# Search Results Diversification
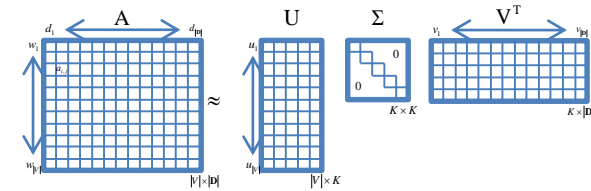
**Kuan-Yu Chen (陳冠宇)**

2020/11/20 @ TR-313, NTUST

# Review.

- Latent Semantic Analysis
  - Vector representation for word is $\Sigma u_i^{\mathrm{T}}$
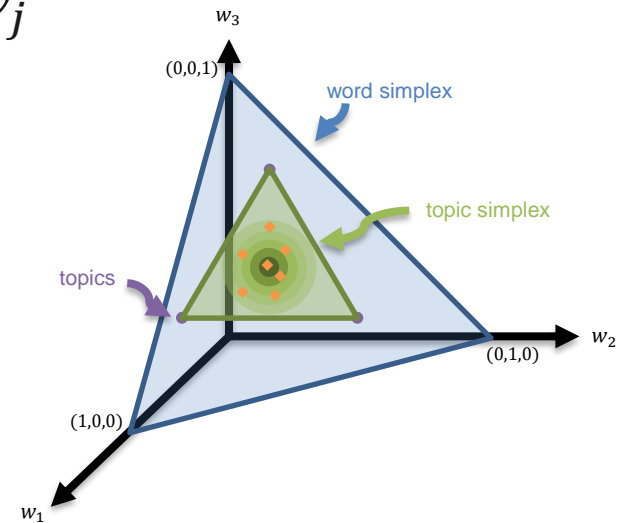  - Vector representation for document is $\Sigma v_j^{\mathrm{T}}$

- Statistical Topic Models
  - Probabilistic Latent Semantic Analysis
    - $\mathcal{L} = \prod_{w_i \in V} \prod_{d_j \in \mathbf{D}} P(w_i, d_j)^{c(w_i, d_j)}$
  - Latent Dirichlet Allocation
    - $\mathcal{L} = \prod_{d_j \in \mathbf{D}} \int P(\theta_{d_j}|\alpha) \left( \prod_{i=1}^{|d_j|} \left( \sum_{k=1}^{K} P(w_i|T_k, \beta) P\left(T_k \middle| \theta_{d_j}\right) \right) \right) d\theta_{d_j}$

# Review..

- 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 $d_j$

3

# Review…

- About the M-step

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

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

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

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

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

# Review….

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

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

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

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

## 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)`.

# Homework 4.

- The evaluation measure is MAP@1000
  - The **hard** deadline is 11/26 23:59
  - You point is depended on your performance on the **private** leaderboard!
    - 50 public queries and 50 private queries

$$YourScore = 5 + \frac{YourMAP - BaselineMAP}{HighestMAP - BaselineMAP} \times 8$$

  - Please submit a **report** and your **source codes** to the Moodle system, otherwise you will get 0 point
    - The report will be judged by TA, and the score is either 1 or 2

# Homework 4..

| # | Team Name |
|---|-----------|
| 1 | 宮澤賢治 |
| 2 | 康帕內魯拉 |
| 📍 | FYI: with 32 topics |
| 3 | Test |
| 4 | M10815048_張晏銘 |
| 5 | Enilesab |
| 6 | TESTT1235 |
| 7 | B10615034_黃柏翰 |
| 8 | Alice |
| 9 | 80847002S_羅天宏 |
| 10 | M10915100_郭智威 |
| 📍 | FYI: with 8 topics |
| 11 | 快樂ㄟ甘蔗man |
| 12 | hongyun |
| 📍 | Baseline 0.4P(w|D) 0.6P(w|BG) |
| 13 | 在小世界裡畫出最耀眼的大幸福 |
| 14 | trytrysee |
| 15 | 騙人的吧 |

| # | △pub | Team Name |
|---|------|-----------|
| 1 | ▲1 | 康帕內魯拉 |
| 2 | ▲3 | Enilesab |
| 3 | — | Test |
| 4 | — | M10815048_張晏銘 |
| 📍 | | FYI: with 32 topics |
| 5 | ▼4 | 宮澤賢治 |
| 6 | ▲1 | B10615034_黃柏翰 |
| 7 | ▲3 | M10915100_郭智威 |
| 8 | ▼2 | TESTT1235 |
| 9 | — | 80847002S_羅天宏 |
| 10 | ▼2 | Alice |
| 📍 | | FYI: with 8 topics |
| 11 | ▲1 | hongyun |
| 12 | ▼1 | 快樂ㄟ甘蔗man |
| 📍 | | Baseline 0.4P(w|D) 0.6P(w|BG) |
| 13 | — | 在小世界裡畫出最耀眼的大幸福 |
| 14 | — | trytrysee |
| 15 | — | 騙人的吧 |

# 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 HW6 and/or final project

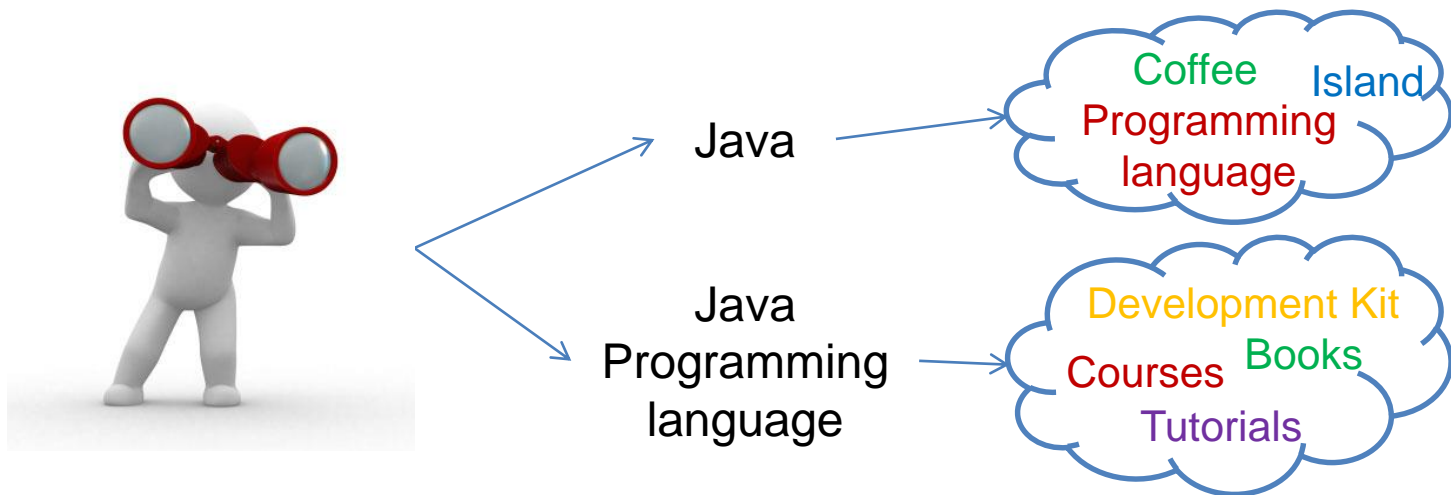| Date | Syllabus | Homework |
|------|----------|----------|
| 9/18 | Course Overview | |
| 9/25 | Break for Rocling2020 | |
| 10/2 | Holiday for Moon Festival | |
| 10/9 | Holiday for National Day | |
| 10/16 | Classic Models | Homework-1(deadline: 10/29 23:59) |
| 10/23 | Extended Probabilistic Models | Homework-2 (deadline: 11/5 23:59) |
| 10/30 | Evaluation & Benchmark Collections | Homework-3 (deadline: 11/12 23:59) |
| 11/6 | Latent Semantic Analysis | |
| 11/13 | Statistical Topic Models | Homework-4 (deadline: 11/26 23:59) |
| 11/20 | Search Results Diversification | |
| 11/27 | Pseudo-Relevance Feedback & Query Models | Homework-5 (deadline: 12/10 23:59) |
| 12/4 | Talk | Submit Your Member List! |
| 12/11 | Representation Learning for Information Retrieval | |
| 12/18 | Supervised Retrieval Models & Information Retrieval in Practice | Homework-6 (deadline: 12/31 23:59) & Submit Your Paper Title! |
| 12/25 | Break for Your Final Project | |
| 1/1 | Holiday for Founding Anniversary | |
| 1/8 | Presentation-1 | |
| 1/15 | Presentation-2 | |

# **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

Java

Coffee    Island
Programming
language

Java
Programming
language

Development Kit
Courses    Books
Tutorials

# 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 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

# Example.



綜合報導

化學定義

周刊科普

相關新聞

# Example..

# 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), $q = \{a_1^q, \cdots, a_K^q\}$ |
| $K$ | Number of sub-queries |
| $R$ | The user's information need |
| $\mathbf{D}$ | A set of documents, $\mathbf{D} = \{d_1, \cdots, d_{|\mathbf{D}|}\}$ |
| $\widetilde{\mathbf{D}}$ | A subset of documents which already selected by new method, $\widetilde{\mathbf{D}} = \{\tilde{d}_1, \cdots, \tilde{d}_{|\widetilde{\mathbf{D}}|}\}$ |

# 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" cab 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

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

$$score(d, q) = \lambda \cdot Rel(d, q) - (1 - \lambda) \cdot \max_{\tilde{d} \in \tilde{\mathbf{D}}} sim(d, \tilde{d})$$

# Explicit MMR – xMMR

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



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

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

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

$$score(d, q) = \lambda \cdot Rel(d, q) - (1 - \lambda) \cdot \max_{\tilde{d} \in \tilde{\mathbf{D}}} sim(d, \tilde{d})$$

# Simple Mixture Model – SMM

- Given the observed new document, we estimate the mixing weight for the background model $\theta_{BG}$ and the previous topic model $\theta_T$

  - 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

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

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

$$score(d, q) = \lambda \cdot Rel(d, q) + (1-\lambda) \cdot \beta$$

# WUME – Motivation



- There are three sub-queries under the given query $q = \{a_1^q, a_2^q, a_3^q\}$, and web documents $\mathbf{D} = \{d_1, \cdots, d_8\}$

- 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:

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

Google Insights for Search or Wikipedia

  - Finally, the ranking function becomes:

$$score(d, q) = \lambda \cdot Rel(d, q) + (1 - \lambda) \cdot \sum_{k=1}^{K} P\left(d\middle|a_k^q, q\right)P\left(a_k^q\middle|q\right)$$

# eXplicit Query Aspect Diversification

- xQuAD: eXplicit Query Aspect Diversification
  - When given an ambiguous query, xQuAD estimates the ranking score by:

$$score(d, q) = \lambda \cdot P(d|q) + (1 - \lambda) \cdot P\left(d, \overline{\overline{\mathbf{D}}}|q\right)$$

  - $P(d|q)$ is the likelihood of document $d$ being observed given the initial query

    The probability can be regarded as modeling *relevance*

  - $P\left(d, \overline{\overline{\mathbf{D}}}|q\right)$ is the likelihood of observing this document but not the documents already in $\widetilde{\mathbf{D}}$

    The probability can be regarded as modeling *diversity*

# xQuAD – 1

- In order to derive $P\left(d, \overline{\overline{\mathbf{D}}}|q\right)$, 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 $P\left(d, \overline{\overline{\mathbf{D}}}|q\right)$ by considering sub-queries:

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

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

$$P\left(d, \overline{\overline{\mathbf{D}}}|a_k^q\right) = \underbrace{P(d|a_k^q)}_{coverage} \underbrace{P\left(\overline{\overline{\mathbf{D}}}|a_k^q\right)}_{novelty}$$

# xQuAD – 2

– For $P\left(\bar{\bar{\mathbf{D}}}|a_k^q\right)$, xQuAD assumes that the relevance of each document in $\widetilde{\mathbf{D}}$ to a given sub-query $a_k^q$ is independent

$$P\left(\bar{\bar{\mathbf{D}}}|a_k^q\right) = P\left(\bar{\bar{d}}_1, \cdots, \bar{\bar{d}}_{|\widetilde{\mathbf{D}}|}\Big|a_k^q\right) = \prod_{\tilde{d}_n \in \widetilde{\mathbf{D}}} P(\bar{\bar{d}}_n|a_k^q) = \prod_{\tilde{d}_n \in \widetilde{\mathbf{D}}} (1 - P(\tilde{d}_n|a_k^q))$$

– To sum up, xQuAD suggests that:

$$\begin{aligned}
P\left(d, \bar{\bar{\mathbf{D}}}|q\right) &= \sum_{k=1}^{K} P\left(d, \bar{\bar{\mathbf{D}}}|a_k^q\right) P(a_k^q|q) \\
&= \sum_{k=1}^{K} P(a_k^q|q) P\left(d|a_k^q\right) P\left(\bar{\bar{\mathbf{D}}}|a_k^q\right) \\
&= \sum_{k=1}^{K} P(a_k^q|q) P(d|a_k^q) \prod_{\tilde{d}_n \in \widetilde{\mathbf{D}}} (1 - P(\tilde{d}_n|a_k^q))
\end{aligned}$$

# xQuAD – 3

- The final score for each document is determined by:

$$score(d, q) = \lambda \cdot P(d|q) + (1 - \lambda) \cdot P\left(d, \bar{\bar{\mathbf{D}}}|q\right)$$

$$= \lambda \cdot P(d|q) + (1 - \lambda) \cdot \sum_{k=1}^{K} \underline{P(a_k^q|q)} \underline{P(d|a_k^q)} \prod_{\tilde{d}_n \in \tilde{\mathbf{D}}} \underline{(1 - P(\tilde{d}_n|a_k^q))}$$

**the importance of $a_k^q$**

**the relevance of $d$ to $a_k^q$**

**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:

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

# Experimental Results

| | TREC09 result | | TREC10 result | |
|---|---|---|---|---|
| | $\alpha$-nDCG@20 | $\alpha$-nDCG@100 | $\alpha$-nDCG@20 | $\alpha$-nDCG@100 |
| $MMR^*$ | 0.365 | 0.427 | 0.344 | 0.415 |
| $WUME^*$ | 0.479 | 0.546 | 0.579 | 0.630 |
| $xQuAD^*$ | 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?

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

- Beyond relevance? Another relevance?

# Questions?



**kychen@mail.ntust.edu.tw**