

Language Models

What's wrong with VSMs?

Vector Space Models work reasonably well, but have a few problems:

- They are based on bag-of-words, so they ignore grammatical context and suffer from term mismatch.
- They don't adapt to the user or collection, but ideal term weights are user- and domain-specific.
- They are heuristic-based, and don't have much explanatory power.

Probabilistic Modeling

We can address these problems by moving to probabilistic models, such as language models:

- We can take grammatical context into account, and trade off between using more context and performing faster inference.
- The model can be trained from a particular collection, or conditioned based on user- and domain-specific features.
- The model is interpretable, and makes concrete predictions about query and document relevance.

In this Module...

1. Ranking as a probabilistic classification task
2. Some specific probabilistic models for classification
3. Smoothing: estimating model parameters from sparse data
4. A probabilistic approach to pseudo-relevance feedback

Ranking with Probabilistic Models

Imagine we have a function that gives us the probability that a document D is relevant to a query Q , $P(R=1|D, Q)$. We call this function a *probabilistic model*, and can rank documents by decreasing probability of relevance.

There are many useful models, which differ by things like:

- Sensitivity to different document properties, like grammatical context
- Amount of training data needed to train the model parameters
- Ability to handle noise in document data or relevance labels

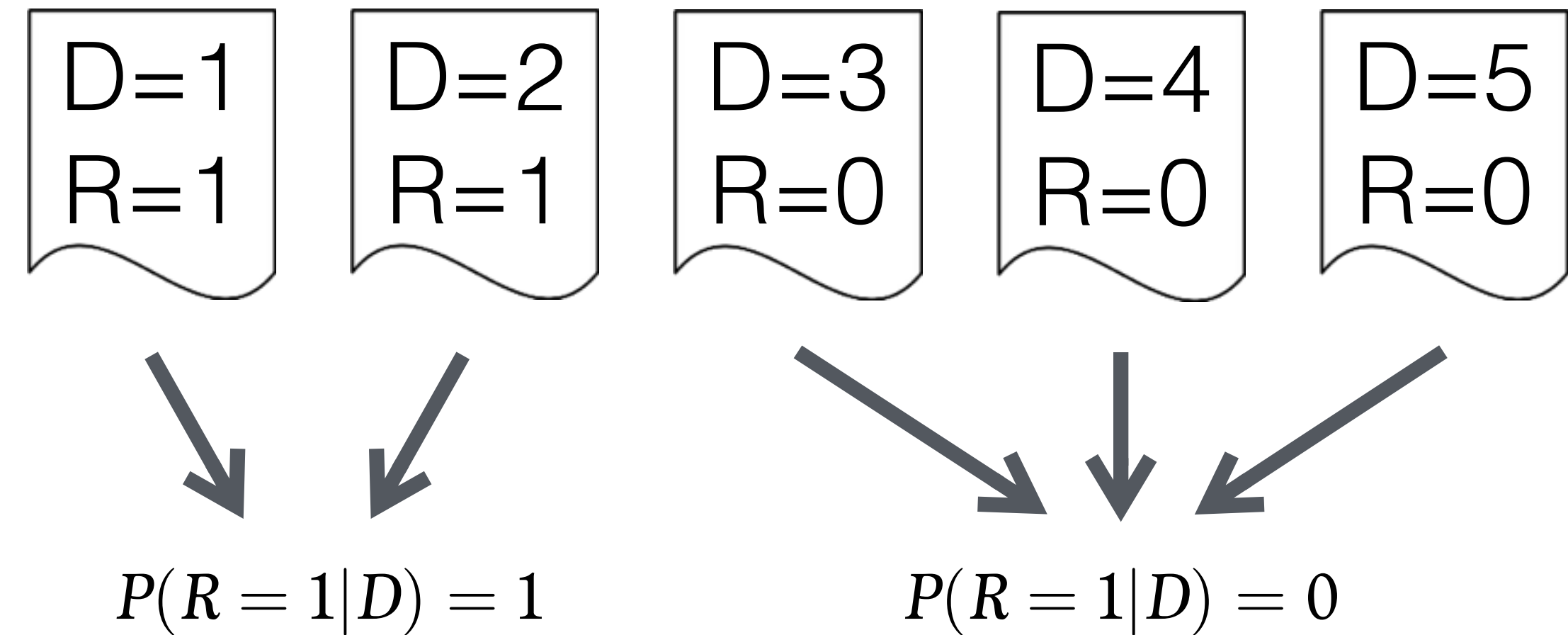
For simplicity here, we will hold the query constant and consider $P(R=1|D)$.

The Flaw in our Plan

Suppose we have documents and relevance labels, and we want to empirically measure $P(R=1|D)$.

Each document has only one relevance label, so every probability is either 0 or 1. Worse, there is no way to generalize to new documents.

Instead, we estimate the probability of documents given relevance labels, $P(D|R=1)$.



	D=1	D=2	D=3	D=4	D=5
$P(D R=1)$	1/2	1/2	0	0	0
$P(D R=0)$	0	0	1/3	1/3	1/3

Language model for doc = prob dist. over words/terms

Recap probab's basics, expectations, estimation, Bayes, R.V.
cond prob, KL
(P/Q)

Bayes' Rule

$R = \text{binary}$

$\text{doc} = (d^1, d^2, \dots, d^T)$
 ↑ prob values

We can estimate $P(D|R=1)$, not $P(R=1|D)$,
 so we apply Bayes' Rule to estimate
 document relevance.

- $P(D|R=1)$ gives the probability that a relevant document would have the properties encoded by the random variable D .
- $P(R=1)$ is the probability that a randomly-selected document is relevant.

$$P(R = 1|D) = \frac{P(D|R = 1)P(R = 1)}{P(D)}$$

$$= \frac{P(D|R = 1)P(R = 1)}{\sum_r P(D|R = r)P(R = r)}$$

$$P(R|D) = \frac{P(D|R) \cdot P(R)}{P(D)}$$

independent
 $P(A, B, C) = P(A) \cdot P(B) \cdot P(C)$

Bayesian Classification

Starting from Bayes' Rule, we can easily build a classifier to tell us whether documents are relevant. We will say a document is relevant if:

want prob(relevance | doc)

$$\begin{aligned} P(R = 1|D) &> P(R = 0|D) \\ \implies \frac{P(D|R = 1)P(R = 1)}{P(D)} &> \frac{P(D|R = 0)P(R = 0)}{P(D)} \\ \implies \frac{P(D|R = 1)}{P(D|R = 0)} &> \frac{P(R = 0)}{P(R = 1)} \end{aligned}$$

We can estimate $P(D|R=1)$ and $P(D|R=0)$ using a language model, and $P(R=0)$ and $P(R=1)$ based on the query, or using a constant. Note that for large web collections, $P(R=1)$ is very small for virtually any query.

$R = \text{relevance} = \text{relevance for a particular query.}$

Unigram Language Model

In order to put this together, we need a language model to estimate $P(D|R)$.

Let's start with a model based on the bag-of-words assumption. We'll represent a document as a collection of independent words ("unigrams").

prob(D|Relev for q) ← $P(D|R)$ *be words*

$$D = (w_1, w_2, \dots, w_n)$$

joint

$$P(D|R) = P(w_1, w_2, \dots, w_n | R)$$

prob that doc d is about query q = R

$$= P(w_1 | R) P(w_2 | R, w_1) P(w_3 | R, w_1, w_2) \dots P(w_n | R, w_1, \dots, w_{n-1})$$

impractical

$$= P(w_1 | R) P(w_2 | R) \dots P(w_n | R)$$

check naive independence term correlations / Co-occurrence

always true

naive assume independent

$$= \prod_{i=1}^n P(w_i | R)$$

Example

Let's consider querying a collection of five short documents with a simplified vocabulary: the only words are apple, baker, and crab.

Document	Rel?	apple?	baker?	crab?	Term	# Rel	# Non Rel	$P(w R=1)$	$P(w R=0)$
apple apple crab	1	1	0	1	apple	2	1	2/2	1/3
crab baker crab	0	0	1	1	baker	1	2	1/2	2/3
apple baker baker	1	1	1	0	crab	1	3	1/2	3/3
crab crab apple	0	1	0	1					
baker baker crab	0	0	1	1					

$P(R = 1) = 2/5$
 $P(R = 0) = 3/5$

Example

Is “apple baker crab” relevant?

$$\frac{P(D|R=1)}{P(D|R=0)} \stackrel{?}{>} \frac{P(R=0)}{P(R=1)}$$

$$\frac{\prod_i P(w_i|R=1)}{\prod_i P(w_i|R=0)} \stackrel{?}{>} \frac{P(R=0)}{P(R=1)}$$

$$\frac{P(\text{apple} = 1|R=1)P(\text{baker} = 1|R=1)P(\text{crab} = 1|R=1)}{P(\text{apple} = 1|R=0)P(\text{baker} = 1|R=0)P(\text{crab} = 1|R=0)} \stackrel{?}{>} \frac{0.6}{0.4}$$

$$\frac{1 \cdot 0.5 \cdot 0.5}{0.3 \cdot 0.6 \cdot 1} \stackrel{?}{>} \frac{0.6}{0.4}$$

$$1.125 < 1.5$$

Term	$P(w R=1)$	$P(w R=0)$
apple	1	1/3
baker	1/2	2/3
crab	1/2	1

$$P(R=1) = 2/5$$

$$P(R=0) = 3/5$$

Retrieval With Language Models

So far, we've focused on language models like $P(D = w_1, w_2, \dots, w_n)$. Where's the query?

Remember the key insight from vector space models: we want to represent queries and documents in the same way. The query is just a "short document:" a sequence of words. There are three obvious approaches we can use for ranking:

1. Query likelihood: Train a language model on a document, and estimate the query's probability.
2. Document likelihood: Train a language model on the query, and estimate the document's probability.
3. Model divergence: Train language models on the document and the query, and compare them.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

ranking
A-s for fix B

$P(A|B) \stackrel{\text{rank}}{\sim} P(B|A) \cdot P(A)$
NOT EXACT, not a prob
proportion ✓
rank ✓

Query Likelihood Retrieval

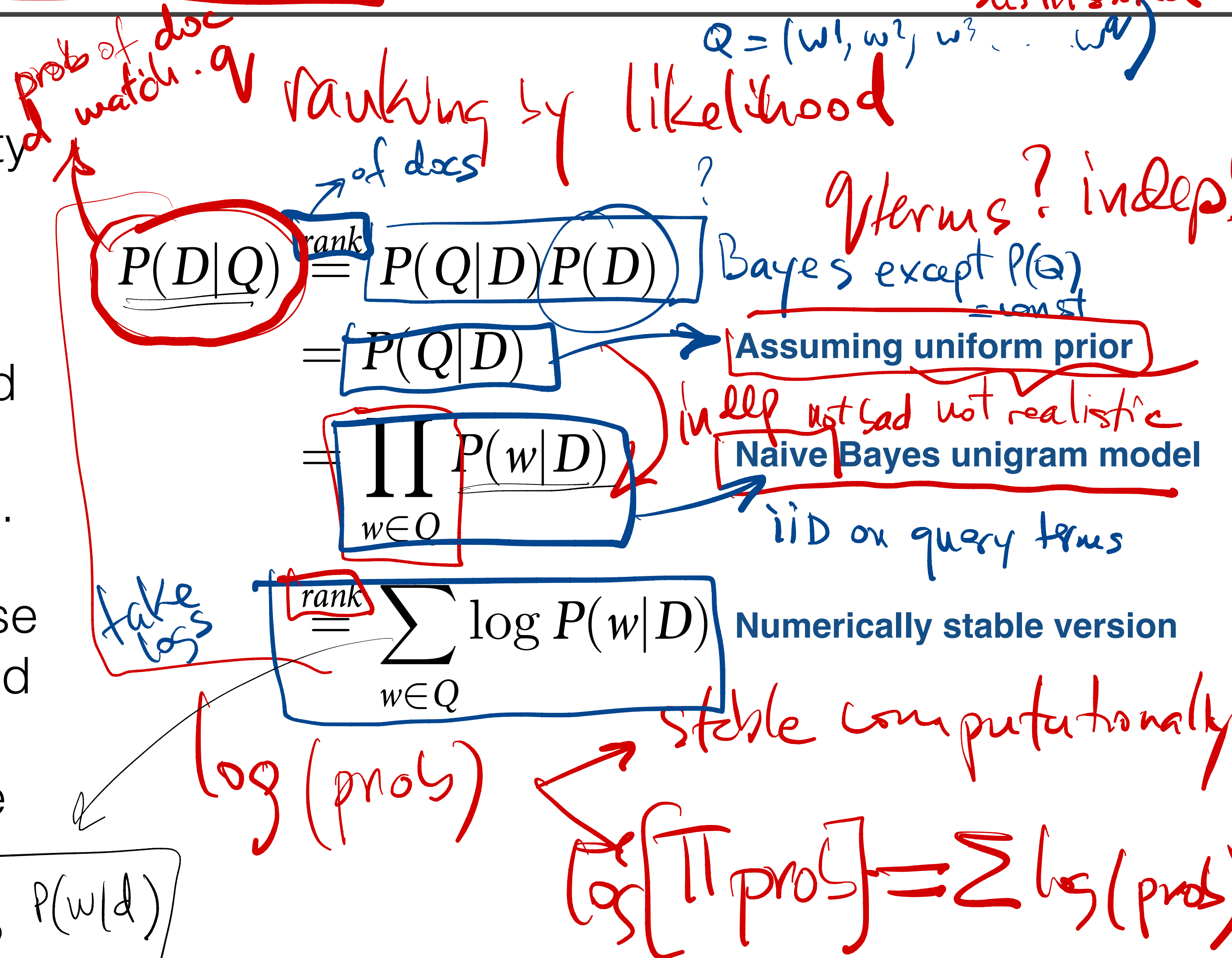
important / HWL
 iid = identical & independent distributed

$$Q = (w^1, w^2, w^3, \dots, w^N)$$

Suppose that the query specifies a topic. We want to know the probability of a document being generated from that topic, or $P(D|Q)$.

However, the query is very small, and documents are long: document language models have less variance.

In the *Query Likelihood Model*, we use Bayes' Rule to rank documents based on the probability of generating the query from the documents' language models.



$$\text{score}(d) \approx \sum_{w \in Q} \log P(w|d)$$

pseudocode
prob type

init: $score(d) = 0 \forall d$

For $w \in q$

For each doc $d \in \text{urlist}(w)$

$score(d) = score(d) + match(w, d)$

skip doc $d \neq w$? incorrect

term w does not appear in doc $d \Rightarrow$ penalty

$\Rightarrow \log[P(w|d)]$

small \rightarrow deep negative

• Can I implement Lang. Model Query Likelihood this way?

— problem: $match \approx \log(prob) \in (-\infty, 0]$

max match value is $0 = \log(1)$

max prob = 1

$w \notin d \Rightarrow$ match(w, d) negative \Rightarrow we cannot simply skip doc.

Solution for L.M.
change init function
init $score(d) = ?$

Example: Query Likelihood

Wikipedia: WWI

World War I (**WWI** or **WW1** or **World War One**), also known as the **First World War** or the **Great War**, was a [global war](#) centred in Europe that began on 28 July 1914 and lasted until 11 November 1918. More than 9 million [combatants](#) and 7 million [civilians died as a result of the war](#), a casualty rate exacerbated by the belligerents' technological and industrial sophistication, and tactical stalemate. It was [one of the deadliest conflicts in history](#), paving the way for major political changes, including revolutions in many of the nations involved.

Query: "deadliest war in history"

Term	P(w D)	log P(w D)
deadliest	1/94 = 0.011	-1.973
war	6/94 = 0.063	-1.195
in	3/94 = 0.032	-1.496
history	1/94 = 0.011	-1.973
	$\Pi = 2.30e-7$	$\Sigma = -6.637$

Keyword based query 34 terms
=> Q.L. model

Example: Query Likelihood

Wikipedia: WWI

World War I (**WWI** or **WW1** or **World War One**), also known as the **First World War** or the **Great War**, was a [global war](#) centred in Europe that began on 28 July 1914 and lasted until 11 November 1918. More than 9 million [combatants](#) and 7 million [civilians died as a result of the war](#), a casualty rate exacerbated by the belligerents' technological and industrial sophistication, and tactical stalemate. It was [one of the deadliest conflicts in history](#), paving the way for major political changes, including revolutions in many of the nations involved.

Query: “deadliest war in history”

Term	P(w D)	log P(w D)
deadliest	1/94 = 0.011	-1.973
war	6/94 = 0.063	-1.195
in	3/94 = 0.032	-1.496
history	1/94 = 0.011	-1.973
	$\Pi = 2.30e-7$	$\Sigma = -6.637$

Example: Query Likelihood

Wikipedia: Taiping Rebellion

The **Taiping Rebellion** was a massive [civil war](#) in [southern China](#) from 1850 to 1864, against the ruling [Manchu Qing dynasty](#). It was a [millenarian movement](#) led by [Hong Xiuquan](#), who announced that he had received visions, in which he learned that he was the younger brother of [Jesus](#). At least 20 million people died, mainly civilians, in one of the [deadliest military conflicts](#) in history.

Query: “deadliest war in history”

Term	P(w D)	log P(w D)
deadliest	1/56 = 0.017	-1.748
war	1/56 = 0.017	-1.748
in	2/56 = 0.035	-1.447
history	1/56 = 0.017	-1.748
	$\Pi = 2.56e-8$	$\Sigma = -6.691$

Summary: Language Model

There are many ways to move beyond this basic model.

- Use n-gram or skip-gram probabilities, instead of unigrams.
- Model document probabilities $P(D)$ based on length, authority, genre, etc. instead of assuming a uniform probability.
- Use the tools from the VSM slides: stemming, stopping, etc.

Next, we'll see how to fix a major issue with our probability estimates: what happens if a query term doesn't appear in the document?

Retrieval With Language Models

There are three obvious ways to perform retrieval using language models:

1. **Query Likelihood Retrieval** ^{HWI} trains a model on the document and estimates the query's likelihood. We've focused on these so far.
2. **Document Likelihood Retrieval** → *not good except in special cases.* trains a model on the query and estimates the document's likelihood. Queries are very short, so these seem less promising.
3. **Model Divergence Retrieval** → *used a lot in practice (not req HW1)* trains models on both the document and the query, and compares them.

Comparing Distributions

The most common way to compare probability distributions is with Kullback-Liebler ("KL") Divergence.

This is a measure from Information Theory which can be interpreted as the expected number of bits you would waste if you compressed data distributed along p as if it was distributed along q .

If $p = q$, $D_{KL}(p||q) = 0$.

measures similarity (closeness) between 2 distributions

$$E_p(v) = \sum_{x \in \Omega} p(x) \cdot v(x)$$

↓
weighted by p .

$$D_{KL}(p||q) = \sum_e p(e) \log \frac{p(e)}{q(e)}$$

$$= E_p \left[\log \frac{p}{q} \right]$$

Bits required to express diff $\frac{p}{q}$

$\log = \log_2$ unless specified

Divergence-based Retrieval

Model Divergence Retrieval works as follows:

1. Choose a language model for the query, $p(w|q)$. M_q = dist over terms
2. Choose a language model for the document, $p(w|d)$. M_d = dist over terms
3. Rank by $-D_{KL}(p(w|q) || p(w|d))$ – more divergence means a worse match.

This can be simplified to a cross-entropy calculation, as shown to the right.

model comparison (M_q, M_d) by compar. distributions

$$\begin{aligned} & D_{KL}(p(w|q) || p(w|d)) \\ &= \sum_w p(w|q) \log \frac{p(w|q)}{p(w|d)} \\ &= \sum_w p(w|q) \log p(w|q) - \sum_w p(w|q) \log p(w|d) \end{aligned}$$

$$\stackrel{\text{rank}}{=} - \sum_w p(w|q) \log p(w|d)$$

both $p(w|q)$ are estimated from large enough data-text \Rightarrow KL dist / dis. makes

discrete

$$p = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8} \right\}$$

$$q = \left\{ \frac{1}{4}, \frac{1}{8}, \frac{1}{8}, \frac{1}{2} \right\}$$

$\Omega = \{x_1, x_2, x_3, x_4\}$

$$KL(p||q) = \sum_{x \in \Omega} p(x) \cdot \log \left(\frac{p(x)}{q(x)} \right)$$

$$= \frac{1}{2} \cdot \log \left(\frac{1/2}{1/4} \right) + \frac{1}{4} \log \left(\frac{1/4}{1/8} \right) +$$

$$+ \frac{1}{8} \log \left(\frac{1/8}{1/8} \right) + \frac{1}{8} \log \left(\frac{1/8}{1/2} \right) =$$

$$= \frac{1}{2} \cdot \log(2) + \frac{1}{4} \log(2) + \frac{1}{8} \log(1) + \frac{1}{8} \log\left(\frac{1}{4}\right)$$

-2

continuous $\sum \rightarrow \int$

Exercise • $KL(p||q) \geq 0$

• $\min KL(p||q) = 0 \iff p(x) = q(x) \forall x$

• KL NOT SYMMETRIC (as usually distance)
 $KL(p||q) \neq KL(q||p)$

Retrieval Flexibility

Model Divergence Retrieval generalizes the Query and Document Likelihood models, and is the most flexible of the three.

Any language model can be used for the query or document. They don't have to be the same. It can help to smooth or normalize them differently.

If you pick the maximum likelihood model for the query, this is equivalent to the query likelihood model.

$$\text{Pick } p(w|q) := \frac{tf_{w,q}}{|q|} = \frac{1}{|q|}$$

$$\begin{aligned} D_{KL}(p(w|q) || p(w|d)) &\stackrel{\text{rank}}{=} - \sum_w p(w|q) \log p(w|d) \\ &= - \sum_w \frac{1}{|q|} \log p(w|d) \end{aligned}$$

Equivalence to Query Likelihood Model

Example: Model Divergence Retrieval

We make the following model choices:

- $p(w|q)$ is Dirichlet-smoothed with a background of words used in historical queries.
- $p(w|d)$ is Dirichlet-smoothed with a background of words used in documents from the corpus.
- $\sum_w qf_w = 500,000$
- $\sum_w cf_w = 1,000,000,000$

Let $qf_w := \text{count}(\text{word } w \text{ in query log})$

$$p(w|q, \mu = 2) = \frac{tf_{w,q} + 2 \times \frac{qf_w}{\sum_w qf_w}}{|q| + 2}$$

$$p(w|d, \mu = 2000) = \frac{tf_{w,d} + 2,000 \times \frac{cf_w}{\sum_w cf_w}}{|d| + 2,000}$$

$$\begin{aligned} D_{KL}(p(w|q) || p(w|d)) &\stackrel{\text{rank}}{=} - \sum_w p(w|q) \log p(w|d) \\ &= - \sum_w \frac{tf_{w,q} + 2 \times \frac{qf_w}{\sum_w qf_w}}{|q| + 2} \log \frac{tf_{w,d} + 2,000 \times \frac{cf_w}{\sum_w cf_w}}{|d| + 2,000} \end{aligned}$$

Ranking by (negative) KL-Divergence provides a very flexible and theoretically-sound retrieval system.

Example: Model Divergence Retrieval

$$\sum_w \frac{tf_{w,q} + 2 \times \frac{qf_w}{\sum_w qf_w}}{|q| + 2} \log \frac{tf_{w,d} + 2,000 \times \frac{cf_w}{\sum_w cf_w}}{|d| + 2,000}$$

Wikipedia: WWI

World War I (**WWI** or **WW1** or **World War One**), also known as the **First World War** or the **Great War**, was a [global war](#) centred in Europe that began on 28 July 1914 and lasted until 11 November 1918. More than 9 million [combatants](#) and 7 million [civilians died as a result of the war](#), a casualty rate exacerbated by the belligerents' technological and industrial sophistication, and tactical stalemate. It was [one of the](#)

Query: "world war one"

	qf	cf	p(w q)	p(w d)	Score
world	2,500	90,000	0.202	0.002	-1.891
war	2,000	35,000	0.202	0.003	-1.700
one	6,000	5E+07	0.205	0.049	-0.893
					-4.484

Example: Model Divergence Retrieval

$$\sum_w \frac{tf_{w,q} + 2 \times \frac{qf_w}{\sum_w qf_w}}{|q| + 2} \log \frac{tf_{w,d} + 2,000 \times \frac{cf_w}{\sum_w cf_w}}{|d| + 2,000}$$

Wikipedia: Taiping Rebellion

The **Taiping Rebellion** was a massive [civil war](#) in [southern China](#) from 1850 to 1864, against the ruling [Manchu Qing dynasty](#). It was a [millenarian movement](#) led by [Hong Xiuquan](#), who announced that he had received visions, in which he learned that he was the younger brother of [Jesus](#). At least 20 million people died, mainly civilians, in one of

Query: "world war one"

	qf	cf	p(w q)	p(w d)	Score
world	2,500	90,000	0.202	8.75E-05	-2.723
war	2,000	35,000	0.202	0.001	-2.199
one	6,000	5E+07	0.205	0.049	-0.890
					-5.812

Modeling Language

Although the bag of words model works very well for text classification, it is intuitively unsatisfying – it assumes the words in a document are independent, given the relevance label, and nobody believes this.

What could we replace it with?

- A “bag of paragraphs” wouldn’t work – too many paragraphs are unique in the collection, so we can’t do meaningful statistics without subdividing them.
- A “bag of sentences” is better, but not much – many sentences are unique, and two documents expressing the same thought are unlikely to choose exactly the same sentence. We need similar documents to have similar features.
- We’ll use sets of words, called *n-grams*, and consider sets of different sizes to balance between good probability estimates (for small *n*) and semantic nuance (for large *n*).

$prob(w|d) = ?$ naive prob = $\frac{count(w,d)}{D}$

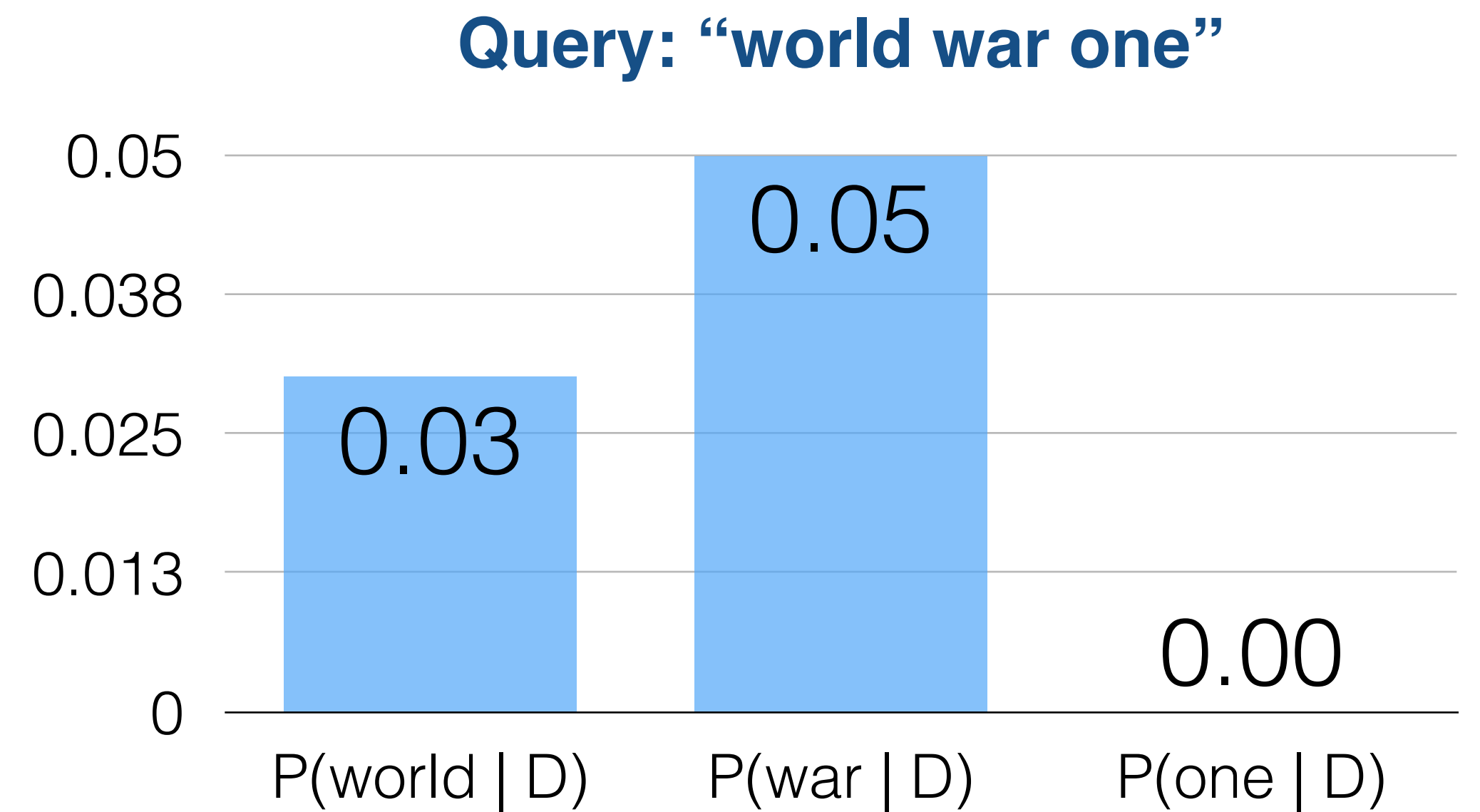
Probability Estimation

\gg naive

Maximum likelihood probability estimates assign zero probability to terms missing from the training data.

This is catastrophic for a Naive Bayes retrieval model: any document that doesn't contain all query terms will get a matching score of zero.

Many other probabilistic models have similar problems. Only truly impossible events should have zero probability.



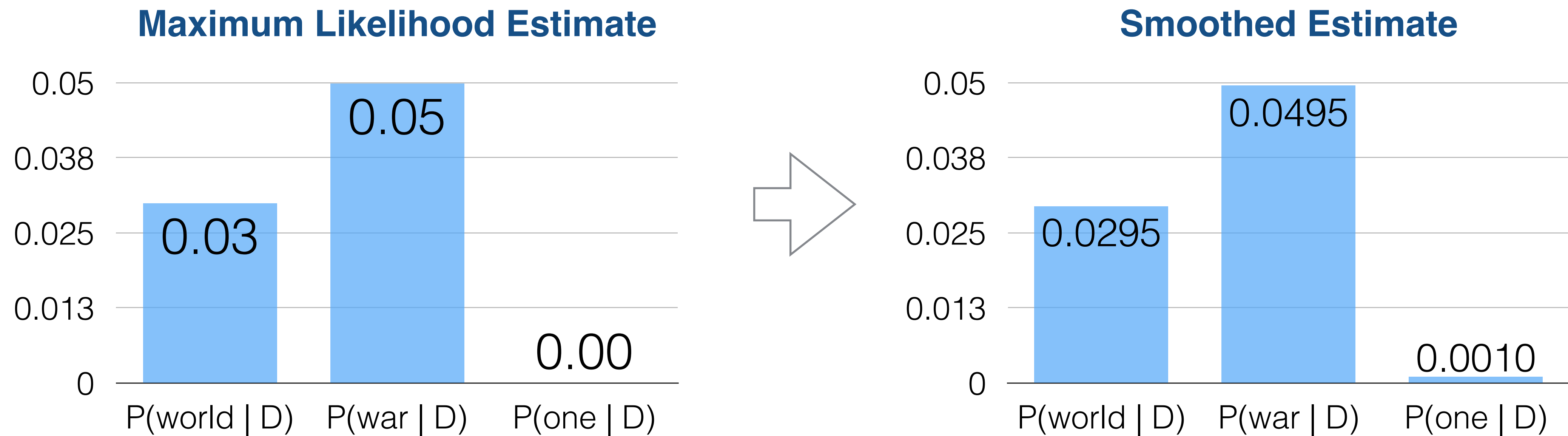
Query Likelihood Model

$$P(D|Q) \stackrel{rank}{=} \prod_{w \in Q} P(w|D)$$

$$= 0.03 \cdot 0.05 \cdot 0$$

Smoothing

The solution is to adjust our probability estimates by taking some probability away from the most-likely events, and moving it to the less-likely events.



This makes the probability distribution less spiky, or “smoother.” The probabilities all move just a little toward the mean.

Smoothing

Smoothing is important for many reasons.

- Assigning zero probability to possible events is incorrect.
- Maximum likelihood estimates from your data don't generalize perfectly to new data, so a Bayesian update from some kind of prior works better.

However, *uniform* smoothing doesn't work very well for language modeling. Next, we'll see why that is, and how we can do better.

Chengxiang Zhai and John Lafferty. 2004. A study of smoothing methods for language models applied to information retrieval.

Laplace Smoothing

Laplace Smoothing, aka “add-one smoothing,” smooths maximum likelihood estimates by adding one count to each event.

$$P(e) = \frac{\text{count}(e) + 1}{\sum_{e \in \text{events}} (\text{count}(e) + 1)}$$

$$P(w|d) = \frac{tf_{w,d} + 1}{|d| + |V|}$$

This is equivalent to a Bayesian posterior with a uniform prior, as we'll see.



Pierre-Simon Laplace (1745-1827)

Image from Wikipedia

Deriving Laplace Smoothing

If we assume nothing about a document's vocabulary distribution, we will use uniform probabilities for all terms.

When we observe the terms in a document, the Bayesian update of these probabilities yields Laplace smoothing.

This Bayesian posterior is our smoothed estimate of the vocabulary distribution for the document's topic.

$$P(\pi|a) \text{ is } \textit{Dirichlet}(\pi|a_1, \dots, a_V) \propto \prod_{i=1}^V \pi_i^{a_i-1}$$

$$P(d|\pi) \text{ is } \textit{Multinomial}(\pi) \propto \prod_{i=1}^V \pi_i^{tf_{i,d}}$$

$$P(w|d) \propto P(d|\pi)P(\pi|a) = \prod_{i=1}^V \pi_i^{tf_{i,d}+a_i-1}$$

is *Dirichlet*($\pi|a_1 + tf_{1,d}, \dots, a_V + tf_{V,d}$)

$$\mathbb{E}[P(w|d)|a = 1] = \frac{tf_{1,d} + 1}{|d| + V}$$

Add- α Smoothing

Laplace smoothing can be generalized from add-one smoothing to add- α smoothing, for $\alpha \in (0, 1]$.

This lets you tune the amount of smoothing you want to use: smaller values of α are closer to the maximum likelihood estimate.

$$P(e) = \frac{\text{count}(e) + a}{\sum_{e \in \text{events}} (\text{count}(e) + a)}$$

$$P(w|d) = \frac{tf_{w,d} + a}{|d| + a|V|}$$

Limits of Uniform Smoothing

Uniform smoothing assigns the same probability to all unseen words, which isn't realistic. This is easiest to see for n-gram models:

$$P(\textit{house}|\textit{the}, \textit{white}) > P(\textit{effortless}|\textit{the}, \textit{white})$$

We strongly believe that “house” is more likely to follow “the white” than “effortless” is, even if neither trigram appears in our training data.

Our bigram counts should help: “white house” probably appears more often than “white effortless.” We can use bigram probabilities as a *background distribution* to help smooth our trigram probabilities.

Jelinek-Mercer Smoothing

One way to combine foreground and background distributions is to take their linear combination. This is the simplest form of Jelinek-Mercer Smoothing.

$$\hat{p}(e) = \lambda p_{fg}(e) + (1 - \lambda) p_{bg}(e), 0 < \lambda < 1$$

For instance, you can smooth n-grams with (n-1)-gram probabilities.

$$\hat{p}(w_n | w_1, \dots, w_{n-1}) = \lambda p(w_n | w_1, \dots, w_{n-1}) + (1 - \lambda) p(w_n | w_2, \dots, w_{n-1})$$

You can also smooth document estimates with corpus-wide estimates.

$$\hat{p}(w|d) = \lambda \frac{tf_{w,d}}{|d|} + (1 - \lambda) \frac{cf_w}{\sum_w cf_w}$$

Relationship to Laplace Smoothing

Most smoothing techniques amount to finding a particular value for λ in Jelinek-Mercer smoothing.

For instance, add-one smoothing is Jelinek-Mercer smoothing with a uniform background distribution and a particular value of λ .

$$\begin{aligned}\text{Pick } \lambda &= \frac{|d|}{|d| + |V|} \\ \hat{p}(w|d) &= \lambda \frac{tf_{w,d}}{|d|} + (1 - \lambda) \frac{1}{|V|} \\ &= \left(\frac{|d|}{|d| + |V|} \right) \frac{tf_{w,d}}{|d|} + \left(\frac{|V|}{|d| + |V|} \right) \frac{1}{|V|} \\ &= \frac{tf_{w,d}}{|d| + |V|} + \frac{1}{|d| + |V|} \\ &= \frac{tf_{w,d} + 1}{|d| + |V|}\end{aligned}$$

Relationship to TF-IDF

TF-IDF is also closely related to Jelinek-Mercer smoothing.

If you smooth the query likelihood model with a corpus-wide background probability, the resulting scoring function is proportional to TF and inversely proportional to DF.

$$\begin{aligned}\log P(q|d) &= \sum_{w \in q} \log \left(\lambda \frac{tf_{w,d}}{|d|} + (1 - \lambda) \frac{df_w}{|c|} \right) \\ &= \sum_{w \in q: tf_{w,d} > 0} \log \left(\lambda \frac{tf_{w,d}}{|d|} + (1 - \lambda) \frac{df_w}{|c|} \right) + \sum_{w \in q: tf_{w,d} = 0} \log(1 - \lambda) \frac{df_w}{|c|} \\ &= \sum_{w \in q: tf_{w,d} > 0} \log \left(\frac{\lambda \frac{tf_{w,d}}{|d|} + (1 - \lambda) \frac{df_w}{|c|}}{(1 - \lambda) \frac{df_w}{|c|}} \right) + \sum_{w \in q} \log(1 - \lambda) \frac{df_w}{|c|} \\ &\stackrel{\text{rank}}{=} \sum_{w \in q: tf_{w,d} > 0} \log \left(\frac{\lambda \frac{tf_{w,d}}{|d|}}{(1 - \lambda) \frac{df_w}{|c|}} + 1 \right)\end{aligned}$$

Dirichlet Smoothing

Dirichlet Smoothing is the same as Jelinek-Mercer smoothing, picking λ based on document length and a parameter μ – an estimate of the average doc length.

$$\lambda = 1 - \frac{\mu}{|d| + \mu}$$

The scoring function to the right is the Bayesian posterior using a Dirichlet prior with parameters:

$$\left(\mu \frac{cf_{w_1}}{\sum_w cf_w}, \dots, \mu \frac{cf_{w_n}}{\sum_w cf_w} \right)$$

$$\hat{p}(w|d) = \frac{tf_{w,d} + \mu \frac{cf_w}{\sum_w cf_w}}{|d| + \mu}$$

$$\log p(q|d) = \sum_{w \in q} \log \frac{tf_{w,d} + \mu \frac{cf_w}{\sum_w cf_w}}{|d| + \mu}$$

Example: Dirichlet Smoothing

Query: “president lincoln”

tf	15
cf	160,000
tf	25
cf	2,400
 d 	1,800
Σ	10
μ	2,000

$$\begin{aligned}\log p(q|d) &= \sum_{w \in q} \log \frac{tf_{w,d} + \mu \frac{cf_w}{\sum_w cf_w}}{|d| + \mu} \\ &= \log \frac{15 + 2,000 \times (160,000/10^9)}{1,800 + 2,000} \\ &\quad + \log \frac{25 + 2,000 \times (2,400/10^9)}{1,800 + 2,000} \\ &= \log(15.32/3,800) + \log(25.005/3,800) \\ &= -5.51 + -5.02 \\ &= -10.53\end{aligned}$$

Effect of Dirichlet Smoothing

Dirichlet Smoothing is a good choice for many IR tasks.

- As with all smoothing techniques, it never assigns zero probability to a term.
- It is a Bayesian posterior which considers how the document differs from the corpus.
- It normalizes by document length, so estimates from short documents and long documents are comparable.
- It runs quickly, compared to many more exotic smoothing techniques.

tf	tf	ML Score	Smoothed Score
15	25	-3.937	-10.53
15	1	-5.334	-13.75
15	0	N/A	-19.05
1	25	-5.113	-12.99
0	25	N/A	-14.4

Witten-Bell Smoothing

Dirichlet Smoothing is the same as Jelinek-Mercer smoothing, picking λ based on

- * doc length $|d|$

- * doc vocabulary $|V|$ (number of unique terms in document)

$$\lambda = \frac{|d|}{|d| + |V|}$$

N-grams and Skip-grams

An *n-gram* is an ordered set of n contiguous words, usually found within a single sentence. Special cases are $n = 1$ (unigrams), $n = 2$ (bigrams), and $n = 3$ (trigrams).

Skip-grams are more “relaxed” – they can appear in any order, and need not be adjacent. They are an unordered set of n words that appear within a fixed window of k words.

Sentence

The quick brown fox jumped over the lazy dog.

Trigrams ($n = 3$)

the quick brown
quick brown fox
brown fox jumped

...

Skip-grams ($n = 3, k = 5$)

quick brown fox
fox jumped quick
lazy dog jumped

...

Markov Chains

We typically construct a generative model of n-grams using *Markov chains* – what is the probability distribution over the next word in the n-gram, given the $n - 1$ words we've seen so far?

$$P(w_n | w_1, w_2, \dots, w_{n-1})$$

This assumes that words are independent, given the relevance label *and the preceding $n - 1$ words*.

We use a special token, like \$, for words “before” the beginning of the sentence.

Sentence

The quick brown fox jumped over the lazy dog.

Trigram Sentence Probability

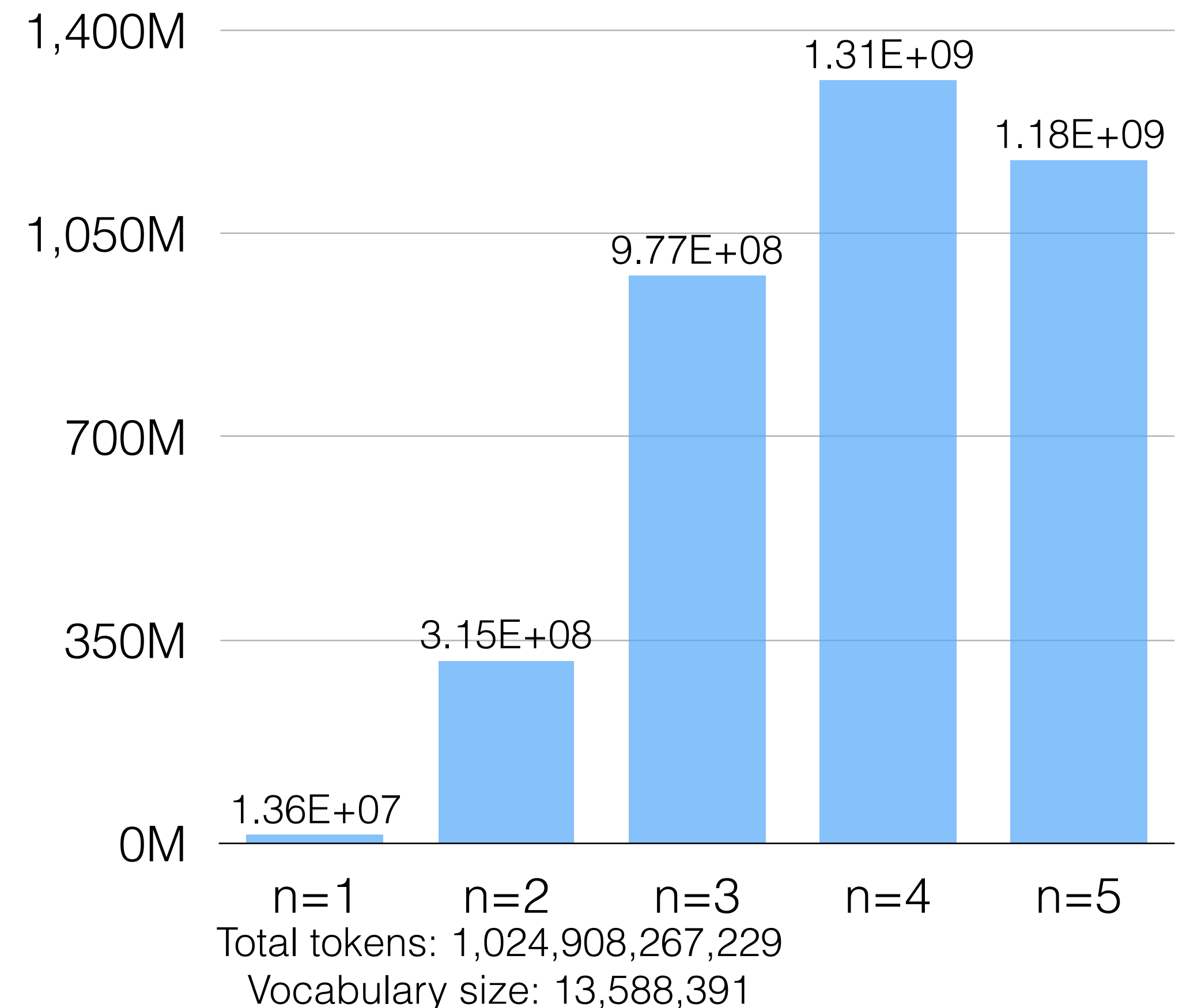
$$\begin{aligned} &P(\textit{the} | \$, \$) \cdot P(\textit{quick} | \$, \textit{the}) \cdot P(\textit{brown} | \textit{the}, \textit{quick}) \\ &\cdot P(\textit{fox} | \textit{quick}, \textit{brown}) \cdot P(\textit{jumped} | \textit{brown}, \textit{fox}) \\ &\cdot P(\textit{over} | \textit{fox}, \textit{jumped}) \cdot P(\textit{the} | \textit{jumped}, \textit{over}) \\ &\cdot P(\textit{lazy} | \textit{over}, \textit{the}) \cdot P(\textit{dog} | \textit{the}, \textit{lazy}) \end{aligned}$$

Number of n-grams in a Corpus

How many n-grams do we expect to see, as a function of the vocabulary size v and n-gram size n ?

- At first glance, you'd expect to see $\binom{v}{n} = O(v^n)$
- However, most possible n-grams will never appear (like “correct horse battery staple?”), and n-grams are limited by typical sentence lengths.
- As n increases, the number of distinct observed n-grams peaks around $n = 4$ and then decreases.

Web 1T 5-gram Corpus



Choosing n-gram Size

The best n-gram size to use depends on a variance-bias tradeoff:

- Smaller values of n have more training data: infrequent n-grams will appear more often, reducing the *variance* of your probability estimates.
- Larger values of n take more context into account: they have more semantic information, reducing the *bias* of your probability estimates.

The best n-gram size is the largest value your data will support. Common choices are $n = 3$ for millions of words, or $n = 2$ for smaller corpora.

Wrapping Up

Using n-grams and skip-grams allows us to include some linguistic context in our retrieval models. This helps disambiguate word senses and improve retrieval performance.

Larger values of n are beneficial, if you have the data to support them. The number of n-grams does not grow exponentially in n , so the index size can be manageable.

Next, we'll see how to use an n-gram language model for retrieval.



language models for retrieval

many slides courtesy James Allan@umass Amherst
some slides courtesy Christopher Manning and Prabhakar Raghavan @ Stanford



what is a retrieval model?

- Model is an idealization or abstraction of an actual process
- Mathematical models are used to study the properties of the process, draw conclusions, make predictions
- Conclusions derived from a model depend on whether the model is a good approximation of the actual situation
- Statistical models represent repetitive processes, make predictions about frequencies of interesting events
- Retrieval models can describe the computational process
 - e.g. how documents are ranked
 - Note that how documents or indexes are *stored* is implementation
- Retrieval models can attempt to describe the human process
 - e.g. the information need, interaction
 - Few do so meaningfully
- Retrieval models have an explicit or implicit definition of relevance



retrieval models

- boolean
- vector space
- latent semantic indexing

today

- statistical language
- inference network
- hyperlink based



language models not in MIR

1. J. M. Ponte and W. B. Croft. A language modeling approach to information retrieval. Proceedings of ACM-SIGIR 1998, pages 275-281.
2. J. M. Ponte. A language modeling approach to information retrieval. Phd dissertation, University of Massachusetts, Amherst, MA, September 1998.
3. D. Hiemstra. Using Language Models for Information Retrieval. PhD dissertation, University of Twente, Enschede, The Netherlands, January 2001.
4. D. R. H. Miller, T. Leek, and R. M. Schwartz. A hidden Markov model information retrieval system. Proceedings of ACM-SIGIR 1999, pages 214-221.
5. F. Song and W. B. Croft. A general language model for information retrieval. In Proceedings of Eighth International Conference on Information and Knowledge Management (CIKM 1999)
6. S. F. Chen and J. T. Goodman. An empirical study of smoothing techniques for language modeling. In Proceedings of the 34th Annual Meeting of the ACL, 1996.
7. C. Zhai and J. Lafferty. A study of smoothing methods for language models applied to ad hoc information retrieval. Proceedings of the ACM-SIGIR 2001, pages 334-342.
8. V. Lavrenko and W. B. Croft. Relevance-based language models. Proceedings of the ACM SIGIR 2001, pages 120-127.
9. V. Lavrenko and W. B. Croft, Relevance Models in Information Retrieval, in Language Modeling for Information Retrieval, W. Bruce Croft and John Lafferty, ed., Kluwer Academic Publishers, chapter 2.



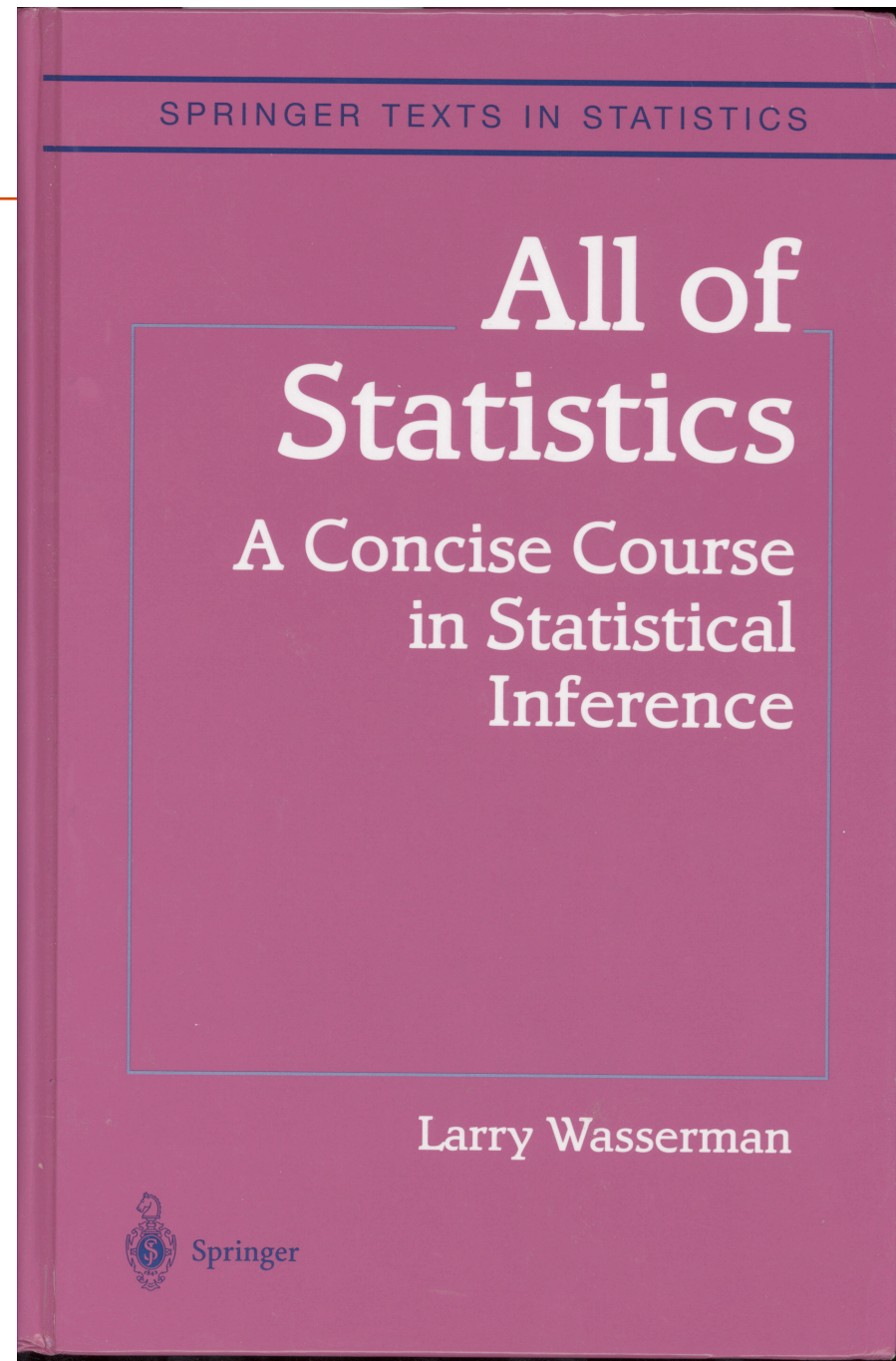
outline

- review: probabilities
- language model
- similarity, ranking in LM
- probability estimation
- smoothing methods
- examples



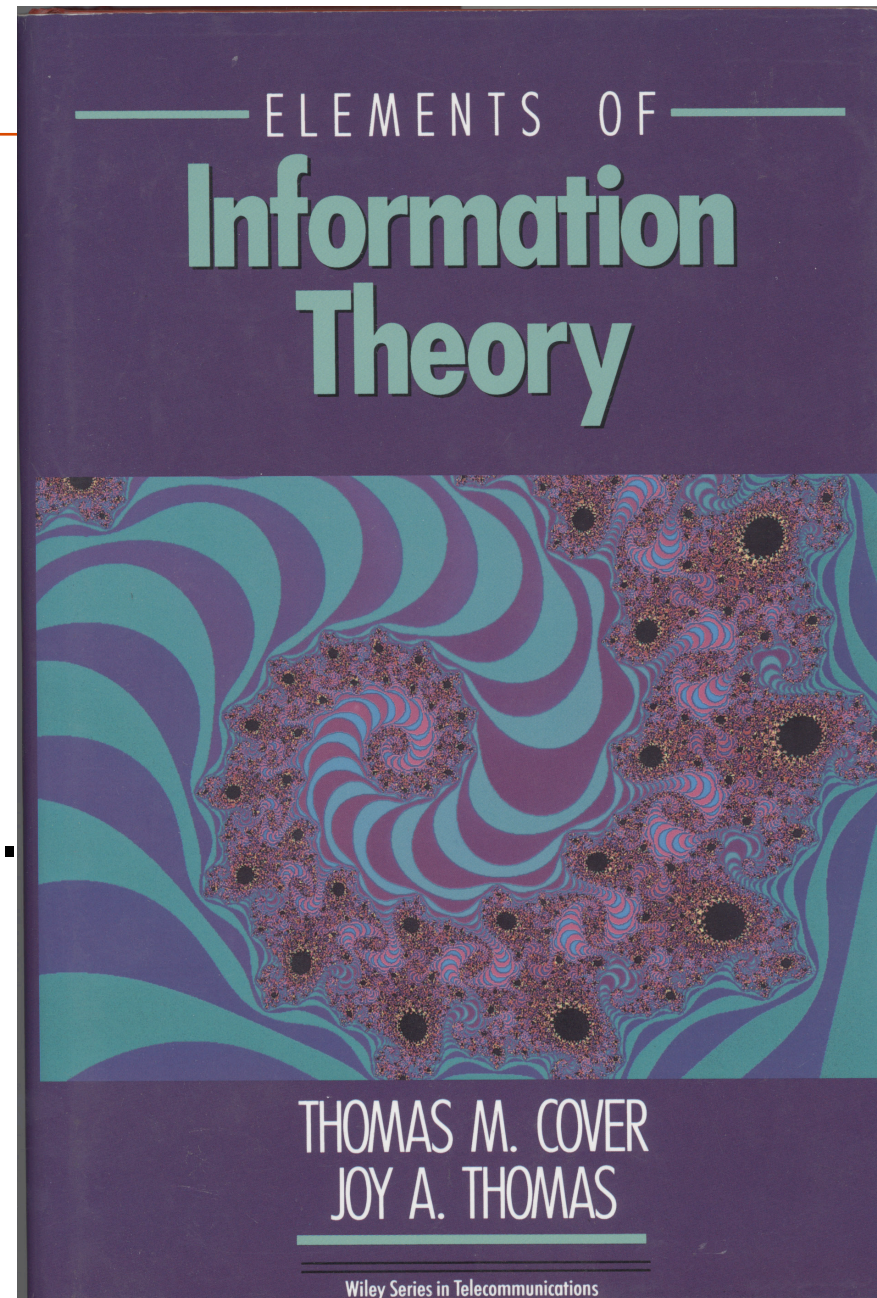
probabilities

- sample space
- probability
- independent events
- cond. probability
- Bayes theorem
- distributions



information theory, coding

- entropy
- joint entropy
- cond. entropy
- relative entropy
- convexity, Jensen ineq.
- optimal coding
- Fano's ineq.





outline

- review: probabilities
- language model
- similarity, ranking in LM
- probability estimation
- smoothing methods
- examples



what is a language model ?

- Probability distribution over strings of text
 - how likely is a given string (observation) in a given “language”
 - for example, consider probability for the following four strings

$p_1 = P(\text{“a quick brown dog”})$

$p_2 = P(\text{“dog quick a brown”})$

$p_3 = P(\text{“быстрая brown dog”})$

$p_4 = P(\text{“быстрая собака”})$

- English: $p_1 > p_2 > p_3 > p_4$

- ... depends on what “language” we are modeling
 - In most of IR, assume that $p_1 == p_2$



lang modeling for IR

- Every document in a collection defines a “language”
 - consider all possible sentences (strings) that author could have written down when creating some given document
 - some are perhaps more likely to occur than others
- subject to topic, writing style, language ...
 - $P(s|MD)$ = probability that author would write down string “s”
- think of writing a billion variations of a document and counting how many time we get “s”
- Now suppose “Q” is the user’s query
 - what is the probability that author would write down “q” ?
- Rank documents D in the collection by $P(Q|MD)$
 - probability of observing “Q” during random sampling from the language model of document D

language models

- estimate probabilities of certain "events" in the text

- based on these probabilities, use likelihood as similarity

generative model

"generator"

task: write an automatic text generator.

- language model based on

- letters?

- words?

- phrases?

which unit for probab?

⇒ prosa stilistic generator



statistical text generation

1. *Zero-order approximation.* (The symbols are independent and equiprobable.) $p(\text{letters}) = \text{uniform}$

XFOML RXKHRJFFJUJ ZLPWCFWKCYJ
FFJEYVKCQSGXYD QPAAMKBZAACIBZLHJQD

~~28 probs~~ 6 probs

2. *First-order approximation.* (The symbols are independent. Frequency of letters matches English text.) $P(\text{letter}) = \text{english probs}$

OCRO HLI RGWR NMIELWIS EU LL NBNESEBYA TH EEI
ALHENHTTPA OOBTTVA NAH BRL

28 probs

3. *Second-order approximation.* (The frequency of pairs of letters matches English text.) $\text{prob}("xy") = \text{english prob}$

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY
ACHIN D ILONASIVE TUCOOWE AT TEASONARE FUSO
TIZIN ANDY TOBE SEACE CTISBE

4. *Third-order approximation.* (The frequency of triplets of letters matches English text.) $\text{prob}("xyz") = \text{same as in english}$

IN NO IST LAT WHEY CRATICT FROURE BERS GROCID
PONDENOME OF DEMONSTURES OF THE REPTAGIN IS
REGOACTIONA OF CRE



5. *Fourth-order approximation.* (The frequency of quadruplets of letters matches English text. Each letter depends on the previous three letters. This sentence is from Lucky's book, *Silicon Dreams* [183].)

$$P(\text{"XYZT"}) = \text{correct}$$

THE GENERATED JOB PROVIDUAL BETTER TRAND THE
DISPLAYED CODE, ABOVERY UPONDULTS WELL THE
CODERST IN THESTICAL IT DO HOCK BOTHE MERG.

(INSTATES CONS ERATION. NEVER ANY OF PUBLE AND TO
THEORY. EVENTIAL CALLEGAND TO ELAST BENERATED IN
WITH PIES AS IS WITH THE)

Instead of continuing with the letter models, we jump to word models.

↳ jump

6. *First-order word model.* (The words are chosen independently but with frequencies as in English.)

$$P(\text{words}) = \text{correct}$$

REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME
CAN DIFFERENT NATURAL HERE HE THE A IN CAME THE TO
OF TO EXPERT GRAY COME TO FURNISHES THE LINE
MESSAGE HAD BE THESE.

7. *Second-order word model.* (The word transition probabilities match English text.)

$$P(\text{"XYZ ABCD"}) = \text{correct}$$

THE HEAD AND INFRONTAL ATTACK ON AN ENGLISH
WRITER THAT THE CHARACTER OF THIS POINT IS
THEREFORE ANOTHER METHOD FOR THE LETTERS THAT THE
TIME OF WHO EVER TOLD THE PROBLEM FOR AN
UNEXPECTED

what prob
Model?

- letters?

- words?

unigrams

- phrases

8. word-trigram

9. phrases

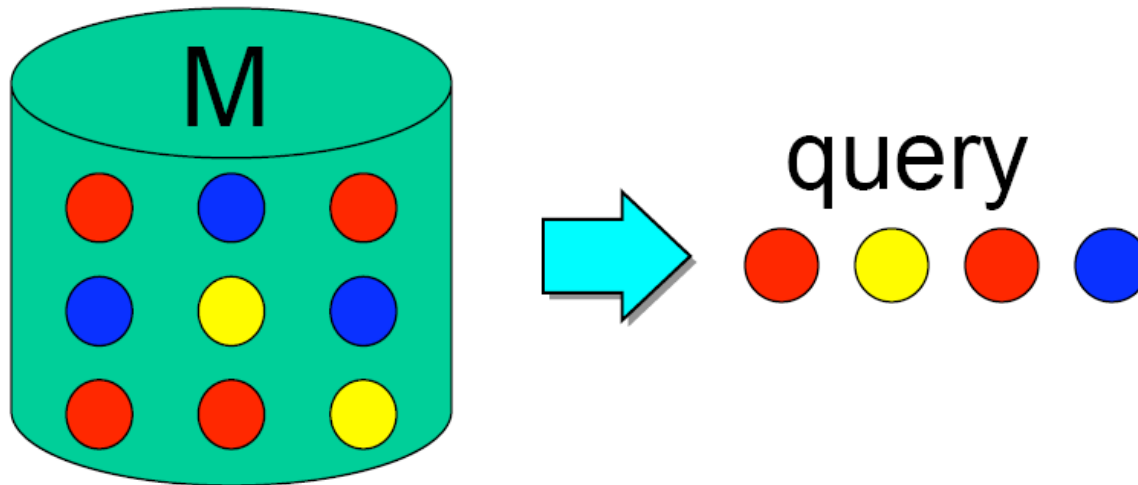


LM choices

- What kind of language model should we use?
 - Unigram or higher-order models?
 - Multinomial or multiple-Bernoulli?
- How can we estimate model parameters?
 - Basic models
 - Translation models
 - Aspect models
 - non-parametric models
- How can we use the model for ranking?
 - Query-likelihood
 - Document-likelihood
 - Likelihood Ratio
 - Divergence of query and document models

unigram LM

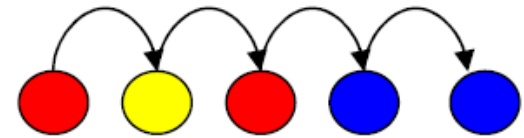
- words are sampled independently, with replacement
- order of the words is lost (no phrases)



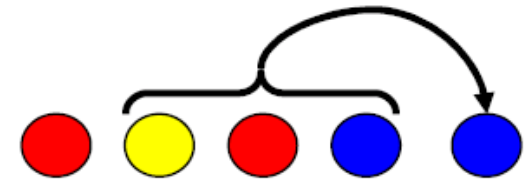
$$\begin{aligned} P(\text{red, yellow, red, blue}) &= P(\text{red}) P(\text{yellow}) P(\text{red}) P(\text{blue}) \\ &= 4/9 * 2/9 * 4/9 * 3/9 \end{aligned}$$

higher-order LM

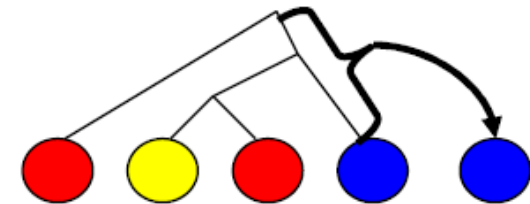
- Unigram model assumes word independence
 - cannot capture surface form: $P(\text{"brown dog"}) == P(\text{"dog brown"})$
- Higher-order models
 - n-gram: condition on preceding words



- cache: condition on a window (cache)



- grammar: condition on parse tree



- Are they useful?
 - no improvements from n-gram, grammar-based models
 - some research on cache-like models (proximity, passages, etc.)
 - parameter estimation is prohibitively expensive

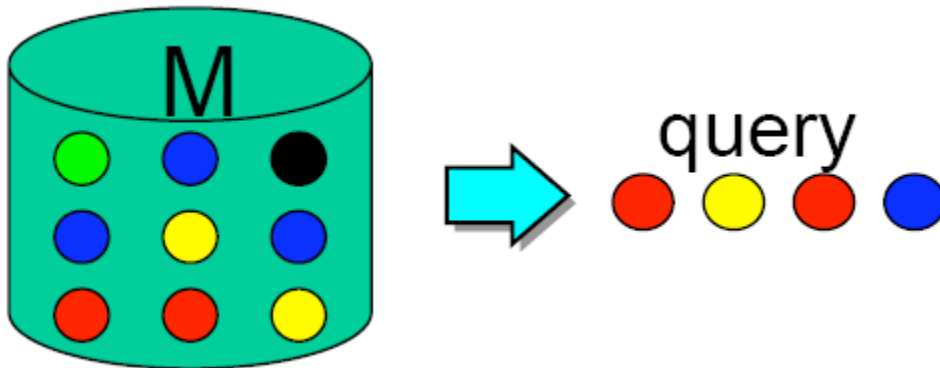


outline

- review: probabilities
- language model
- similarity, ranking in LM
- probability estimation
- smoothing methods
- examples

multinomial similarity

- Predominant model
- Fundamental event:
what is the identity of the i 'th query token?
- observation is a sequence of events, one for each query token

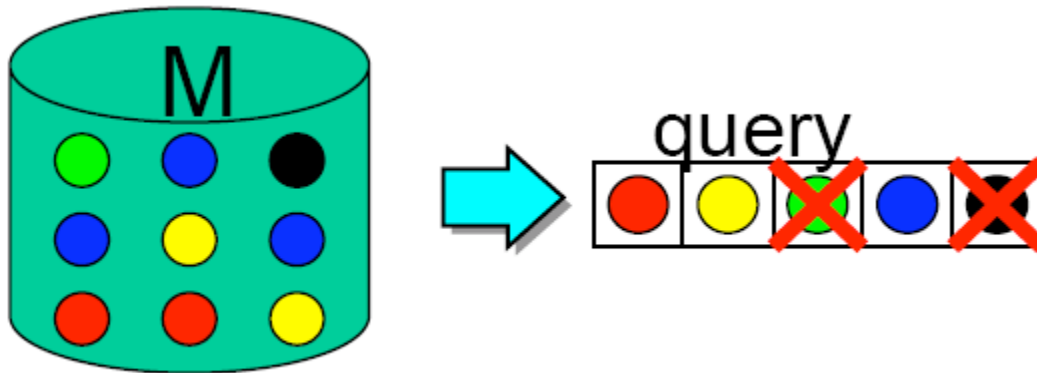


$$P(q_1 \dots q_k | M) = \prod_{i=1}^k P(q_i | M)$$



multiple-Bernoulli similarity

- Original model
- fundamental event: *does the word w occur in the query?*
- observation is a vector of binary events, one for each possible word



$$P(q_1 \dots q_k | M) = \prod_{w \in q_1 \dots q_k} P(w | M) \prod_{w \notin q_1 \dots q_k} [1 - P(w | M)]$$



score, ranking in LM

- what is the probability to generate the given query, given a language model?
- what is the probability to generate the given document, given a language model?
- how "close" are 2 statistical models?



score: query likelihood

- Standard approach: query-likelihood
 - estimate a language model M_D for every document D in the collection
 - rank docs by the probability of “generating” the query

$$P(q_1 \dots q_k \mid M_D) = \prod_{i=1}^k P(q_i \mid M_D)$$

- |
 - no notion of relevance in the model: everything is random sampling
 - user feedback / query expansion not part of the model
 - examples of relevant documents cannot help us improve the language model M_D
 - does not directly allow weighted or structured queries



score: document likelihood

- Flip the direction of the query-likelihood approach
 - estimate a language model M_Q for the query Q
 - rank docs D by the likelihood of being a random sample from M_Q
 - M_Q expected to “predict” a typical relevant document
- Problems:
 - different doc lengths, probabilities not comparable
 - favors documents that contain frequent (low content) words

$$P(D | M_Q) = \prod_{w \in D} P(w | M_Q)$$



score: likelihood ratio

- Try to fix document likelihood:
 - Bayes' likelihood that M_q was the source, given that we observed D
 - related to Probability Ranking Principle: $P(D|R) / P(D|N)$
 - allows relevance feedback, query expansion, etc.
 - can benefit from complex estimation of the query model MQ

$$P(M_q | D) = \frac{P(M_q)P(D | M_q)}{P(D)} \approx \frac{c \prod_{w \in D} P(w | M_q)}{\prod_{w \in D} P(w | GE)}$$

score: model comparison

- Combine advantages of two ranking methods
 - estimate a model of both the query MQ and the document MD
 - directly compare similarity of the two models
 - natural measure of similarity is cross-entropy (others exist):

$$H(M_Q || M_D) = - \sum_w P(w | M_Q) \log P(w | M_D)$$

- number of bits we would need to "encode" MQ using MD
 - equivalent to Kullback-Leibler divergence
 - equivalent to query-likelihood if MQ is simply counts of words in Q
- Cross-entropy is not symmetric: use H (MQ || MD)
 - reverse works consistently worse, favors different document
 - use reverse if ranking multiple queries w.r.t. one document

Q.L. model

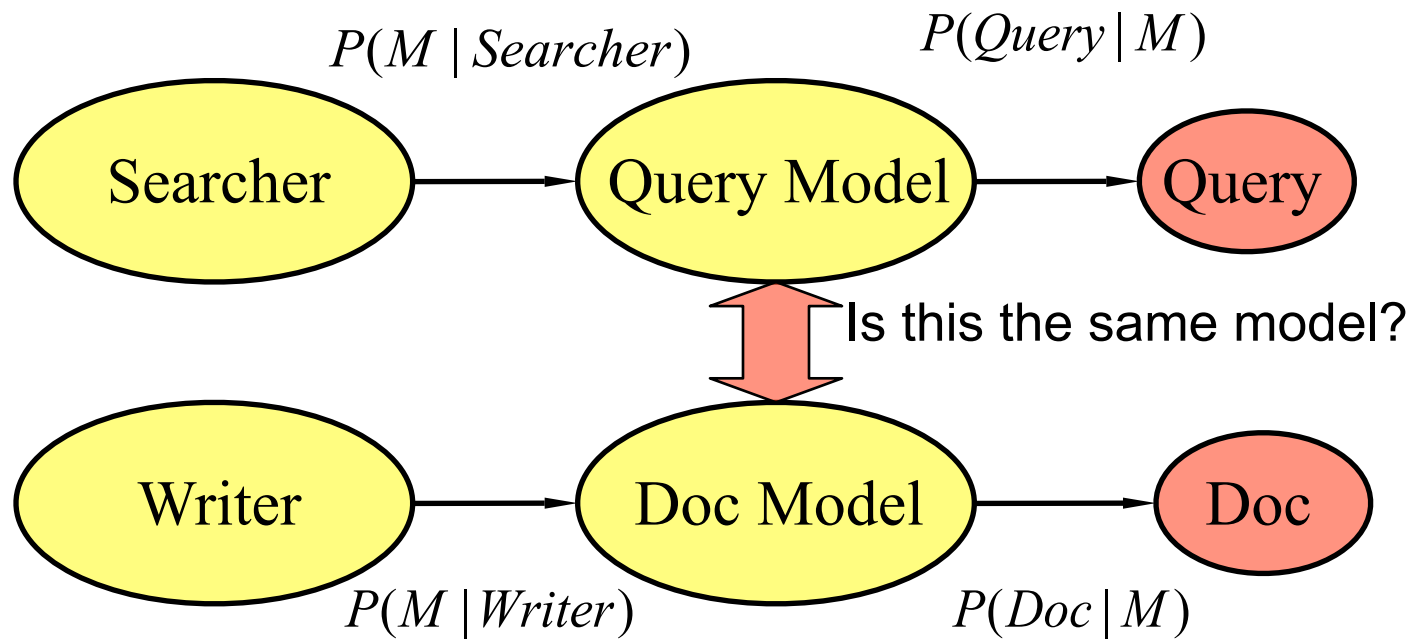
$$\text{score}(d) \approx \sum_{w \in q} \log [P(w|d)]$$

basic estim $w \in q$

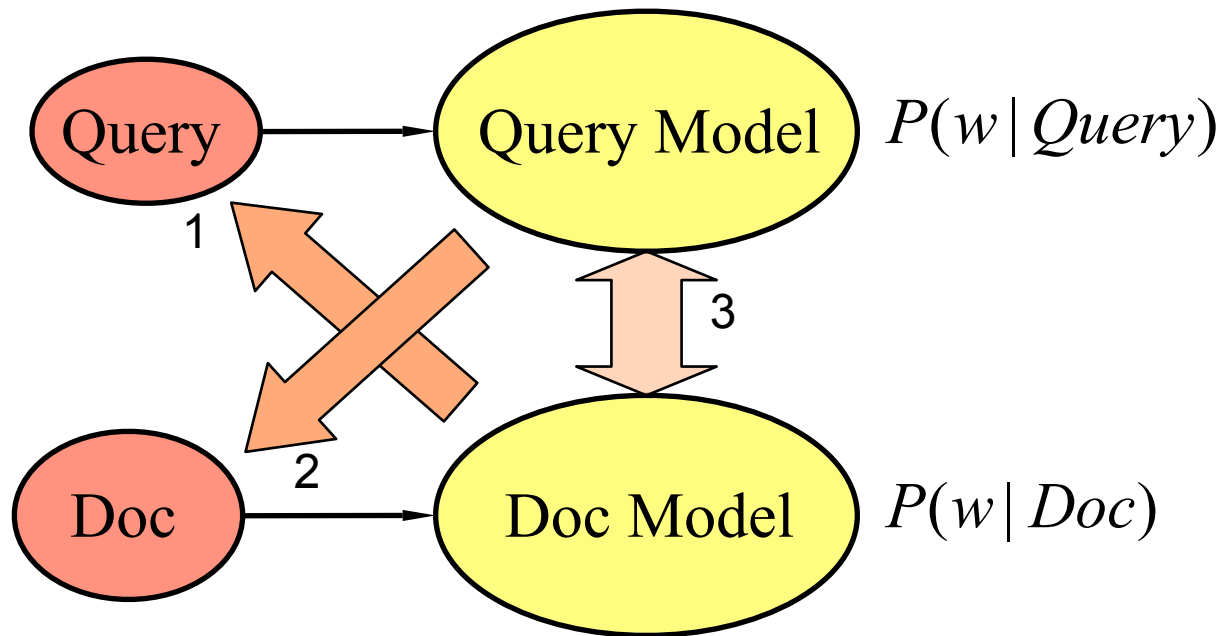
$$w \notin d \Rightarrow P(w|d) = 0 \Rightarrow \log(P(w|d)) = -\infty$$

$\frac{TF(w,d)}{DL(d)}$

Models of Text Generation



Retrieval with Language Models



Retrieval:

- Query likelihood (1)
- Document likelihood (2)
- Model comparison (3)



LM: popular choices

- Use Unigram models
 - no consistent benefit from using higher order models
 - estimation is much more complex (e.g. bi-gram from a 3-word query)
- Use Multinomial models
 - well-studied, consistent with other fields that use LMs
 - extend multiple-Bernoulli model to non-binary events?
- Use Model Comparison for ranking
 - allows feedback, expansion, etc. through estimation of MQ and MD
 - use $KL(MQ || MD)$ for ranking multiple documents against a query
- Estimation of MQ and MD is a crucial step
 - very significant impact on performance (more than other choices)
 - key to cross-language, cross-media and other applications



Translation model (Berger and Lafferty)

- Basic LMs do not address issues of synonymy.
 - Or any deviation in expression of information need from language of documents
- A translation model lets you generate query words not in document via “translation” to synonyms etc.
 - Or to do cross-language IR, or multimedia IR

$$P(\vec{q} | M) = \prod_i \sum_{v \in \text{Lexicon}} \underbrace{P(v | M)}_{\text{Basic LM}} \underbrace{T(q_i | v)}_{\text{Translation}}$$

- Need to learn a translation model (using a dictionary or via statistical machine translation)

outline



SMOOTHING { ① estimation with very few words (short DL)
② $w \notin d \Rightarrow P(w|d) = 0$ breaks the whole score ZERO FREQ

- review: probabilities
- language model
- similarity, ranking in LM

- probability estimation
- smoothing methods

- examples

prob(words)

too few trials
good when it fails?

coin bias

basic estimation from trials

flip coin N times

$$p(\text{head}) = \frac{\# \text{heads}}{\# \text{trials}} = \frac{n}{N}$$

prob($w|D$) = ?
model - prob of term in doc?
basic prob($w|D$) = $\frac{\#(w, D)}{DL(D)}$



$$N=3 \\ n=2 \Rightarrow p(\text{head}) = \frac{2}{3}$$

estimation

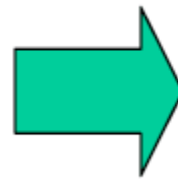
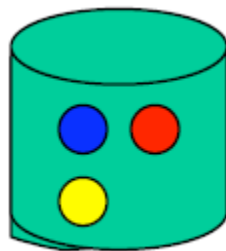
same problem for very short docs

- Want to estimate MQ and/or MD from Q and/or D
- General problem:
 - given a string of text S ($= Q$ or D), estimate its language model MS
 - S is commonly assumed to be an i.i.d. random sample from MS
 - Independent and identically distributed
- Basic Language Models
 - maximum-likelihood estimator and the zero frequency problem
 - discounting, interpolation techniques
 - Bayesian estimation

maximum likelihood

- count relative frequencies of words in S
- maximum-likelihood property:
 - assigns highest possible likelihood to the observation
- unbiased estimator:
 - if we repeat estimation an infinite number of times with different starting points S , we will get correct probabilities (on average)
 - this is not very useful...

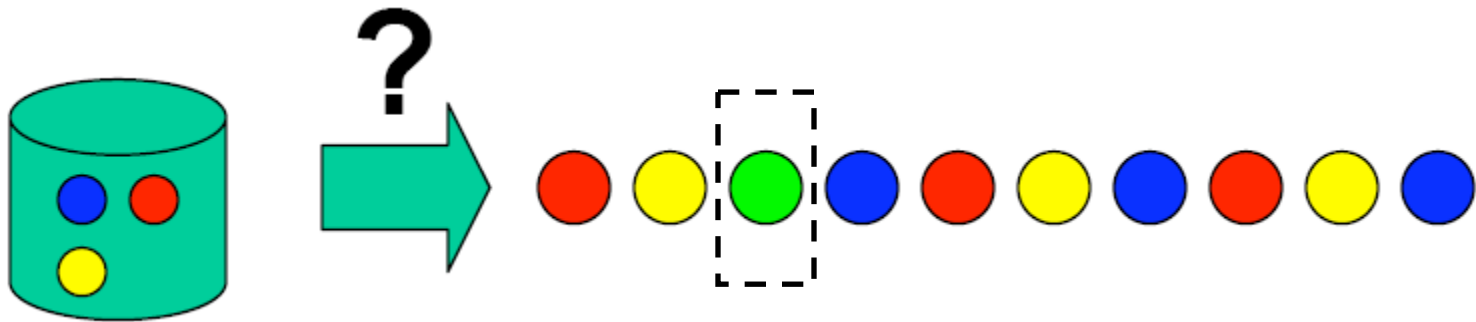
$$P_{ml}(w|M_S) = \#(w, S) / |S|$$



$$\begin{aligned} P(\bullet) &= 1/3 \\ P(\bullet) &= 1/3 \\ P(\bullet) &= 1/3 \\ P(\bullet) &= 0 \\ P(\bullet) &= 0 \end{aligned}$$

zero-frequency problem

- Suppose some event not in our observation S
 - Model will assign zero probability to that event
 - And to any set of events involving the unseen event
- Happens very frequently with language
- It is incorrect to infer zero probabilities
 - especially when creating a model from short samples



Discount-smooth

Laplace smoothing

derived very complicated (with integrals)

$K = \# \text{ possible} = 2$
wins Laplace
 N trial $H+T=N$
heads tails

- count events in observed data
- add 1 to every count
- renormalize to obtain probabilities
- it corresponds to uniform priors

Laplace estimate

$$\left[\frac{H+1}{N+2}, \frac{T+1}{N+2} \right]$$

N trials

- if event counts are (m_1, m_2, \dots, m_k) with $\sum_i m_i = N$ then

max likelihood estimates are $(\frac{m_1}{N}, \frac{m_2}{N}, \dots, \frac{m_k}{N})$

laplace estimates are $(\frac{m_1+1}{N+k}, \frac{m_2+1}{N+k}, \dots, \frac{m_k+1}{N+k})$

basic estimates = Max likelihood

laplace estimates

$$\frac{m_1}{N} \quad \frac{m_2}{N} \quad \dots \quad \frac{m_k}{N} \quad \longrightarrow \quad \frac{m_1+1}{N+k} \quad \frac{m_2+1}{N+k} \quad \dots \quad \frac{m_k+1}{N+k}$$

$\sum_i \frac{m_i}{N} = 1$? yes $\sum_i \frac{m_i+1}{N+k} = \frac{\sum_i m_i + k}{N+k} = \frac{N+k}{N+k} = 1$

Laplace $\frac{deg}{vol}$

$$N = 1000$$

$$w_1 = \#1$$

$$w_2 = \#2$$

$$w_6 = \#6$$

$$k=6$$

direct / (basi) / Max-Likeli /

Laplace

$$\frac{w_1}{N}$$

$$\frac{w_2}{N}$$

$$\frac{w_6}{N}$$

$$\frac{w_1+1}{N+6}$$

$$\frac{w_2+1}{N+6}$$

$$\frac{w_6+1}{N+6}$$

discounting methods



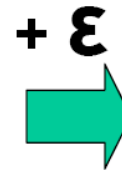
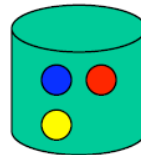
$$\sum_i \frac{w_i + 1}{N + k} = \frac{\sum w_i + \sum 1}{N + k} = \frac{N + k}{N + k} = 1$$

→ valid distribution

- Laplace smoothing

- Lindstone correction

- add ϵ to all count, renormalize



$$\begin{aligned}
 P(\bullet) &= (1 + \epsilon) / (3 + 5\epsilon) \\
 P(\bullet) &= (1 + \epsilon) / (3 + 5\epsilon) \\
 P(\bullet) &= (1 + \epsilon) / (3 + 5\epsilon) \\
 P(\bullet) &= (0 + \epsilon) / (3 + 5\epsilon) \\
 P(\bullet) &= (0 + \epsilon) / (3 + 5\epsilon)
 \end{aligned}$$

- absolute discounting

- subtract ϵ , redistribute probab mass

$$\frac{w_1}{N} \quad \frac{w_2}{N} \quad \dots \quad \frac{w_k}{N}$$

basic

$$\left(\frac{w_1 + \epsilon}{N + k\epsilon} \right) \left(\frac{w_2 + \epsilon}{N + k\epsilon} \right) \dots \left(\frac{w_k + \epsilon}{N + k\epsilon} \right)$$

$\frac{1}{k} = \text{uniform}$ general Laplace (ϵ) $\sum 1 = 1$

ϵ high \Rightarrow strong uniform prior

discounting methods

ϵ low ($\epsilon = 0.0001$) \Rightarrow estimates \approx basic

- Held-out estimation

- Divide data into training and held-out sections
- In training data, count N_r , the number of words occurring r times
- In held-out data, count T_r , the number of times those words occur
- $r^* = T_r/N_r$ is adjusted count (equals r if training matches held-out)
- Use r^*/N as estimate for words that occur r times

- Deleted estimation (cross-validation)

- Same idea, but break data into K sections
- Use each in turn as held-out data, to calculate $T_r(k)$ and $N_r(k)$
- Estimate for words that occur r times is average of each

- Good-Turing estimation

- From previous, $P(w|M) = r^* / N$ if word w occurs r times in sample
- In Good-Turing, steal total probability mass from next most frequent word
- Provides probability mass for words that occur $r=0$ times
- Take what's leftover from $r>0$ to ensure adds to one

$$TPM(r+1) = N_{r+1} \cdot \frac{r+1}{N}$$

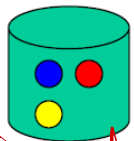
$$\begin{aligned} P(w_r|M) &= TPM(r+1)/N_r \\ &= \frac{N_{r+1}}{N_r} \cdot \frac{r+1}{N} \end{aligned}$$

interpolation methods

discounting

- Problem with all discounting methods:
 - discounting treats unseen words equally (add or subtract ϵ)
 - some words are more frequent than others
- Idea: use background probabilities
 - “interpolate” ML estimates with General English expectations (computed as relative frequency of a word in a large collection)
 - reflects expected frequency of events

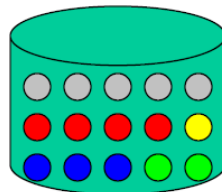
ML estimate



direct estimate
= basic estimate
(from doc)

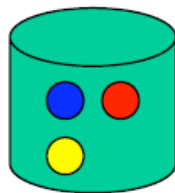
final estimate =

background probability

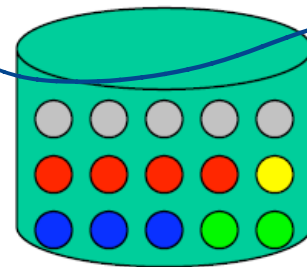


estimate from large
data (from collection,
English prior)

λ



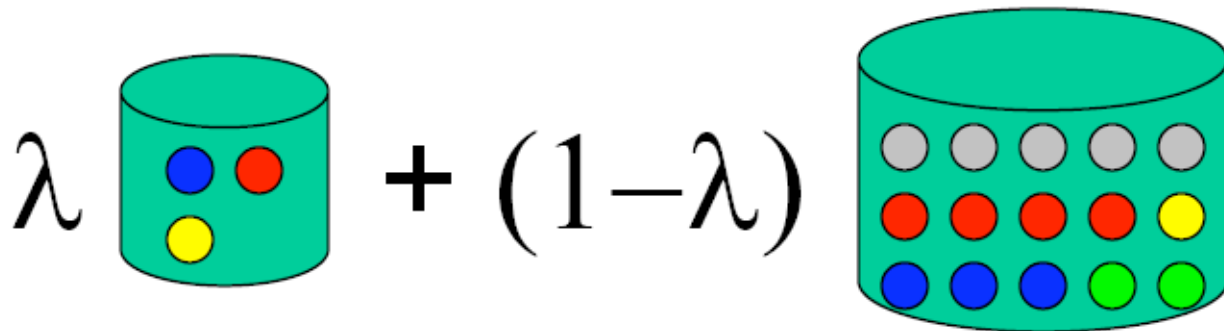
+ (1- λ)





Jelinek Mercer smoothing

- Correctly setting λ is very important
- Start simple
 - set λ to be a constant, independent of document, query
- Tune to optimize retrieval performance
 - optimal value of λ varies with different databases, query sets, etc.



Dirichlet smoothing

- Problem with Jelinek-Mercer:
 - longer documents provide better estimates
 - could get by with less smoothing
- Make smoothing depend on sample size
- N is length of sample = document length
- μ is a constant

$$\underbrace{N / (N + \mu)}_{\lambda} \text{ } \text{cylinder with 3 dots} + \underbrace{\mu / (N + \mu)}_{(1-\lambda)} \text{ } \text{cylinder with 15 dots}$$

The diagram illustrates the Dirichlet smoothing formula. It shows two terms being added together. The first term is $N / (N + \mu)$, which is labeled as λ . This term is represented by a small green cylinder containing three colored dots (blue, red, and yellow). The second term is $\mu / (N + \mu)$, which is labeled as $(1-\lambda)$. This term is represented by a larger green cylinder containing fifteen colored dots (red, blue, green, and yellow).

Witten-Bell smoothing

- A step further:
 - condition smoothing on “redundancy” of the example
 - long, redundant example requires little smoothing
 - short, sparse example requires a lot of smoothing
- Derived by considering the proportion of new events as we walk through example
 - N is total number of events = document length
 - V is number of unique events = number of unique terms in doc

$$\underbrace{N / (N + V)}_{\lambda} \text{ } \text{cylinder with 3 colored balls} + \underbrace{V / (N + V)}_{(1-\lambda)} \text{ } \text{cylinder with 15 colored balls}$$

The diagram illustrates the Witten-Bell smoothing formula. It shows two cylinders representing event distributions. The first cylinder, representing the document's event distribution, contains 3 balls (blue, red, yellow) and is associated with the fraction $N / (N + V)$, which is labeled as λ . The second cylinder, representing the uniform distribution over the vocabulary, contains 15 balls (5 grey, 5 red, 5 blue) and is associated with the fraction $V / (N + V)$, which is labeled as $(1-\lambda)$. The two cylinders are separated by a plus sign, indicating their combination in the smoothing process.



interpolation vs back-off

- Two possible approaches to smoothing
- Interpolation:
 - Adjust probabilities for all events, both seen and unseen
- Back-off:
 - Adjust probabilities only for unseen events
 - Leave non-zero probabilities as they are
 - Rescale everything to sum to one: rescales “seen” probabilities by a constant
- Interpolation tends to work better
 - And has a cleaner probabilistic interpretation

Two-stage smoothing



Query = "the algorithms for data mining"

d1:	0.04	0.001	0.02	0.002	0.003
d2:	0.02	0.001	0.01	0.003	0.004

Two-stage smoothing



Query = "the algorithms for data mining"

d1:	0.04	0.001	0.02	0.002	0.003
d2:	0.02	0.001	0.01	0.003	0.004

$p(\text{"algorithms"}|d1) = p(\text{"algorithm"}|d2)$

$p(\text{"data"}|d1) < p(\text{"data"}|d2)$

$p(\text{"mining"}|d1) < p(\text{"mining"}|d2)$

But $p(q|d1) > p(q|d2)$!

Two-stage smoothing



Query = "the algorithms for data mining"

d1:	0.04	0.001	0.02	0.002	0.003
d2:	0.02	0.001	0.01	0.003	0.004

$p(\text{"algorithms"}|d1) = p(\text{"algorithm"}|d2)$

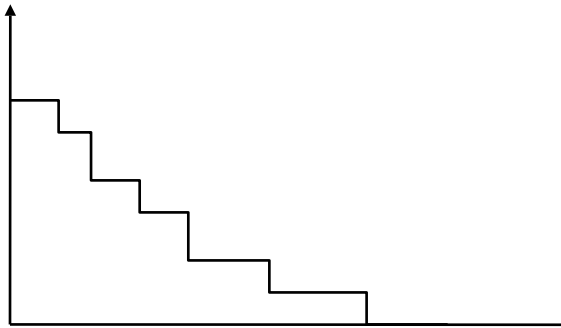
$p(\text{"data"}|d1) < p(\text{"data"}|d2)$

$p(\text{"mining"}|d1) < p(\text{"mining"}|d2)$

But $p(q|d1) > p(q|d2)!$

We should make $p(\text{"the"})$ and $p(\text{"for"})$ **less different** for all docs.

Two-stage smoothing

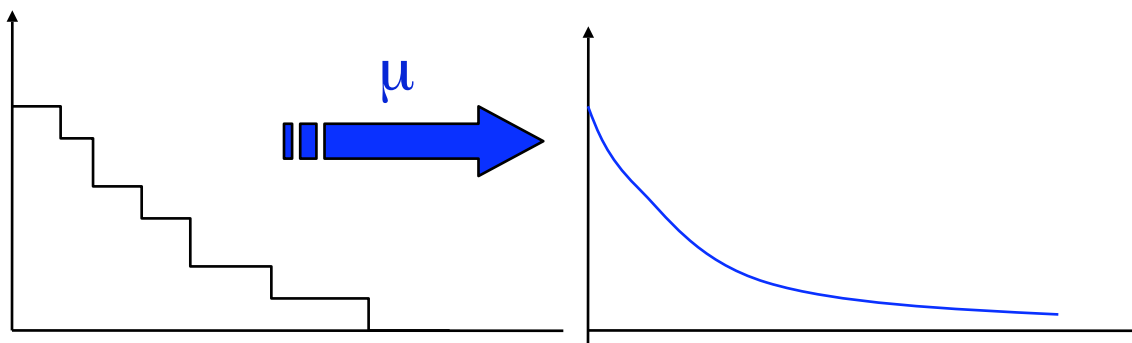


$$P(w|d) = \frac{c(w,d)}{|d|}$$

Two-stage smoothing

Stage-1

- Explain unseen words
- Dirichlet prior(Bayesian)



$$P(w|d) = \frac{c(w,d) + \mu p(w|C)}{|d| + \mu}$$

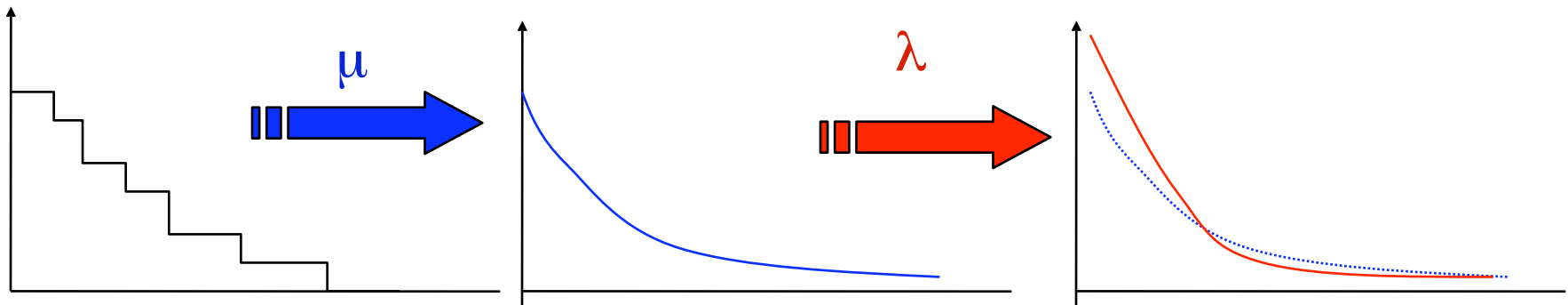
Two-stage smoothing

Stage-1

- Explain unseen words
- Dirichlet prior(Bayesian)

Stage-2

- Explain noise in query
- 2-component mixture

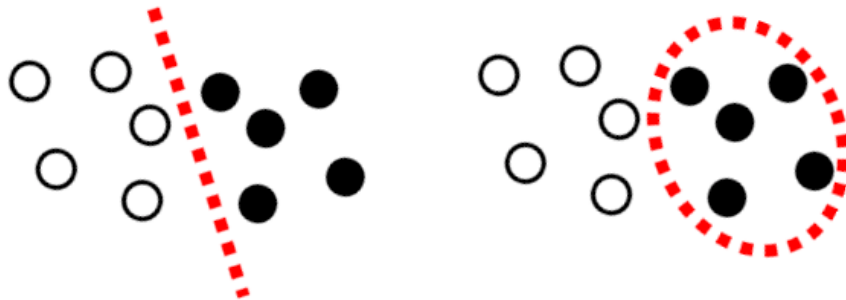


$$P(w|d) = (1-\lambda) \frac{c(w,d) + \mu p(w|C)}{|d| + \mu} + \lambda p(w|U)$$



LM are generative techniques

- How do we determine if a given model is a LM?
- LM is generative
 - at some level, a language model can be used to generate text
 - explicitly computes probability of observing a string of text
 - Ex: probability of observing a query string from a document model
probability of observing an answer from a question model
 - model an entire population
- Discriminative approaches
 - model just the decision boundary
 - Ex: is this document relevant?
does it belong to class X or Y



- have a lot of advantages,
- but these are not generative approaches



LM: summary

- Goal: estimate a model M from a sample text S
- Use maximum-likelihood estimator
 - count the number of times each word occurs in S , divide by length
- Smoothing to avoid zero frequencies
 - discounting methods: add or subtract a constant, redistribute mass
 - better: interpolate with background probability of a word
 - smoothing has a role similar to IDF in classical models
- Smoothing parameters very important
 - Dirichlet works well for short queries (need to tune the parameter)
 - Jelinek-Mercer works well for longer queries (also needs tuning)
 - Lots of other ideas being worked on



Language models: pro & con

- Novel way of looking at the problem of text retrieval based on probabilistic language modeling
 - Conceptually simple and explanatory
 - Formal mathematical model
 - Natural use of collection statistics, not heuristics (almost...)
- LMs provide effective retrieval and can be improved to the extent that the following conditions can be met
 - Our language models are accurate representations of the data.
 - Users have some sense of term distribution.



Comparison With Vector Space

- There's some relation to traditional tf.idf models:
 - (unscaled) term frequency is directly in model
 - the probabilities do length normalization of term frequencies
 - the effect of doing a mixture with overall collection frequencies is a little like idf: terms rare in the general collection but common in some documents will have a greater influence on the ranking



Comparison With Vector Space

- Similar in some ways
 - Term weights based on frequency
 - Terms often used as if they were independent
 - Inverse document/collection frequency used
 - Some form of length normalization useful
- Different in others
 - Based on probability rather than similarity
 - Intuitions are probabilistic rather than geometric
 - Details of use of document length and term, document, and collection frequency differ