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

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

Jamie Callan
Carnegie Mellon University
callan@cs.cmu.edu

1

Outline

Statistical language models

- Introduction
- Query likelihood model
- Kullback-Leibler (KL) Divergence
- Indri

HW2 implementation (see notes posted online)

- Indri default beliefs
- Window operator

2

© 2020 Jamie Callan

2

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

3

© 2020 Jamie Callan

3

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

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

4

© 2020 Jamie Callan

4

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
- A language model 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

5

© 2020 Jamie Callan

5

Terminology

d: A document

θ_d : The language model for document d

q: A query

θ_q : The language model for query q

d and θ_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 $p(d | \theta_q)$

6

© 2020 Jamie Callan

6

Retrieval Model Based Upon Statistical Language Models

A document d defines a probability distribution θ_d over index terms

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

t	$p(t)$
apple	0.0204
banana	0.0001
campaign	0.0034
dog	0.0102
:	:

A query q also defines a probability distribution θ_q

- A sparse distribution (more on this later)

Two common methods of using language modeling to rank document d

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

7

© 2020 Jamie Callan

7

Outline

Statistical language models

- Introduction
- Query likelihood model
- Kullback-Leibler (KL) Divergence
- Indri

HW2 implementation (see notes posted online)

- Indri default beliefs
- Window operator

8

© 2020 Jamie Callan

8

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 little frequency information ($qtf=1$ for most terms)

Query: digital camera

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

Solution

$$p(d|q) = \frac{p(q|d)p(d)}{p(q)} \quad \text{Bayes rule}$$

$$\propto p(q|d)p(d) \quad \text{Drop document-independent term}$$

Key issues

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

9

© 2020 Jamie Callan

9

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...

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

$$\propto p(q|d) \quad \text{Drop the constant (uniform probability)}$$

10

© 2020 Jamie Callan

10

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 ...

$$p(q|d) = \prod_{q_i \in q} p(q_i|d) \quad q_i: \text{ term in query } q$$

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?

11

© 2020 Jamie Callan

11

Estimating $p(q_i|d)$

Maximum likelihood estimation (MLE) is a simple approach

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

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

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

- This is exact-match AND
- We would prefer best-match AND

12

© 2020 Jamie Callan

12

Estimating $p(q_i|d)$: Smoothing

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

- Imprecise probability estimates from MLE
 - E.g., for infrequent terms (few observations)
 - E.g., for short documents (small sample)
- 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?

13

© 2020 Jamie Callan

13

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

Linear interpolation with a reference language model

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

Average an estimate from a small sample (document) with an estimate from a large sample (collection)

Smoothing decreases as $\lambda \rightarrow 0$

What value of λ is best?

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

14

© 2020 Jamie Callan

14

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

Method #2: Bayesian smoothing using Dirichlet priors

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

$$p_{MLE}(q_i|C) = \frac{ctf_{q_i}}{length_{tokens}(C)}$$

Increase tf of terms that are expected to be frequent in d (small sample) because they are frequent in C (large sample)

What value of μ is best?

- μ in [1,000-10,000] seems best
- μ might be approximately related to average document length

15

© 2020 Jamie Callan

15

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 harder to see

- Smoothing provides an idf-like effect
- Next slides...

16

© 2020 Jamie Callan

16

No Smoothing Example:

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

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

Document d_1

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

$(2/50) \times$

$(3/50) =$

$0.04 \times 0.06 =$

$p(q|d_1) = 0.0024$

Document d_2

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

$(3/50) \times$

$(2/50) =$

$0.06 \times 0.04 =$

$p(q|d_2) = 0.0024$

The model does not distinguish between frequent and rare terms

17

© 2020 Jamie Callan

17

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

Document d_1

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

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

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

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

$(0.6 \times 0.06 + 0.0004) =$

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

$p(q|d_1) = 0.001019$

Document d_2

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

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

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

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

$(0.6 \times 0.04 + 0.0004) =$

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

$p(q|d_2) = 0.000976$

The model does distinguish between frequent and rare terms

18

© 2020 Jamie Callan

18

Dirichlet Smoothing Example:

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

Two query terms: $p(\text{apple} | C) = 0.01$, $p(\text{ipad} | C) = 0.001$
 $\mu = 2000$ “frequent” “rare”

Document d_1

doclen = 50, $tf_{\text{apple}} = 2$, $tf_{\text{ipad}} = 3$
 $(2 + 2000 \times 0.01) / (50 + 2000) \times$
 $(3 + 2000 \times 0.001) / (50 + 2000) =$

$(2 + 20) / 2050 \times$
 $(3 + 2) / 2050 =$

$0.0107 \times 0.0024 =$

$p(q|d_1) = 0.000026$

Document d_2

doclen = 50, $tf_{\text{apple}} = 3$, $tf_{\text{ipad}} = 2$
 $(3 + 2000 \times 0.01) / (50 + 2000) \times$
 $(2 + 2000 \times 0.001) / (50 + 2000) =$

$(3 + 20) / 2050 \times$
 $(2 + 2) / 2050 =$

$0.0112 \times 0.0020 =$

$p(q|d_2) = 0.000022$

The model does distinguish between frequent and rare terms

19

© 2020 Jamie Callan

19

Why Does Smoothing Work?

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

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

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

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

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

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

$p_{\delta}(\text{ipad} | d) = 0.6 \times (1/50) = 0.012$

– 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

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

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

20

© 2020 Jamie Callan

20

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

How Jelenick-Mercer smoothing has an idf-like effect

- Several steps of math manipulation are omitted for clarity

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

21

© 2020 Jamie Callan

21

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 λ approaches 0?

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

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

22

© 2020 Jamie Callan

22

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 effect” is less important, so small λ is best
- Long queries
 - Many query terms, so (usually) most query terms must match
 - “idf effect” is more important
 - Give priority to the less common terms, so larger λ is best

23

© 2020 Jamie Callan

23

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 μ approaches 0?

- There is no smoothing (maximum likelihood only)

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

24

© 2020 Jamie Callan

24

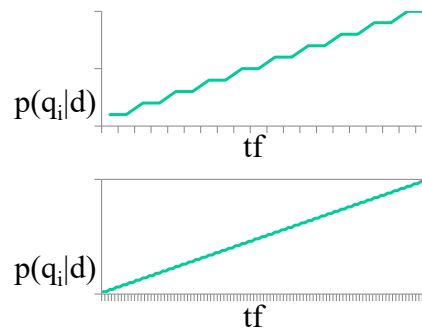
Estimating $p(q_i|d)$: Bayesian Smoothing With Dirichlet Priors

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 μ is more important
- Long documents
 - Probabilities are more smooth
 - Larger μ is less important



25

© 2020 Jamie Callan

25

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:** Compensates for different word importance
 - » “idf effect”
 - **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

© 2020 Jamie Callan

26

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

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

- $tf_{q,d}$:** 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

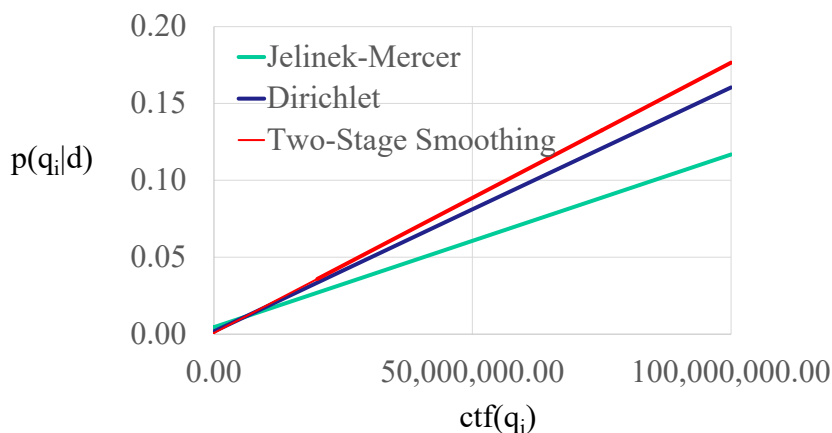
27

© 2020 Jamie Callan

27

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

How do different smoothing methods treat different types of terms?



Assumptions

- tf : 5
- Doc len: 450
- Coll len: 530M
- λ : 0.6
- μ : 2,500

**λ and μ control the slope
(more smoothing)**

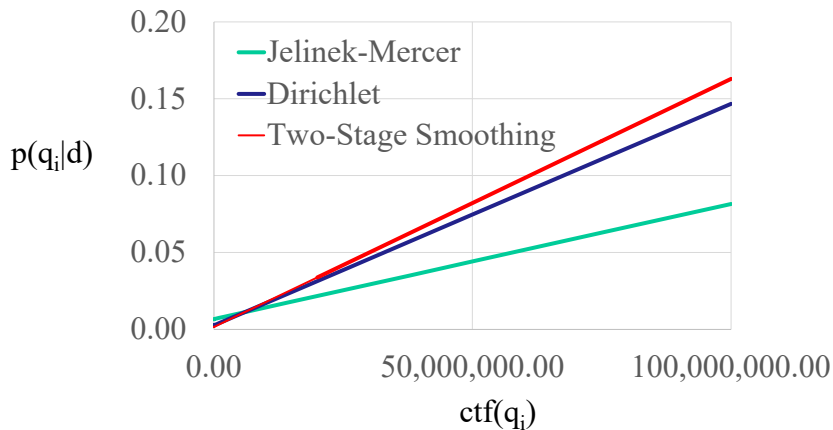
28

© 2020 Jamie Callan

28

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

How do different smoothing methods treat different types of terms?



Assumptions

- tf: 5
- Doc len: 450
- Coll len: 530M
- λ : 0.4
- μ : 1,500

λ and μ control the slope
(less smoothing)

29

© 2020 Jamie Callan

29

Outline

Statistical language models

- Introduction
- Query likelihood model
- Kullback-Leibler (KL) Divergence
- Indri

HW2 implementation (see notes posted online)

- Indri default beliefs
- Window operator

30

© 2020 Jamie Callan

30

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 the similarity of θ_d and θ_q “KL divergence”

In some cases, these approaches are equivalent...

31

© 2020 Jamie Callan

31

KL Divergence: Measuring the Similarity of Language Models

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

$$KL(a||b) = \sum_{x \in X} a(x) \log \frac{a(x)}{b(x)}$$

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?

32

© 2020 Jamie Callan

32

Measuring the Similarity of Statistical Language Models

Ranking by (negative) KL-divergence

$$-KL(q||d) = - \sum_{w \in V} p(w|q) \log \frac{p(w|q)}{p(w|d)}$$

33

© 2020 Jamie Callan

33

Measuring the Similarity of Statistical Language Models

Ranking by (negative) KL-divergence

$$\begin{aligned} -KL(q||d) &= - \sum_{w \in V} p(w|q) \log \frac{p(w|q)}{p(w|d)} \\ &= - \sum_{q_i \in q} p(q_i|q) \log \frac{p(q_i|q)}{p(q_i|d)} \end{aligned}$$

P(w|q) = 0 when w ∉ q

34

© 2020 Jamie Callan

34

Measuring the Similarity of Statistical Language Models

Ranking by (negative) KL-divergence

$$\begin{aligned}
 -KL(q \| d) &= - \sum_{w \in V} p(w | q) \log \frac{p(w | q)}{p(w | d)} \\
 &= - \sum_{q_i \in q} p(q_i | q) \log \frac{p(q_i | q)}{p(q_i | d)} \quad \mathbf{P(w|q) = 0 \text{ when } w \notin q} \\
 &= - \left(\sum_{q_i \in q} p(q_i | q) \log p(q_i | q) - \sum_{q_i \in q} p(q_i | q) \log p(q_i | d) \right)
 \end{aligned}$$

35

© 2020 Jamie Callan

35

Measuring the Similarity of Statistical Language Models

Ranking by (negative) KL-divergence

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

36

© 2020 Jamie Callan

36

Measuring the Similarity of Statistical Language Models

Ranking by (negative) KL-divergence

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

37

© 2020 Jamie Callan

37

Measuring the Similarity of Statistical Language Models

Ranking by (negative) KL-divergence

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

38

© 2020 Jamie Callan

38

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)

39

© 2020 Jamie Callan

39

Outline

Statistical language models

- Introduction
- Query likelihood model
- Kullback-Leibler (KL) Divergence
- Indri

HW2 implementation (see notes posted online)

- Indri default beliefs
- Window operator

40

© 2020 Jamie Callan

40

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

41

© 2020 Jamie Callan

41

Inference Nets + Language Models: Indri

Documents

....

d_i

....

42

© 2020 Jamie Callan

42

Inference Nets + Language Models: Indri

Documents

.... d_i

Representation Nodes
(terms, #near/n, ...)

r_1 r_2 r_3 r_4 r_j

43

© 2020 Jamie Callan

43

Inference Nets + Language Models: Indri

Smoothing parameters

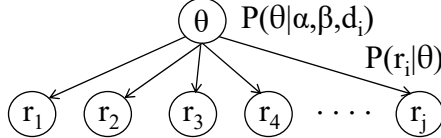
α, β

Documents

.... d_i

Language Model

Representation Nodes
(terms, #near/n, ...)



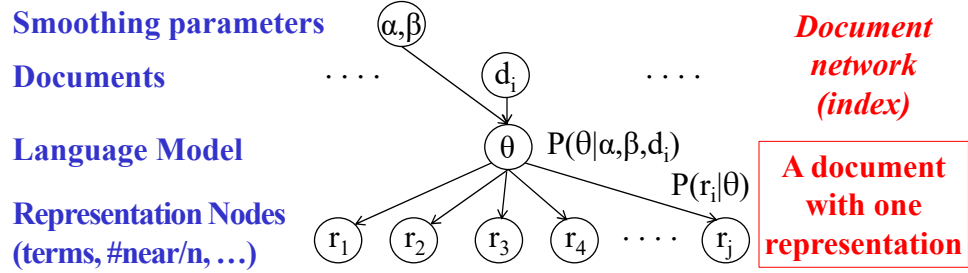
A document
with one
representation

44

© 2020 Jamie Callan

44

Inference Nets + Language Models: Indri

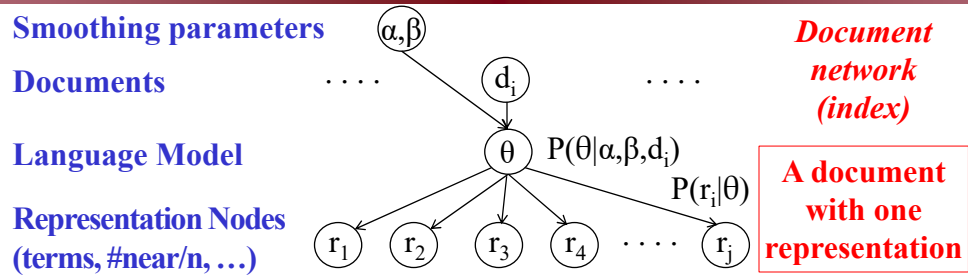


45

© 2020 Jamie Callan

45

Inference Nets + Language Models: Indri



Information Need

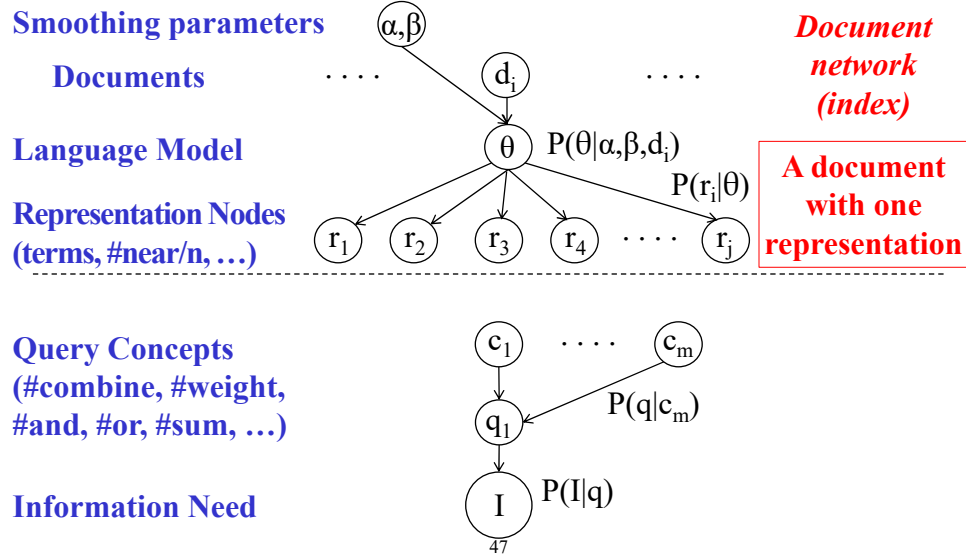


46

© 2020 Jamie Callan

46

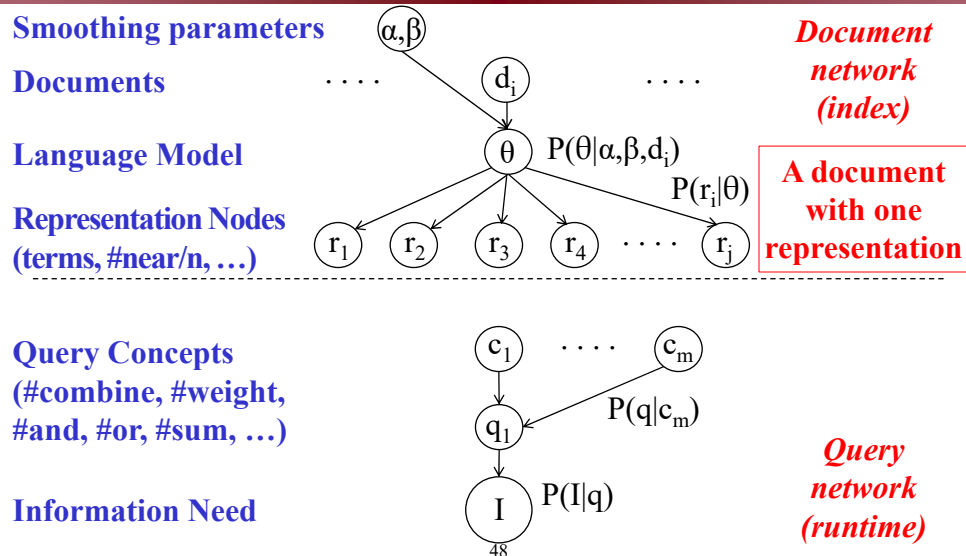
Inference Nets + Language Models: Indri



© 2020 Jamie Callan

47

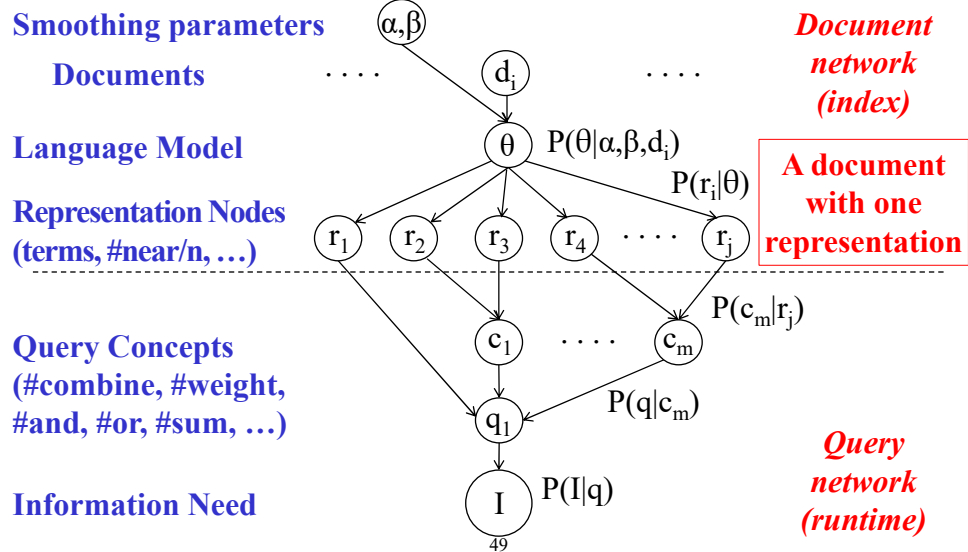
Inference Nets + Language Models: Indri



© 2020 Jamie Callan

48

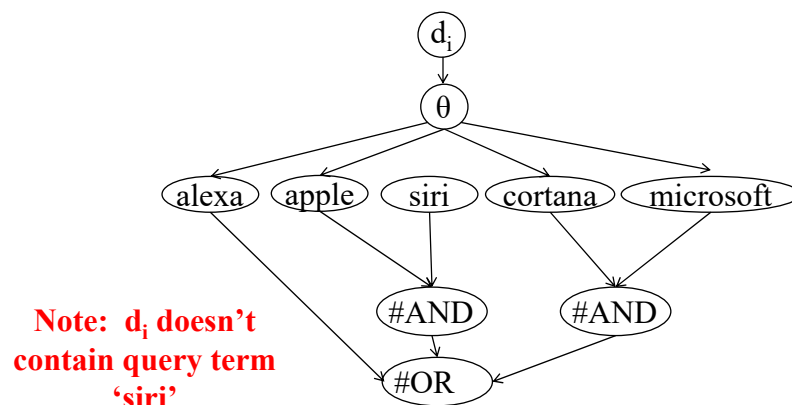
Inference Nets + Language Models: Indri



49

Inference Nets + Language Models: Indri

This isn't as complicated as it looks



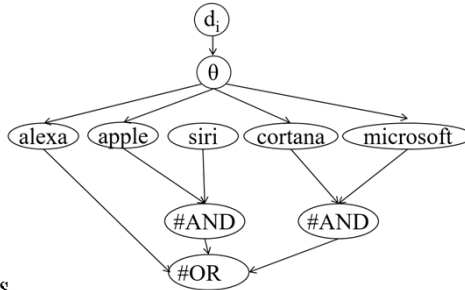
50

50

Inference Nets + Language Models: Indri

This isn't as complicated as it looks

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



51

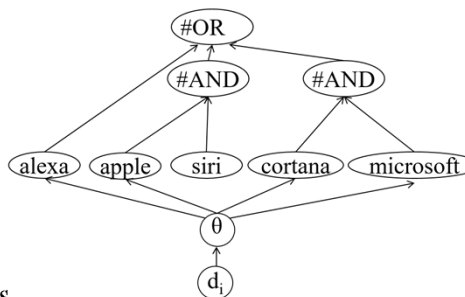
© 2020 Jamie Callan

51

Inference Nets + Language Models: Indri

This isn't as complicated as it looks

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

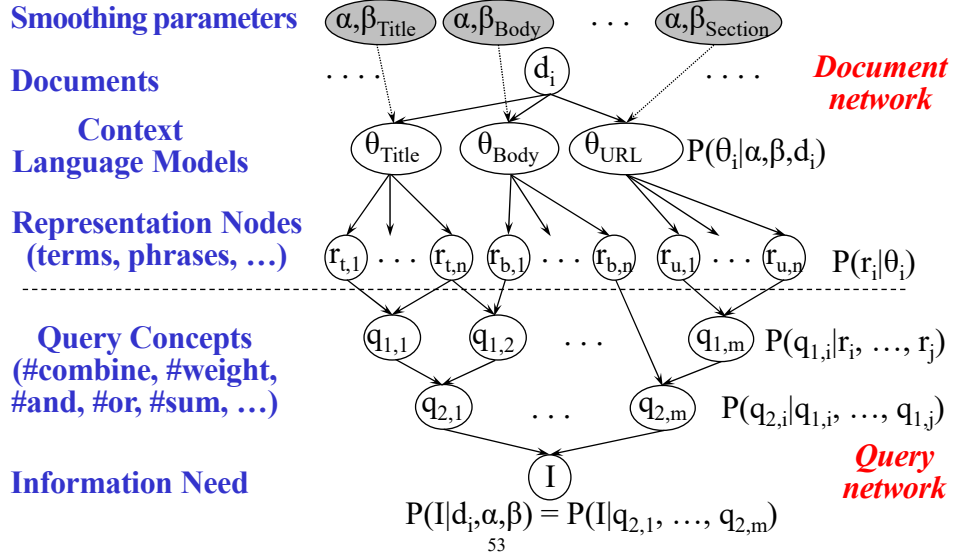


52

© 2020 Jamie Callan

52

Inference Nets + Language Models: Indri



53

© 2020 Jamie Callan

53

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

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

54

© 2020 Jamie Callan

54

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

What about operators that combine scores?

- E.g., ScoreList+ \rightarrow ScoreList
- E.g., #AND, #OR, ...
- Indri views these as operators that combine evidence
 - “evidence” means “conditional probabilities”

55

© 2020 Jamie Callan

55

Indri Query Operators

AND is the default query operator for most language modeling systems

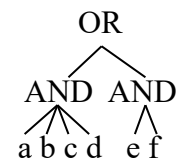
- Typically implemented as the product of the argument weights

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

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

Indri's AND operator uses the geometric mean

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



Which AND operator has more impact on the score of the OR operator?

56

© 2020 Jamie Callan

56

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

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

57

© 2020 Jamie Callan

57

Indri Query Operators

Indri provides an OR operator

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

58

© 2020 Jamie Callan

58

Indri Query Operators

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

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

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

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

59

© 2020 Jamie Callan

59

Indri Query Operators

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

60

© 2020 Jamie Callan

60

Outline

Statistical language models

- Introduction
- Query likelihood model
- Kullback-Leibler (KL) Divergence
- Indri

HW2 implementation (see notes posted online)

- Indri default beliefs
- Window operator

61

© 2020 Jamie Callan

61

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?

62

© 2020 Jamie Callan

62

Retrieval Models Summary

Important differences among the models that you should know

- Unranked vs. ranked retrieval
- Exact-match vs. best-match (partial-match) 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

63

© 2020 Jamie Callan

63

Retrieval Models Summary

All of these retrieval models are used widely

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

64

© 2020 Jamie Callan

64