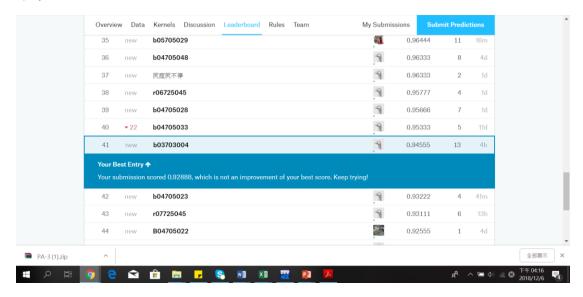
B03703004 財金五陳冠宇

Introduction to Information Retrieval HW03

執行環境: Dev C++ 5.11

使用語言: C++

結果:



參數選擇:

Feature < 500

Chi-square > 20

MI>1.5

每篇總字數最多 45

P(t=1,c=1)>6 or P(t=1,c=0)=0

((P(t=1,c=1)-E[P(t=1,c=1)]) + (P(t=0,c=0)-E[P(t=0,c=0)]) - (P(t=0,c=1)-E[P(t=0,c=1)]) - (P(t=1,c=0)-E[P(t=1,c=0)]) > 22

→number of features 472, score = 0.92666

Feature>500

Chi-square > 0

MI>1.5

P(t=1,c=1)>1

 \rightarrow number of features = 3342, score = 0.94555

各函數介紹:(1. Feature Selection 2. Train model 3. Apply model)

- ✓ Feature Selection
 - 1. Chi-square
 - 2. Pointwise MI
 - 3. (在 class_i 內的 df 值>某數 || 在非 class_i 內的 df 值<某數)
 - 4. Chi-square 表格(1,1)+(2,2)-(1,2)-(2,1)>某數
 - 5. 每個 class 所擁有 feature 數<某數
- ※參數選擇請見第一頁

```
913
      void SelectFeatures(int classes[][16])
914 🖵 {
           map<string, double> dict1;
int total = 0;
915
916
917 =
918 =
           for(int i=0; i<13; i++){
                for(int j=1; j<16; j++){
919
                    string doc = std::to_string(classes[i][j]);
                    ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_OUT\\"+doc+".txt");
920
               string line;
char *t, *s;
921
922
                getline(docM, line);
924
                char delim[] = '
925 🖨
                while (getline(docM, line)){
                    t = strtok ((char*)line.c_str(), delim); //parse with delim
while (t != NULL){
926 T
927 □
928
                        string term(t);
                         s = strtok(NULL, delim); //t pointing to the next delimiter position
930
                         if(dict1.find(term) == dict1.end())
931
                             dict1[term] = 1:
932
                             //cout<<term<<endl;
933
                         else
936
                             dict1[term] += 1;
937
                         t = strtok(NULL, delim);
                    }
938
               }
939
941
           ofstream outfile ("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_Class\\merge.txt");
942
943
           outfile<<std::left<<setw(70)<<"term"<<setw(70)<<"total_num"<<endl;
           map<string, double>::iterator it;
944
945 <del>|</del>
946 |
           for(it = dict1.begin(); it != dict1.end(); it++){
                outfile<<std::left<<setw(70)<<(*it).first;
947
948
               outfile<<std::left<<setw(70)<<(*it).second<<endl:
948
           }
949
950
           for(int i=0; i<13; i++){
951 T
952 च
               map<string, Arr5> score;
               for(int j=1; j<16; j++){
                   string doc = std::to_string(classes[i][j]);
ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_OUT\\"+doc+".txt");
953
954
               string line;
char *t, *s;
getline(docM, line);
955
956
957
958
                char delim[] = "
               while (getline(docM, line)){
959 🗀
                    t = strtok ((char*)line.c_str(), delim); //parse with delim
while (t != NULL){
960
961
962
                       string term(t);
                        s = strtok(NULL, delim); //t pointing to the next delimiter position
963
964
                        if(score.find(term) == score.end())
965
966
                            score[term] = {1,0,0,0,0};
968
969
                        else
970
                            score[term].num[0]++;
971
                        t = strtok(NULL, delim);
972
973
               }
974
975
976
           string now = std::to_string(i+1);
                                         \\Mark\\Desktop\\大五\\資訊檢索\\IRTM_Class\\fs_"+now+".txt");
977
           ofstream outfile ("C:
           map<string, double> D2V;
978
           map<string, double> D2V2;
980
           map<string, Arr5>::iterator it1;
```

```
981 <del>|</del>
982
983
                                   for(it1 = score.begin(); it1 != score.end(); it1++){
   double expected_f[4];
   double chi_score = 0;
                                              double chi_score = 0;
score[(*it1).first].num[1] = 15-score[(*it1).first].num[0];
score[(*it1).first].num[2] = dict1[(*it1).first] - score[(*it1).first].num[0];
score[(*it1).first].num[3] = 180 - score[(*it1).first].num[2];
expected_f[0] = 195.0*(15.0/195)*(dict1[(*it1).first]/195.0);
expected_f[1] = 195.0*(15.0/195)*(dict1[(*it1).first])/195.0);
expected_f[2] = 195.0*(180.0/195)*(dict1[(*it1).first])/195.0);
expected_f[3] = 195.0*(180.0/195)*(dict1[(*it1).first])/195.0);
cxpected_f[3] = 195.0*(180.0/195)*(dict1[(*it1).first])/195.0);
chi_score = pow(score[(*it1).first].num[0]-expected_f[0], 2)/(expected_f[0]) + pow(score[(*it1).first].num[1]-expected_f[1], 2)/(expected_f[2]) + pow(score[(*it1).first].num[3]-expected_f[3]);
double MI = log2((score[(*it1).first].num[0]/195.0)/((15.0/195)*(dict1[(*it1).first]/195)));
//MI = MIC/-log2((score[(*it1).first].num[0]/195.0))/((15.0/195)*(dict1[(*it1).first]/195)));
    984
    986
    987
   988
989
    990
   991
992
    993
   994
995
                                               //MI = MI/(-log2((score[(*it1).first].num[0]/195.0)));

double EMI = score[(*it1).first].num[0]*log2(score((*it1).first].num[0]/(15*dict1[(*it1).first]))
+score[(*it1).first].num[1]*log2(score[(*it1).first].num[1]/(15*(195.0-dict1[(*it1).first])))
+score[(*it1).first].num[2]*log2(score[(*it1).first].num[2]/(180.0*dict1[(*it1).first])))
+score[(*it1).first].num[3]*log2(score[(*it1).first].num[3]/(180.0*dict1[(*it1).first])));
if((score[(*it1).first].num[0])>[|score[(*it1).first].num[2]/(180.0*dict1[(*it1).first])));
if((score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).fir
    996
   997
998
   999
1000
1001
                                                else
1002
                                               1003
1005
1006
1008
1009
                                    //&& ((score[(*it1).first].num[1]-expected_f[1])+(score[(*it1).first].num[2]-expected_f[2]))<0 && (expected_f[0]>1 || chi_score * expected_f[0]>20)
//&&((score[(*it1).first].num[0]-expected_f[0])+(score[(*it1).first].num[3]-expected_f[3]))>0 && ((score[(*it1).first].num[1]-expected_f[1])+(score[
1011
1012
                                    //aa((store[('tt)/_trst].num|0)-expected_f[0])*(store[('t
vector/capricstring,double) > vec (D2V.begin(), D2V.end());
sort(vec.begin(),vec.end(),comp_by_value);
1015
                                                  vector<pair<string,double> >::iterator it5;
1016
                                                 int count = 0;
1017 <del>|</del>
1018 <del>|</del>
                                                  for(it5 = vec.begin(); it5!= vec.end() && count<45; ++it5){</pre>
                                                                  if(D2V[(*it5).first]>20){
1019
                                                                                  // && (D2V2[(*it5).first] > 2) && count<50
                                                                                  outfile<<std::left<<setw(70)<<(*it5).first;
1020
                                                                                  outfile<<std::left<<setw(70)<<D2V[(*it5).first]<<endl;
1021
1022
                                                                                   count++;
1023
                                                                                    total++;
1024
                                                                  }
1025
1026
1027
                                                  //printVec(vec);
                                                  //cout<<endl<<endl;
1028
1029
1030
                                                  cout<<total<<endl:
1031
```

outfile.close();

✓ Train model (Multinomial)

利用 13*15 個 Training documents 為每一個 class 計算 add-one smoothing 的 $P(X=t_k|c)$ 機率

```
732
       void TrainMultinomialNB(int classes[][16], int num[13])
733 🖵 {
734
           map<string, double> total;
735 <del>|</del>
           for(int i=0; i<13; i++){
               for(int j=1; j<16; j++){
                    string doc = std::to_string(classes[i][j]);
737
                    ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_OUT\\"+doc+".txt");
738
739
                    string line;
                   char *t, *s;
740
741
                   getline(docM, line);
742
                   char delim[] = "
743 <del>|</del>
                   while (getline(docM, line)){
                   t = strtok ((char*)line.c_str(), delim); //parse with delim
745 🖵
                   while (t != NULL){
746
                        string term(t);
747
                        s = strtok(NULL, delim); //t pointing to the next delimiter position
748
                        double ans = atof(s);
749
                        if(total.find(term) == total.end())
750
                            total[term] = 1;
                        t = strtok(NULL, delim);
751
752
753
754
               }
755
           }
756
757 🗀
           for(int i=0; i<13; i++){
758
               map<string, double> dict1;
               int distinct_num = 0, total_num = 0;
string clas = std::to_string(i+1);
759
760
               ofstream outfile ("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_Class\\"+clas+".txt");
761
762
763
               map<string, double> featuredict;
764
               string doc = std::to_string(i+1);
```

```
ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM CLASS\\fs "+doc+".txt");
765
              string line;
766
767
               char *t, *s;
               getline(docM, line);
768
769
               char delim[] = "
770 🗀
               while (getline(docM, line)){
771
               t = strtok ((char*)line.c_str(), delim); //parse with delim
772 🗀
               while (t != NULL){
                  string term(t);
773
                  s = strtok(NULL, delim); //t pointing to the next delimiter position
774
775
                  double ans = atof(s);
                  featuredict[term] = ans;
776
777
                  t = strtok(NULL, delim);
778
779
780
781
782
               for(int j=1; j<16; j++){
                  string doc = std::to_string(classes[i][j]);
783
                  ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_OUT\\"+doc+".txt");
784
785
                  string line;
786
                  char *t, *s;
                  getline(docM, line);
787
788
                  char delim[] = "
                  while (getline(docM, line)){
789 🗀
                  t = strtok ((char*)line.c_str(), delim); //parse with delim
790
791 🖨
                  while (t != NULL){
792
                      string term(t);
                       s = strtok(NULL, delim); //t pointing to the next delimiter position
793
                      double ans = atof(s);
if(dict1.find(term) == dict1.end())
794
795
796
797
                           dict1[term] = ans;
798
                           distinct_num++;
```

```
799
                                            total_num++;
 800
                                     else
 801
 802
                                           dict1[term] += ans;
total_num++;
 803
 804
                                     t = strtok(NULL, delim);
 806
 807
 808
 809
                        }
outfile<<std::left<<setw(70)<<"term"<<setw(70)<<"prob"<<endl;
map<string, double>::iterator it;
for(it = dict1.begin(); it != dict1.end(); it++){
    if(featuredict.find((*it).first) != featuredict.end())
 810
 811
 812
 813
813
                                     outfile<<std::left<<setw(70)<<(*it).first;
outfile<<std::left<<setw(70)<<((*it).second+1)/(total_num+total.size()-1)<<endl;</pre>
 815
 816
                               }
 817
 818
                        outfile.close();
num[i] = total_num+total.size()-1;
 819
 820
 821
 822
823
```

✓ Apply model

利用各 testing documents 的 df dictionary 做 log 機率連加,若沒有出現在 training documents 中的 term 就直接略過,最大的即是分類答案

```
int ApplyMultinomialNB(string d, int num[13])
825 🗏 {
               int result = -9999999;
               int class_result = 0;
map<string, int> dict1;
827
828
829
               ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_OUT\\"+d);
               string line;
char *m, *n;
getline(docM, line);
char delim[] = " ";
while (getline(docM, line)){
831
832
833
834
835 <del>|</del>
               m=strtok ((char*)line.c_str(), delim); //parse with delim
while (m != NULL){
836 T
837 =
                    string term(m);

n = strtok(NULL, delim); //t pointing to the next delimiter position
838
839
                    double ans = atoi(n);
dict1[term] = ans;
840
841
842
                    m = strtok(NULL, delim);
844
               map<string, double> dictT;
for(int i=1; i<=13; i++)</pre>
846
848
                     double temp = 0;
string doc = std::to_string(i);
ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_Class\\"+doc+".txt");
849
850
                    string line;
char *t, *s;
getline(docM, line);
852
853
854
                    char delim[] = "";
while (getline(docM, line)){
t = strtok ((char*)line.c_str(), delim); //parse with delim
855
856
857
858 🗀
                    while (t != NULL){
859
                          string term(t);
s = strtok(NULL, delim); //t pointing to the next delimiter position
860
861
                          double ans = atof(s);
dictT[term] = ans;
862
                          t = strtok(NULL, delim);
863
865
              }
867
               for(int i=1; i<=13; i++)
869
870
                     map<string, double> dict2;
                    double temp = 0;

string doc = std::to_string(i);

ifstream docM("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_Class\\"+doc+".txt");
871
872
873
874
875
                    string line;
char *t, *s;
                    char = C, s;
getline(docM, line);
char delim[] = " ";
while (getline(docM, line)){
    t = strtok ((char*)line.c_str(), delim);  //parse with delim
while (t != NULL){
876
877
878 🛱
879 T
880 F
                          string term(t);
                          s = strok(NULL, delim); //t pointing to the next delimiter position
double ans = atof(s);
dict2[term] = ans;
882
883
884
885
                           t = strtok(NULL, delim);
886
887
                    map<string, int>::iterator it;
for(it = dict1.begin(); it != dict1.end(); it++){
    if(dict2.find((*it).first) == dict2.end() && (dictT.find((*it).first) != dict2.end()))
    temp = temp + (*it).second*log(1.0/num[i-1]);
888
889
890
891
                        else if (dict2.find((*it).first) != dict2.end() )
  temp = temp + (*it).second*log((dict2[(*it).first]));
 894
                         895
 898
                   }
//if(d=="127.txt")
 899
 999
900
901
902
903 =
904
905
906 -
907
                    //cout<<temp*(1.0/13)<<" "<<result<<endl;
                    if(temp*(1.0/13)>result)
                        result = temp*(1.0/13);
class_result = i;
                    s
//cout<<"文章"<<d<<" "<<i<<" "<<result<<endl;
908
909
910
911
              return class_result;
```

✓ 主程式

Feature Selection -> Train Multinomial Model -> Apply the model and return most possible class

```
224
            int classes[13][16];
 225
            int testing[195];
 226
            int num[13];
 227
            int count = 0;
 228
            char *kk;
            int class_num = 0;
 229
 230
            int p = 0;
 231
            string line2;
            ifstream doc2("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\class.txt"); while (getline(doc2, line2)){
 232
 233 -
 234
                p = 0;
kk = strtok ((char*)line2.c_str(), " "); //parse with delim
 235
                while (kk != NULL){
    classes[class_num][p] = atoi(kk);
 236
                                                           //copy the string pointed to char array 'target'
 237
                    if(p!=0)
 238
 239
 240
                        testing[count] = atoi(kk);
 241
                        count++;
 242
 243
                    //cout<<p<<" "<<(classes[class_num][p])<<" ";
 244
                     //cout<<"count"<<" "<<count<<" ";
 245
                    kk = strtok(NULL, " "); //t pointing to the next delimiter position
 246
  247
 248
                 /cout<<endl;
 249
                class num++;
 250
  251
            SelectFeatures(classes);
 252
            TrainMultinomialNB(classes, num);
 253
            char InputPathA[65535] = "C:\\Users\\Mark\\Desktop\\大五\\資訊檢索\\IRTM_OUT"; //放要讀取檔案的資料來路徑到InputPath字串裡 ofstream myfile ("C:\\Users\\Mark\\Desktop\\大五\\資訊檢索");
 254
 255
            myfile.open ("result.csv");
myfile <<"Id"<<","<"Value"<<"\n";</pre>
 256
 257
            char szDirA[65535];
258
259
            char dirA[65535];
260
            int startid = 0;
261
            WIN32_FIND_DATA AFileData;
            HANDLE AhList;
sprintf(szDirA, "%s\\*", InputPathA );
262
263
264
            int f=0:
            if ( (AhList = FindFirstFile(szDirA, &AFileData))==INVALID_HANDLE_VALUE )
265
                printf("No files be found.\n\n");
266
267 <del>|</del>
268 <del>|</del>
269 <del>|</del>
            else {
                while (1) {
                    if (!FindNextFile(AhList, &AFileData)) {
270
                         if (GetLastError() == ERROR_NO_MORE_FILES)
271
                             break;
272
273 🖃
274
                         int is_test = 0;
275
                         f++;
276
                         string x = AFileData.cFileName;
                         for(int i=0; i<195; i++)
277
278
                              if(atoi(x.substr(0, x.find(".txt", 0)).c str())==testing[i])
279
280
281
                                  is test = 1;
282
                                  break:
283
284
                         if (startid>0 && is_test!=1)
285
286
                              287
288
                         startid ++;
289
290
291
292
            }
293
294
295
            return 0;
296
```