

Introduction to
Assignment Project Exam Help
Informa |
<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
Lecture 9: Probabilistic M guage Model

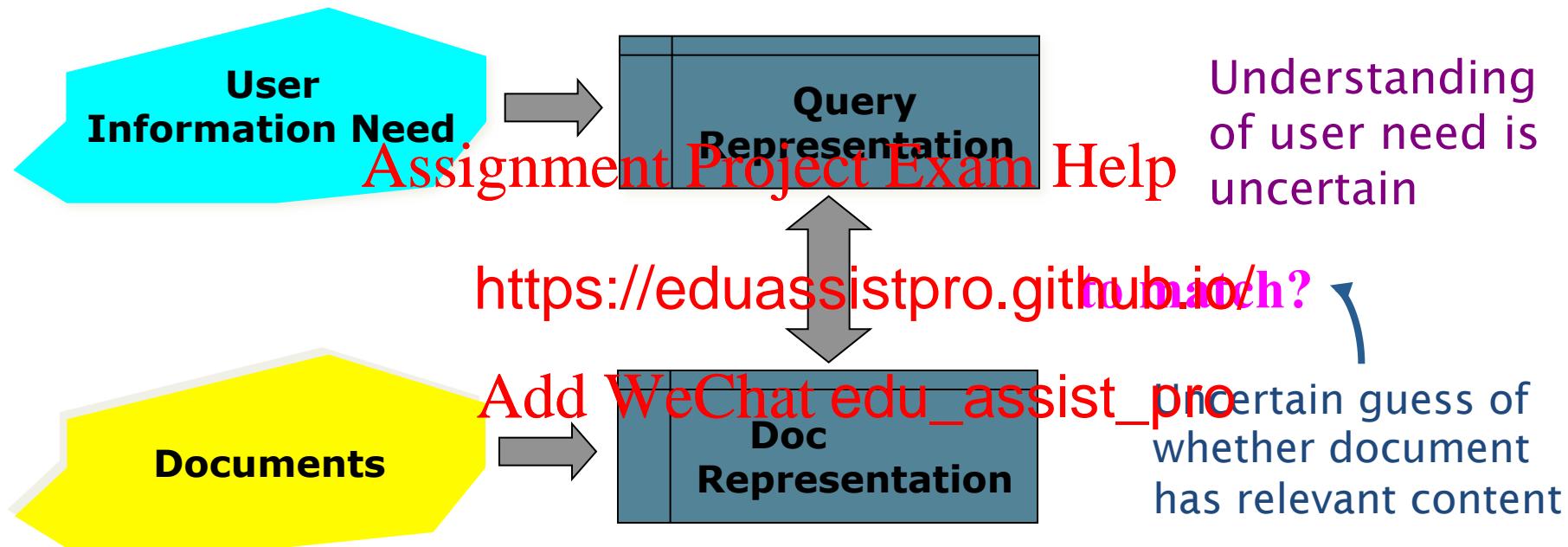
Recap of the last lecture

- Improving search results
 - Especially for high recall. E.g., searching for *aircraft* so it matches with *plane*; *thermodynamic* with *heat*
- Options for i
 - Global meth <https://eduassistpro.github.io/>
 - Query expansion
 - Thesauri
 - Automatic thesaurus generation
 - Global indirect relevance feedback
 - Local methods
 - Relevance feedback
 - Pseudo relevance feedback

Probabilistic relevance feedback

- Rather than reweighting in a vector space...
- If user has told us some relevant and some irrelevant documents, then we can proceed to build a probabilistic c
<https://eduassistpro.github.io/> ve Bayes model:
 - $P(t_k|R) = |\mathbf{D}_{rk}| / |\mathbf{D}_r|$
 - $P(t_k|NR) = |\mathbf{D}_{nrk}| / |\mathbf{D}_{nr}|$
 - t_k is a term; \mathbf{D}_r is the set of known relevant documents; \mathbf{D}_{rk} is the subset that contain t_k ; \mathbf{D}_{nr} is the set of known irrelevant documents; \mathbf{D}_{nrk} is the subset that contain t_k .

Why probabilities in IR?



In traditional IR systems, matching between each document and query is attempted in a semantically imprecise space of index terms.

Probabilities provide a principled foundation for uncertain reasoning

Can we use probabilities to quantify our uncertainties?

Probabilistic IR topics

- Classical probabilistic retrieval model
 - Probability ranking principle, etc.
- (Naïve) Baye
- Bayesian ne
- Language model
- *Probabilistic methods are one of the oldest but also one of the currently hottest topics in IR.*
 - *Traditionally: neat ideas, but they've never won on performance. It may be different now.*

The document ranking problem

- We have a collection of documents
- User issues a query
- A list of documents is returned
- Ranking method <https://eduassistpro.github.io/>
- In what order do we present the results to the user?
Add WeChat edu_assist_pro
- We want the “best” document to be first, second best, second, etc....
- Idea: Rank by probability of relevance of the document w.r.t. information need
 - $P(\text{relevant} | \text{document}_i, \text{query})$

Recall a few probability basics

- For events a and b :
- Bayes' Rule

$$p(a, b) = p(a \text{ Assignment Project Exam Help}) = p(b | a)p(a)$$

$$p(\bar{a} | b)p(b) \quad \text{https://eduassistpro.github.io/}$$

$$p(a | b) = \frac{p(b | a)p(a)}{p(b)} = \frac{\cancel{p(b | a)p(a)}}{\sum_{x=a, \bar{a}} p(b | x)p(x)}$$

Posterior

Prior

- Odds:

$$O(a) = \frac{p(a)}{p(\bar{a})} = \frac{p(a)}{1 - p(a)}$$

The Probability Ranking Principle

"If a reference retrieval system's response to each request is a ranking of the documents in the collection in order of decreasing probability of relevance to the user who submitted the request, where the probability is based on what is available to the system for this purpose, then the system will be the best that is obtainable given those data."

- [1960s/1970s] S. Robertson, W.S. Cooper, M.E. Maron; van Rijsbergen (1979:113); Manning & Schütze (1999:538)

Probability Ranking Principle

Let x be a document in the collection.

Let R represent relevance of a document w.r.t. given (fixed) query and let NR represent non-relevance.

$R=\{0,1\}$ vs. NR/R

Need to find $p(\text{document } x \text{ is relevant})$

$$p(R | x) = \frac{p(x | R)p(R)}{p(x)}$$

Add WeChat edu_assist_pro prior probability of retrieving a (non) relevant document

$$p(NR | x) = \frac{p(x | NR)p(NR)}{p(x)}$$

$$p(R | x) + p(NR | x) = 1$$

$p(x|R)$, $p(x|NR)$ - probability that if a relevant (non-relevant) document is retrieved, it is x .

Probability Ranking Principle (PRP)

- Simple case: no selection costs or other utility concerns that would differentially weight errors
Assignment Project Exam Help
- ***Bayes' Optim***
 - x is relevant <https://eduassistpro.github.io/>
- PRP in action: Rank all documents $p(R|x)$
Add WeChat edu_assist_pro
- Theorem:
 - Using the PRP is optimal, in that it minimizes the loss (Bayes risk) under 1/0 loss
 - Provable if all probabilities correct, etc. [e.g., Ripley 1996]

Probability Ranking Principle

- More complex case: retrieval costs.
 - Let d be a document
 - C - cost of retrieving document
 - C' - cost of retrieving next document
- Probability of ranking d first given query q :
$$C \cdot p(R|d) + C' \cdot (1 - p(R|d)) \leq C \cdot p(R|d') + C' \cdot (1 - p(R|d'))$$
for all d' not yet retrieved, then d is the next document to be retrieved
- We won't further consider loss/utility from now on

Probability Ranking Principle

- How do we compute all those probabilities?
 - Do not know exact probabilities, have to use estimates
 - **Binary Independence Retrieval (BIR)** which we discuss later
- Questionable <https://eduassistpro.github.io/>
 - “Relevance” of each document independent of relevance of other documents
 - Really, it’s bad to keep on returning **duplicates**
 - Boolean model of relevance
 - That one has a single step information need
 - Seeing a range of results might let user refine query

$$MMR \stackrel{\text{def}}{=} \operatorname{Arg} \max_{D_i \in R \setminus S} \left[\lambda(Sim_1(D_i, Q) - (1-\lambda) \max_{D_j \in S} Sim_2(D_i, D_j)) \right]$$

Probabilistic Retrieval Strategy

- Estimate how terms contribute to relevance
 - How do things like tf, df, and length influence your judgments about document relevance?
 - One answer <https://eduassistpro.github.io/>
- Combine to find document probability
 - Add WeChat `edu_assist_pro`
- Order documents by decreasing probability

Probabilistic Ranking

Basic concept:

"For a given query, *if* we know some documents that are relevant, terms that occur in those documents should be given greater w <https://eduassistpro.github.io/> other relevant documents.

Add WeChat [edu_assist_pro](#)
By making assumptions about the distribution of terms and applying Bayes Theorem, it is possible to derive weights theoretically."

Van Rijsbergen

Binary Independence Model

- Traditionally used in conjunction with PRP
- “**Binary**” = **Boolean**: documents are represented as binary incidence vectors of terms (cf. lecture 1):
■ $\vec{x} = (x_1, \dots,$
■ $x_i = 1$ iff te <https://eduassistpro.github.io/>
- “**Independence**”: terms occur independently
- Different documents can be modelled
- Bernoulli Naive Bayes model (cf. text categorization!)

Binary Independence Model

- Queries: binary term incidence vectors
- Given query q ,
 - for each document d need to compute $p(R | q, d)$.
 - replace with $\frac{p(R | q, \vec{x})}{p(NR | q, \vec{x})}$ where \vec{x} is binary term incidence vector
- Will use odds and Bayes' Rule

$$O(R | q, \vec{x}) = \frac{p(R | q, \vec{x})}{p(NR | q, \vec{x})} = \frac{\frac{p(R | q) p(\vec{x} | R, q)}{p(\vec{x} | q)}}{\frac{p(NR | q) p(\vec{x} | NR, q)}{p(\vec{x} | q)}}$$

Binary Independence Model

$$O(R | q, \vec{x}) = \frac{p(R | q, \vec{x})}{p(NR | q, \vec{x})} = \frac{p(R | q)}{p(NR | q)} \cdot \frac{p(\vec{x} | R, q)}{p(\vec{x} | NR, q)}$$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Needs estimation

- Using ~~Independent~~ Add WeChat ~~Assum~~ edu_assist_pro

$$\frac{p(\vec{x} | R, q)}{p(\vec{x} | NR, q)} = \prod_{i=1}^n \frac{p(x_i | R, q)}{p(x_i | NR, q)}$$

$$\text{So : } O(R | q, d) = O(R | q) \cdot \prod_{i=1}^n \frac{p(x_i | R, q)}{p(x_i | NR, q)}$$

Binary Independence Model

$$O(R | q, d) = O(R | q) \cdot \prod_{i=1}^n \frac{p(x_i | R, q)}{p(x_i | NR, q)}$$

Assignment Project Exam Help

- Since x_i is e

$$O(R | q, d) = O \frac{\frac{p(x_i = 1 | R, q)}{p(x_i = 0 | R, q)}}{\frac{p(x_i = 1 | NR, q)}{p(x_i = 0 | NR, q)}}$$

https://eduassistpro.github.io/
Add WeChat edu_assist_pro

- Let $p_i = p(x_i = 1 | R, q)$; $r_i = p(x_i = 1 | NR, q)$;

- **Assume**, for all terms ***not occurring*** in the query ($q_i = 0$) $p_i = r_i$

Then...

This can be
changed (e.g., in
relevance feedback)

Binary Independence Model

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{\substack{x_i = q_i = 1 \\ r_i}} \frac{p_i}{r_i} \cdot \prod_{\substack{x_i = 0 \\ q_i = 1 \\ r_i}} \frac{1 - p_i}{1 - r_i}$$

Assignment Project Exam Help

All match

Non-matching
query terms

<https://eduassistpro.github.io/>

$$= O(R | q) \cdot \prod_{x_i = q_i = 1} \frac{p_i}{r_i} (1 - p_i) \prod_{q_i = 1} \frac{1 - p_i}{1 - r_i}$$

All matching terms

All query terms

Binary Independence Model

$$O(R | q, \vec{x}) = O(R | q) \cdot \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} \cdot \prod_{q_i=1} \frac{1-p_i}{1-r_i}$$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`
to be estimated
for rankings

- Retrieval Status Value:

$$RSV = \log \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} = \sum_{x_i=q_i=1} \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

$$= \sum_{x_i=q_i=1} (\log(odds(p_i)) - \log(odds(r_i))) = \sum_{x_i=q_i=1} (\text{logit}(p_i) - \text{logit}(r_i))$$

Binary Independence Model

- All boils down to computing RSV.

$$RSV = \log \prod_{x_i = q_i = 1} p_i (1 - r_i)$$

$$RSV = \sum_i c_i; \quad \frac{\log \frac{p_i (1 - r_i)}{1 - p_i}}{\text{Assignment Project Exam Help}}$$

<https://eduassistpro.github.io/>
Add WeChat r_i $c_i = \text{it}(p_i) - \log \text{it}(r_i)$
`edu_assist_pro`

So, how do we compute c_i 's from our data ?

Binary Independence Model

- Estimating RSV coefficients.
- For each term i look at this table of document counts:

	Assignment	Project	Exam	Help
Documents	Relevant	Non-Relevant	Total	
$X_i = 1$	https://eduassistpro.github.io/			
$X_i = 0$	n N-n			
Total	Add WeChat edu_assist_pro			

• Estimates: $p_i \approx \frac{s}{S}$ $r_i \approx \frac{(n-s)}{(N-S)}$

$$c_i \approx K(N, n, S, s) = \log \frac{s/(S-s)}{(n-s)/(N-n-S+s)}$$

However, these estimates could be 0.

Add $\frac{1}{2}$ Smoothing

- Add $\frac{1}{2}$ to each of the center four cells.

	Assignment	Project	Exam	Help
Documents	Relevant	Non-Relevant	Total	
$X_i = 1$			$n+1$	
$X_i = 0$			$N-n+1$	
Total	Add WeChat $\sum_{i=1}^{n+1}$	$\sum_{i=1}^{N-n+1}$	$\sum_{i=1}^{N+1}$	$\sum_{i=1}^{N+1}$

$$c_i \approx K(N, n, S, s) = \log \frac{(s + 1/2)/(S - s + 1/2)}{(n - s + 1/2)/(N - n - S + s + 1/2)}$$

Example /1

- Query = $\{x_1, x_2\}$
- $O(R=1 | D_3, q)$

Doc	Judgment	x_1	x_2	x_3
D_1	R	1	1	1
D_2	R	0	0	1
D_3	R	1	0	0
D_4	NR	1	0	1
D_5	NR	0	1	1

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Example /2

- Estimate p_i and r_i

Doc	Judgment	x_1	x_2	x_3
D_1	R	1	1	1
D_2	R	0	0	1
D_3	R	1	0	0
D_4	NR	1	0	1
D_5	NR	0	1	1

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

in fact, tf-idf can be deemed as the cross-entropy

Estimation – key challenge

- If non-relevant documents are approximated by the whole collection, then r_i (prob. of occurrence in non-relevant documents for query) is n/N and
 - $\log(1 - r_i)/r$ IDF!
- p_i (probabilistic weight of relevant documents) can be estimated in various ways:
 - from relevant documents if know some
 - Relevance weighting can be used in feedback loop
 - constant (Croft and Harper combination match – 0.5) – then just get idf weighting of terms
 - proportional to prob. of occurrence in collection
 - more accurately, to log of this (Greiff, SIGIR 1998)

Iteratively estimating p_i

1. Assume that p_i constant over all x_i in query
 - $p_i = 0.5$ (even odds) for any given doc
2. Determine ~~Assignment Project Exam Help~~ set:
 - V is fixed set of documents on this model (not <https://eduassistpro.github.io/>)
3. We need to improve our ~~Add WeChat~~ $\text{edu_assist}_{\text{pro}}$ and r_i , so
 - Use distribution of x_i in docs in V . Let V_i be set of documents containing x_i
 - $p_i = |V_i| / |V|$
 - Assume if not retrieved then not relevant
 - $r_i = (n_i - |V_i|) / (N - |V|)$
4. Go to 2. until converges then return ranking

Probabilistic Relevance Feedback

1. Guess a preliminary probabilistic description of R and use it to retrieve a first set of documents V , as above. [Assignment](#) [Project](#) [Exam](#) [Help](#)
2. Interact with <https://eduassistpro.github.io/> description: learn some d and NR
[Add WeChat](#) [edu_assist_pro](#)
3. Reestimate p_i and r_i on the hese
 - Or can combine new information with original guess (use Bayesian prior):
$$p_i^{(2)} = \frac{|V_i| + \kappa p_i^{(1)}}{|V| + \kappa}$$

κ is prior weight
4. Repeat, thus generating a succession of approximations to R .

PRP and BIR

- Getting reasonable approximations of probabilities is possible.
Assignment Project Exam Help
- Requires r :
 - *term ind* <https://eduassistpro.github.io/>
 - *terms not in query don't* [Add WeChat](#) [edu_assist_pro](#)
 - *boolean representation* [ts/queries/relevance](#)
 - *document relevance values are independent*
- Some of these assumptions can be removed
- Problem: either require partial relevance information or only can derive somewhat inferior term weights

Okapi BM25

- Heuristically extend the BIR to include information of term frequencies, document length, etc.

Assignment Project Exam Help

contribution of tf

$$RSVd = \sum_{t \in q} \left(\log \frac{1}{df_t} \right) \frac{(k_3 + 1)tf_{t,q}}{k_1 \left((1 - b) \frac{\text{Add WeChat edu_assist_pro}}{L_{ave}} \right)} \cdot \frac{k_3 + tf_{t,d}}{k_3 + tf_{t,q}}$$

idf

Normalized term
freq (doc)

Normalized term
freq (query)

- Typically, $k_1, k_3 \in [1.2, 2.0], b = 0.75$

Good and Bad News

- Standard Vector Space Model
 - Empirical for the most part; success measured by results
 - Few properties provable
- Probabilistic Mo
 - Based on a firm <https://eduassistpro.github.io/>
 - Theoretically justified optimal rank
- Disadvantages
 - Making the initial guess to get V
 - Binary word-in-doc weights (not using term frequencies)
 - Independence of terms (can be alleviated)
 - Amount of computation
 - Has never worked convincingly better in practice

Resources

- S. E. Robertson and K. Spärck Jones. 1976. Relevance Weighting of Search Terms. *Journal of the American Society for Information Sciences* 27(3): 129–146.
- C. J. van Rijsbergen. 1975. *Information Retrieval*. 2nd ed. London: Butterworths, h] http://www.dcs.gla.ac.uk/~hjw/Assignment_Project_Exam_Help
www.dcs.gla.ac.uk/~hjw/Assignment_Project_Exam_Help
<https://eduassistpro.github.io/>
- N. Fuhr. 1992. Probabilistic Models in Information Retrieval. *The Computer Journal*, 35(3), 243–255. [Easiest]
- F. Crestani, M. Lalmas, C. J. van Rijsbergen, and I. Campbell. 1998. Is This Document Relevant? ... Probably: A Survey of Probabilistic Models in Information Retrieval. *ACM Computing Surveys* 30(4): 528–552.
<http://www.acm.org/pubs/citations/journals/surveys/1998-30-4/p528-crestani/>
[Adds very little material that isn't in van Rijsbergen or Fuhr]

Resources

- H.R. Turtle and W.B. Croft. 1990. Inference Networks for Document Retrieval.
Proc. ACM SIGIR: 1-24.
- E. Charniak. Bayesian Nets without Tears. *AI Magazine* 12(4):50-63 (1991). <http://www.aaai.org/Library/1991/aaai91-0001.pdf>
- D. Heckerman. 1995 <https://eduassistpro.github.io/> Bayesian Networks. Microsoft Technical Report MSR-TR-95-06
<http://www.research.microsoft.com/~hecker/> Add WeChat edu_assist_pro
- N. Fuhr. 2000. Probabilistic Datalog: Implementing Logical Information Retrieval for Advanced Applications. *Journal of the American Society for Information Science* 51(2): 95–110.
- R. K. Belew. 2001. *Finding Out About: A Cognitive Perspective on Search Engine Technology and the WWW*. Cambridge UP 2001.
- MIR 2.5.4, 2.8

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

LANGUAGE MODEL

Today

- The Language Model Approach to IR

- Basic query generation model

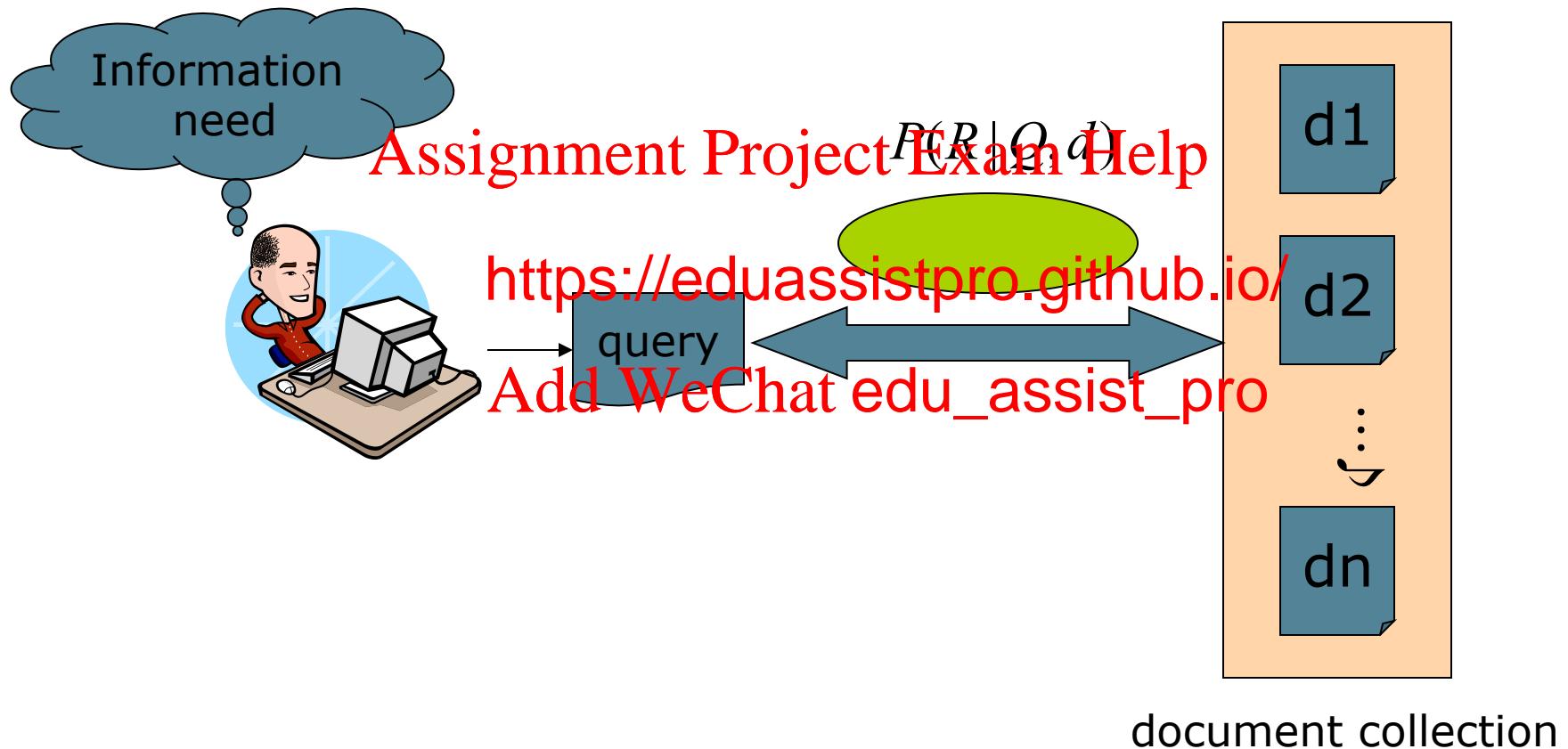
- Alternative models

Assignment Project Exam Help

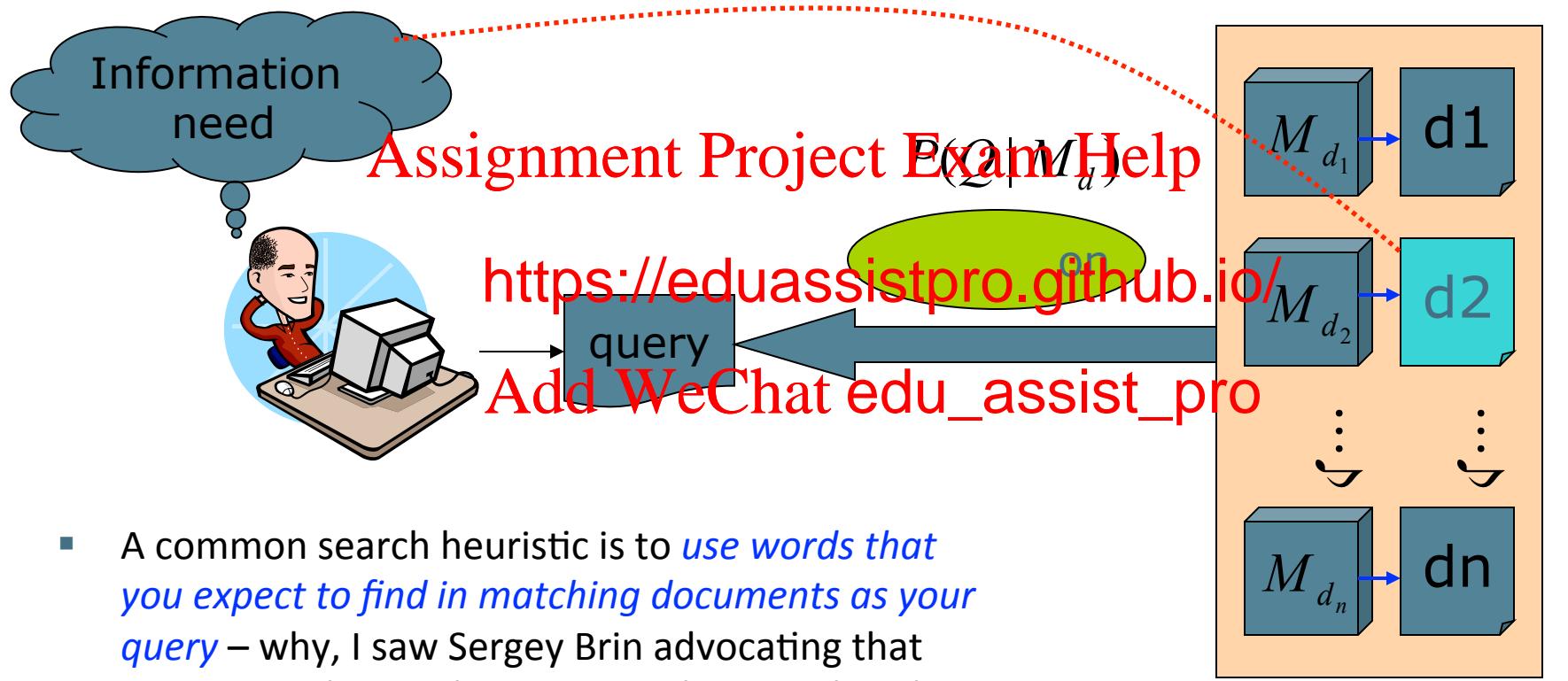
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Standard Probabilistic IR



IR based on Language Model (LM)



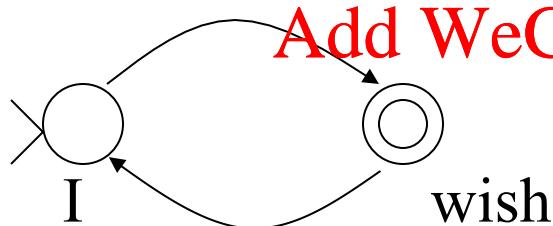
- A common search heuristic is to *use words that you expect to find in matching documents as your query* – why, I saw Sergey Brin advocating that strategy on late night TV one night in my hotel room, so it must be good!
- The LM approach directly exploits that idea!
- See later slides for a more formal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

document collection

Formal Language (Model)

- Traditional generative model: generates strings
 - Finite state machines or regular grammars, etc.
- Example: **Assignment Project Exam Help**

<https://eduassistpro.github.io/>



I wish

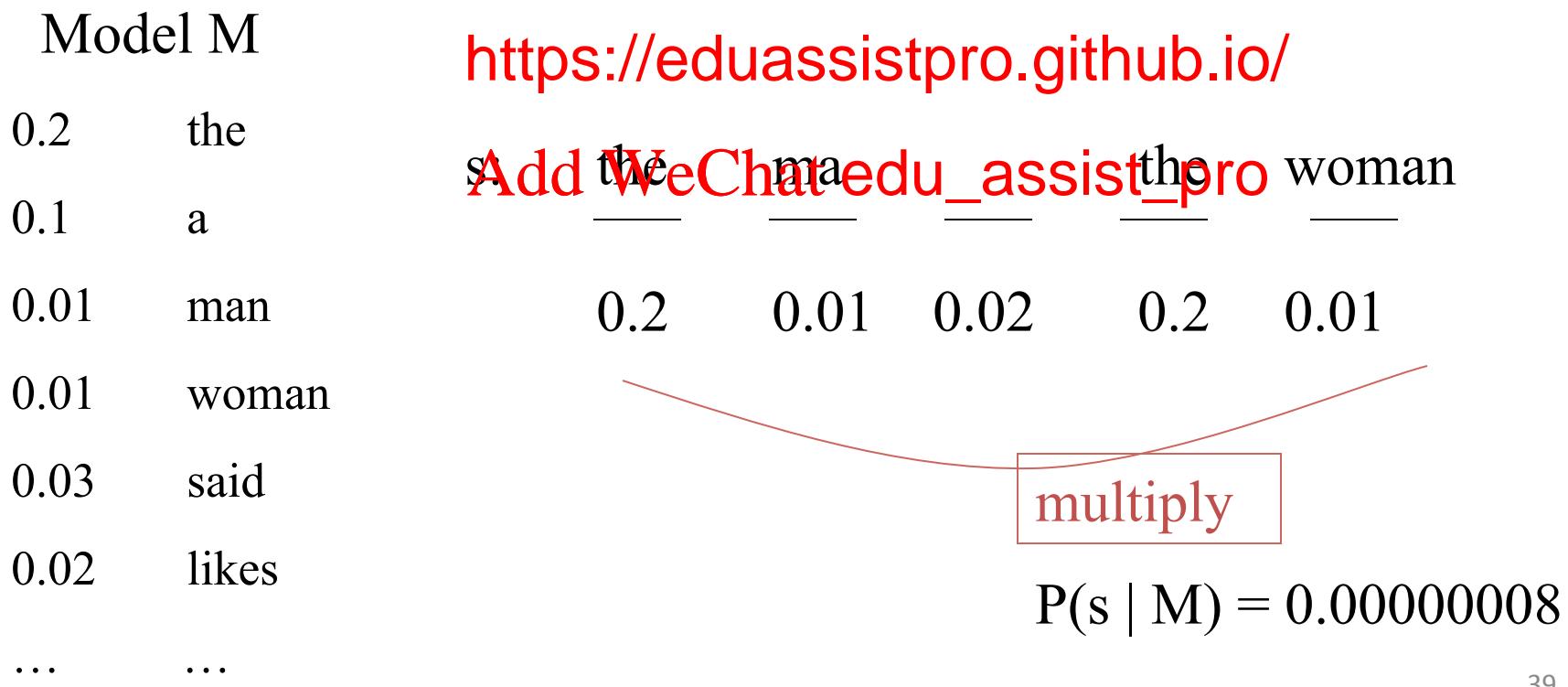
I wish I wish

I wish I wish I wish I wish

...

Stochastic Language Models

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



Stochastic Language Models

- Model *probability* of generating any string

Assignment Project Exam Help

Model M1

0.2	the
0.01	class
0.0001	sayst
0.0001	pleaseth
0.0001	yon
0.0005	maiden
0.01	woman

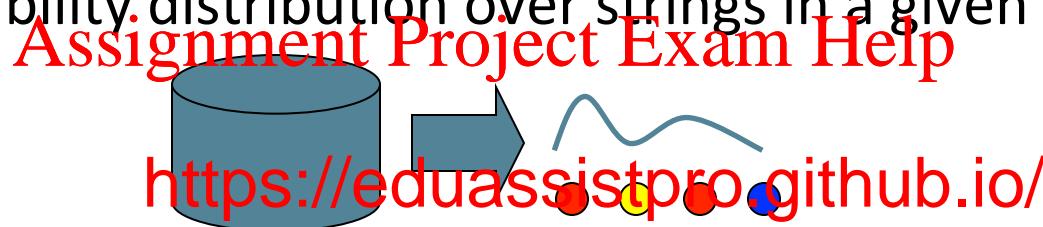
<https://eduassistpro.github.io/>

0.2	the	0.0001	class	Add WeChat	t	0.2	0.01	0.0001	0.0001	0.0005
0.03	sayst	0.02	pleaseth	edu_assist_pro	—	0.2	0.0001	0.02	0.1	0.01
0.1	yon	0.01	maiden		—	0.01	0.0001	0.02	0.1	0.01
0.0001	woman				—	0.0001	0.02	0.1	0.01	0.0005

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

Stochastic Language Models

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



$$\begin{aligned} P(\bullet \bullet \bullet | M) &= P(\bullet \\ &\quad P(\bullet | M, \bullet) \\ &\quad P(\bullet | M, \bullet \bullet) \\ &\quad P(\bullet | M, \bullet \bullet \bullet) \end{aligned}$$

Unigram and higher-order models

$$P(\bullet \bullet \bullet \bullet)$$

$$= P(\bullet) P(\bullet | \bullet) P(\bullet | \bullet \bullet) P(\bullet | \bullet \bullet \bullet)$$

Assignment Project Exam Help

- Unigram Language Models <https://eduassistpro.github.io/>

$$P(\bullet) P(\bullet) P(\bullet) \quad \bullet$$

Add WeChat edu_assist_pro

Easy.
Effective!

- Bigram (generally, n -gram) Language Models

$$P(\bullet) P(\bullet | \bullet) P(\bullet | \bullet) P(\bullet | \bullet)$$

- 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 $P(d | q)$
- $P(d | q) = P(q | \text{https://eduassistpro.github.io/}) P(d | \text{https://eduassistpro.github.io/})$
 - $P(q)$ is the same for all documents
 - $P(d)$ [the prior] is often treated as a constant 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

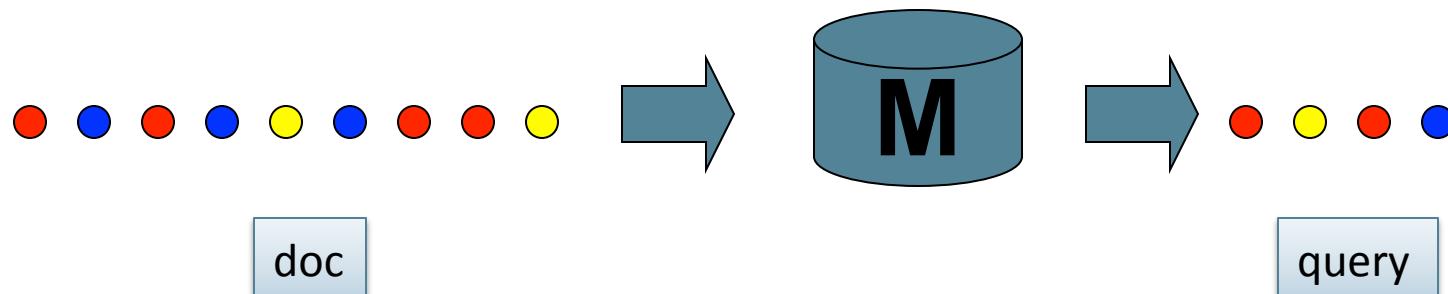
The fundamental problem of LMs

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

Assignment Project Exam Help

P (● ○ ● ● https://eduassistpro.github.io/)

- Estimate a language model
 - 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 ob le from the respective d
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- 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 the query according to each of the documents.
 - Rank the documents according to their probabilities.
 - Usually a unigram estimate or bigram estimate
 - Some work on bigrams, parallelizing van Rijsbergen's algorithm.

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 these documents from others
- 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

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

Assignment Project Exam Help

- The probability of document d using <https://eduassistpro.github.io/> the language model of

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

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

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 <https://eduassistpro.github.io/>
 - A non-occurring term is considered more likely than would be expected by chance.
 - 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 α
- There's a wide range of smoothing techniques to handle this problem, such as adding 1, $\frac{1}{2}$ or ϵ to the counts, Laplace priors, discounting, and interpolation
 - [See FSNLP ch. 6 or CS224N if you want more]
- A simple idea that works well in practice is to use a mixture between the document multinomial and the collection multinomial distribution

Mixture model

- Jelinek-Mercer method
 - $P(w|d) = \lambda P_{mle}(w|M_d) + (1 - \lambda)P_{mle}(w|M_c)$
- Mixes the probability from the document with the general collection <https://eduassistpro.github.io/>
- Correctly setting λ is very important
- A high value of lambda makes it more “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 (cf. Dirichlet prior or Witten-Bell smoothing)

Basic mixture model summary

- General formulation of the LM for IR

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

collection/background <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- 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.

Relationship to idf

Note here (i.e., [CMS09]) λ is multiplied to the background model.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

proportional to the tf, inversely proportional to the cf

$f_{qi,D} = 0 \rightarrow$ query word that does not occur in the doc

Add contributions from $i: f_{qi,D} > 0$

Becomes a constant | Q, C

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 under $\text{ts}; \lambda = \frac{1}{2}$
- 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: $d_1 > d_2$

Ponte and Croft Experiments

- Data
 - TREC topics 202-250 on TREC disks 2 and 3
 - Natural language queries consisting of one sentence each
 - TREC topics 51-100 on TREC disk 3 using the concept fields
 - Lists of good <https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

<num>
<dom> ational Economics
<title>Topic: Satellite Launch Contracts
<desc>Description:
... </desc>

<con>Concept(s):
1. Contract, agreement
2. Launch vehicle, rocket, payload, satellite
3. Launch services, ... </con>

Precision/recall results 202-250

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Precision/recall results 51-100♪

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Language models: pro & con

- Novel way of looking at the problem of text retrieval based on probabilistic language modeling
 - Conceptually simple and explanatory
 - Formal mat
 - Natural use <https://eduassistpro.github.io/>
- LMs provide be improved to the extent that the following conditions 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 frequencies is a little like idf in the general collection but common in so ts will have a greater influence on the ranking

<https://eduassistpro.github.io/>

all collection

Add WeChat edu_assist_pro

in the general

ts will have a

Comparison With Vector Space

- Similar in some ways
 - Term weights based on frequency
 - Terms often used as if they were independent
 - Inverse document used
 - Some form of full
- 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. 560–584.
- A. Berger and J. Lafferty. stical translation. *SIGIR 22*, pp. 222–229.
- D.R.H. Miller, T. Leek, an retrieval system. *SIGIR 22*, pp. 214–221.
- [Several relevant newer papers at *SIGIR 23–25*
- 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.