

Московский авиационный институт
(национальный исследовательский университет)

Институт информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Курсовой проект по курсу «Дискретный анализ»
«Автоматическая классификация документов»

Студент: И. П. Моисеенков
Преподаватель: С. А. Сорокин
Группа: М8О-208Б-19
Дата: 29.06.2021
Оценка: хорошо
Подпись:

Москва, 2021

1 Задание

Требуется реализовать систему, которая на основе базы вопросов и тегов к ним будет предлагать варианты тегов, которые подходят к новым вопросам.

Формат запуска программы в режиме обучения:

```
./prog learn -input <input file> -output <stats file>
```

- -input - входной файл с вопросами
- -output - выходной файл с рассчитанной статистикой

Формат запуска программы в режиме классификации:

```
./prog classify -stats <stats file> -input <input file> -output <output file>
```

- -stats - файл со статистикой, полученной на предыдущем этапе
- -input - входной файл с вопросами
- -output - выходной файл с тегами к вопросам

Формат входных файлов при обучении:

```
<Количество строк в вопросе [n]>  
<Тег 1>,<Тег 2>,...,<Тег m>  
<Заголовок вопроса>  
<Текст вопроса [n строк]>
```

Формат входных файлов при запросах:

```
<Количество строк в вопросе [n]>  
<Заголовок вопроса>  
<Текст вопроса [n строк]>
```

Формат выходного файла: для каждого запроса в отдельной строке выводится предполагаемый набор тегов, через запятую.

2 Описание

Для решения этой задачи я использовал наивный байесовский классификатор. Это простой вероятностный классификатор, который работает на основе формулы Байеса:

$$P(A|B) = \frac{P(B|A)*P(A)}{P(B)}$$

Классификатор является «наивным», т.к. он работает со строгим предположением о независимости классов. Таким образом, для вычисления вероятности принадлежности текста к определенному классу используется формула:

$$P(class|word_1, \dots, word_n) = \frac{P(word_1|class)*\dots*P(word_n|class)*P(class)}{P(word_1)*\dots*P(word_n)}$$

Необходимо рассчитать вероятность принадлежности к каждому классу, применить функцию Softmax и выбрать классы с наибольшими вероятностями (в моей реализации это все классы с вероятностями больше $1/|classes|$). Именно к этим классам и будет относиться рассматриваемый вопрос.

Можно заметить, что знаменатель остается неизменным при рассмотрении одного вопроса, поэтому он не будет влиять на итоговое распределение вероятностей. Это значит, что его можно не вычислять.

Во избежание переполнения типа double в классификаторе вычисляется не произведение вероятностей, а сумма их логарифмов.

Отдельно стоит упомянуть следующую проблему. Если в тексте встретилось слово, которого не было в обучающей выборке, то по формуле получается нулевая вероятность принадлежности любому классу. Поэтому при вычислении вероятностей применяется сглаживание Лапласа:

$$P(word|class) = \frac{P(word)+\alpha}{P(class)+\alpha*|classes|}$$

В моей реализации по умолчанию $\alpha = 1$.

При запуске программы в режиме обучения классификатор собирает статистику - сколько всего слов в выборке, сколько слов встретилось при определенном классе, сколько раз встретилось каждое слово при определенном классе и сколько раз встретилось каждое слово в общем. Эта же информация и сохраняется в файл со статистикой.

3 Исходный код

Класс NaiveBayesClassifier

```
1  #include <iostream>
2  #include <vector>
3  #include <unordered_map>
4
5  using TText = std::vector<std::string>;
6
7  struct data {
8      TText tag;
9      TText doc;
10 };
11
12 class NaiveBayesClassifier {
13 public:
14     void fit(std::vector<data> &dataset) {
15         for (const auto &data : dataset) {
16             TText tags = data.tag;
17             TText doc = data.doc;
18             for (auto &tag: tags) {
19                 for (auto &word : doc) {
20                     add(tag, word);
21                 }
22             }
23             dataset_size += doc.size();
24         }
25     }
26
27     std::unordered_map<std::string, double> predict_proba(TText &doc) {
28         std::unordered_map<std::string, double> probabilities;
29         for (auto &data : tag_summary) {
30             std::string cur_tag = data.first;
31             probabilities[cur_tag] = probability(cur_tag, doc);
32         }
33
34         return probabilities;
35     }
36
37     void save_stats(std::ofstream &out) {
38         out << dataset_size << "\n";
39         for (auto &tag_data : tag_summary) {
40             std::string tag = tag_data.first;
41             auto data = tag_data.second;
42             out << tag << " " << data.size() << " ";
43             for (auto &word_count : data) {
44                 out << word_count.first << " " << word_count.second << " ";
45             }
46         }
47     }
48 }
```

```

46         out << "\n";
47     }
48 }
49
50 void load_stats(std::ifstream &in) {
51     in >> dataset_size;
52     std::string tag;
53     while (in >> tag) {
54         int amount_of_words;
55         in >> amount_of_words;
56         for (int i = 0; i < amount_of_words; ++i) {
57             std::string word;
58             int count;
59             in >> word >> count;
60             add(tag, word, count);
61         }
62     }
63 }
64
65 int get_number_of_tags() {
66     return tag_sizes.size();
67 }
68
69 private:
70     std::unordered_map<std::string, int> tag_sizes; // amounts of words with each tag
71     std::unordered_map<std::string, std::unordered_map<std::string, int>> tag_summary;
72     // amounts of each word with each tag
73     std::unordered_map<std::string, int> words_summary; // amounts of each word
74     int dataset_size = 0; // total amount of words in dataset
75
76     double probability(std::string &tag, TText &doc, double alpha = 1) { // P(tag | doc)
77         // the Bayes rule
78
79         double prob = 0;
80         for (auto word : doc) {
81             prob += log(probability(word, tag, alpha));
82         }
83
84         prob += log(probability(tag, alpha));
85
86         return prob;
87     }
88
89     double probability(std::string &tag, double alpha = 1) { // P(tag)
90         // applying Laplace smoothing
91         return (1. * tag_sizes[tag] + alpha) / (dataset_size + alpha * dataset_size);
92     }

```

```

93 | double probability(std::string &word, std::string &tag, double alpha = 1) { // P(
    |     word / tag)
94 |     // applying Laplace smoothing
95 |     return (1. * tag_summary[tag][word] + alpha) / (words_summary[word] + alpha *
    |         tag_summary.size());
96 | }
97 |
98 | void add(std::string &tag, std::string &word) {
99 |     if (tag_summary[tag][word]) {
100 |         ++tag_summary[tag][word];
101 |     } else {
102 |         tag_summary[tag][word] = 1;
103 |     }
104 |
105 |     if (words_summary[word]) {
106 |         ++words_summary[word];
107 |     } else {
108 |         words_summary[word] = 1;
109 |     }
110 |
111 |     if (tag_sizes[tag]) {
112 |         tag_sizes[tag] = 1;
113 |     } else {
114 |         ++tag_sizes[tag];
115 |     }
116 | }
117 |
118 | void add(std::string &tag, std::string &word, int count) {
119 |     tag_summary[tag][word] = count;
120 |
121 |     if (words_summary[word]) {
122 |         words_summary[word] += count;
123 |     } else {
124 |         words_summary[word] = count;
125 |     }
126 |
127 |     if (tag_sizes[tag]) {
128 |         tag_sizes[tag] += count;
129 |     } else {
130 |         tag_sizes[tag] = count;
131 |     }
132 | }
133 | };

```

main.cpp

```

1 | #include <cstring>
2 | #include <fstream>

```

```

3 #include <unordered_set>
4 #include <cmath>
5
6 #include "NaiveBayesClassifier.h"
7
8 TText convert(std::string &text); // convert text to vector of words
9 void print_error();
10 std::unordered_map<std::string, double> softmax(std::unordered_map<std::string, double
    >& preds);
11
12 int main(int argc, char *argv[]) {
13     std::string input_file_name;
14     std::string stats_file_name;
15     std::string output_file_name;
16
17     if (argc < 6) {
18         print_error();
19         return 1;
20     }
21
22     if (!strcmp(argv[1], "learn")) {
23         if (argc != 6) {
24             print_error();
25             return 1;
26         }
27
28         for (int i = 2; i < argc; ++i) {
29             if (!strcmp(argv[i], "--input") && input_file_name.empty()) {
30                 input_file_name = argv[++i];
31             } else if (!strcmp(argv[i], "--output") && stats_file_name.empty()) {
32                 stats_file_name = argv[++i];
33             } else {
34                 print_error();
35                 return 1;
36             }
37         }
38
39         if (input_file_name.empty() || stats_file_name.empty()) {
40             print_error();
41             return 1;
42         }
43
44         std::ifstream input_file(input_file_name);
45         std::ofstream stats_file(stats_file_name);
46         std::vector<data> dataset;
47
48         int lines;
49         while (input_file >> lines) {
50             std::string tags;

```

```

51         std::string text;
52
53         input_file.ignore();
54         getline(input_file, tags);
55         for (int i = 0; i < lines + 1; ++i) {
56             std::string cur_line;
57             getline(input_file, cur_line);
58             text += " " + cur_line;
59         }
60         dataset.push_back({convert(tags), convert(text)});
61     }
62
63     NaiveBayesClassifier NB;
64     NB.fit(dataset);
65     NB.save_stats(stats_file);
66 } else if (!strcmp(argv[1], "classify")) {
67     if (argc != 8) {
68         print_error();
69         return 1;
70     }
71
72     for (int i = 2; i < argc; ++i) {
73         if (!strcmp(argv[i], "--input") && input_file_name.empty()) {
74             input_file_name = argv[++i];
75         } else if (!strcmp(argv[i], "--output") && output_file_name.empty()) {
76             output_file_name = argv[++i];
77         } else if (!strcmp(argv[i], "--stats") && stats_file_name.empty()) {
78             stats_file_name = argv[++i];
79         } else {
80             print_error();
81             return 1;
82         }
83     }
84
85     if (input_file_name.empty() || stats_file_name.empty() || output_file_name.
86         empty()) {
87         print_error();
88         return 1;
89     }
90
91     std::ifstream input_file(input_file_name);
92     std::ifstream stats_file(stats_file_name);
93     std::ofstream output_file(output_file_name);
94
95     NaiveBayesClassifier NB;
96     NB.load_stats(stats_file);
97
98     int lines;
99     while (input_file >> lines) {

```



```

99         std::string text;
100
101         input_file.ignore();
102         for (int i = 0; i < lines + 1; ++i) {
103             std::string cur_line;
104             getline(input_file, cur_line);
105             text += " " + cur_line;
106         }
107         auto doc = convert(text);
108         auto probas = NB.predict_proba(doc);
109         auto preds = softmax(probas);
110
111         bool need_comma = false;
112         double threshold = 1. / NB.get_number_of_tags();
113         for (const auto& pred : preds) {
114             if (pred.second > threshold) {
115                 if (need_comma) {
116                     output_file << ", ";
117                 }
118                 output_file << pred.first;
119                 need_comma = true;
120             }
121         }
122         output_file << "\n";
123     }
124
125     } else {
126         print_error();
127         return 1;
128     }
129 }
130
131 TText convert(std::string &text) { // convert text to vector of words
132     std::vector<std::string> words;
133     int i = 0;
134     std::string cur_word;
135     while (i < text.size()) {
136         if (isalpha(text[i])) {
137             cur_word += tolower(text[i]);
138         } else if (cur_word.length()) {
139             words.push_back(cur_word);
140             cur_word = "";
141         }
142         if (i + 1 == text.size() && cur_word.length()) {
143             words.push_back(cur_word);
144         }
145         ++i;
146     }
147     return words;

```

```

148 }
149
150 void print_error() {
151     std::cerr << "Incorrect syntax!" << std::endl <<
152         "./classifier learn --input <input file> --output <output file>" << std::
153             endl <<
154             "./classifier classify --stats <stats file> --input <input file> --output
155                 <output file>" << std::endl;
156 }
157
158 std::unordered_map<std::string, double> softmax(std::unordered_map<std::string, double
159     >& probs) {
160     double max = -1e10;
161     for (const auto& prob : probs) {
162         if (prob.second > max) {
163             max = prob.second;
164         }
165     }
166
167     double sum = 0;
168     for (auto& prob : probs) {
169         sum += exp(prob.second - max);
170     }
171
172     double constant = max + log(sum);
173     std::unordered_map<std::string, double> res;
174     for (auto prob : probs) {
175         res[prob.first] = exp(prob.second - constant);
176     }
177     return res;
178 }

```

4 Демонстрация работы

Продemonстрирую работоспособность классификатора на простом датасете. Будем классифицировать вопросы на три группы - про кошек, собак или хомяков.

Обучающие данные

1
Dogs
Buying a dog
Where can I buy a dog?
1
Dogs
Puppies
How many puppies are there in his house?
1
Dogs
Cute animals
What are the most beautiful breeds of dogs?
2
Dogs
Lost my dog
I've lost my dog.
How can I find it now?
1
Cats
A name for a cat
What are the best names for cats?
1
Cats
Mother cat with her kittens
My cat has many kittens. How should I take care of them?
1
Cats
Different cats
How many different kinds of cats exist?
1
Cats
Black cat
Is it true that a black cat causes problems?
1

Hamsters
Hamster's age
How many years do hamsters live?
1
Hamsters
Very small animals
What is the smallest hamster in the world?
1
Hamsters
Feeding hamsters
I have two hamsters. How should I feed them?
1
Hamsters
Hamsters living together
Can two small hamsters live together in one cage?
1
Cats, dogs
Friendship between cats and dogs
Can cats and dogs live together?
1
Dogs, hamsters
Animals
Can my puppy make friends with a new little hamster?

Тестовые данные

1
Feeding
How to feed my little friend who lives in a cage?
1
Puppies and kittens
Who is better - a puppy or a kitten?
1
Many animals!
Is it ok, that a little hamster, two black cats and a puppy live together with me?

Результат

hamsters
dogs, cats
dogs, cats

5 Выводы

Выполнив данное задание, я познакомился с наивным байесовским классификатором. Это очень простая, но в то же время эффективная модель для классификации текстов. Байесовский классификатор используется, например, для фильтра спама и определения тегов к текстам.

Я реализовал свой байесовский классификатор с нуля, при этом я познакомился с различными техниками для улучшения его работоспособности (например, сглаживание Лапласа).

При написании кода я использовал знания, полученные на курсах дискретного анализа и объектно-ориентированного программирования.

Для тестирования классификатора я создал несколько датасетов. Первый - с вопросами о животных, второй - с электронными письмами (классификация на спам и нормальные письма). На втором датасете я измерил точность работы классификатора. Полученный результат: $accuracy \approx 75\%$

Список литературы

- [1] *Наивный классификатор Байеса*
URL: <https://towardsdatascience.com/naive-bayes-classifier-81d512f50a7c>.
- [2] *Наивный классификатор Байеса*
URL: <https://vk.cc/c3sewd>.
- [3] *Сглаживание Лапласа*
URL: <https://vk.cc/c3sexP>.
- [4] *Функция Softmax*
URL: <https://slaystudy.com/implementation-of-softmax-activation-function-in-c-c/>.