

Introduction to Information Retrieval

4. seminar

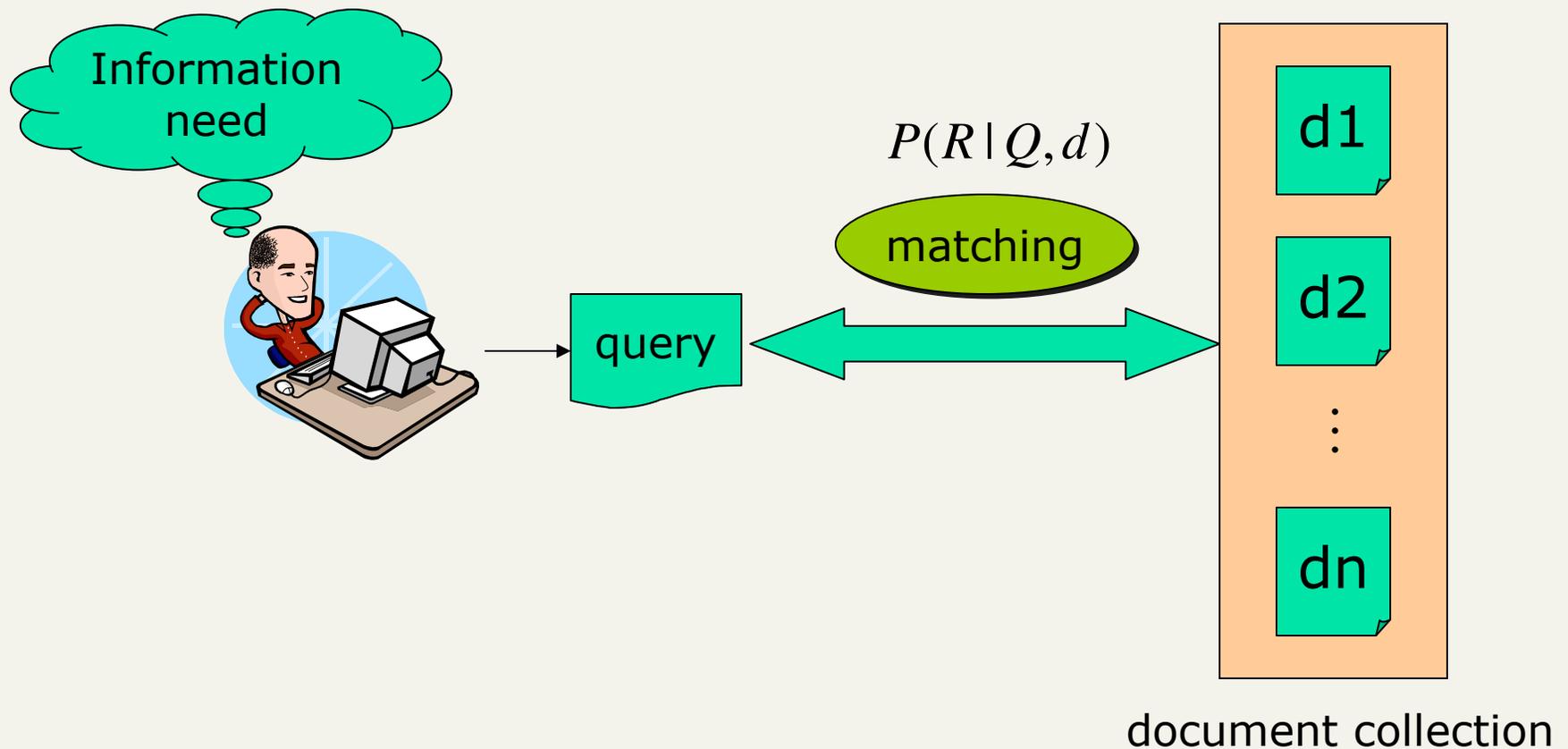
Language models for information retrieval

[Borrowed slides from the Stanford CS276 class on Text Retrieval and Mining, which borrowed slides from Viktor Lavrenko and Chengxiang Zhai]

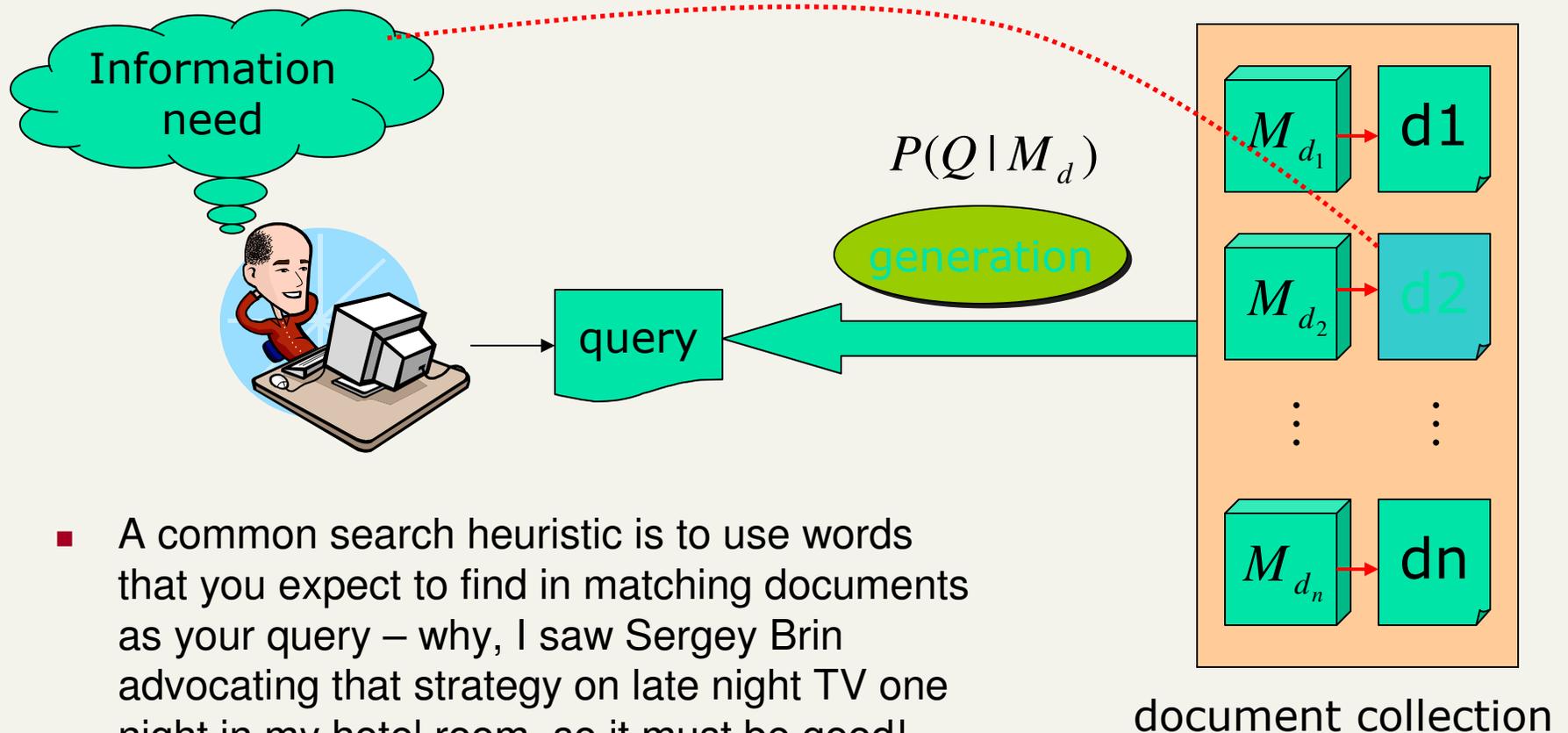
University of Pannonia

Tamás Kiezer, Miklós Erdélyi

Standard Probabilistic IR



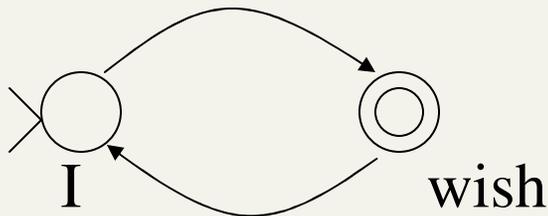
IR based on Language Model (LM)



- A common search heuristic is to use words that you expect to find in matching documents as your query – why, I saw Sergey Brin advocating that strategy on late night TV one night in my hotel room, so it must be good!
- The LM approach directly exploits that idea!

Formal Language (Model)

- Traditional generative model: generates strings
 - Finite state machines or regular grammars, etc.
- Example:



I wish

I wish I wish

I wish I wish I wish

I wish I wish I wish I wish

...

*wish I wish

Stochastic Language Models

- Models *probability* of generating strings in the language (commonly all strings over alphabet Σ)

Model M

0.2	the					
		the	man	likes	the	woman
0.1	a	—	—	—	—	—
0.01	man	0.2	0.01	0.02	0.2	0.01
0.01	woman					
0.03	said					
0.02	likes					
...						

multiply

$$P(s | M) = 0.00000008$$

Stochastic Language Models

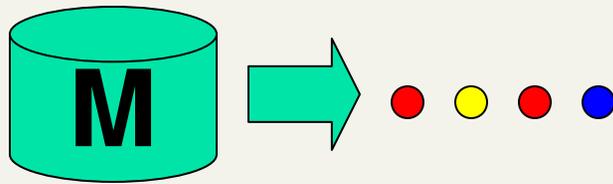
- Model *probability* of generating any string

Model M1		Model M2						
0.2	the	0.2	the	the	class	pleaseth	yon	maiden
0.01	class	0.0001	class	_____	_____	_____	_____	_____
0.0001	sayst	0.03	sayst	0.2	0.01	0.0001	0.0001	0.0005
0.0001	pleaseth	0.02	pleaseth	0.2	0.0001	0.02	0.1	0.01
0.0001	yon	0.1	yon					
0.0005	maiden	0.01	maiden					
0.01	woman	0.0001	woman					

$P(s|M2) > P(s|M1)$

Stochastic Language Models

- A statistical model for generating text
 - Probability distribution over strings in a given language



$$\begin{aligned} P(\text{red yellow red blue} \mid M) &= P(\text{red} \mid M) \\ &\quad P(\text{yellow} \mid M, \text{red}) \\ &\quad P(\text{red} \mid M, \text{red yellow}) \\ &\quad P(\text{blue} \mid M, \text{red yellow red}) \end{aligned}$$

Unigram and higher-order models

$$P(\text{red yellow red blue})$$

$$= P(\text{red}) P(\text{yellow} | \text{red}) P(\text{red} | \text{red yellow}) P(\text{blue} | \text{red yellow red})$$

- Unigram Language Models

$$P(\text{red}) P(\text{yellow}) P(\text{red}) P(\text{blue})$$



- Bigram (generally, n -gram) Language Models

$$P(\text{red}) P(\text{yellow} | \text{red}) P(\text{red} | \text{yellow}) P(\text{blue} | \text{red})$$

- Other Language Models

- Grammar-based models (PCFGs), etc.
 - Probably not the first thing to try in IR

Using Language Models in IR

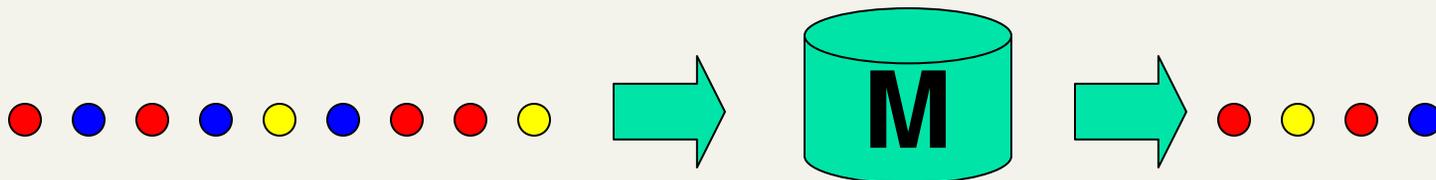
- Treat each document as the basis for a model (e.g., unigram sufficient statistics)
- Rank document d based on $P(d | q)$
- $P(d | q) = P(q | d) \times P(d) / P(q)$
 - $P(q)$ is the same for all documents, so ignore
 - $P(d)$ [the prior] is often treated as the same for all d
 - But we could use criteria like authority, length, genre
 - $P(q | d)$ is the probability of q given d 's model
- Very general formal approach

The fundamental problem of LMs

- Usually we don't know the model **M**
 - But have a sample of text representative of that model

$$P(\text{red yellow red blue} \mid M(\text{red blue red blue yellow blue red red yellow}))$$

- Estimate a language model from a sample
- Then compute the observation probability



Language Models for IR

- Language Modeling Approaches
 - Attempt to model query generation process
 - Documents are ranked by the probability that a query would be observed as a random sample from the respective document model
- Multinomial approach

$$P(Q|M_D) = \prod_w P(w|M_D)^{q_w}$$

Retrieval based on probabilistic LM

- Treat the generation of queries as a random process.
- Approach
 - Infer a language model for each document.
 - Estimate the probability of generating the query according to each of these models.
 - Rank the documents according to these probabilities.
 - Usually a unigram estimate of words is used.

Retrieval based on probabilistic LM

- Intuition
 - Users ...
 - Have a reasonable idea of terms that are likely to occur in documents of interest.
 - They will choose query terms that distinguish these documents from others in the collection.
- Collection statistics ...
 - Are integral parts of the language model.
 - Are not used heuristically as in many other approaches.
 - In theory. In practice, there's usually some wiggle room for empirically set parameters

Query generation probability (1)

- Ranking formula

$$p(Q, d) = p(d) p(Q | d)$$

$$\approx p(d) p(Q | M_d)$$

- The probability of producing the query given the language model of document d using MLE is:

$$\hat{p}(Q | M_d) = \prod_{t \in Q} \hat{p}_{ml}(t | M_d)$$

$$= \prod_{t \in Q} \frac{tf_{(t,d)}}{dl_d}$$

Unigram assumption:
Given a particular language model,
the query terms occur independently

M_d : language model of document d

$tf_{(t,d)}$: raw tf of term t in document d

dl_d : total number of tokens in document d

Insufficient data

- Zero probability $p(t | M_d) = 0$
 - May not wish to assign a probability of zero to a document that is missing one or more of the query terms [gives conjunction semantics]

- General approach

- A non-occurring term is possible, but no more likely than would be expected by chance in the collection.

- If $tf_{(t,d)} = 0$ $p(t | M_d) = \frac{cf_t}{cs}$

cf_t : raw count of term t in the collection

cs : raw collection size (total number of tokens in the collection)

Insufficient data

- Zero probabilities spell disaster
 - We need to smooth probabilities
 - Discount nonzero probabilities
 - Give some probability mass to unseen things
- There's a wide space of approaches to smoothing probability distributions to deal with this problem, such as adding 1, $\frac{1}{2}$ or ϵ to counts, Dirichlet priors, discounting, and interpolation
- A simple idea that works well in practice is to use a mixture between the document multinomial and the collection multinomial distribution

Mixture model

- $P(w|d) = \lambda P_{\text{mle}}(w|M_d) + (1 - \lambda)P_{\text{mle}}(w|M_c)$
- Mixes the probability from the document with the general collection frequency of the word.
- Correctly setting λ is very important
- A high value of lambda makes the search “conjunctive-like” – suitable for short queries
- A low value is more suitable for long queries
- Can tune λ to optimize performance
 - Perhaps make it dependent on document size

Basic mixture model summary

- General formulation of the LM for IR

$$p(Q, d) = p(d) \prod_{t \in Q} ((1 - \lambda) p(t) + \lambda p(t | M_d))$$

general language model

individual-document model

- The user has a document in mind, and generates the query from this document.
- The equation represents the probability that the document that the user had in mind was in fact this one.

Example

- Document collection (2 documents)
 - d_1 : Xerox reports a profit but revenue is down
 - d_2 : Lucent narrows quarter loss but revenue decreases further
- Model: MLE unigram from documents; $\lambda = 1/2$
- Query: *revenue down*
 - $P(Q|d_1) = [(1/8 + 2/16)/2] \times [(1/8 + 1/16)/2]$
 $= 1/8 \times 3/32 = 3/256$
 - $P(Q|d_2) = [(1/8 + 2/16)/2] \times [(0 + 1/16)/2]$
 $= 1/8 \times 1/32 = 1/256$
- Ranking: $d_1 > d_2$

Vector Space (tf-idf) vs. LM

Rec.	precision			significant?
	tf-idf	LM	%chg	
0.0	0.7439	0.7590	+2.0	
0.1	0.4521	0.4910	+8.6	
0.2	0.3514	0.4045	+15.1	*
0.4	0.2093	0.2572	+22.9	*
0.6	0.1024	0.1405	+37.1	*
0.8	0.0160	0.0432	+169.6	*
1.0	0.0028	0.0050	+76.9	
11-point average	0.1868	0.2233	+19.6	*

- The language modeling approach always does better in these experiments.
 - Note that where the approach shows significant gains is at higher levels of recall.

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)

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.*
 - *Or we get more sophisticated with translation model

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

Resources

- J.M. Ponte and W.B. Croft. 1998. A language modelling approach to information retrieval. In *SIGIR 21*.
- D. Hiemstra. 1998. A linguistically motivated probabilistic model of information retrieval. *ECDL 2*, pp. 569–584.
- A. Berger and J. Lafferty. 1999. Information retrieval as statistical translation. *SIGIR 22*, pp. 222–229.
- Workshop on Language Modeling and Information Retrieval, CMU 2001.
<http://la.lti.cs.cmu.edu/callan/Workshops/lmir01/> .
- The Lemur Toolkit for Language Modeling and Information Retrieval.
<http://www-2.cs.cmu.edu/~lemur/> . CMU/Umass LM and IR system in C(++), currently actively developed.