

資料工程 Final Project

404410030 資工四 鄭光宇

系統需求

要執行這支程式，系統必須具備：

- gcc, g++
- make
- apache2.0
- php7.0
- php7.0-mbstring

編譯

```
1 | make
```

執行

先把 ettoday 資料放在 makefile 同目錄下 data 資料夾

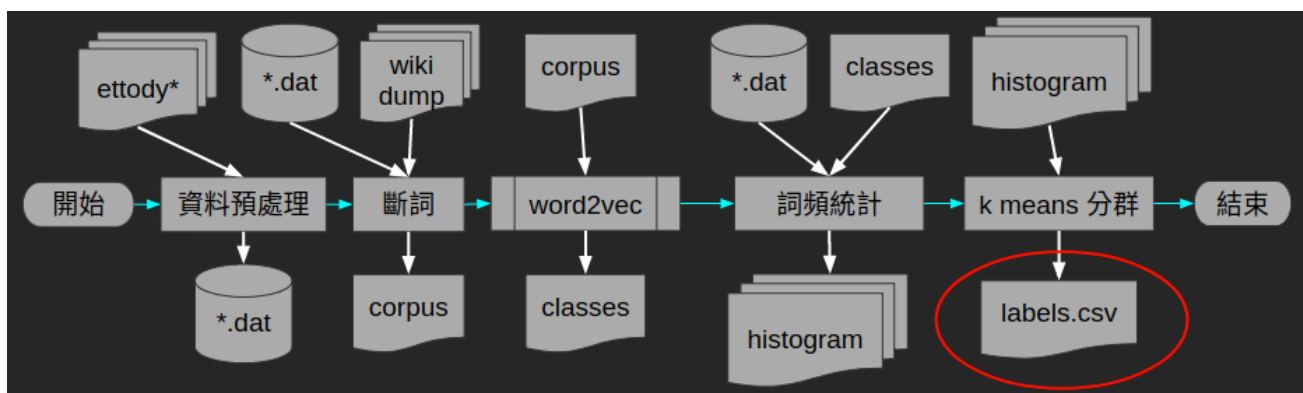
wiki dump 放在 wiki_data 資料夾

```
1 | make # 編譯
2 | make data_cleaning # 對 ettoday, wiki dump資料清理
3 | make run_preprocessing # 訓練詞向量，得到 100 類相似詞
4 | make compile_kmeans # 編譯 kmeans 程式
5 | ./word_count # 得到每篇新聞內文詞頻
6 | ./feature_extractor # 產生 kmeans 訓練資料
7 | # 分成 13 群，容忍值 1e-4，最大迭代次數 200 次
8 | # 用不同初始值跑 3 次，選較好的 centroids
9 | ../kmeans/kmeans ../.db/word_count/features.bin 13 1e-4 200 3
```

詞類別會被存放在 `../.db/word2vec/classes.txt`

文章分群類別會被存放在 `./labels.csv`

全部文章的 histogram 存放在 `../.db/word_count/features.bin`

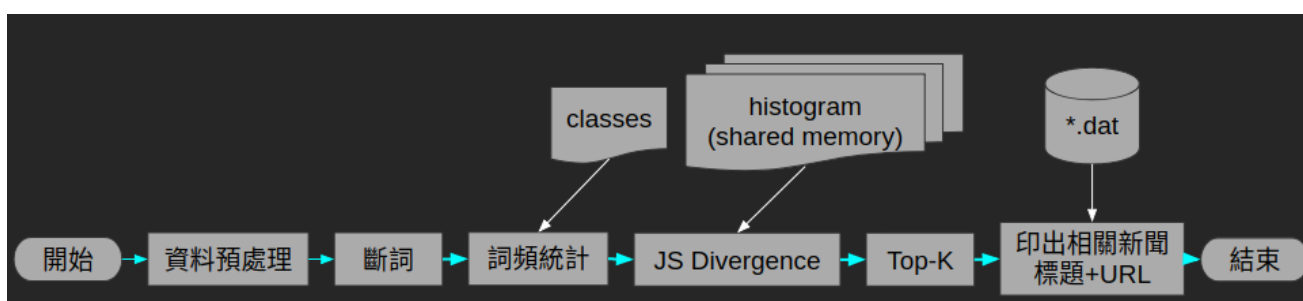


把專案目錄軟連結到 Apache 能識別到的地方，

之後執行 `./loadData2Shm`

將文章 histogram 的 binary 載入 shared memory，方便快速查詢

再使用瀏覽器開啟 `index.php`，輸入文章查詢類似文章。



資料

使用 `ettoday0~ettoday5`

並且以：

，。；：「」等標點符號斷句。

實作細節

專案目錄結構

```

1 FINAL_PJ
2 |— closeShm.c
3 |— compute_copy_range.cpp
4 |— data_cleaning.sh
5 |— feature_extractor.cpp
6 |— find_news_by_id.cpp
7 |— flat_hash_map.hpp
8 |— gen_corpus.cpp
9 |— gen_wiki_corpus.cpp
10 |— index.php
11 |— jieba_word_count.hpp

```

```
12 | ├── kmeans
13 | |   ├── kmeans.c
14 | |   ├── kmeans.h
15 | |   ├── main.c
16 | |   └── makefile
17 | ├── loadData2Shm.c
18 | ├── makefile
19 | ├── ngram_word_count.hpp
20 | ├── parser.c
21 | ├── put_data_here
22 | ├── read_histogram.hpp
23 | ├── record_structure.h
24 | ├── record_structure_io.c
25 | ├── retrieve_topN.cpp
26 | ├── Trie.hpp.patch
27 | ├── word_count.cpp
28 | ├── wstringcv.h
29 | ├── 資料工程 Final Project.md
30 | └── 資料工程 Final Project.pdf
31 |
32 | 1 directory, 28 files
```

實作細節

除了 word2vec、斷詞以外，其他都自己使用 C/C++ 實作，

盡量使用 OpenMP 平行化。

網頁部份使用 Apache2.0 + PHP7.0，

呼叫編譯好的 C++ 程式做新聞預處理、比較、查詢。

Word2Vec

使用 Google 提供的 C 版本 word2vec，800 多行的程式碼，很簡潔很強大。

<https://github.com/tmikolov/word2vec>

斷詞

使用結巴 (jieba) 斷詞，這裡使用 runtime 較高效的 `cppjieba`

<https://github.com/yanyiwu/cppjieba>

histogram 差異計算

使用 Jensen-Shannon divergence 的平方計算兩 histogram 分佈差異

若要對文章長度變化不敏感，只在乎詞頻的相對關係，可以選擇對詞頻 histogram normalize。

```
1 double hist_intersection(const float *P, const float *Q, unsigned int cols) {
2     double P_M = 0.0;
3     double Q_M = 0.0;
4     double JSD = 0.0;
```

```

5     for (unsigned int i=0; i<cols; ++i) {
6         float M_i = (P[i]+Q[i]) / 2.0 + 1e-6;
7         P_M += P[i]*log((P[i]+1e-6) / M_i);
8         Q_M += Q[i]*log((Q[i]+1e-6) / M_i);
9     }
10    JSD = (P_M+Q_M) / 2.0;
11    return JSD*JSD;
12 }
13
14 double hist_intersection_normalized(const float *P, const float *Q, unsigned int
cols) {
15     double P_M = 0.0;
16     double Q_M = 0.0;
17     double JSD = 0.0;
18     double P_S = 1e-6;
19     double Q_S = 1e-6;
20     for (unsigned int i=0; i<cols; ++i) {
21         P_S += P[i];
22         Q_S += Q[i];
23     }
24     for (unsigned int i=0; i<cols; ++i) {
25         float p=P[i]/P_S, q=Q[i]/Q_S;
26         float M_i = (p+q) / 2.0 + 1e-6;
27         P_M += p*log((p+1e-6) / M_i);
28         Q_M += q*log((q+1e-6) / M_i);
29     }
30     JSD = (P_M+Q_M) / 2.0;
31     return JSD*JSD;
32 }

```

K-means 程式碼節錄

1. 隨機選取 K 個初始 centroids
2. 分配每個資料點到最近的 centroid，形成 k 個 cluster（很容易平行化）
3. 計算每個 cluster 內資料點平均值，得到新的 centroids
4. 重複 2, 3 直到收斂，或達到指定迭代次數

複雜度為 $O(tNDK)$, N 為資料筆數, D 為資料維度, K 為 cluster 數, t 為迭代次數。實際上有可能跑到不好的局部最佳解，所以通常會選擇不同的初始值多跑幾次。

```

1  double kmeans_intersec_int(unsigned int **data, unsigned int **return_labels,
double ***return_centroid, int rows, int cols, int K, double tol, int max_iter,
char verbose) {
2      double mean_centroid_d = DBL_MAX;
3      double **centroids = NULL;
4      double **new_centroids = NULL;
5      unsigned int *lab_counts=NULL;
6      unsigned int iter_counter=0;
7      int *labels=NULL;
8      centroids = (double**)malloc(sizeof(double*)*K);
9      new_centroids = (double**)malloc(sizeof(double*)*K);
10     labels = (int*)malloc(sizeof(int)*rows);

```

```

11     lab_counts = (int*)malloc(sizeof(int)*K);
12
13     for (int i=0; i<K; ++i) centroids[i] = (double*)malloc(sizeof(double)*cols);
14     for (int i=0; i<K; ++i) new_centroids[i] =
15     (double*)malloc(sizeof(double)*cols);
16
17     // initialize (randomly pick k samples wo replacement)
18     for (int i=0; i<K; ++i) {
19         int l=rand()%rows;
20         // check repetition
21         char fail;
22         do {
23             fail=0;
24             for (int j=0; j<i; ++j)
25                 if (lab_counts[j]==1) {
26                     fail=1;
27                     l=rand()%rows; // pick another
28                     break;
29                 }
30             } while(fail);
31         lab_counts[i] = 1;
32         for (int j=0; j<cols; ++j) {
33             centroids[i][j] = (double)data[l][j];
34         }
35     }
36
37     iter_counter=0;
38     while(iter_counter<max_iter && mean_centroid_d>tol) {
39         // determine labels
40         #pragma omp parallel for shared(data, centroids, labels) schedule(static,1)
41         for (int i=0; i<rows; ++i) {
42             unsigned int best_l=0;
43             double min_d=DBL_MAX;
44             for (int k=0; k<K; ++k) {
45                 double d = hist_intersection(data[i], centroids[k], cols);
46                 if (d<min_d) {
47                     min_d = d;
48                     best_l = k;
49                 }
50             }
51             labels[i] = best_l;
52         }
53         // determine new centroids
54         for (int k=0; k<K; ++k) memset(new_centroids[k], 0x00,
55         sizeof(double)*cols);
56         memset(lab_counts, 0x00, sizeof(int)*K);
57         for (int i=0; i<rows; ++i) {
58             unsigned int l = labels[i];
59             ++lab_counts[l];
60             for (int j=0; j<cols; ++j) {
61                 new_centroids[l][j] += data[i][j]; // sum

```

```

62     for (int k=0; k<K; ++k) {
63         for (int j=0; j<cols; ++j) new_centroids[k][j] /= (double)
(lab_counts[k]+1e-8); // mean
64     }
65     mean_centroid_d = 0;
66     for (int k=0; k<K; ++k) {
67         mean_centroid_d += hist_intersection_f(centroids[k], new_centroids[k],
cols);
68     }
69     mean_centroid_d /= (double)K;
70     // assign new centroids to centroids
71     for (int k=0; k<K; ++k) {
72         double *ptr = centroids[k];
73         centroids[k] = new_centroids[k];
74         new_centroids[k] = ptr;
75     }
76     ++iter_counter;
77     if (verbose) fprintf(stderr, "[%d/%d] d:%.4lf\n", iter_counter, max_iter,
mean_centroid_d);
78 }
79
80 double *intra_distance = NULL;
81 intra_distance = (double*)malloc(sizeof(double)*K);
82 memset(intra_distance, 0x00, sizeof(double)*K);
83
84 for (int i=0; i<rows; ++i) {
85     unsigned int k = labels[i];
86     intra_distance[k] += hist_intersection(data[i], centroids[k], cols);
87 }
88
89 double mean_intra_distance = 0.0;
90 double max_intra_distance = 0.0;
91 for (int k=0; k<K; ++k) {
92     if (intra_distance[k]>max_intra_distance) max_intra_distance =
intra_distance[k];
93     mean_intra_distance += intra_distance[k];
94     intra_distance[k] /= (double)(lab_counts[k]+1e-8);
95 }
96 mean_intra_distance /= (double)rows;
97 if (verbose) {
98     fprintf(stderr, "mean data-centroid distance:\n");
99     for (int k=0; k<K; ++k) {
100         fprintf(stderr, "%3d: %.4f\n", k, intra_distance[k]);
101     }
102     fprintf(stderr, "average: %.4f\n", mean_intra_distance);
103     fprintf(stderr, "maximum: %.4f\n", max_intra_distance);
104     fprintf(stderr, "each class count:\n");
105     for (int k=0; k<K; ++k) {
106         fprintf(stderr, "%3d: %10d\n", k, lab_counts[k]);
107     }
108 }
109 free(intra_distance); intra_distance=NULL;
110 if (verbose) {

```

```

111     fprintf(stderr, "centroid-centroid distance:\n");
112     fprintf(stderr, "          ");
113     for (int i=0; i<K; ++i) fprintf(stderr, "%12d", i);
114     fprintf(stderr, "\n");
115     for (int i=0; i<K-1; ++i) {
116         // print upper triangle matrix
117         fprintf(stderr, "%12d", i);
118         for (int j=0; j<=i; ++j) fprintf(stderr, "          ");
119         for (int j=i+1; j<K; ++j) fprintf(stderr, "%12.3f",
hist_intersection_f(centroids[i], centroids[j], cols));
120         fprintf(stderr, "\n");
121     }
122     fprintf(stderr, "\n");
123 }
124
125 for (int k=0; k<K; ++k) free(new_centroids[k]);
126
127 free(lab_counts);
128 free(new_centroids);
129
130 if(return_centroid!=NULL) *return_centroid = centroids;
131 else free(centroids);
132 if(return_labels!=NULL)  *return_labels  = labels;
133 else free(labels);
134
135 return mean_intra_distance;
136 }

```

Top K

1. 假設資料為 $A[N]$ ，要找到 top K 個最小元素
2. 初始化容量為 K 的 max-heap H
3. push $A[0...K-1]$ -> H
4. 依序檢查 $A[K...N-1]$ 中的資料 a_i , $i: K \rightarrow N-1$
5. 如果 $a_i < H.top$, $H.pop()$, $H.push(a_i)$
6. 最後 H 中會留下 top K 個最小元素

複雜度大約 $O(N \log K)$

```

1  for (int i=0; i<topN; ++i) max_heap.push({distances[i], i});
2  for (unsigned int i=topN; i<n_rows; ++i) {
3      if (distances[i]<max_heap.top().first) {
4          max_heap.pop();
5          max_heap.push({distances[i], i});
6      }
7  }
8  topN = max_heap.size();
9  for (int i=topN-1; i>=0; --i) {
10     topN_id[i] = max_heap.top().second;
11     max_heap.pop();
12 }

```

Copy-append model

有 A, B 兩個檔案，我們要找到一連串操作 δ ，使得：

$A + \delta \rightarrow B$

δ 允許的操作有：從A複製貼上、加上新內容

使用 Copy-append model 去偵測兩篇文章整段相同的部份

虛擬碼：

```
1  k=4 # or 5 or 6, 7, 8, ...
2  # S 為 A 的 k-gram index, N 為 B 的 string length
3  i=0
4  while i<N:
5      b = B[i:min(i+K, N)] # B 的 K-gram
6      if b not in S: # binary search / hash table
7          put(append,b)
8          i+=b.length
9      else:
10         (l,m) = longest_match(A,i,S,B)
11         put(copy,l,m)
12         i+=l
```

C++ 實現，longest match 的部份主要使用 Z function

還沒有與 K-gram index + Boyer Moore 的方法比較速度

但實際跑起來蠻快的

```
1  inline std::pair<int,int> find_longest_match(const unsigned char *fileA, const
std::vector<int> &shortcut_index, const unsigned char *fileB, int A_len, int B_len,
int *z) {
2      using std::min;
3      using std::max;
4      fileA += shortcut_index[0];
5      A_len -= shortcut_index[0];
6      memset(z, 0x00, sizeof(int)*(A_len+B_len));
7      #define s(i) (i<B_len?fileB[i]:fileA[i-B_len])
8      // From Eddy's codebook:
9      int l=0, r=0;
10     z[0]=A_len+B_len;
11     for (int i=1; i<A_len+B_len; ++i) {
12         int j = max(min(z[i-1],r-i),0);
13         while(i+j<A_len+B_len&&s(i+j)==s(j)) ++j;
14         z[i] = j;
15         if (i+z[i]>r) r=i+z[i], l=i;
16     }
17     #undef s
18     int M=B_len;
19     int L=z[M];
20     for (auto v : shortcut_index) {
21         int m = v+B_len-shortcut_index[0];
```



```

22     int l = z[m];
23     if (l>L) {
24         L = l;
25         M = m;
26     }
27 }
28 if (L>B_len) L = B_len; // Although this happen, the value of M is correct.
29 M = M-B_len+shortcut_index[0];
30 return {M,L};
31 }
32
33 std::string copy_append_encoder(const unsigned char *fileA, const unsigned char
                                *fileB, int A_len, int B_len, unsigned int K_size)
34 {
35     using namespace std;
36     int B_index=0;
37     ska::flat_hash_map<string, vector<int> > kgramA;
38     string delta;
39     int *z_buffer = new int[A_len+B_len+5];
40
41     /* Generate Kgram index for fileA */
42     for (int i=0; i<A_len; ++i) {
43         int l = min((int)(i+K_size), A_len) - i;
44         string ss((char*)(fileA+i), l);
45         if(kgramA.count(ss)==0) kgramA.insert({ss, vector<int>(1, i)});
46         else kgramA[ss].push_back(i);
47     }
48
49     string lastKgramB;
50     int last_msg_length = 0;
51
52     while (B_index<B_len) {
53         int l = min((int)(B_index+K_size), B_len) - B_index;
54         int m;
55         string kgramB((char*)(fileB+B_index), l);
56         if (kgramA.count(kgramB)==0) {
57             string append_msg;
58             if (last_msg_length>0) delta.resize(delta.length()-last_msg_length);
59             lastKgramB += kgramB;
60             append_msg += "a " + to_string(lastKgramB.length());
61             append_msg += "\n" + lastKgramB;
62             delta += append_msg;
63             last_msg_length = append_msg.length();
64         } else {
65             /* Find longest match between A and B[B_index:] */
66             pair<int,int> ML = find_longest_match(fileA, kgramA[kgramB],
67             fileB+B_index, A_len, B_len-B_index, z_buffer);
68             m = ML.first;
69             l = ML.second;
70             string copy_msg;
71             copy_msg += "c " + to_string(m);
72             copy_msg += ", " + to_string(l);
73             copy_msg += "\n";

```

```

72         delta += copy_msg;
73         lastKgramB.clear();
74         last_msg_length = 0;
75     }
76     B_index += 1;
77 }
78 lastKgramB.clear();
79 if (z_buffer!=NULL) delete[] z_buffer; z_buffer=NULL;
80 return delta;
81 }

```

實作上將連續的 append 操作變成同一個操作，避免浪費字元數。

考慮到使用者可能比較在意文章「出自」那一些段落，所以在伺服器上改為只偵測複製的操作，並找出複製的 unique 區段。

```

1  std::vector< std::pair<int,int> > copy_detect(const unsigned char *fileA, const
    unsigned char *fileB, int A_len, int B_len, unsigned int K_size) {
2      using namespace std;
3      int B_index=0;
4      ska::flat_hash_map<string, vector<int> > kgramA;
5      vector< pair<int,int> > ret;
6      int *z_buffer = new int[A_len+B_len+5];
7
8      /* Generate Kgram index for fileA */
9      for (int i=0; i<A_len; ++i) {
10         int l = min((int)(i+K_size), A_len) - i;
11         string ss((char*)(fileA+i), l);
12         if(kgramA.count(ss)==0) kgramA.insert({ss, vector<int>(1, i)});
13         else kgramA[ss].push_back(i);
14     }
15
16     while (B_index<B_len) {
17         int l = min((int)(B_index+K_size), B_len) - B_index;
18         int m;
19         string kgramB((char*)(fileB+B_index), l);
20         if (kgramA.count(kgramB)>0) {
21             /* Find longest match between A and B[B_index:] */
22             pair<int,int> ML = find_longest_match(fileA, kgramA[kgramB],
fileB+B_index, A_len, B_len-B_index, z_buffer);
23             m = ML.first;
24             l = ML.second;
25             ret.push_back({m,m+l}); // [, )
26         }
27         B_index += 1;
28     }
29     if (z_buffer!=NULL) delete[] z_buffer; z_buffer=NULL;
30     return ret;
31 }
32
33 inline std::vector< std::pair<int,int> > summary_ranges(std::vector<
    std::pair<int,int> > &ranges) {
34     using namespace std;

```

```
35     vector< pair<int,int> > ret;
36     sort(ranges.begin(), ranges.end());
37     ranges.push_back({INT_MAX,INT_MAX}); // 哨兵
38     int l=ranges[0].first, r=ranges[0].second;
39     for (int i=1; i<ranges.size(); ++i) {
40         if (ranges[i].first<=r) r = max(r, ranges[i].second);
41         else {
42             ret.push_back({l,r});
43             l = ranges[i].first;
44             r = ranges[i].second;
45         }
46     }
47     ranges.pop_back();
48     return ret;
49 }
```

實驗結果

實驗環境：

- Ubuntu 16.04 64-bit
- 32GB DDR4 (2400MHz)
- AMD Ryzen 5 1600 (6C12T)
- 5400rpm HDD

分群效果

資料距離群心的平均距離 (JSD^2)

	data-centroid	count
0	3152.1657	57238
1	1422.4215	42418
2	1422.4490	22251
3	111797.1769	157
4	1189.1763	85518
5	2667.7802	53586
6	9599.0289	2698
7	12069.3628	5496
8	4196.4556	9458
9	99240.8710	107
10	13561.2694	1829
11	10511.0250	5945
12	2159.9833	243577

群心間距離

centroid-centroid distance:												
	0	1	2	3	4	5	6	7	8	9	10	11
0		1726.177	2635.470	574997.263	2847.557	779.596	9773.077	29298.529	13485.311	994810.870	15664.413	13011.640
1			230.995	714038.255	4045.732	1718.443	26484.825	36188.242	22616.659	1063529.988	31050.319	20550.180
2				619361.984	7325.204	1875.415	20703.235	25512.895	25759.213	927342.939	25631.757	10833.333
3					867918.455	514452.040	171303.461	134522.862	583342.601	102481.399	166575.413	294752.677
4						5195.936	43592.105	79396.880	4924.774	1358128.557	55902.481	44777.198
5							10833.258	16775.143	18346.265	873512.894	21049.818	13267.452
6								7570.148	32179.436	478678.231	7548.072	13524.838
7									81353.877	204036.603	10266.032	4717.326
8										1129539.647	53287.469	60263.092
9											459350.232	422082.587
10												9483.434
11												

觀察

肉眼能觀察出來的大概有以下幾類：

財經新聞 (0)

體育新聞 (1,2)

長文 (3)

股市相關公告-短文 (4)

股市相關公告-長文 (8)

政治新聞 (5)

甄嬛傳系列文章/原創長文 (9)

3C新聞 (10)

*(num) 代表某個 cluster

有些 cluster 特別容易混淆，例如 財經新聞(0) 與 政治新聞 (5)

有些離群值，其群心與其他類別隔非常遠，肉眼閱讀文章也能發現用詞有極大差異，例如 (9)

而有些 cluster 應該處於同一類，例如 (1) (2) 同樣是體育新聞，可能是因為 cluster 數量設太大。

文章查詢效果

記者洪聖壹／西班牙巴塞隆納報導

BlackBerry 今天（2/25）上午在 MWC 2014 會場當中宣布與鴻海集團結盟，確認將於今年四月份，在印尼推出首款合作新機 BlackBerry Z3，未來將由鴻海集團負責硬體生產，黑莓則會把資源放在軟件研發，雙方還將推出代號為「Classic」的 Q20 黑莓機。

BlackBerry 執行長程守宗與鴻海董事長郭台銘兩人今天共同宣布合作成果，並發表了 BlackBerry Z3，配置了全觸控螢幕機身，搭載 BlackBerry 10.2.1作業系統，但是空機售價如同今年一月份傳出的一樣，將以不到 200 美元的價格，於四月份先在印尼推出。

顯示前幾筆： ☒ 比較文章長度

標題	差異度	重複的內文
MWC2014／黑莓、鴻海聯手推合作新機BlackBerryZ3	163.92	[1] 記者洪聖壹／西班牙巴塞隆納報導 [2] 2/25) 上午在 [3] 會場當中 [4] 集團結盟，確認將於今年四月份，在印尼推出首款合作新機 [5] Z3，未來將由鴻海集團負責硬體生產，黑莓則會把資源放在軟件研發，雙方還將推出代號為「Classic」的 [6] 行長程守宗與鴻海董事長 [7] 今天共同宣 [8] 果，並發表了 [9] Z3，配置了全觸控螢幕機身，搭載 [10] 10.2.1作業系統，但是空機售價如同今年一月份傳出的一樣，將以不到 [11] 的價格，於四月份先在印尼推出
LG第二代彎曲手機GFlex2傳將於2015年初亮相	185.72	[1] 記者洪聖壹／
蝴蝶飛向歐洲！HTC將於IFA展出8核手機與二代蝴蝶機	189.36	[1] 記者洪聖壹／
快訊／Google宣佈推出搭載原生Android4.2的GALAXYS4	204.39	[1] 記者洪聖壹／
前臺灣三星總經理杜偉昱將出任德誼數位轉戰蘋果通路	221.48	[1] 記者洪聖壹／
終於要來了！HTC自拍神器re將於12月2日登場	233.50	[1] 記者洪聖壹／

如上圖，可以輸入文章查詢類似文章，輸出前 N 筆最像的文章標題與 URL，並顯示它與欲查詢文章的 JSD² 作為差異度。最後列出與查詢的內文有整段重複的部份。

在我的電腦上大約 1~2 秒一個查詢

在系上工作站大約 3~6 秒一筆查詢

GitHub

程式碼：

https://github.com/peter0749/Data-Engineering/tree/master/FINAL_PJ

https://github.com/peter0749/Data-Engineering/tree/master/copy_append_model