

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

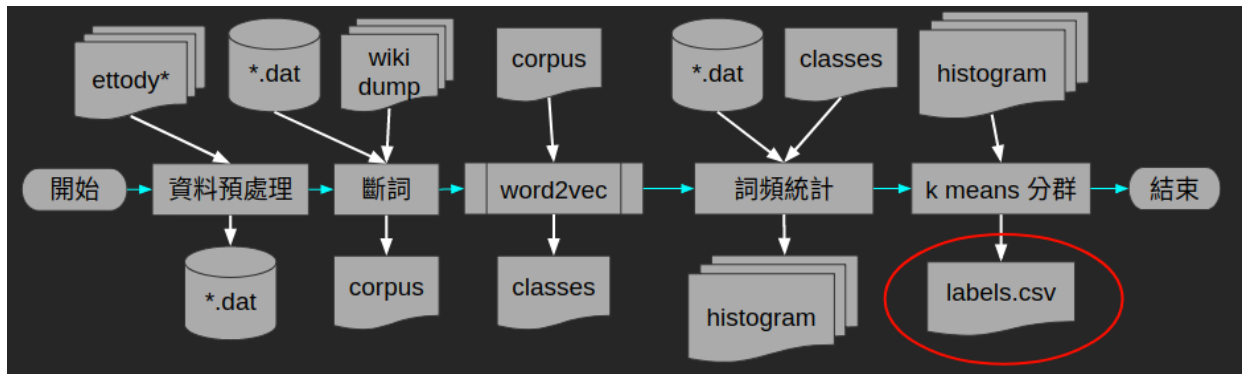
```
make
```

執行

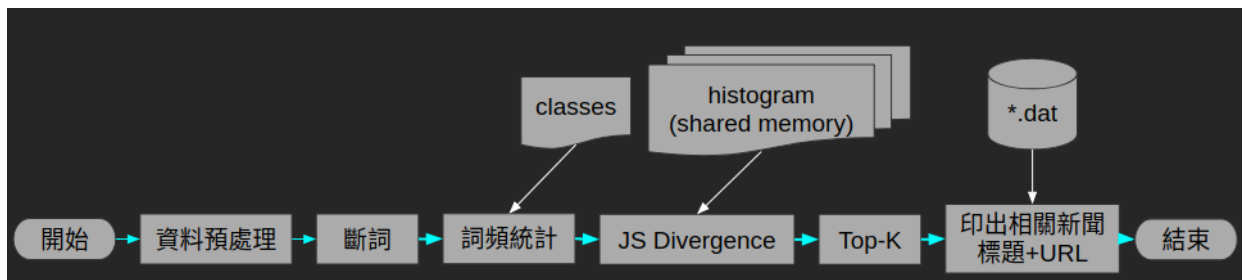
先把 ettoday 資料放在 makefile 同目錄下 data 資料夾 wiki dump 放在 wiki_data 資料夾

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

```
FINAL_PJ
├─ closeShm.c
├─ data_cleaning.sh
├─ feature_extractor.cpp
├─ find_news_by_id.cpp
├─ gen_corpus.cpp
├─ gen_wiki_corpus.cpp
├─ index.php
├─ jieba_word_count.hpp
├─ kmeans
│   └─ kmeans.c
│   └─ kmeans.h
│   └─ main.c
│   └─ makefile
├─ loadData2Shm.c
└─ makefile
```

```
├─ ngram_word_count.hpp
├─ parser.c
├─ put_data_here
├─ read_histogram.hpp
├─ record_structure.h
├─ record_structure_io.c
├─ retrieve_topN.cpp
├─ Trie.hpp.patch
├─ word_count.cpp
└─ wstringcvt.hpp
```

1 directory, 24 files

實作細節

除了 word2vec、斷詞以外，其他都自己使用 C/C++ 實作，儘量使用 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。

```
double hist_intersection(const float *P, const float *Q, unsigned int cols)
{
    double P_M = 0.0;
    double Q_M = 0.0;
    double JSD = 0.0;
    for (unsigned int i=0; i<cols; ++i) {
        float M_i = (P[i]+Q[i]) / 2.0 + 1e-6;
        P_M += P[i]*log((P[i]+1e-6) / M_i);
        Q_M += Q[i]*log((Q[i]+1e-6) / M_i);
    }
    JSD = (P_M+Q_M) / 2.0;
    return JSD*JSD;
}
```

```

double hist_intersection_normalized(const float *P, const float *Q,
unsigned int cols) {
    double P_M = 0.0;
    double Q_M = 0.0;
    double JSD = 0.0;
    double P_S = 1e-6;
    double Q_S = 1e-6;
    for (unsigned int i=0; i<cols; ++i) {
        P_S += P[i];
        Q_S += Q[i];
    }
    for (unsigned int i=0; i<cols; ++i) {
        float p=P[i]/P_S, q=Q[i]/Q_S;
        float M_i = (p+q) / 2.0 + 1e-6;
        P_M += p*log((p+1e-6) / M_i);
        Q_M += q*log((q+1e-6) / M_i);
    }
    JSD = (P_M+Q_M) / 2.0;
    return JSD*JSD;
}

```

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。

```

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) {
    double mean_centroid_d = DBL_MAX;
    double **centroids = NULL;
    double **new_centroids = NULL;
    unsigned int *lab_counts=NULL;
    unsigned int iter_counter=0;
    int *labels=NULL;
    centroids = (double**)malloc(sizeof(double*)*K);
    new_centroids = (double**)malloc(sizeof(double*)*K);
    labels = (int*)malloc(sizeof(int)*rows);
    lab_counts = (int*)malloc(sizeof(int)*K);

    for (int i=0; i<K; ++i) centroids[i] =
(double*)malloc(sizeof(double)*cols);
    for (int i=0; i<K; ++i) new_centroids[i] =
(double*)malloc(sizeof(double)*cols);

```

```

// initialize (randomly pick k samples wo replacement)
for (int i=0; i<K; ++i) {
    int l=rand()%rows;
    // check repetition
    char fail;
    do {
        fail=0;
        for (int j=0; j<i; ++j)
            if (lab_counts[j]==1) {
                fail=1;
                l=rand()%rows; // pick another
                break;
            }
    } while(fail);
    lab_counts[i] = 1;
    for (int j=0; j<cols; ++j) {
        centroids[i][j] = (double)data[l][j];
    }
}

iter_counter=0;
while(iter_counter<max_iter && mean_centroid_d>tol) {
    // determine labels
    #pragma omp parallel for shared(data, centroids, labels)
    schedule(static,1)
    for (int i=0; i<rows; ++i) {
        unsigned int best_l=0;
        double min_d=DBL_MAX;
        for (int k=0; k<K; ++k) {
            double d = hist_intersection(data[i], centroids[k], cols);
            if (d<min_d) {
                min_d = d;
                best_l = k;
            }
        }
        labels[i] = best_l;
    }
    // determine new centroids
    for (int k=0; k<K; ++k) memset(new_centroids[k], 0x00,
sizeof(double)*cols);
    memset(lab_counts, 0x00, sizeof(int)*K);
    for (int i=0; i<rows; ++i) {
        unsigned int l = labels[i];
        ++lab_counts[l];
        for (int j=0; j<cols; ++j) {
            new_centroids[l][j] += data[i][j]; // sum
        }
    }
}

```

```

        for (int k=0; k<K; ++k) {
            for (int j=0; j<cols; ++j) new_centroids[k][j] /= (double)
(lab_counts[k]+1e-8); // mean
        }
        mean_centroid_d = 0;
        for (int k=0; k<K; ++k) {
            mean_centroid_d += hist_intersection_f(centroids[k],
new_centroids[k], cols);
        }
        mean_centroid_d /= (double)K;
        // assign new centroids to centroids
        for (int k=0; k<K; ++k) {
            double *ptr = centroids[k];
            centroids[k] = new_centroids[k];
            new_centroids[k] = ptr;
        }
        ++iter_counter;
        if (verbose) fprintf(stderr, "[%d/%d] d:%.4lf\n", iter_counter,
max_iter, mean_centroid_d);
    }

    double *intra_distance = NULL;
    intra_distance = (double*)malloc(sizeof(double)*K);
    memset(intra_distance, 0x00, sizeof(double)*K);

    for (int i=0; i<rows; ++i) {
        unsigned int k = labels[i];
        intra_distance[k] += hist_intersection(data[i], centroids[k],
cols);
    }

    double mean_intra_distance = 0.0;
    double max_intra_distance = 0.0;
    for (int k=0; k<K; ++k) {
        if (intra_distance[k]>max_intra_distance) max_intra_distance =
intra_distance[k];
        mean_intra_distance += intra_distance[k];
        intra_distance[k] /= (double)(lab_counts[k]+1e-8);
    }
    mean_intra_distance /= (double)rows;
    if (verbose) {
        fprintf(stderr, "mean data-centroid distance:\n");
        for (int k=0; k<K; ++k) {
            fprintf(stderr, "%3d: %.4f\n", k, intra_distance[k]);
        }
        fprintf(stderr, "average: %.4f\n", mean_intra_distance);
        fprintf(stderr, "maximum: %.4f\n", max_intra_distance);
        fprintf(stderr, "each class count:\n");
        for (int k=0; k<K; ++k) {

```

```

        fprintf(stderr, "%3d: %10d\n", k, lab_counts[k]);
    }
}
free(intra_distance); intra_distance=NULL;
if (verbose) {
    fprintf(stderr, "centroid-centroid distance:\n");
    fprintf(stderr, "          ");
    for (int i=0; i<K; ++i) fprintf(stderr, "%12d", i);
    fprintf(stderr, "\n");
    for (int i=0; i<K-1; ++i) {
        // print upper triangle matrix
        fprintf(stderr, "%12d", i);
        for (int j=0; j<=i; ++j) fprintf(stderr, "          ");
        for (int j=i+1; j<K; ++j) fprintf(stderr, "%12.3f",
hist_intersection_f(centroids[i], centroids[j], cols));
        fprintf(stderr, "\n");
    }
    fprintf(stderr, "\n");
}

for (int k=0; k<K; ++k) free(new_centroids[k]);

free(lab_counts);
free(new_centroids);

if(return_centroid!=NULL) *return_centroid = centroids;
else free(centroids);
if(return_labels!=NULL) *return_labels = labels;
else free(labels);

return mean_intra_distance;
}

```

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. $push\ A[0...K-1] \rightarrow H$
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 $O(N\log K)$

```
for (int i=0; i<topN; ++i) max_heap.push({distances[i], i});
for (unsigned int i=topN; i<n_rows; ++i) {
    if (distances[i]<max_heap.top().first) {
        max_heap.pop();
        max_heap.push({distances[i], i});
    }
}
topN = max_heap.size();
for (int i=topN-1; i>=0; --i) {
    topN_id[i] = max_heap.top().second;
    max_heap.pop();
}
```

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	12
0		1726.177	2635.470	574997.263	2847.557	779.596	9773.077	29298.529	13485.311	994810.870	15664.413	13011.640	885.770
1			230.995	714038.255	4045.732	1718.443	26484.825	36188.242	22616.659	1063529.988	31050.319	20550.180	951.453
2				619361.984	7325.204	1875.415	20703.235	25512.895	25759.213	927342.939	25631.757	16833.333	1290.238
3					867918.455	514452.040	171383.461	134522.862	583342.691	102481.399	166575.413	294752.677	659091.197
4						5195.936	43592.105	79396.880	4924.774	1358128.557	58902.481	44777.198	4659.396
5							10833.258	16775.143	18346.265	873512.894	21049.818	13267.452	491.653
6								7570.148	32179.436	478678.231	7548.072	13524.838	21390.045
7									81353.877	204036.603	10266.032	4717.326	22604.782
8										1129539.647	53287.469	60263.092	25981.084
9											459350.232	422082.587	950550.873
10												9483.434	24666.488
11													10903.596

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

2連敗。(圖／達志影像／美聯社)最近10場比賽拿下7勝、狀況回穩的塞爾提克，與公牛的比賽在上半場結束就以49比38領先。易籃再戰，公牛吹起反攻的號角，第三節不只將落後的比數追平甚至超前，塞爾提克浪費前兩節11分的領先，到了第四節差距更逐漸被拉大，無力追趕的波士頓，終場就以7分差輸掉比賽。全場得到最高26分的丹恩光是下半場就砍下18分，是率領球隊逆轉戰局的大功臣，諾亞(JoakimNoah)則繳出17分9籃板的成績，布瑟(CarlosBoozer)12分14籃板有雙十表現。反觀塞爾提克除了皮爾斯(PaulPierce)有22分8籃板4助攻，賈奈特(KevinGarnett)12分14籃板和朗多(RajonRondo)10分12助攻都有雙十演出外，全隊共有5人得分上雙，但最後仍不敵下半場大爆發的公牛，皮爾斯在決勝節更只靠罰球拿下1分。風城當家球星羅斯(DerrickRose)今天連續第12場缺賽，MVP缺陣期間球隊仍得到8勝4敗的成績，總教練席波杜(TomThibodeau)認為「我們應該專注於我們現有的。我們球隊不管少掉哪個人，都還是會維持高昂的鬥志，我們相信自己能贏，整個球季都會如此。」

顯示前幾筆： ☒ 比較文章長度

Go!

標題	差異度
NBA／丹恩下半場大爆發 率公牛撞倒綠衫軍逆轉戰局	0.00
NBA／「火車男孩」拜能不受傷勢影響 率湖人撞翻快艇	138.82
NBA季後賽／手感冰冷 綠衫軍遭逼進第七戰	154.15
NBA／安東尼負傷上陣 絕殺三分巫師無計可施	159.80
NBA／巫師末節逆襲 湖人苦吞二連敗	165.14
NBA／塞爾提克踰走山貓 排名追上76人並列第一	178.69
NBA季後賽／要命進攻犯規 綠衫軍系列賽遭扳平	186.64
NBA／「西河」掛傷號 「雷帝」杜蘭特撐大局砍34分	194.00

搜尋時間：1.05 秒

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 JSD^2 作為差異度。

在我的電腦上大約 1~2 秒一個查詢 在系上工作站也許是因為 libgomp 遺失了，只能編譯單執行緒程式，大約 6~7 秒一筆查詢

GitHub

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