

Assignment Project Exam Help

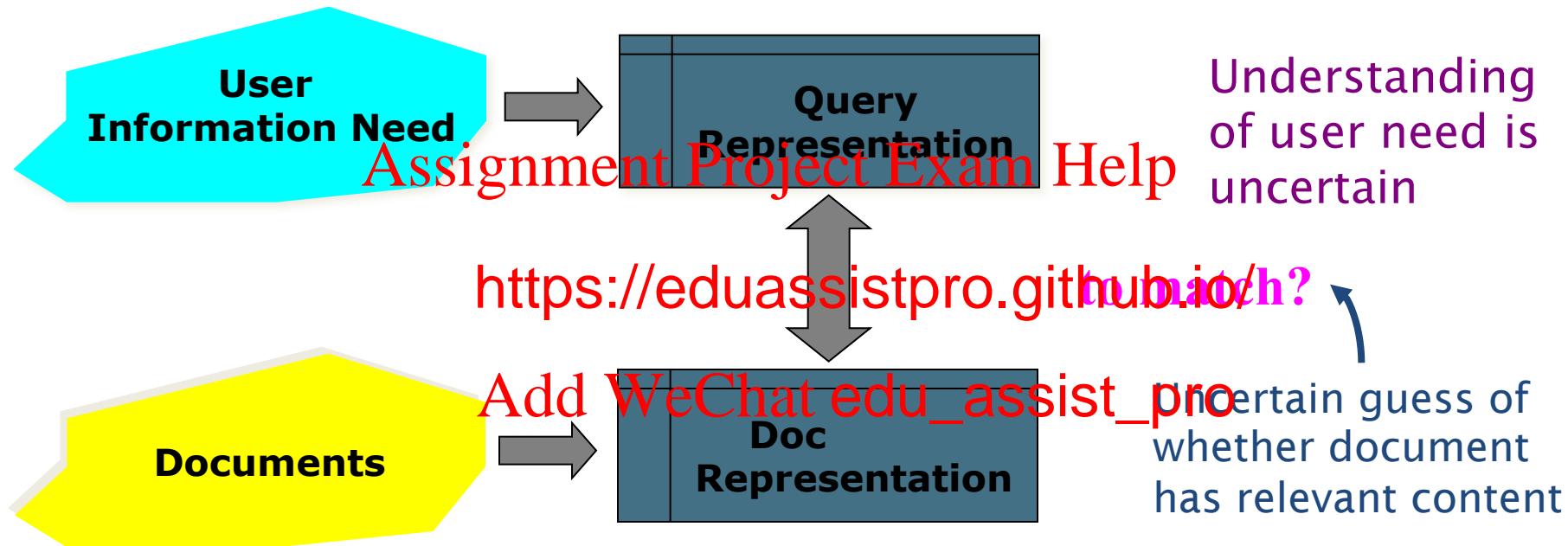
Add WeChat `edu_assist_pro`

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

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

Assignment Project Exam Help

Why probabilities?



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?

Assignment Project Exam Help

Probabilistic IR

Add WeChat to top
edu_assist_pro

- Classical probabilistic retrieval model
 - Probability ranking principle, etc.
- (Naïve) Bayes
- Bayesian neural networks <https://eduassistpro.github.io/>
- Language models
 - Add WeChat to top
edu_assist_pro
 - An important emphasis in recent work
- *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.*

Assignment Project Exam Help

The document ran Add WeChat edu_assist_pro blem

- 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 them to the user?
 - We want the “best” document, then second best, etc....
- Idea: Rank by probability of relevance of the document w.r.t. information need
 - $P(\text{relevant} | \text{document}_i, \text{query})$

Assignment Project Exam Help

Recall a few Add WeChat edu_assist_pro

- For events a and b :
- Bayes' Rule

$$p(a, b) = p(a \mid b) p(b) = p(b \mid a) p(a)$$

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

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

↑ Prior

↑ Posterior

- Odds:

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

Assignment Project Exam Help

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 of relevance is based on what is available to the system for this purpose, then that system to its user will be the best that is obtainable from those data."

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

Assignment Project Exam Help

Probability Ranking

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 $\frac{\text{document } x \text{ is relevant}}{\text{document } x \text{ is non-relevant}}$

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

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 .

Assignment Project Exam Help

Probability Ranking P^{Add WeChat edu_assist pro}(PRP)

- Simple case: no selection costs or other utility concerns that would differentially weight errors
- *Bayes' Optim*
 - x is relevant <https://eduassistpro.github.io/>
- PRP in action: Rank all docu $p(R|x)$
- 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]

Assignment Project Exam Help

Probability Ranking

- More complex case: retrieval costs.

- Let d be a document
- C - cost of retrieving document
- C' - cost of retrieving next document

- Probability Ranking

$$C \cdot p(R|d) + C' \cdot (1 - p(R|d)) \leq C \cdot p(R|d') + C' \cdot (1 - p(R|d'))$$

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

Assignment Project Exam Help

Probability Ranking

- 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

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

Assignment Project Exam Help

ProbabilisticRetrievAdd WeChat edu_assist_pro

- 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
- Order documents by decreasing probability

Assignment Project Exam Help

Probabilistic Ranking

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

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](https://eduassistpro.github.io/)
By making assumptions about the distribution of terms and applying Bayes Theorem, it is possible to derive weights theoretically."

Van Rijsbergen

Assignment Project Exam Help

Binary Independence

- 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 modeled
- Bernoulli Naive Bayes model (cf. text categorization!)

Assignment Project Exam Help

Binary Independence

- Queries: binary term incidence vectors
- Given query q ,
 - for each document d need to compute $p(R|q,d)$.
 - replace with \vec{x} where $x_i = 1$ if t_i is binary term in document d
- Will use odds and Bayes' Rule

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

Assignment Project Exam Help

Binary Independence

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

Assignment Project Exam Help

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

Needs estimation

- Using Independence Assumption

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

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

Assignment Project Exam Help

Binary Independence

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

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

Assignment Project Exam Help

- Since x_i is e

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

- 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)

Assignment Project Exam Help

Binary Independence

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

~~Add WeChat~~ ~~edu_assist~~ ~~pro~~

All match

Non-matching
query terms

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

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

~~Add WeChat~~ ~~edu_assist~~ ~~pro~~

All matching terms

All query terms

Assignment Project Exam Help

Binary Independence

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

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

Assignment Project Exam Help

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

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)
to be estimated
for rankings

- Retrieval Status Value:

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

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

Assignment Project Exam Help

Binary Independence

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

- All boils down to computing RSV.

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

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

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

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

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

Assignment Project Exam Help

Binary Independence

Add WeChat edu_assist_pro

- 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/	n	N-n
$X_i = 0$			
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.

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

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

Documents	Relevant	Non-Relevant	Total
$X_i = 1$			$n+1$
$X_i = 0$			$N-n+1$
Total	$\Sigma+1$	Add WeChat edu_assist+2pro	

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

Assignment Project Exam Help

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	1	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

Assignment Project Exam Help

Example /2

- Estimate p_i and r_i

Doc	Judgment	x_1	x_2	x_3
D ₁	R		1	1
D ₂	R		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

Estimation

Add WeChat [edu_assist_pro](#)

- 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 (probability of relevant documents) Add WeChat [edu_assist_pro](#) 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 Add WeChat estimate

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 document set:
 - V is fixed set of documents on this model (not <https://eduassistpro.github.io/>)
3. We need to improve our guess p_i 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 Relevancy 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
3. Reestimate p_i and r_i on the ese
 - Or can combine new information with original guess (use Bayesian prior):
4. Repeat, thus generating a succession of approximations to R .

$$p_i^{(2)} = \frac{|V_i| + \kappa p_i^{(1)}}{|V| + \kappa}$$

κ is prior weight

Assignment Project Exam Help

PRP and BIR

Add WeChat edu_assist_pro

- Getting reasonable approximations of probabilities is possible.
- Requires r
 - *term ind* <https://eduassistpro.github.io/>
 - *terms not in query don't come*
 - *boolean representation*
 - *documents/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

Assignment Project Exam Help

Okapi BM₂₅^{Add WeChat edu_assist_pro}

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

Assignment Project Exam Help

caps the contribution of tf

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

idf

Normalized term
freq (doc)

Normalized term
freq (query)

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

Assignment Project Exam Help

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 ran
- 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

Assignment Project Exam Help

Resources

Add WeChat edu_assist_pro

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. 1979. *Information Retrieval*. 2nd ed. London: Butterworths, [th]

<http://www.dcs.bbk.ac.uk/~cjb/pubs/IR.pdf>

N. Fuhr. 1992. Probabilistic Models in Information Retrieval. *The Computer Journal*, 35(3), 243–255. [Easiest]

Add WeChat edu_assist_pro

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]

Assignment Project Exam Help

Resources

Add WeChat [edu_assist_pro](#)

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/Magazine/Vol12/12-04/vol12-04.html>

D. Heckerman. 1995. <https://eduassistpro.github.io/> etworks. Microsoft
Technical Report MSR-TR-95-06
<http://www.research.microsoft.com/~heckerm/>

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

Add WeChat `edu_assist_pro`

Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

LANGUAGE MODEL

Assignment Project Exam Help

Today

Add WeChat `edu_assist_pro`

- The Language Model Approach to IR

- Basic query generation model
- Alternative models

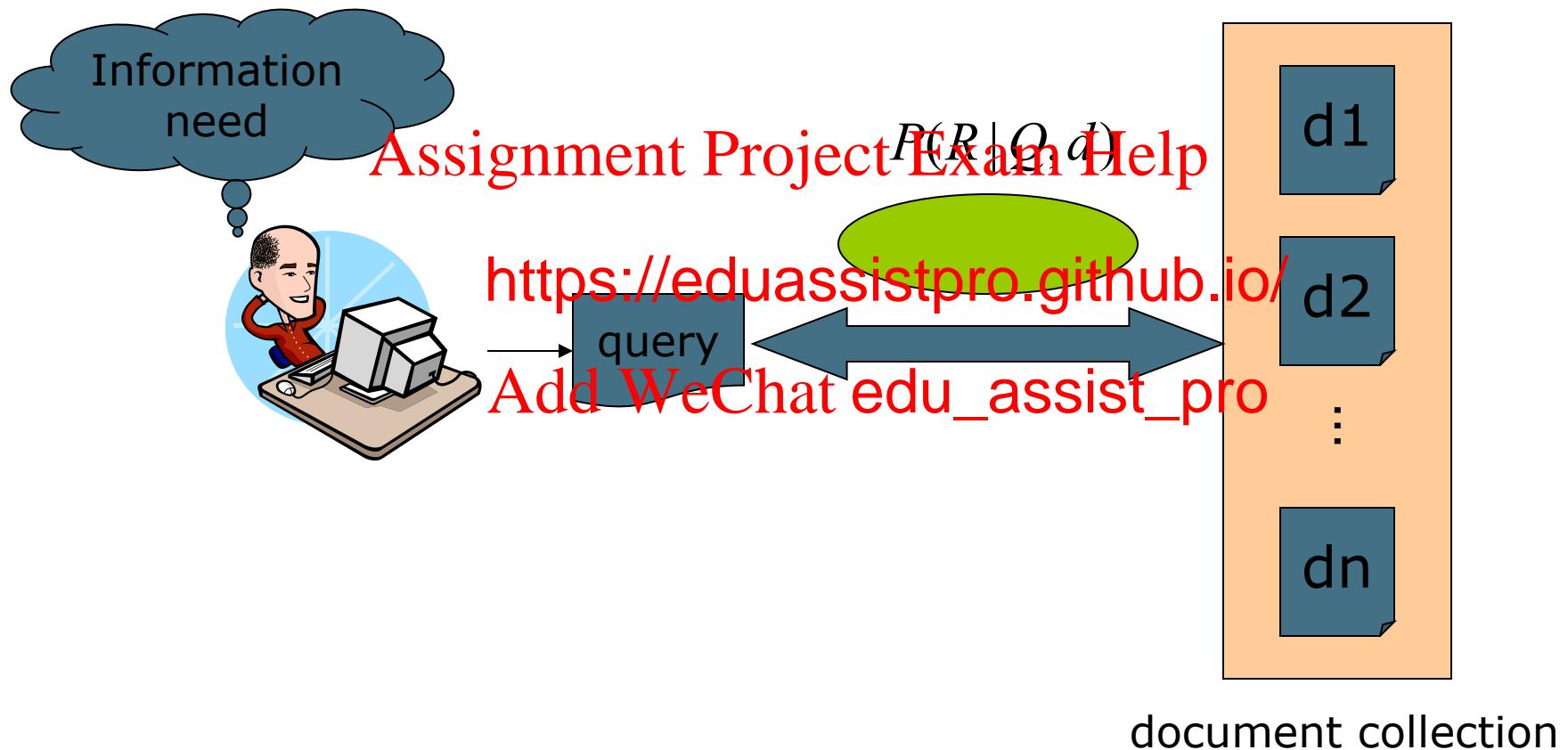
~~Assignment Project Exam Help~~

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

Add WeChat `edu_assist_pro`

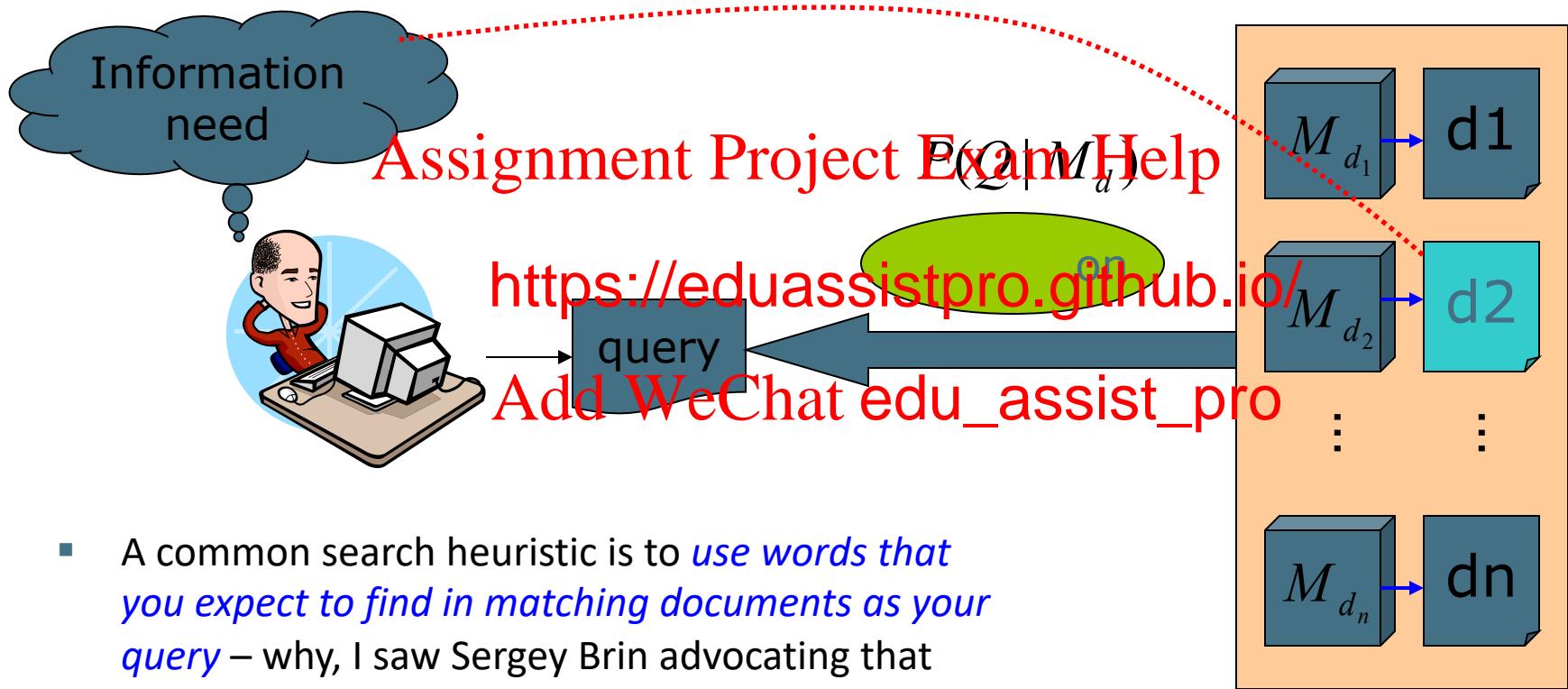
Assignment Project Exam Help

Standard Probability



Assignment Project Exam Help

IR based on Language (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

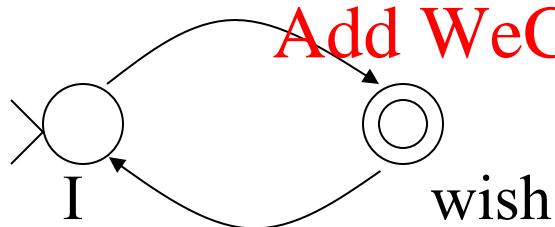
document collection

Assignment Project Exam Help

Formal Language (~~Add WeChat~~ ~~edu_assist_pro~~)

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

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



Add WeChat ~~edu_assist_pro~~

I wish

I wish I wish

I wish I wish I wish I wish

...

Assignment Project Exam Help

Stochastic Language Model

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

Model M

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

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

s the ma the woman

0.2 0.01 0.02 0.2 0.01

multiply

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

Assignment Project Exam Help

Stochastic Language Model

- 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/	
0.0001	Add WeChat
0.001	class
0.03	sayst
0.02	pleaseth
0.1	yon
0.01	maiden
0.0001	woman

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

Assignment Project Exam Help

Stochastic Language

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

Assignment Project Exam Help



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

$$P(\bullet \circ | M, \bullet)$$

$$P(\bullet \circ \bullet | M, \bullet \circ)$$

$$P(\bullet \circ \bullet \circ | M, \bullet \circ \bullet)$$

Assignment Project Exam Help

Unigram and higher-order Language Models

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

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

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

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

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

Assignment Project Exam Help

Using Language Model

- Treat each document as the basis for a model (e.g., unigram sufficient statistics)
- Rank document d based on $P(d | q)$
- $P(d | q) = P(q | \text{https://eduassistpro.github.io/}) P(d | \text{prior})$
 - $P(q)$ is the same for all documents
 - $P(d)$ [the prior] is often treated the same for all documents
 - 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

Assignment Project Exam Help

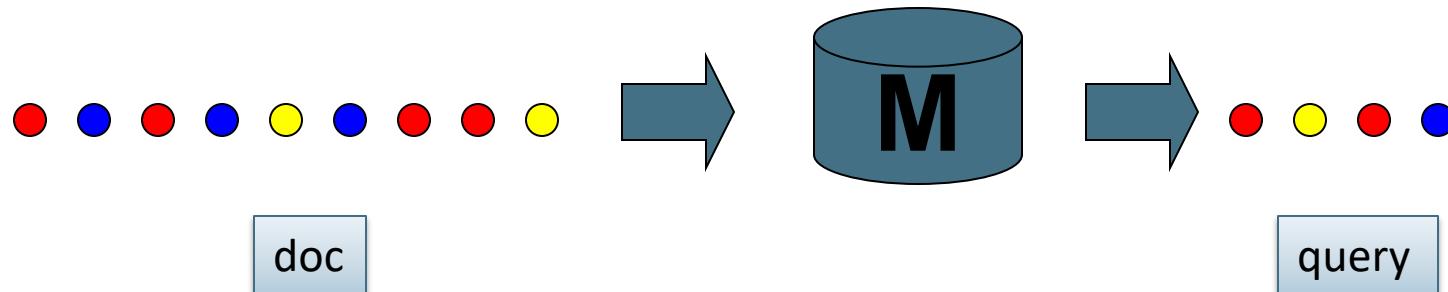
The fundamental problem

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

Assignment Project Exam Help

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

- Estimate a language model
- Then compute the observation probability



Assignment Project Exam Help

Language Models

- Language Modeling Approaches
 - Attempt to model query generation process
 - Documents are ranked by the probability that a query would be obe from the <https://eduassistpro.github.io/> respective do
- Multinomial approach

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

Assignment Project Exam Help

Retrieval based on ~~Add WeChat~~ $\text{prob}_{\text{edu_assist_pro}}$

- 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 <https://eduassistpro.github.io/>
 - Rank the documents according to probabilities.
 - Usually a unigram estimate or sed
 - Some work on bigrams, paralleling van Rijsbergen

Assignment Project Exam Help

Retrieval based on [Add WeChat edu_assist_pro](#)

- Intuition

- Users ...

- Have a reasonable idea of terms that are likely to occur in documents of interest.
 - They will [Assignment Project Exam Help](#) [https://eduassistpro.github.io/](#) from others

- Collection statistics [Add WeChat edu_assist_pro](#)

- 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

Assignment Project Exam Help

Query generation

- 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

Assignment Project Exam Help

Insufficient data

- Zero probability

$$p(t \mid M_d) = 0$$

- May not wish to assign a probability of zero to a document that is missing one or more of the query terms [gives conjunction sema]

- LM-based sm <https://eduassistpro.github.io/>

- A non-occurring term is ~~more~~ likely than would be expected by chance

- Naïve Idea: if $tf_{(t,d)} = 0$ then $p(t \mid M_d) \stackrel{?}{=} \frac{cf_t}{cs}$

- Need to work on the maths so that

$$\sum_{t \in V} p(t \mid M_d) = 1$$

cf_t : raw count of term t in the collection

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

Assignment Project Exam Help

Insufficient data

- Zero probabilities spell disaster
 - We need to smooth probabilities
 - Discount nonzero probabilities
 - Give some
- There's a wide range of smoothing techniques to address this problem, such as adding 1, $\frac{1}{2}$ or ϵ to the document multinomial prior, 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

Assignment Project Exam Help

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 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 (cf. Dirichlet prior or Witten-Bell smoothing)

Assignment Project Exam Help

Basic mixture mod

- 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.

Assignment Project Exam Help

n: # terms in Q

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

Relationship to tf

Add WeChat edu_assist_pro

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

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

proportional to the tf, inversely proportional to the cf

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

Becomes a constant | Q, C

Assignment Project Exam Help

Example Add WeChat edu_assist_pro

- 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 $\hat{\theta}_1 = 1/8$, $\hat{\theta}_2 = 2/16$, $\hat{\theta}_3 = 0$, $\hat{\theta}_4 = 1/16$, $\lambda = 1/2$
- 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: $d_1 > d_2$

Assignment Project Exam Help

Ponte and Croft Ex

- 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

<num> 1
<dom> www.sciencedirect.com
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>

Assignment Project Exam Help

Precision/recall res

Add WeChat edu_assist_pro

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

Add WeChat edu_assist_pro

Assignment Project Exam Help

Precision/recall res Add WeChat edu_assist_pro 100

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Assignment Project Exam Help

Language models: Add WeChat edu_assist_pro

- 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 metrics (almost...) be improved to the extent that the following 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

Assignment Project Exam Help

Comparison With ~~Add WeChat edu_assist_pro~~ pace

- 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 <https://eduassistpro.github.io/> all collection
 - the effect of frequencies is ~~Add WeChat edu_assist_pro~~ in the general collection but common in some will have a greater influence on the ranking

Assignment Project Exam Help

Comparison With ~~Add WeChat edu_assist_pro~~ pace

- Similar in some ways
 - Term weights based on frequency
 - Terms often used as if they were independent
 - Inverse document frequency 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

Assignment Project Exam Help

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. statistical translation. *SIGIR* 22, pp. 222–229.

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

D.R.H. Miller, T. Leek, an markov model information 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.