are more smooth
  - Larger  $\mu$  is less important



25

© 2017 Jamie Callan

## Estimating $p(q_i|d)$ : Two-Stage Smoothing

Mixture modeling and Bayesian smoothing with Dirichlet priors are both effective and common ... which should you use?

- They do different things
  - **Mixture model:** Compensate for differences in word importance
    - » “idf effects”
  - **Dirichlet prior:** Improves the estimate of the document model
    - » “small sample”, “unseen words”
- Two-stage smoothing gives the best effects of both methods

$$p(q_i | d) = (1 - \lambda) \frac{tf_{q_i,d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu} + \lambda p_{MLE}(q_i | C)$$

26

© 2017 Jamie Callan

## Query Likelihood With Two-Stage Smoothing: Putting it All Together

$$\begin{aligned} p(q|d) &= \prod_{q_i \in q} p(q_i|d) \\ &= \prod_{q_i \in q} \left( (1-\lambda) \frac{tf_{q_i,d} + \mu p_{MLE}(q_i|C)}{length(d) + \mu} + \lambda p_{MLE}(q_i|C) \right) \\ &= \prod_{q_i \in q} \left( (1-\lambda) \frac{tf_{q_i,d} + \mu \frac{ctf(q_i)}{length(c)}}{length(d) + \mu} + \lambda \frac{ctf(q_i)}{length(c)} \right) \end{aligned}$$

**tf<sub>q,d</sub>:** Term frequency of term q in document d  
**ctf (q):** Term frequency of term q in the entire collection  
**length (d):** Total number of word occurrences in document d  
**length (c):** Total number of word occurrences in collection c

27

© 2017 Jamie Callan

## Outline

### Statistical language models

- Introduction to language models
- Query likelihood
- Kullback-Leibler (KL) Divergence
- Indri

28

© 2017 Jamie Callan

## Retrieval Model Based Upon Statistical Language Models

### How is a language model used to rank document d?

- Rank d by  $p(d \mid \theta_q)$  “query likelihood”
- Rank d by  $p(q \mid \theta_d)$
- Rank d by the similarity of  $\theta_d$  and  $\theta_q$  “KL divergence”

In some cases, some of these approaches turn out to be equivalent...

29

© 2017 Jamie Callan

## KL Divergence: Measuring the Similarity of Language Models

**Kullback-Leibler divergence measures the relative entropy between two probability mass functions**

$$KL(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

**Language models are probability mass functions**

- The probability of observing a term in a sample of text
- Why not use KL divergence to measure the similarity of q and d?

30

© 2017 Jamie Callan

## Measuring the Similarity of Statistical Language Models

### Ranking by (negative) KL-divergence

$$\begin{aligned}
 -KL(q \| d) &= - \sum_{w \in V} p(w | q) \log \frac{p(w | q)}{p(w | d)} \\
 &= - \sum_{q_i \in q} p(q_i | q) \log \frac{p(q_i | q)}{p(q_i | d)} && \mathbf{P(w|q) = 0 \text{ when } w \notin q} \\
 &= - \left( \sum_{q_i \in q} p(q_i | q) \log p(q_i | q) - \sum_{q_i \in q} p(q_i | q) \log p(q_i | d) \right) \\
 &\propto \sum_{q_i \in q} p(q_i | q) \log p(q_i | d) && \mathbf{Drop \ constant \ term} \\
 &\propto \sum_{q_i \in q} \frac{1}{|q|} \log p(q_i | d) && \mathbf{Uniform \ probabilities \ for \ query \ terms} \\
 &\propto \sum_{q_i \in q} \log p(q_i | d) && \mathbf{Drop \ constant \ term}
 \end{aligned}$$

31

© 2017 Jamie Callan

## Two Different Paths to a Common Destination

Query likelihood ranks by  $p(q | d) = \prod_{q_i \in q} p(q_i | d)$

KL diverge ranks by  $\sum_{q_i \in q} \log p(q_i | d)$

These are rank equivalent  $p(q | d) = \prod_{q_i \in q} p(q_i | d)$

$$\propto \sum_{q_i \in q} \log p(q_i | d)$$

So...the query likelihood and KL divergence approaches are equivalent (when document priors are uniform)

32

© 2017 Jamie Callan



## Outline

### Statistical language models

- Introduction to language models
- Query likelihood
- Kullback-Leibler (KL) Divergence
- Indri

33

© 2017 Jamie Callan

## Inference Nets + Language Models: Indri

The Indri retrieval model combines statistical language models with Bayesian inference networks

- A probabilistic retrieval model
- Structured queries
- Documents with multiple representations
- Documents with structure (a later lecture)

**The theory is sophisticated, but the underlying ideas are familiar**

34

© 2017 Jamie Callan

## Inference Nets + Language Models: Indri

**Documents**

....

$d_i$

....

35

© 2017 Jamie Callan

## Inference Nets + Language Models: Indri

**Documents**

....

$d_i$

....

**Representation Nodes**  
(terms, phrases, ...)

$r_1$

$r_2$

$r_3$

$r_4$

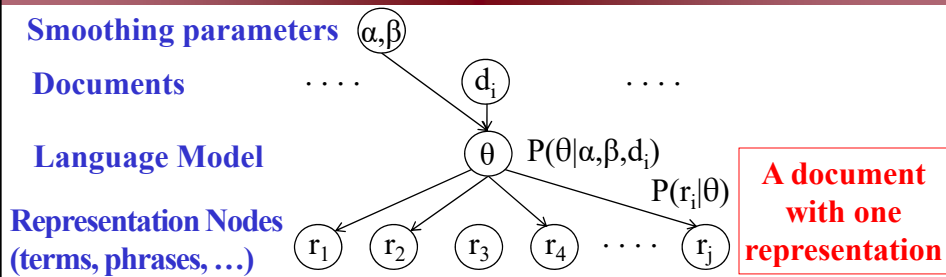
....

$r_j$

36

© 2017 Jamie Callan

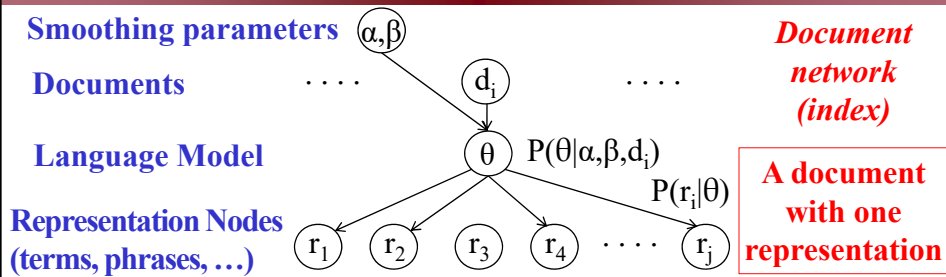
## Inference Nets + Language Models: Indri



37

© 2017 Jamie Callan

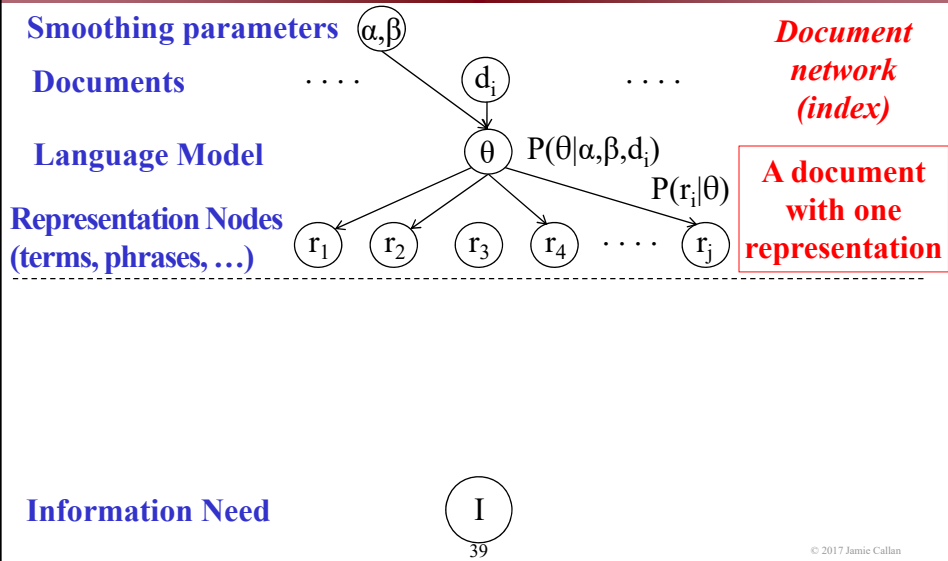
## Inference Nets + Language Models: Indri



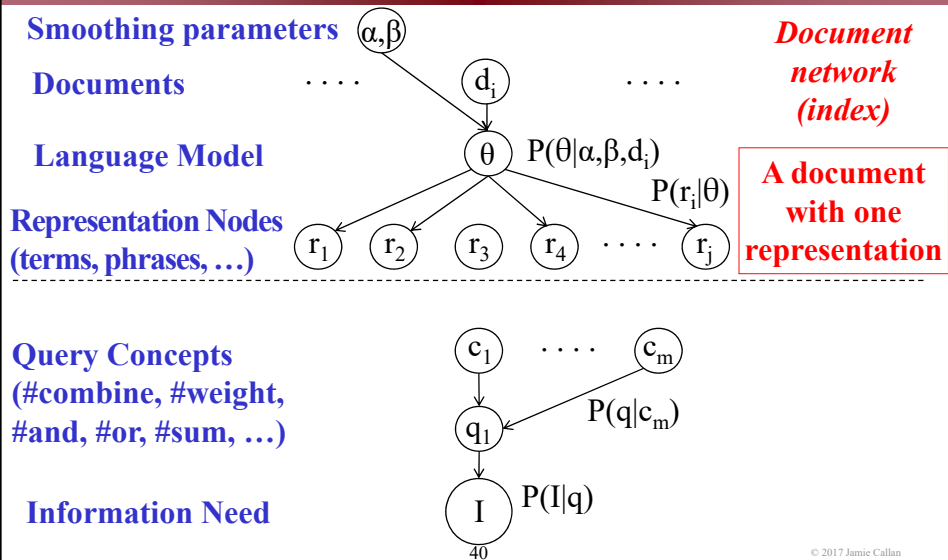
38

© 2017 Jamie Callan

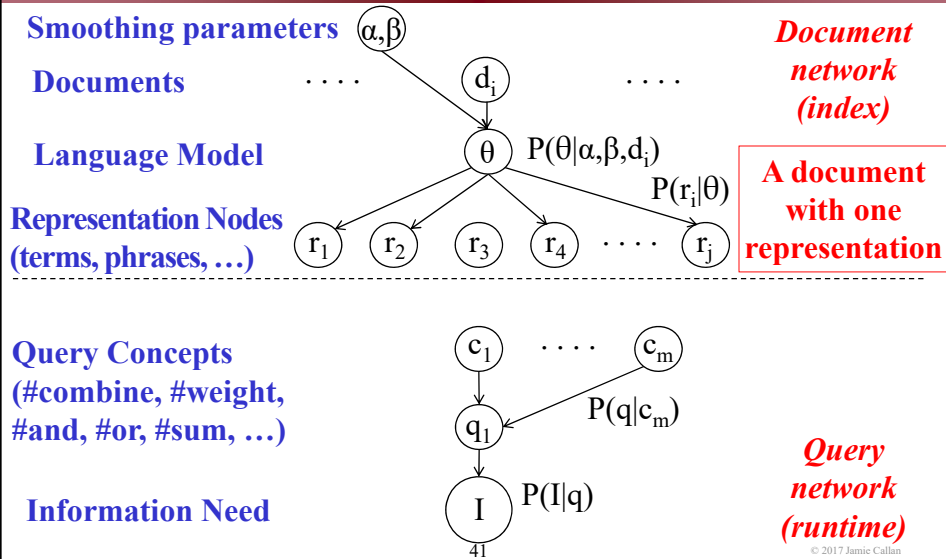
## Inference Nets + Language Models: Indri



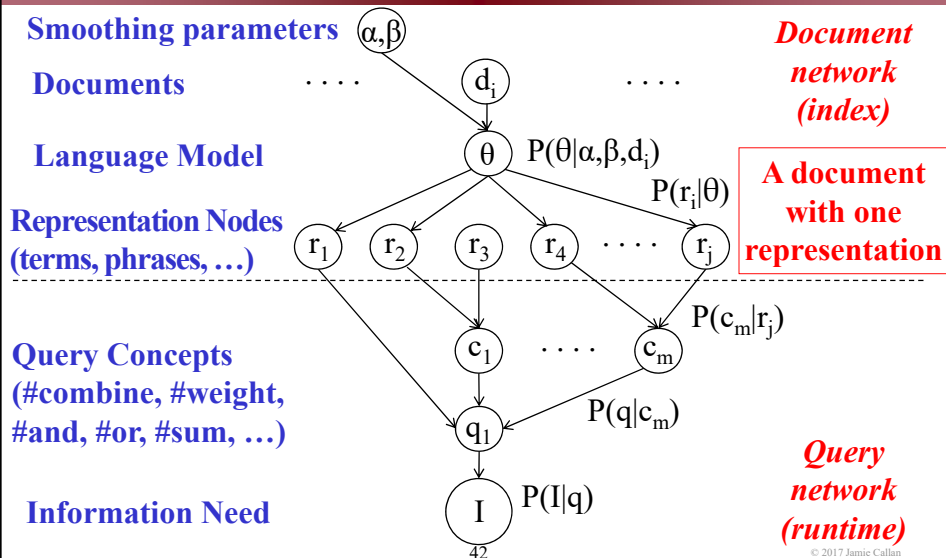
## Inference Nets + Language Models: Indri



## Inference Nets + Language Models: Indri



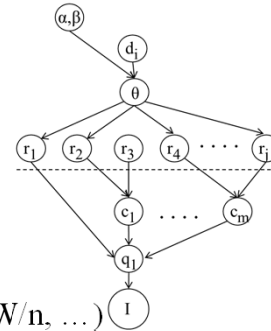
## Inference Nets + Language Models: Indri



## Inference Nets + Language Models: Indri

### This isn't as complicated as it looks

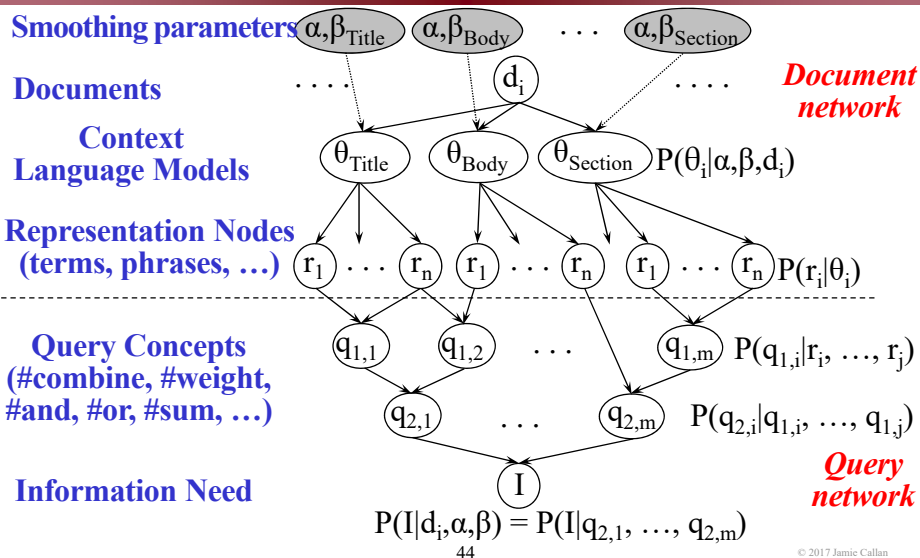
- Document + smoothing parameters  $(\alpha, \beta)$   
→ language model  $(\theta)$
- The language model vocabulary is defined by representation nodes  $(r_i)$ 
  - Terms stored in the index
  - Query operators that create index terms  
» QryIop ( $\#SYN$ ,  $\#NEAR/n$ ,  $\#WINDOW/n$ , ...)
- Information needs are represented by queries
- Queries are composed of query operators that combine evidence  
» QrySop ( $\#AND$ ,  $\#OR$ ,  $\#SUM$ ,  $\#WSUM$ , ...)



43

© 2017 Jamie Callan

## Inference Nets + Language Models: Indri



44

© 2017 Jamie Callan

## Indri Term Weights

**Indri can use any language modeling method to estimate  $p(r_i|\theta)$**

- This weight is equivalent to  $p(q_i|d)$  in most language models
- Dirichlet smoothing is the most common choice

$$p(q_i | d) = \frac{tf_{q_i,d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu}$$

- Two-stage smoothing is also used

$$p(q_i | d) = (1 - \lambda) \frac{tf_{q_i,d} + \mu p_{MLE}(q_i | C)}{length(d) + \mu} + \lambda p_{MLE}(q_i | C)$$

45

© 2017 Jamie Callan

## Indri Query Operators

**Operators that map inverted lists to an inverted list are easy**

- E.g., InvertedList+  $\rightarrow$  InvertedList
- E.g., #NEAR/n, #WINDOW/n, #SYN, ...
- To the language model these look just like ordinary terms

**Operators that combine scores are somewhat less clear**

- E.g., ScoreList+  $\rightarrow$  ScoreList
- E.g., #AND, #OR, ...
- How should each query operator combine evidence?

46

© 2017 Jamie Callan

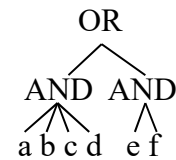
## Indri Query Operators

### AND is the default query operator for most language modeling systems

- Typically implemented as the product of the argument weights

$$p_{and}(q | d) = \prod_{q_i \in q} p(q_i | d)$$

- This is pleasing theoretically, but it has a problem
  - Typically  $\#AND(a\ b\ c\ d) < \#AND(a\ b)$



### Indri's AND operator uses the geometric mean

$$p_{and}(q | d) = \prod_{q_i \in q} p(q_i | d)^{\frac{1}{|q|}}$$

47

© 2017 Jamie Callan

## Indri Query Operators

### We may want to give different weight to different evidence

$\#wand ($

0.7  $\#and$  (time traveler wife)

0.2  $\#and$  ( $\#near/1$  (time traveler)  $\#near/1$  (traveler wife))

0.1  $\#and$  ( $\#window/8$  (time traveler)  $\#window/8$  (traveler wife)))

### Indri's WAND operator (also called WEIGHT) gives this control

- WAND: weighted AND

$$p_{wand}(q | d) = \prod_{q_i \in q} p(q_i | d)^{\frac{w_i}{\sum w_i}}$$

48

© 2017 Jamie Callan



## Indri Query Operators

### Indri provides a NOT operator

$$p_{not}(q | d) = 1 - p(q | d)$$

- NOT operators aren't used a lot in probabilistic systems

### Indri provides an OR operator

$$p_{or}(q | d) = 1 - \prod_{q_i \in q} (1 - p(q_i | d))$$

49

© 2017 Jamie Callan

## Indri Query Operators

### The AND operator assumes that its arguments are estimates of independent probabilities

- E.g., #AND (buy ipod)

### The WSUM operator assumes that its arguments are different ways of estimating the same probability

- E.g., the probability that this document is about “apple”) #WSUM (0.3 apple.title 0.1 apple.url 0.6 apple.body)
- WSUM takes a weighted average of the estimates
  - It should be called WAVG – the name is a historical artifact

$$p_{wsum}(q | d) = \sum_{q_i \in q} \frac{w_i}{\sum w_i} p(q_i | d)$$

50

© 2017 Jamie Callan

## Indri Query Operators

$$\begin{aligned}\text{AND, COMBINE: } p_{and}(q|d) &= \prod_{q_i \in q} p(q_i|d)^{\frac{1}{|q|}} \\ \text{WAND, WEIGHT: } p_{wand}(q|d) &= \prod_{q_i \in q} p(q_i|d)^{\frac{w_i}{w}}, \quad w = \sum w_i \\ \text{OR: } p_{or}(q|d) &= 1 - \prod_{q_i \in q} (1 - p(q_i|d)) \\ \text{WSUM: } p_{wsum}(q|d) &= \sum_{q_i \in q} \frac{w_i}{w} p(q_i|d), \quad w = \sum w_i \\ \text{NOT: } p_{not}(q|d) &= 1 - p(q|d)\end{aligned}$$

51

© 2017 Jamie Callan

## Indri Implementation

**Mostly the Indri query operators are easy to implement**

- See the preceding slide for the score calculations
- **Calculate scores only for documents that contain a query term**
  - Use inverted or score lists – similar to HW1
- Use document length, ctf, and corpus length for smoothing
  - Lookup from the index – see the HW2 web page

**But, one aspect is a little tricky to get right...**

52

© 2017 Jamie Callan

## Indri Implementation

**Query:** #or (a #and (b c) )

**Document:** a

**Query terms b and c do not appear in this document**

**... what is the score of the #AND operator?**

- $tf_{a,d} = 1$        $tf_{b,d} = 0$        $tf_{c,d} = 0$
- Do the usual Indri score calculation
  - So, only smoothing scores for b and c

**This is simple conceptually, but how is it implemented?**

- You don't want to calculate #AND scores for every document
  - ... just the documents that have at least one query term

53

© 2017 Jamie Callan

## Indri Implementation

**Query:** #or (a #and (b c) )

**Document:** a

**Add a new method to all QrySop operators**

`double getDefaultScore (RetrievalModel r, long docid)`

**When any QrySop operator calculates scores**

If the  $i^{\text{th}}$  query argument contains document d  
then call the  $i^{\text{th}}$  query argument's `getScore` method  
else call the  $i^{\text{th}}$  query argument's `getDefaultScore` method

54

© 2017 Jamie Callan

## Indri Implementation

**Query:** #or (a #and (b c) )

**Document:** a

**QrySopScore.getDefaultScore (RetrievalModel r, long docid)**

- The standard Indri SCORE calculation done with  $tf=0$

If  $r == \text{RetrievalModel.Indri}$

$$p_{scoreDefault}(t | docid) = (1 - \lambda) \frac{0 + \mu p_{MLE}(t | C)}{length(docid) + \mu} + \lambda p_{MLE}(t | C)$$

**This is the only difference.  
Do the usual calculation, but with  $tf=0$ .**

55

© 2017 Jamie Callan

## Indri Implementation

**Query:** #or (a #and (b c) )

**Document:** a

**QrySopAnd.getDefaultScore (RetrievalModel r, long docid)**

- The standard Indri AND calculation done on the default score of each argument

If  $r == \text{RetrievalModel.Indri}$

$$p_{andDefault}(q | d) = \prod_{q_i \in q} p_{q_i default}^{1/|q|}(q_i | d)$$

**This is the only difference.  
Call the  $i^{\text{th}}$  query argument's getDefaultScore method.**

56

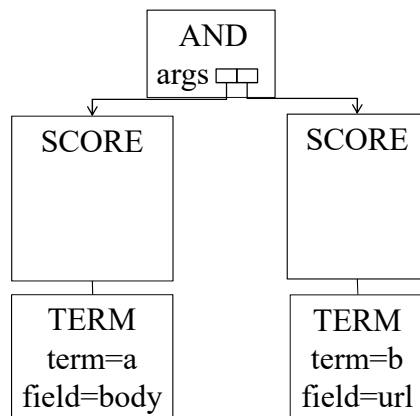
© 2017 Jamie Callan

## QryEval Example

57

© 2017 Jamie Callan

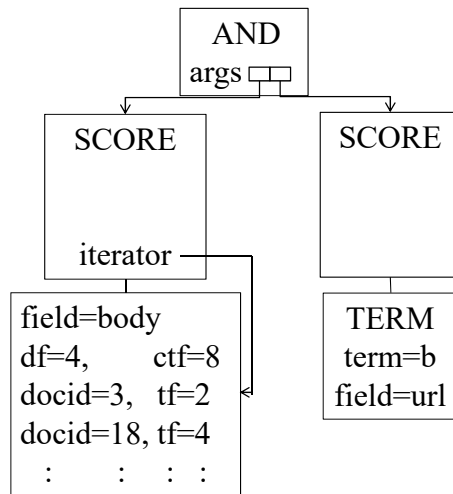
## The initial query: #AND (a.body b.url)



58

© 2017 Jamie Callan

## Query Initialization



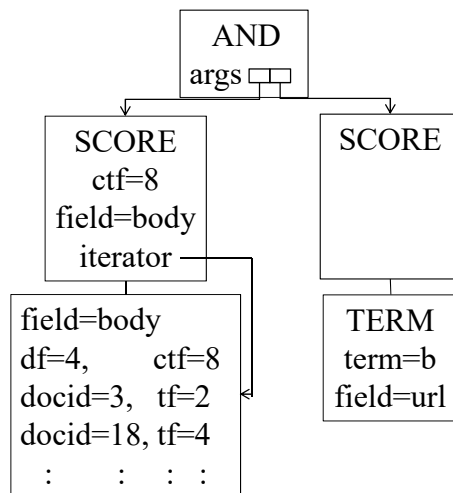
### AND operator initialization

AND initializes its first arg.  
SCORE initializes its arg.  
The result is an inverted list.

59

© 2017 Jamie Callan

## Query Initialization



### AND operator initialization

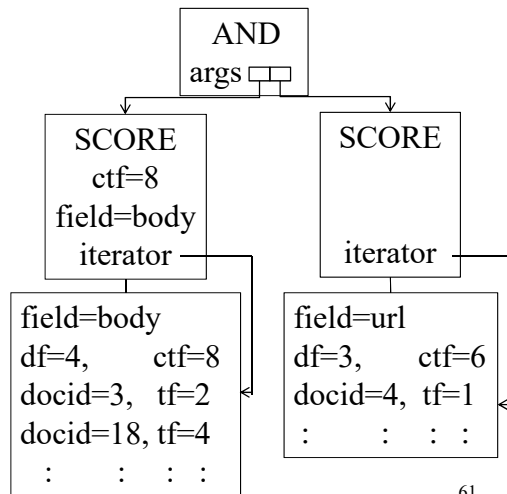
AND initializes its first arg.  
SCORE initializes its arg.  
The result is an inverted list.

The SCORE operator  
caches information that it  
will need later.

60

© 2017 Jamie Callan

## Query Initialization

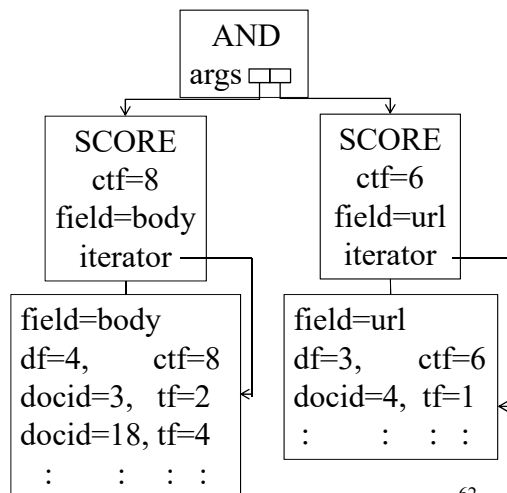


### AND operator initialization

AND initializes its second arg.  
SCORE initializes its arg.  
The result is an inverted list.

© 2017 Jamie Callan

## Query Initialization



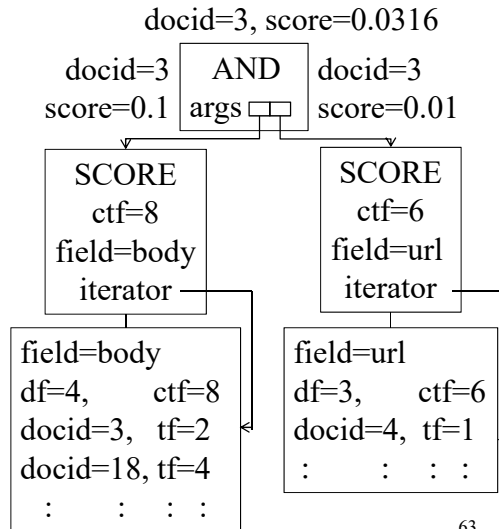
### AND operator initialization

AND initializes its second arg.  
SCORE initializes its arg.  
The result is an inverted list.

The SCORE operator  
caches information that it  
will need later.

© 2017 Jamie Callan

## Call to docIteratorHasMatch & getScore (First Time)



### AND Operator Evaluation

Min document is 3.

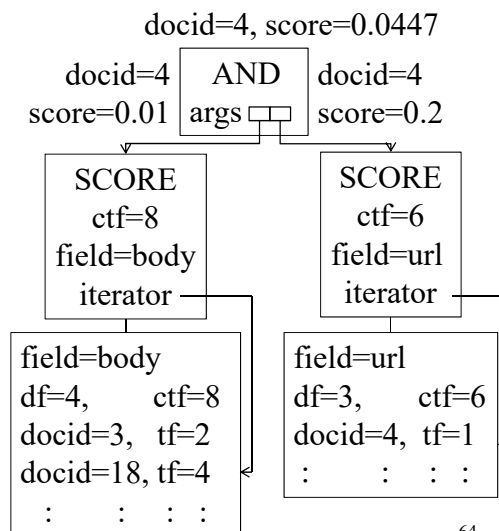
args[0] matches, so call  
args[0].getScore ().  
Suppose the result is 0.1.

args[1] does not match, so call  
args[1].getDefaultScore(3).  
Suppose the result is 0.01.

$$\text{Score}_{\text{AND}}(3) = (0.1^{0.5} \times 0.01^{0.5})$$

© 2017 Jamie Callan

## Call to docIteratorHasMatch & getScore (Second Time)



### AND Operator Evaluation

Min document is 4.

args[0] does not match, so call  
args[0].getDefaultScore (4).  
Suppose the result is 0.01.

args[1] matches, so call  
args[1].getScore().  
Suppose the result is 0.2.

$$\text{Score}_{\text{AND}}(4) = (0.01^{0.5} \times 0.2^{0.5})$$

© 2017 Jamie Callan



## Default Belief Scores Are Only a Small Complication

### When evaluating a query argument

- If it matches the current document
  - Ask the query argument to calculate the document score for the current document
  - Else ask the query argument to calculate a default score for the current document
    - » E.g., a SCORE operator (the example given)
    - » E.g., an OR operator (similar logic)

65

© 2017 Jamie Callan

## Default Belief Scores Are Only a Small Complication

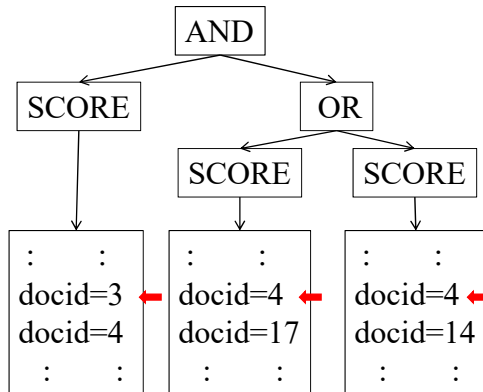
### Which types of query operators calculate default scores?

- If an operator calculates scores
  - ... it also calculates default scores
- QrySop operators calculate default scores
- QryIop operators do not calculate default scores

66

© 2017 Jamie Callan

## Call to docIteratorHasMatch & getScore (First Time)



### AND Operator Evaluation

Min document is 3.

args[0] matches, so call  
args[0].getScore ().

Suppose the result is 0.3.

args[1] does not match, so call  
args[1].getDefaultScore(3).

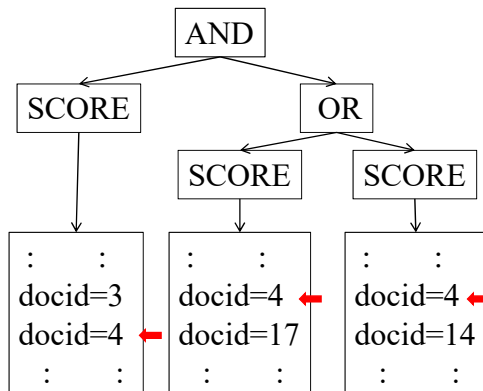
OR calls getDefaultScore(3)  
for all of its args and computes  
a score. Suppose it is 0.01.

$$\text{Score}_{\text{AND}}(3) = (0.3^{0.5} \times 0.01^{0.5})$$

67

© 2017 Jamie Callan

## Call to docIteratorHasMatch & getScore (Second Time)



### AND Operator Evaluation

Min document is 4.

args[0] matches, so call  
args[0].getScore ().  
Suppose the result is 0.3.

args[1] matches, so call  
args[1].getScore().  
Suppose the score is 0.2.

$$\text{Score}_{\text{AND}}(4) = (0.3^{0.5} \times 0.2^{0.5})$$

68

© 2017 Jamie Callan

## Using Default Belief Scores Properly Requires Two Components

### Add a new method to all QrySop operators

double getDefaultScore (RetrievalModel r, long docid)

- QrySopScore.getDefaultScore calculates a score for a term
- QrySop <other>. getDefaultScore combines scores for  $n$  terms

### When any QrySop operator calculates scores

If the  $i^{\text{th}}$  query argument contains document  $d$

then read its score from the  $i^{\text{th}}$  score list

else call the  $i^{\text{th}}$  query argument's getDefaultScore method

**It may sound complicated now, but actually it is very easy**

69

© 2017 Jamie Callan

## Outline

### Statistical language models

- Introduction to language models
- Query likelihood
- Kullback-Leibler (KL) Divergence
- Indri

70

© 2017 Jamie Callan

## Retrieval Models Summary

### We have discussed the following retrieval models

- Unranked Boolean
- Ranked Boolean, using tf scoring
- Vector Space, using Inc.ltc scoring
- Okapi BM25
- Language Models
  - Query likelihood, with two-stage smoothing
  - KL Divergence and Jensen-Shannon Divergence
  - Indri

### That's a lot of material ... what's important?

71

© 2017 Jamie Callan

## Retrieval Models Summary

### Important differences among the models that you should know

- Boolean vs. ranked retrieval
- The kinds of statistics used by most ranking functions
  - The theories are different, but the stats are mostly the same
- How well the different models support query operators
  - Inverted-list operators vs. score operators
  - Strict Boolean vs. probabilistic Boolean
- Which models you could use (or not use) for different types of tasks

72

© 2017 Jamie Callan

## Retrieval Models Summary

---

**All of these retrieval models are used widely**

- There must be a reason why ... you should know what it is