

Search Results Diversification

Kuan-Yu Chen (陳冠宇)

2019/11/01 @ TR-514, NTUST

Review

- Latent Semantic Analysis
- Topic Models
 - Probabilistic Latent Semantic Analysis
 - Latent Dirichlet Allocation

About HW3.

- The Expectation-Maximization algorithm
 - E-step

$$P(T_k | w_i, d_j) = \frac{P(w_i | T_k) P(T_k | d_j)}{\sum_{k=1}^K P(w_i | T_k) P(T_k | d_j)}$$

- M-step

$$P(w_i | T_k) = \frac{\sum_{d_j \in \mathbf{D}} c(w_i, d_j) P(T_k | w_i, d_j)}{\sum_{i'=1}^{|V|} \sum_{d_j \in \mathbf{D}} c(w_{i'}, d_j) P(T_k | w_{i'}, d_j)}$$

$$P(T_k | d_j) = \frac{\sum_{i=1}^{|V|} c(w_i, d_j) P(T_k | w_i, d_j)}{\sum_{i'=1}^{|V|} c(w_{i'}, d_j)} = \frac{\sum_{i=1}^{|V|} c(w_i, d_j) P(T_k | w_i, d_j)}{|d_j|}$$

- $P(w_i | T_k)$ and $P(T_k | d_j)$ are random initial with two constrains
 - $\sum_{w_i \in V} P(w_i | T_k) = 1$, for every topic T_k
 - $\sum_{k=1}^K P(T_k | d_j) = 1$, for every document d_j

About HW3..

- The probability $P(q|d_j)$ should be calculated in log domain

$$P(q|d_j) = \prod_{i=1}^{|q|} \left[\alpha \cdot P(w_i|d_j) + \beta \cdot \left(\sum_{k=1}^K P(w_i|T_k)P(T_k|d_j) \right) + (1 - \alpha - \beta) \cdot P_{BG}(w_i) \right]$$

$$\begin{aligned} \log P(q|d_j) &= \sum_{i=1}^{|q|} \log \left[\alpha \cdot P(w_i|d_j) + \beta \cdot \left(\sum_{k=1}^K P(w_i|T_k)P(T_k|d_j) \right) + (1 - \alpha - \beta) \cdot P_{BG}(w_i) \right] \\ &= \sum_{i=1}^{|q|} \left\{ [\log \alpha + \log P(w_i|d_j)] \oplus \left[\log \beta + \log \left(\sum_{k=1}^K P(w_i|T_k)P(T_k|d_j) \right) \right] \oplus [\log(1 - \alpha - \beta) + \log P_{BG}(w_i)] \right\} \end{aligned}$$

numpy.logaddexp

```
numpy.logaddexp(x1, x2, /, out=None, *, where=True, casting='same_kind', order='K', dtype=None, subok=True[, signature, extobj]) = <ufunc 'logaddexp'>
```

Logarithm of the sum of exponentiations of the inputs.

Calculates $\log(\exp(x1) + \exp(x2))$. This function is useful in statistics where the calculated probabilities of events may be so small as to exceed the range of normal floating point numbers. In such cases the logarithm of the calculated probability is stored. This function allows adding probabilities stored in such a fashion.








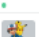
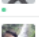
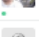



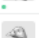




Parameters:

- x1, x2** : array_like
Input values.
- out** : ndarray, None, or tuple of ndarray and None, optional
A location into which the result is stored. If provided, it must have a shape that the inputs broadcast to. If not provided or None, a freshly-allocated array is returned. A tuple (possible only as a keyword argument) must have length equal to the number of outputs.
- where** : array_like, optional
Values of True indicate to calculate the ufunc at that position, values of False indicate to leave the value in the output alone.
- **kwargs**
For other keyword-only arguments, see the [ufunc docs](#).

Returns:

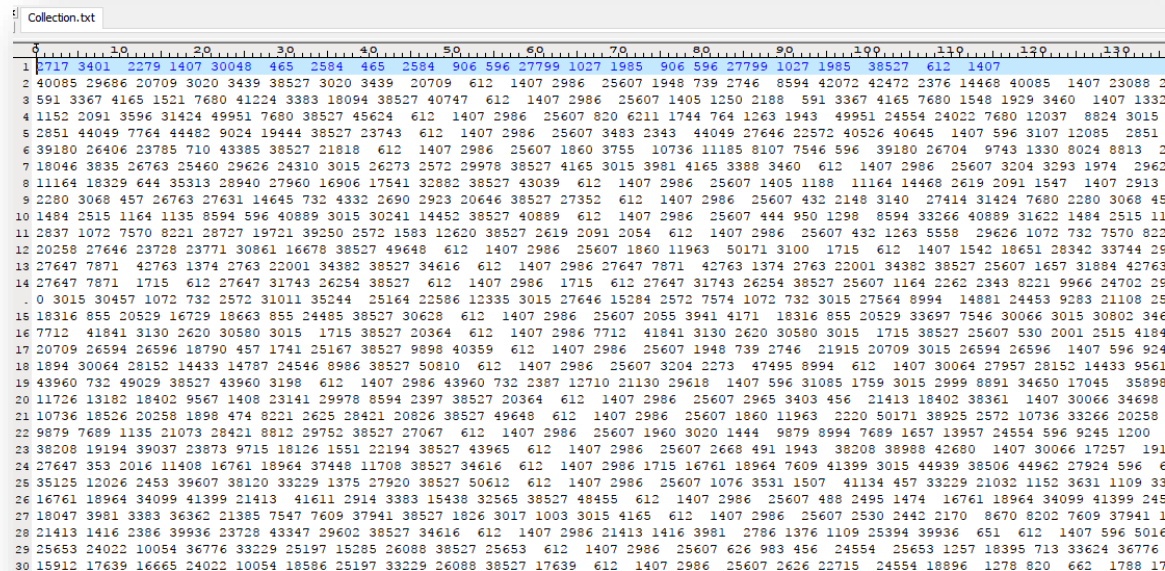
- result** : ndarray
Logarithm of $\exp(x1) + \exp(x2)$.

About HW3...

#	Team Name	Notebook	Team Members	Score ?	Entries
1	B			0.52528	55
2	ye			0.51220	43
3	M			0.50242	55
📍	ULM			0.50229	
4	te			0.49764	9
5	B			0.49679	47
6	M			0.49538	27
7	Pe			0.49281	20
8	M			0.49250	1
9	B			0.49096	16
10	M			0.45864	75
11	M			0.44721	61
📍	VSM			0.26967	
12	M			0.22473	15
13	M			0.22344	7
14	M			0.19460	6
15	te			0.10610	3
16	M			0.05606	8
17	al			0.03981	13
18	M			0.02019	1

About HW3...

- A Collection.txt is released
 - Each line is a document
 - About 18 thousand documents in total
- Train a PLSA model on the corpus with 128 topics, we can achieve 0.53x by tuning the parameters easily
 - But the fold-in strategy is needed!



The screenshot shows a text file named 'Collection.txt' with a large number of lines, each containing a sequence of integers. The numbers are arranged in a dense grid, with some lines starting with a small number (likely a document ID) followed by a series of topic or term values. The values range from 0 to 1000, with some lines showing values up to 1591. The file appears to be a large dataset of documents, each represented by a line of numbers.

About Final Project

- Group your team!
 - 2~4 team members
 - Choose a paper
- Do you have GPU units?
 - We have to make sure you can do HW5 and final project

Resource

- Conferences

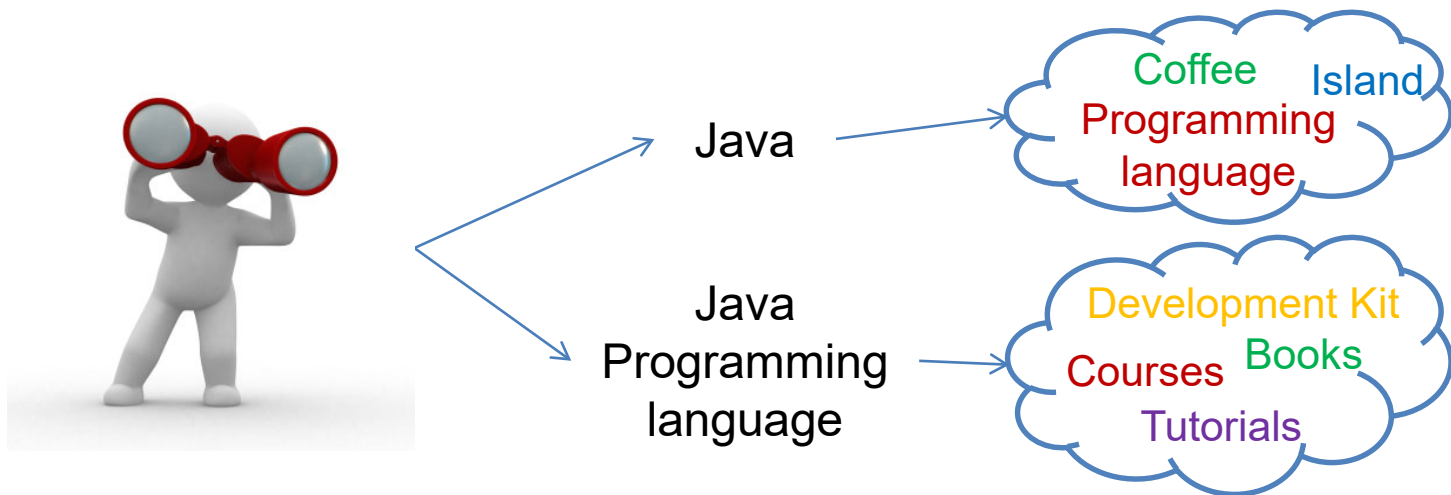
- ACM Annual International Conference on Research and Development in Information Retrieval (SIGIR)
- International Joint Conferences on Artificial Intelligence (IJCAI)
- ACM Conference on Information Knowledge Management (CIKM)
- Annual Meeting of the Association for Computational Linguistics (ACL)
- International Conference on Learning Representations (ICLR)

- Journals

- Journal of the American Society for Information Science (JASIS)
- ACM Transactions on Information Systems (TOIS)
- Information Processing and Management (IP&M)
- ACM Transactions on Asian Language Information Processing (TALIP)
- Information Retrieval Journal (IRJ)

Introduction – What's going on?

- Traditional retrieval functions ignore the relations among returned documents
 - Top ranked documents may contain relevant yet **redundant information**
 - In order to maximize the satisfaction of different search users, it is necessary to diversify search results
 - Search results diversification can play an initial step for many search system

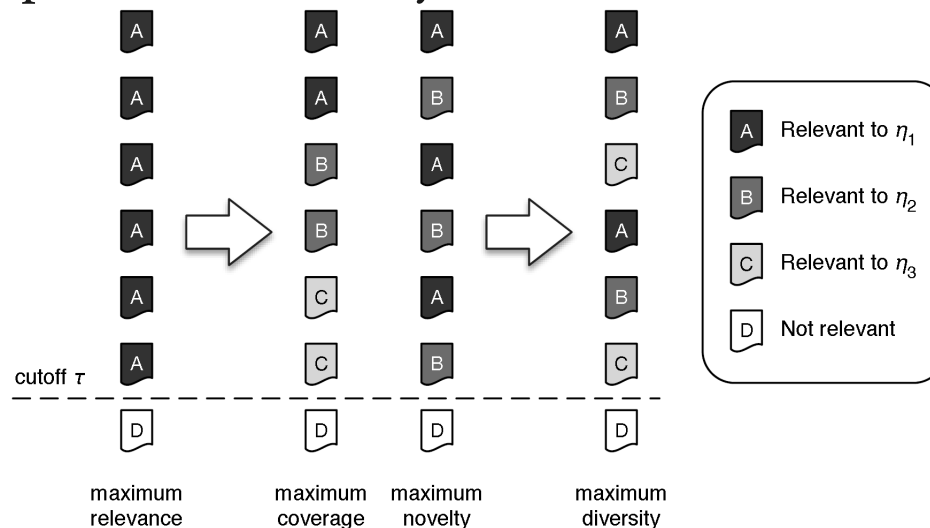


Relevance, Coverage, Novelty, & Diversity

- Most of the retrieval models assume that the relevance of a document can be estimated with **certainty** and **independently** of the estimation of the other retrieved documents
 - Ambiguous queries
 - Ensuring a high **coverage** of the possible information needs
 - Redundancy results
 - Ensuring the retrieved documents provide a high **novelty**

Relevance, Coverage, Novelty, & Diversity

- Coverage and novelty can also be conflicting objectives
 - A ranking with maximum coverage may not attain maximum novelty
 - Although covering all information needs, the ranking may place all documents covering a particular need ahead than others
 - A ranking with maximum novelty may not attain maximum coverage
 - Although covering each need as early as possible in the ranking, not all possible needs may be covered



Introduction – Various Modeling

- Many diversification methods have been proposed
 - balance the relevance and the redundancy: MMR
 - distinguish previous topics and new coming: SMM
 - language modeling approach: WUME
 - probabilistic framework: xQuAD
- These methods mainly differ in **diversity modeling**
 - **Implicitly**: The diversity is implicitly modeled through document similarities
 - **Explicitly**: It can be explicitly modeled through the coverage of query subtopics, and document dependency

Introduction – Notations

Symbol	Description
q	A given query
a_k^q	Sub-queries (aspect), $q = \{a_1^q, \dots, a_K^q\}$
K	Number of sub-queries
R	The user's information need
\mathbf{D}	A set of documents, $\mathbf{D} = \{d_1, \dots, d_{ \mathbf{D} }\}$
$\tilde{\mathbf{D}}$	A subset of documents which already selected by new method, $\tilde{\mathbf{D}} = \{\tilde{d}_1, \dots, \tilde{d}_{ \tilde{\mathbf{D}} }\}$

Maximal Marginal Relevance – MMR

- MMR motivated the need for “relevant novelty” as a potentially superior criterion
 - An approximation to measuring relevant novelty is to **measure relevance and novelty independently**
- “Marginal Relevance” can be regarded as the metric
 - A document has high **marginal relevance** if it is both relevant to the query and contains minimal similarity to previously selected documents

$$Div_{MMR}(d, q) = - \max_{\tilde{d} \in \tilde{\mathbf{D}}} sim(d, \tilde{d})$$

$$d^* = \operatorname{argmax}_{d \in \mathbf{D}} \lambda \cdot Rel(d, q) - (1 - \lambda) \cdot \max_{\tilde{d} \in \tilde{\mathbf{D}}} sim(d, \tilde{d})$$

Simple Mixture Model – SMM

- Given the observed new document, we estimate the mixing weight for the background model θ_{BG} and the previous topic model θ_T
 - The simplest previous topic model can be modeled as:

$$P(w|\theta_T) = \sum_{\tilde{d} \in \tilde{\mathbf{D}}} \frac{1}{N} P(w|\tilde{d})$$

- The mixture weight for the background model can serve as a measure of novelty or redundancy

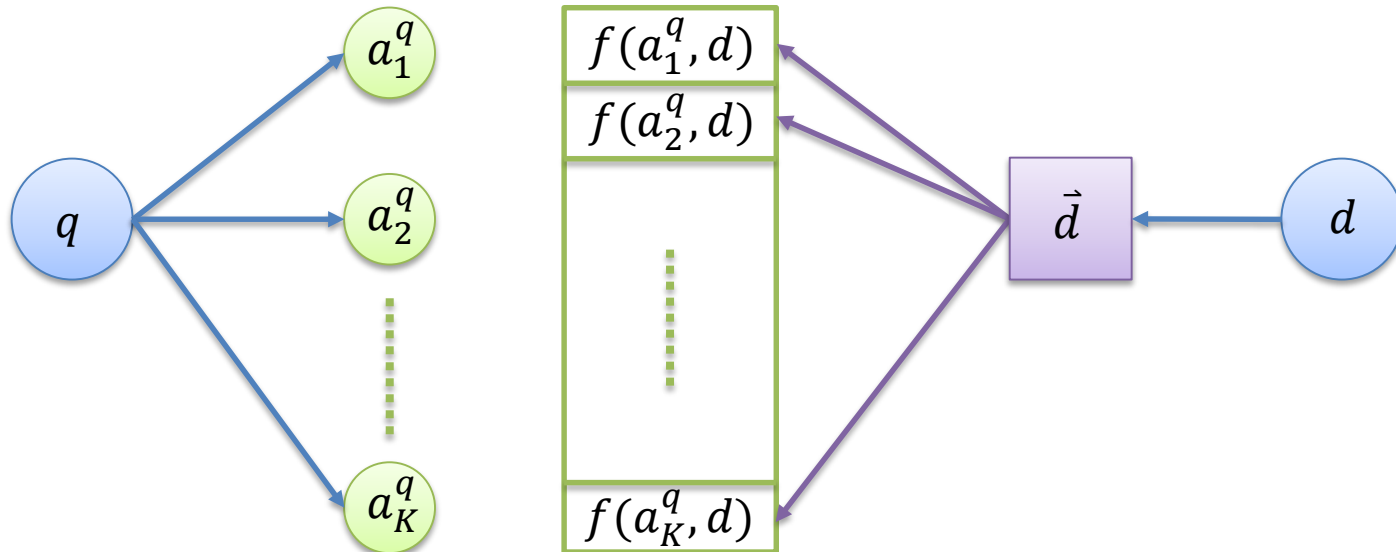
$$L(\beta|d, \theta_{BG}, \theta_T) = \prod_{w \in V} \left((1 - \beta) \cdot P(w|\theta_T) + \beta \cdot P(w|\theta_{BG}) \right)^{c(w,d)}$$

$$d^* = \operatorname{argmax}_{d \in \mathbf{D}} \lambda \cdot \operatorname{Rel}(d, q) + (1 - \lambda) \cdot \operatorname{argmax}_{\beta} L(\beta|d, \theta_{BG}, \theta_T)$$

$$\begin{aligned} 1 - \beta &= P(\theta_T|d) \\ \beta &= P(\theta_{BG}|d) \end{aligned}$$

Explicit MMR – xMMR

- For a given query with its sub-queries, each document can be represented by a K -dimensional vector over sub-queries

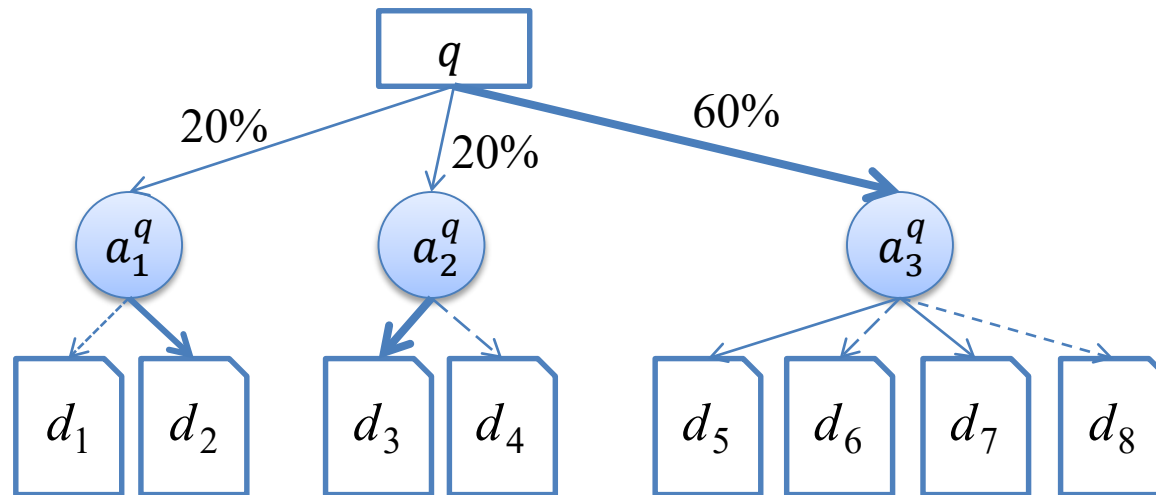


$$f(a^q, d) \equiv P(d|a^q) \quad f(a^q, d) \equiv \cos(a^q, d)$$

- By doing so, the redundancy score can be defined by considering sub-queries

$$d^* = \operatorname{argmax}_{d \in \mathbf{D}} \lambda \cdot \operatorname{Rel}(d, q) - (1 - \lambda) \cdot \max_{\vec{d} \in \tilde{\mathbf{D}}} \operatorname{sim}(\vec{d}, \vec{\tilde{d}})$$

WUME – Motivation



- There are three sub-queries under the given query $q = \{a_1^q, a_2^q, a_3^q\}$, and web documents $\mathbf{D} = \{d_1, \dots, d_8\}$
- Although d_3 is more relevant to one of the sub-query a_2^q than d_5 to a_3^q , given that a_2^q attracts less user interest than a_3^q , d_3 **should** still be ranked lower than d_5

WUME

- WUME formalize the diversification method as:
 - Given a query Q , the probability that a retrieved document meets user's information need R can be written as:

$$P(R|d) = \frac{P(R)P(d|R)}{P(d)} \propto P(d|R)$$

- Take sub-query information into consideration:

$$P(d|R) \approx P(d|q) = \sum_{k=1}^K P(d|a_k^q, q)P(a_k^q|q)$$

Google Insights for
Search or Wikipedia



- Finally, the ranking function becomes:

$$d^* = \operatorname{argmax}_{d \in \mathbf{D}} \lambda \cdot \operatorname{Rel}(d, q) + (1 - \lambda) \cdot \sum_{k=1}^K P(d|a_k^q, q)P(a_k^q|q)$$

eXplicit Query Aspect Diversification

- xQuAD: eXplicit Query Aspect Diversification
 - When given an ambiguous query, xQuAD builds a new ranked list by:

$$d^* = \operatorname{argmax}_{d \in \mathbf{D}} \lambda \cdot P(d|q) + (1 - \lambda) \cdot P(d, \tilde{\mathbf{D}}|q)$$

- $P(d|q)$ is the likelihood of document d being observed given the initial query
 - The probability can be regarded as modeling *relevance*
- $P(d, \tilde{\mathbf{D}}|q)$ is the likelihood of observing this document but not the documents already in $\tilde{\mathbf{D}}$
 - The probability can be regarded as modeling *diversity*

xQuAD – 1

- In order to derive $P(d, \bar{\mathbf{D}}|q)$, xQuAD explicitly consider the possibly several aspects underlying the initial query as a set of sub-queries
 - By assuming $\sum_{k=1}^K P(a_k^q|q) = 1$, xQuAD calculates $P(d, \bar{\mathbf{D}}|q)$ by considering sub-queries:

$$P(d, \bar{\mathbf{D}}|q) = \sum_{k=1}^K P(d, \bar{\mathbf{D}}|a_k^q)P(a_k^q|q)$$

- Further, $P(d, \bar{\mathbf{D}}|a_k^q)$ can be broken down by independent assumption:

$$P(d, \bar{\mathbf{D}}|a_k^q) = P(d|a_k^q)P(\bar{\mathbf{D}}|a_k^q)$$

coverage *novelty*
↓ ↓

xQuAD – 2

- xQuAD also assumes that the relevance of each document in $\tilde{\mathbf{D}}$ to a given sub-query a_k^q is independent

$$P(\tilde{\mathbf{D}}|a_k^q) = P(\tilde{d}_1, \dots, \tilde{d}_{|\tilde{\mathbf{D}}|} | a_k^q) = \prod_{\tilde{d}_n \in \tilde{\mathbf{D}}} P(\tilde{d}_n | a_k^q) = \prod_{\tilde{d}_n \in \tilde{\mathbf{D}}} (1 - P(\tilde{d}_n | a_k^q))$$

- To sum up, xQuAD suggests that:

$$Div_{xQuAD}(d, q) = \sum_{k=1}^K \underbrace{P(a_k^q | q)}_{\text{the importance of } a_k^q} \underbrace{P(d | a_k^q)}_{\text{the relevance of } d \text{ to } a_k^q} \prod_{\tilde{d}_n \in \tilde{\mathbf{D}}} \underbrace{(1 - P(\tilde{d}_n | a_k^q))}_{\text{the satisfaction degree } a_k^q}$$

- Instead of comparing a document d to all documents already selected in $\tilde{\mathbf{D}}$, xQuAD estimates the utility of any document satisfying the sub-query a_k^q , given how well it is already satisfied by the documents in $\tilde{\mathbf{D}}$

Analytical Comparisons

- Diversity Modeling:
 - MMR and SMM **implicitly** model the diversity through document similarities
 - xMMR, WUME and xQuAD **explicitly** model the diversity through the coverage of query subtopics
- Document Dependency:
 - WUME assumes that the diversity score of a document is independent of other documents
 - The other three methods assume that the diversity score depends on the previously selected documents

General Framework

- Most of these methods **iteratively select** the document that is not only **relevant** to the query but also **diversified** to cover more query subtopics, explicitly or implicitly
- All of methods fit into a general framework that iteratively selects with the highest relevance and diversity scores:

$$d^* = \operatorname{argmax}_{d \in \mathbf{D}} \lambda \cdot \operatorname{Rel}(d, q) + (1 - \lambda) \cdot \operatorname{Div}(d, q)$$

Experimental Results

	TREC09 result		TREC10 result	
	α -nDCG@20	α -nDCG@100	α -nDCG@20	α -nDCG@100
<i>MMR</i> *	0.365	0.427	0.344	0.415
<i>WUME</i> *	0.479	0.546	0.579	0.630
<i>xQuAD</i> *	0.482	0.550	0.588	0.636

- All the parameters in each method are set to the optimum values
 - Both xQuAD and WUME perform significantly better than MMR
 - Using explicit sub-queries in diversification is better
 - The performances of xQuAD and WUME are not significantly different

Conclusions

- The experiment result shows that the explicit sub-query modeling and sub-query importance penalization strategies perform better
- It is interesting to find that how the sub-queries affect the overall performance
- Finally, we can think about that what's the difference between sub-queries and latent topics?
 - Supervised v.s. Unsupervised?
- Beyond relevance or another relevance?

$$P(d|R) \approx P(d|q) = \sum_{k=1}^K P(d|a_k^q, q)P(a_k^q|q)$$

Questions?



kychen@mail.ntust.edu.tw