🔒 **MelvinLecoy** / **gitcode**   ( Private )

<> **Code**    ⊙ Issues    ⌥ Pull requests    ▷ Actions    ⊞ Projects    ⊘ Security    ⬚ Insights    ⚙ Settings

⎇ master ▾                                                                                      •••

**gitcode** / **p5-ml** / **main.cpp**

⬤  **Kwan Ting Lau** gitcode                                                              ⟲ History

⚇ **0** contributors

---

223 lines (192 sloc) │ 6.12 KB                                                            •••

```cpp
1    // Project UID db1f506d06d84ab787baf250c265e24e
2    #include "csvstream.h"
3    #include <cmath>
4    #include <cstring>
5    #include <iostream>
6    #include <iterator>
7    #include <map>
8    #include <set>
9    #include <string>
10   #include <utility>
11
12   using namespace std;
13
14   class Piazza {
15   private:
16     int post = 0;
17     int unique_word = 0;
18     int word_post = 0;  // number of post of with w
19     int label_post = 0; // bumber of post of label C
20
21     map<string, int> word_map;              // number of post of one label
22     map<string, int> label_map;             // number of post of one label
23     map<pair<string, string>, int> want_map; // C with w    time
24     // map<string, string> row;
25
26     set<string> unique_words_file; // set of unique word in the file
27     set<string> label_file;        // set of label in the file
28     set<string> given_post_word;   // set of given words
29
30   public:
31     double get_post() { return post; }
32
33     int get_uniq_word() { return unique_word; }
34
35     int get_word_post() { return word_post; }
```

```cpp
36
37      int get_label_post() { return label_post; }
38
39      set<string> get_unique_words_file() {
40        return unique_words_file;
41      } // get the set of unique word
42
43      set<string> get_label_file() {
44        return label_file;
45      } // get the set of label label in file
46
47      set<string> get_givn_post_word() { return given_post_word; }
48
49      int size_of_unique_words_file() { return unique_words_file.size(); }
50
51      int size_of_label_file() { return label_file.size(); }
52
53      int size_of_given_post_word_file() { return given_post_word.size(); }
54
55      double log_prior(string label) {
56        return log(label_map[label] / static_cast<double>(post));
57      }
58
59      double log_likelihood(string label, string unique_word) {
60        pair<string, string> pair_1;
61        pair_1.first = label;
62        pair_1.second = unique_word;
63        if (want_map.find(pair_1) ==
64            want_map.end()) { // doesn't appear in a post with label C
65          if ((word_map.find(unique_word) ==
66              word_map.end())) { // doesn't appear in training set
67            return log(1 / static_cast<double>(post));
68          } else {
69            return log(word_map[unique_word] / static_cast<double>(post));
70          }
71        }
72        return log(want_map[pair_1] / static_cast<double>(label_map[label]));
73      }
74
75      // log probability score, set is the set of words from a post
76      double log_probability_score(string label, set<string> words) {
77        double score = log_prior(label);
78        for (auto word : words) {
79          score += log_likelihood(label, word);
80        }
81        return score;
82      }
83
84      set<string> unique_words(const string &str) {
85        istringstream source(str);
86        set<string> words;
87        string word;
88        // Read word by word from the stringstream and insert into the set
89        while (source >> word) {
```

```cpp
 90              words.insert(word);
 91          }
 92          return words;
 93      }
 94
 95      void train(csvstream &csvin, bool debug) {
 96          map<string, string> row; // store the csv
 97          if (debug) {
 98              cout << "training data:" << endl;
 99          }
100          while (csvin >> row) {
101              if (debug) {
102                  cout << "  label = " << row["tag"] << ", content = " << row["content"]
103                        << endl;
104              }
105              unique_words_file = unique_words(row["content"]);
106              post++;
107              label_map[row["tag"]]++;
108              for (auto b : unique_words_file) {
109                  word_map[b]++;
110                  pair<string, string> pair_2;
111                  pair_2.first = row["tag"];
112                  pair_2.second = b;
113                  ++want_map[pair_2];
114              }
115          }
116          cout << "trained on " << post << " examples" << endl;
117          if (debug) {
118              cout << "vocabulary size = " << word_map.size() << endl << endl;
119              cout << "classes:" << endl;
120              for (auto x : label_map) {
121                  cout << "  " << x.first << ", " << x.second
122                        << " examples, log-prior = " << log_prior(x.first) << endl;
123              }
124              cout << "classifier parameters:" << endl;
125              for (auto y : want_map) {
126                  cout << "  " << y.first.first << ":" << y.first.second
127                        << ", count = " << y.second << ", log-likelihood = "
128                        << log_likelihood(y.first.first, y.first.second) << endl;
129              }
130          }
131
132          cout << endl;
133      } // #post
134      // #label  num
135      // #unique words num
136      // label  word num
137      // sum of unique words word_appear.size()
138
139      pair<string, double> predict(set<string> words) {
140          double score = -INFINITY;
141          string label = "";
142          pair<string, double> pair_3;
143          for (auto element : label_map) {
```

```cpp
144            string label_1 = element.first;
145            double score_1 = log_probability_score(label_1, words);
146            if (score_1 > score) {
147              score = score_1;
148              label = label_1;
149              pair_3.first = label_1;
150              pair_3.second = score_1;
151            }
152          }
153          return pair_3;
154        }

156        void test(csvstream &csvin) {
157          pair<string, double> test_pair;
158          map<string, string> row;
159          int correct_predict = 0;
160          int count = 0;

162          cout << "test data:" << endl;
163          while (csvin >> row) {
164            string tag = row["tag"];
165            string content = row["content"];
166            test_pair = predict(unique_words(content));
167            cout << "  correct = " << tag;
168            cout << ", predicted = " << test_pair.first;
169            cout << ", log-probability score = " << test_pair.second << endl;
170            cout << "  content = " << content << endl << endl;
171            if (test_pair.first == tag) {
172              correct_predict++;
173            }
174            count++;
175          }
176          cout << "performance: " << correct_predict << " / ";
177          cout << count << " posts predicted correctly" << endl;
178        }
179      };

181      int main(int argc, char *argv[]) {
182        cout.precision(3);

184        if ((argc != 3) && (argc != 4)) {
185          cout << "Usage: main.exe TRAIN_FILE TEST_FILE [--debug]" << endl;
186          return -1;
187        }

189        if ((argc == 4) && (string(argv[3]) != "--debug")) {
190          cout << "Usage: main.exe TRAIN_FILE TEST_FILE [--debug]" << endl;
191          return -1;
192        }

194        Piazza piazza;

196        try {
197          string s_1 = argv[1];
```

```cpp
198        csvstream csvin(s_1);
199    } catch (const exception &e) {
200      cout << "Error opening file: " << argv[1] << endl;
201      return -1;
202    }
203
204    try {
205      string s_2 = argv[2];
206      csvstream csvin1(s_2);
207    } catch (const exception &e) {
208      cout << "Error opening file: " << argv[2] << endl;
209      return -1;
210    }
211    string s_1 = argv[1];
212    csvstream csvin(s_1);
213    string s_2 = argv[2];
214    csvstream csvin1(s_2);
215
216    bool debug = false;
217    if (argc == 4 && string(argv[3]) == "--debug") {
218      debug = true;
219    }
220
221    piazza.train(csvin, debug);
222    piazza.test(csvin1);
223  }
```

[Give feedback](#)