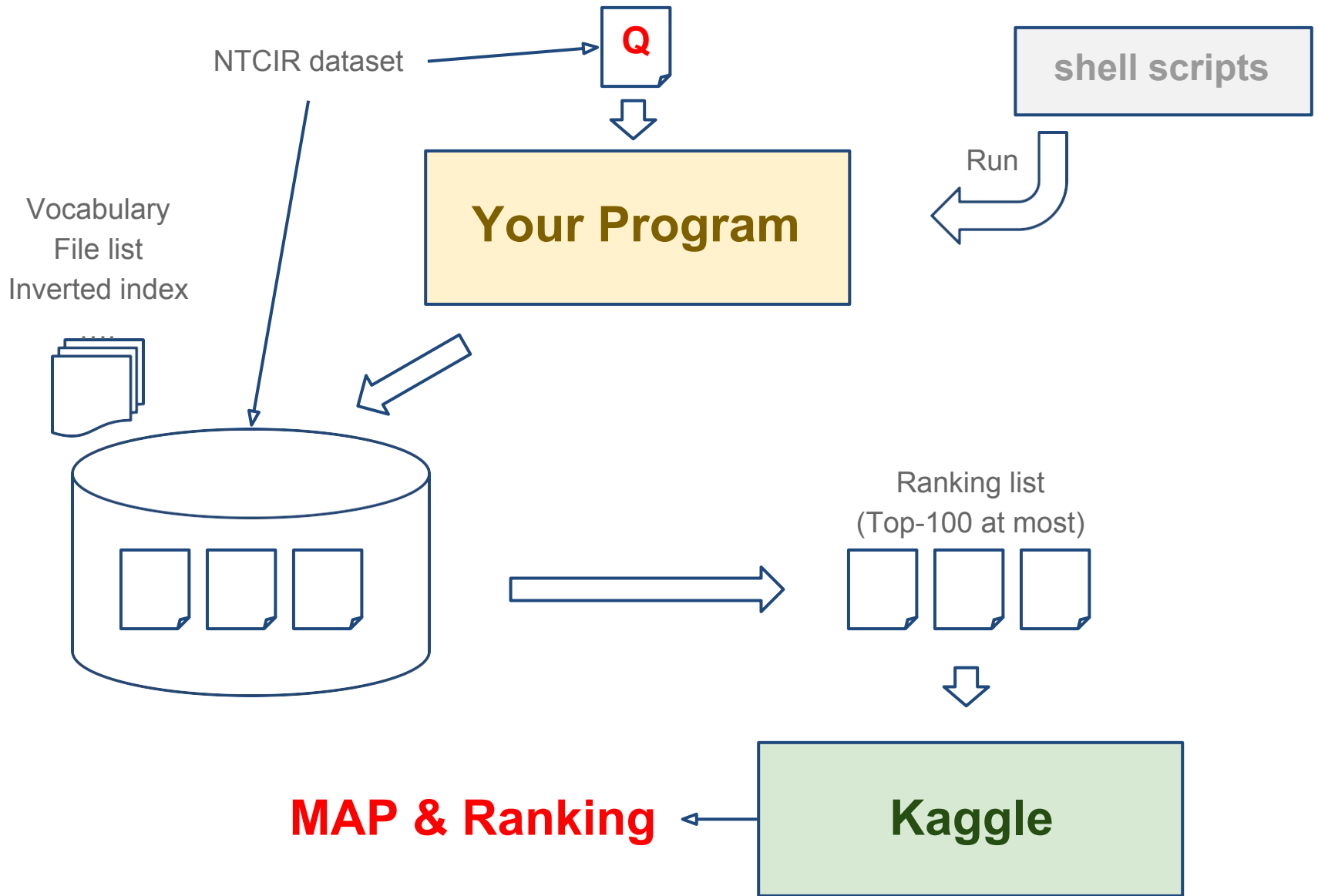# Programming HW 1

## Web Retrieval and Mining spring 2018

# Introduction

- In this homework, you are asked to implement **a small information retrieval system**.

- We will give you a bunch of Chinese news articles and several queries in NTCIR format, and your task is to find the relevant documents among these articles according to the given queries.

- You should implement the retrieval system by **Vector Space Model** (VSM) with **Rocchio Relevance Feedback** (pseudo version).

# Kaggle

[https://www.kaggle.com/c/ntucsie-wm2018-vsm](https://www.kaggle.com/c/ntucsie-wm2018-vsm)

# Data

1. CIRB010.tar.gz
2. model.tar.gz
3. user_aggreement_form_for_ntcir.pdf
4. query_train.xml
5. ans_train.csv
6. query_test.xml
7. script.tar.gz
8. sample_submission.csv

# NTCIR Document Format

- The NTCIR document format conforms to XML 1.0
- The root element is <xml>, it contains exactly one <doc> tag.
- A <doc> tag represents exactly one newswire article, in which several sub-elements are used to specify different type of information:
  - <id>: An unique document ID.
  - <date>: The publication date.
  - <title>: The title of the article.
  - <text>: The content of the article, which may include one or more passages enclosed in <p> tags.

# CIRB010.tar.gz

```xml
CDN_ECO_0000006

1  <?xml version="1.0" encoding="UTF-8"?>
2  <xml>
3  <doc>
4  <id>cdn_eco_0000006</id>
5  <date>1998-05-08</date>
6  <title> 培訓軟體人才實施三年計畫 </title>
7  <text>
8  <p>
9  行政院會昨日通過「加強資訊軟體人才培訓方案」，預定自今（八十七）年七月起至九十年六月等三年度內，培訓二萬二千五百人次投入就業市場，預計培訓人才全部
   就業後，每年將可增加軟體及相關產業產值約十八億美元。
10 </p>
11 <p>
12 這項方案由行政院副院長劉兆玄召集各部會研擬而成，蕭萬長在昨日院會中，肯定此項方案的研擬過程充分發揮行政團隊精神，以快速、實際的行動協助業者解決問題
   ，為行動內閣建立良好的典範。
13 </p>
14 <p>
15 目前各行業及資訊業均面臨資訊軟體人才嚴重短缺的問題，但若由正規教育體系培養供應，估計每年仍不足五千至一萬人，為在短期間協助解決人力不足問題，行政院
   特別由副院長劉兆玄召集「軟體人才培訓專案小組」，以培訓非資訊科系畢業生第二專才之養成訓練為主，培訓在職工作人員之進階訓練為輔，其中在職工作人員不限
   於政府公職人員，也包括私人企業人員。
16 </p>
17 <p>
18 蕭萬長在院會中指出，在硬體方面，我國資訊產品高居世界第三位，但軟體方面尚須加強，他特別舉比爾‧蓋茲為例強調軟體的重要性。蕭萬長指出，在資訊產品硬、
   軟體方面都屬領先的ＩＢＭ公司，在利潤製造上，反而不及比爾‧蓋茲，而比爾‧蓋茲能累積這麼多財富，企業如此成功，最重要就是靠軟體產品。
19 </p>
20 <p>
21 據「加強資訊軟體人才培訓方案」的目標，將於八十八年至九十年度間依序培訓五千人次、七千五百人次、一萬人次，共計培養二萬二千五百人次。在經費方面，八十
   八年度約需新臺幣二億七千萬元，八十九年度需四億元，九十年度約五億四千萬元，養成訓練費用中的四分之三由政府負擔，四分之一由受訓學員負擔。
22 </p>
23 <p>
24 蕭萬長指出，目前由於經濟環境變遷，產業結構調整快速，正規教育體系無法迅速補充產業所需人力，短期內經由轉訓非資訊科系畢業生予以調整，是正確的手段，但
   長期而言，中高級人力正規教育的培育體系仍應與經濟發展的方向互相配合。此培訓計畫將自七月一日起開辦，蕭萬長特地指示新聞局配合宣導，以廣召甫畢業或退伍
   的青年參訓。
25 </p>
26 </text>
```
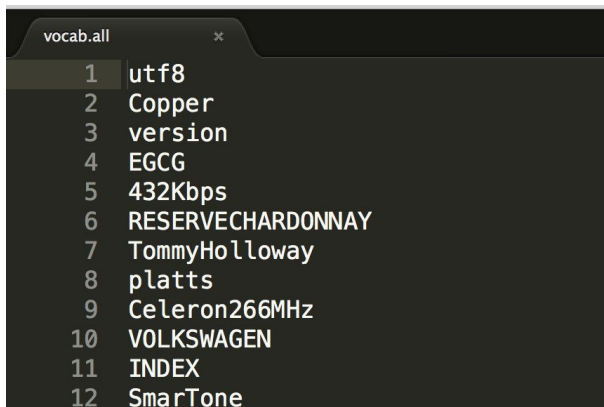
# NTCIR Document Set

- **Please** sign up the user agreement form and hand it to the TAs (at R302 of CSIE dept) in order to use this corpus.
  - Note that you'll get no point if you don't sign up the user agreement form.

- We have indexed the NTCIR documents and produced three model files for you:
  - vocab.all
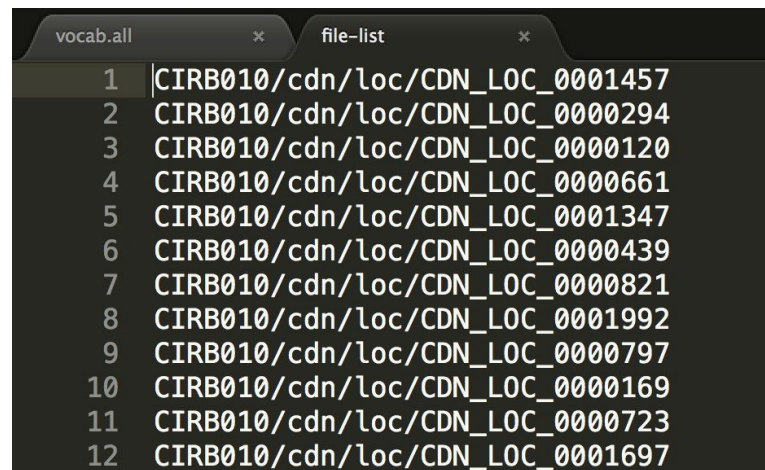  - file-list
  - inverted-file

# Format: *vocab.all*

- The file contains all vocabularies in NTCIR documents.
- The first line is character encoding format.
- Each line of the following is a vocabulary.
- Vocabularies are case-sensitive.
- Each vocabulary will have a **vocab_id** according to its line number(start from 0).
    - E.g. the **vocab_id** of **"Valentine"** is **1**; **Powell** is **2**.

```
vocab.all                    ×
    1   utf8
    2   Copper
    3   version
    4   EGCG
    5   432Kbps
    6   RESERVECHARDONNAY
    7   TommyHolloway
    8   platts
    9   Celeron266MHz
   10   VOLKSWAGEN
   11   INDEX
   12   SmarTone
```

# Format: *file-list*

- This is a list of all NTCIR documents.
- Each line denotes a document which has its line number (start from 0) as its **file_id**.
  - E.g.
    **./CIRB010/cdn/chi/cdn_chi_0001457** has **file_id 0**,
    **./CIRB010/cdn/chi/cdn_chi_0000294** has **file_id 1**.

# Format: *inverted-file*

- **vocab_id** and **file_id** referred from **vocab.all** and **file-list**.
- **vocab_id_1 vocab_id_2** denotes an unigram when **vocab_id_2**==**-1** or a bigram when **vocab_id_2**!=**-1**.
- If there are **N** files containing **vocab_id_1 vocab_id_2**, there will be the number **N** next to **vocab_id_2**, followed by **N** lines that display the counts of this term in each file.

```
vocab_id_1 vocab_id_2 number_of_file_contain_that_unigram_or_bigram
file_id_1 count_of_the_term_in_that_file
file_id_2 count_of_the_term_in_that_file
file_id_3 count_of_the_term_in_that_file

...
vocab_id_3 vocab_id_4 number_of_file_contain_that_unigram_or_bigram

...
```

# *inverted-file*



```
vocab.all              ×    inverted-file              ×

  1    1 -1 2
  2    33689 1
  3    38365 1
  4    2 -1 1
  5    33256 1
  6    2 12371 1
  7    33256 1
  8    3 -1 1
  9    10849 2
 10    3 6756 1
 11    10849 1
 12    3 6850 1
 13    10849 1
 14    4 -1 1
 15    33320 1
 16    4 5600 1
 17    33320 1
 18    5 -1 1
 19    32656 1
 20    5 12374 1
 21    32656 1
 22    6 -1 1
 23    10346 1
```

# Query File Format

- **The** NTCIR topic format conforms to XML 1.0
- The file contains multiple topics, each of them is enclosed in a <topic> tag. In each topic, different types of information are specified by the following tags:
    - <number>: The topic number.
    - <title>: The topic title.
    - <question>: A short description about the query topic.
    - <narrative>: Even more verbose descriptions about the topic.
    - <concepts>: A set of keywords that can be used in retrieval about the topic.

- <span style="color:red">You have to retrieve several relevant documents for each topic.</span>
- All the content of **title**, **question**, **narrative**, and **concepts** can be used as the query of the topic, it's **your own choice** to decide which part(s) you want to use.

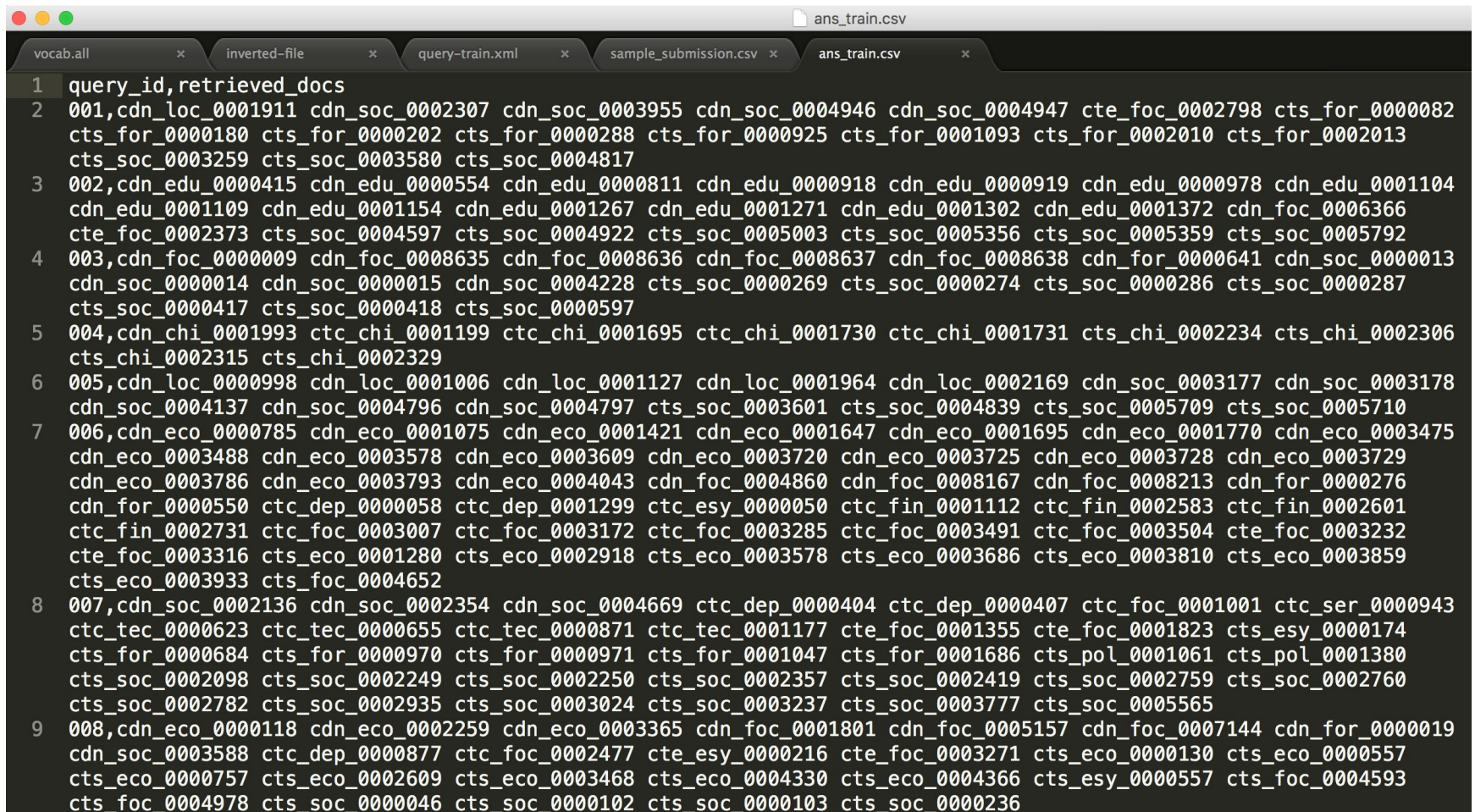# Query File Format

```
1    <?xml version="1.0" encoding="UTF-8"?>
2    <xml>
3    <topic>
4    <number>CIRB010TopicZH001</number>
5    <title>集會遊行法與言論自由</title>
6    <question>
7    查詢集會遊行法中有關主張共產主義或分裂國土規定之修正與討論。
8    </question>
9    <narrative>
10   相關文件內容應敘述集會遊行法原本對主張共產主義或分裂國土之限制，其是否符合憲法中對言論自由等基本人權的保障，大法官對此議題的相關解釋，學者專
     討論與看法，以及集會遊行法條文的修改現況。
11   </narrative>
12   <concepts>
13   集會遊行法、集會遊行、集遊法、憲法、言論自由、保障、共產主義、分裂國土、大法官會議、立法、修正條文。
14   </concepts>
15   </topic>
16
17   <topic>
18   <number>CIRB010TopicZH002</number>
19   <title>新三不政策與台獨</title>
20   <question>
21   查詢美國所提出之對台新三不政策對民進黨台獨主張的影響。
22   </question>
23   <narrative>
24   １９９８年柯江會談中，美方非正式發表了對台新三不政策：不支持台獨、不支持一中一台的主張，不支持台灣以主權國家身分加入國際組織。相關文件中應論
     對民進黨台獨主張的衝擊，包括民進黨內人士對此的反應、看法，該黨人士所發表的相關言論，所引起的黨內爭議及採取的因應策略等。
25   </narrative>
26   <concepts>
27   新三不政策、柯江會談、民進黨、建國黨、許信良、林義雄、施明德、台獨、統獨、獨派、兩岸、中國、美國、中共、大陸、台獨黨綱、主權獨立、轉型。
28   </concepts>
29   </topic>
```

# Ranking List Format

- **Each** line contains a topic and a document relevant to this topic.
- First column: **topic number**, which is the last three digits in <number>…</number> tag in the query xml file.
- Second column: **document_id**, which is the string in <id>…</id> tag in the NTCIR document. Please note it should be in lowercase.
- The two columns should be separated by a space.
- There will be several lines for the same topic if there are many.
- Note that the lines of a same topic should be sorted by their ranks, which means that, take topic **001** for example, **cdn_foc_0004185** should be more relevant to **001** than **cdn_foc_0004196**.
- You can return at most 100 documents for each topic.

# Ranking List Format

# Program IO

- Your program is required to support input of a **query file**, and output a **ranking list**.
- We provide 30 query files for you as inputs.(We provide only 10 answers, 10 queries for public scoreboard and 10 queries for private scoreboard)

- There is no restriction to the programming language you use, but make sure your program is **executable on R217 workstation**.

- Using the third party tools directly for VSM or Relevance Feedback is prohibited.

# Program Execution Details

- You are given two shell scripts to compile and run your program.

- You should edit these two scripts according to how you implement this assignment.

- When testing your program, we will execute the following commands **on R217 workstation**, please make sure your program is executable on the workstation.
  - $./compile.sh
  - $./execute.sh -option1 value1 -option2 value2...

# Program Execution Details

- Here are the required options that must be supported by your program. (Options without default values are guaranteed to be specified when we test your program.)

```
vocab.all          ×    inverted-file      ×    query-train.xml     ×    sample_submission.csv  ×    ans_train.csv      ×    SYNOPSIS:           ●
1  SYNOPSIS:
2      ./execute.sh [-r] [-b] -i query-file -o ranked-list -m model-dir -d NTCIR-dir
3
4  OPTIONS:
5      -r
6          If specified, turn on the relevance feedback on your program.
7
8      -b
9          If specified, run your best version of your program.
10
11     -i  query-file
12         The input query file.
13
14     -o  ranked-list
15         The output ranked list file.
16
17     -m model-dir
18         The input model directory, which includes three files:
19             model-dir/vocab.all
20             model-dir/file-list
21             model-dir/inverted-file
22
23     -d NTCIR-dir
24         The directory of NTCIR documents, which is the path name of CIRB010 directory.
25         Ex. If the directory's pathname is /tmp2/CIRB010, it will be "-d /tmp2/CIRB010".
```

# Evaluation

- We will use the **Mean Average Precision (MAP)** value to evaluate your ranking list.

- We provide an answer **ranking list** for **query-train.xml**.
  - There're two columns in the answer list, first is the **topic number**, and second is the **document_id** relevant to this topic.
  - You can use this answer list to check your system's performance.

# Report (8%)

- Please write your report in a **Report.pdf** and put it into the submission zip. The report should contain the following content:

  - (2%) Describe your VSM (e.g., parameter….)

  - (2%) Describe your Rocchio Relevance Feedback (e.g., how do you define relevant documents, parameter…)

  - (3%) **Results of Experiments**

    - **MAP value under different parameters of VSM**

    - **Feedback vs. no Feedback**

    - Other experiments you tried

  - (1%) Discussion: what you learn in the homework.

# Submission

- Please put report ,scripts and code into the directory named your **student ID**. Package this folder into a <span style="color:red">zip</span> file and submit it to CEIBA, following is the structure and content of the zip:

- For example: R05922XXX.zip

  +---R05922XXX(directory)

     +---**report.pdf**
     +---**compile.sh**

     +---**execute.sh**

     +---<span style="color:blue">source/</span>

     <span style="color:red">(Note that you don't have to submit the model files and NTCIR documents</span>)

# Scoring(15 points)

- 2% for VSM model.

- 2% for Rocchio relevance feedback.

- 8% for your report.

- 3% for performance (2% for simple baseline, 1% for strong baseline on public scoreboard)

- Note that you'll get 0 point if
  1. you don't have any record on kaggle
  2. you don't sign up the user_agreement_form
  3. you get caught cheating by kaggle
  4. TA can't reproduce your work

# Bonus

- Top-3 ranking at <span style="color:red">public scoreboard</span>.

  - 1% for $1^{st}$-$3^{rd}$

- Top-10 ranking at <span style="color:red">private scoreboard</span>.

  - 3% for $1^{st}$-$3^{rd}$

  - 2% for $4^{th}$-$5^{th}$

  - 1% for $6^{th}$-$10^{th}$

# Bonus package limitation

1. Any external data is allowed except for ground truth ! (pre-trained word2vec model)
2. Any package is allowed except for search engine package !
3. You need to describe your method in your report !

# Rules

- Kaggle:
  - display name: 學號_ID (Ex: r05922032_嘿嘿)
  - 5 times submissions a day
  - 2 entries for private scoreboard submissions
- Deadline:
  - Kaggle:    2018/04/13 08:00:00
  - Ceiba:     2018/04/14 23:59:59
- Late policy: 10% per day
- Email to TAs: **ntucsiewm@gmail.com**