# E11 Naive Bayes (C++/Python)

18340215 张天祎

2020 年 11 月 26 日

# 目录

# 1  Datasets

The UCI dataset (`http://archive.ics.uci.edu/ml/index.php`) is the most widely used dataset for machine learning. If you are interested in other datasets in other areas, you can refer to `https://www.zhihu.com/question/63383992/answer/222718972`.

Today's experiment is conducted with the **Adult Data Set** which can be found in `http://archive.ics.uci.edu/ml/datasets/Adult`.

| Data Set Characteristics: | Multivariate | Number of Instances: | 48842 | Area: | Social |
|---|---|---|---|---|---|
| Attribute Characteristics: | Categorical, Integer | Number of Attributes: | 14 | Date Donated | 1996-05-01 |
| Associated Tasks: | Classification | Missing Values? | Yes | Number of Web Hits: | 1305515 |

You can also find 3 related files in the current folder, `adult.name` is the description of **Adult Data Set**, `adult.data` is the training set, and `adult.test` is the testing set. There are 14 attributes in this dataset:

>50K,  <=50K.


1. age: continuous.
2. workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.
3. fnlwgt: continuous.
4. education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 5. 1st-4th, 10th, Doctorate, 5th-6th, Preschool.
5. education-num: continuous.
6. marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.
7. occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.
8. relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.
9. race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.
10. sex: Female, Male.

11. capital−gain: continuous.

12. capital−loss: continuous.

13. hours−per−week: continuous.

14. native−country: United−States, Cambodia, England, Puerto−Rico, Canada, Germany, Outlying−US(Guam−USVI−etc), India, Japan, Greece, South, China, Cuba, Iran, Honduras, Philippines, Italy, Poland, Jamaica, Vietnam, Mexico, Portugal, Ireland, France, Dominican−Republic, Laos, Ecuador, Taiwan, Haiti, Columbia, Hungary, Guatemala, Nicaragua, Scotland, Thailand, Yugoslavia, El−Salvador, Trinadad&Tobago, Peru, Hong, Holand−Netherlands.

**Prediction task is to determine whether a person makes over 50K a year.**

# 2 Naive Bayes

Naive Bayes is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. It is not a single algorithm for training such classifiers, but a family of algorithms based on a common principle: all naive Bayes classifiers assume that **the value of a particular feature is independent of the value of any other feature**, given the class variable.

For example, a fruit may be considered to be an apple if it is red, round, and about 10 cm in diameter. A naive Bayes classifier considers each of these features to contribute independently to the probability that this fruit is an apple, regardless of any possible correlations between the color, roundness, and diameter features.

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of conditional independence between every pair of features given the value of the class variable. Bayes' theorem states the following relationship, given class variable $y$ and dependent feature vector $x_1$ through $x_n$:

$$P(y \mid x_1, ..., x_n) = \frac{P(y)P(x_1, ..., x_n \mid y)}{P(x_1, ..., x_n)}$$

Using the naive conditional independence assumption that

$$P(x_i \mid y, x_1, ..., x_{i-1}, x_{x+1}, ..., x_n) = P(x_i \mid y)$$

, for all $i$, this relationship is simplified to

$$P(y \mid x_1, ..., x_n) = \frac{P(y) \prod_{i=1}^{n} P(x_i \mid y)}{P(x_1, ..., x_n)}$$

Since $P(x_1, ..., x_n)$ is constant given the input, we can use the following classification rule:

$$P(y \mid x_1, ..., x_n) \propto P(y) \prod_{i=1}^{n} P(x_i \mid y)$$

$$\hat{y} = \arg\max_{y} P(y) \prod_{i=1}^{n} P(x_i \mid y),$$

and we can use Maximum A Posteriori (MAP) estimation to estimate $P(y)$ and $P(x_i \mid y)$, the former is then the relative frequency of class $y$ in the training set.

The different naive Bayes classifiers differ mainly by the assumptions they make regarding the distribution of $P(x_i \mid y)$.

- When attribute values are discrete, $P(x_i \mid y)$ can be easily computed according to the training set.

- When attribute values are continuous, an assumption is made that the values associated with each class are distributed according to Gaussian i.e., Normal Distribution. For example, suppose the training data contains a continuous attribute $x$. We first segment the data by the class, and then compute the mean and variance of $x$ in each class. Let $\mu_k$ be the mean of the values in $x$ associated with class $y_k$, and let $\sigma_k^2$ be the variance of the values in $x$ associated with class $y_k$. Suppose we have collected some observation value $x_i$. Then, the probability distribution of $x_i$ given a class $y_k$, $P(x_i \mid y_k)$ can be computed by plugging $x_i$ into the equation for a Normal distribution parameterized by $\mu_k$ and $\sigma_k^2$. That is,

$$P(x = x_i \mid y = y_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}$$

## 3   Task

- Given the training dataset `adult.data` and the testing dataset `adult.test`, please accomplish the prediction task to determine whether a person makes over 50K a year in `adult.test` by using Naive Bayes algorithm (C++ or Python), and compute the accuracy.

- Note: keep an eye on the discrete and continuous attributes.

- Please finish the experimental report named `E11_YourNumber.pdf`, and send it to `ai_2020@foxmail.com`

## 4   Codes and Results

```python
1    import csv
2    import math
3    import numpy as np
4    from collections import Counter
5
6    def load(file_name):
7            with open(file_name, 'r') as file:
8                    reader = csv.reader(file)
9                    rows = []
10                   for row in reader:
11                           if row != [] and len(row) == 15:
12                                   tem = []
13                                   for ind, word in enumerate(row):
14                                           if ind == 0:
15                                                   tem.append(word)
16                                           elif ind == len(row) - 1 and word
                                                   [-1]=='.':
17                                                   tem.append(word[1:-1])
18                                           else:
19                                                   tem.append(word[1:])
20                                   rows.append(tem)
21                   return rows
22
23
24   def load_attributes(file_name):
25           with open(file_name, 'r') as file:
26                   reader = csv.reader(file)
27                   rows = [row for row in reader if row != []]
28                   attributes = {}
29                   con_list = []
30                   dis_map = {}
31                   for row in rows:
32                           if row[0][0] != '|':
33                                   row[-1] = row[-1][:-1]
34                                   if row[0].find(':') == -1:
35                                           labels = row
36                                   else:
```

```
37                                                tem = row [0]. split (':')
38                                                attribute = tem [0]
39                                                row [0] = tem [1]
40                                                attributes [attribute] = [word [1:]
                                                     for word in row]
41                          for ind, (attribute, val) in enumerate(attributes.items()):
42                                  if val == ['continuous']:
43                                          con_list.append(ind)
44                                  else:
45                                          dis_map[ind] = val
46                          return labels, con_list, dis_map
47
48

49    class NB():
50          def __init__(self, labels, con_list, dis_map, train):
51
52                  train_data = [i[:-1] for i in train]
53                  train_label = [i[-1] for i in train]
54
55                  # 处理标签
56                  self.label_cpt = {}
57                  count = Counter(train_label)
58                  data_size = len(train)
59                  for i in count:
60                          self.label_cpt[i] = count[i]/data_size
61
62                  # 处理连续数据
63                  self.con_list = {}
64                  for i in con_list:
65                          for j in self.label_cpt:
66                                  # self.con_list[(i,j)] = (np.mean([float(k[i
                                      ]) for k in train if int(k[i]) != 0 and k
                                      [-1] == j]), math.sqrt(np.var([float(k[i
                                      ]) for k in train if int(k[i]) != 0  and
                                      k[-1] == j])))
67                                  self.con_list[(i,j)] = (np.mean([float(k[i])
                                      for k in train if k[-1] == j]), math.
```

```python
                              sqrt(np.var([float(k[i]) for k in train
                                  if k[-1] == j])))
                  # 处理离散数据
                  self.dis_cpt = {}
                  for i in dis_map:
                      count = dict(Counter([j[i] for j in train_data]))
                      unknow = count.get('?', 0)
                      for j in count:
                          if j != '?':
                              count[j] = count[j] / (data_size -
                                  unknow)
                      for j in dis_map[i]:
                          if count.get(j, "") == "":
                              count[j] = 0
                      if count.get('?', "") != "":
                          count.pop('?')
                      self.dis_cpt[i] = count

      def test(self, test_data):
              rst = []
              # print(self.con_list)
              for test in test_data:
                      ph = {}
                      for label in self.label_cpt:
                              ph[label] = self.label_cpt[label]
                              for ind, i in enumerate(test):
                                      if ind not in [k[0] for k in self.
                                          con_list]:
                                              if i != '?':
                                                      ph[label] *= self.
                                                          dis_cpt[ind][i]
                                              else:
                                                      ph[label] *= max(
                                                          self.dis_cpt[ind
                                                          ].items(), key=
                                                          lambda x:x[1])[1]
                                      else:
```
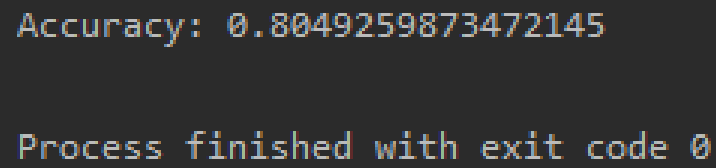
```python
 97                                             if float(i) != 0:
 98                                                 u = self.con_list[(
                                                        ind,label)][0]
 99                                                 sig = self.con_list
                                                        [(ind,label)][1]
100                                                 ph[label] *= np.exp
                                                        (-(float(i) - u)
                                                        ** 2 /(2* sig
                                                        **2))/(math.sqrt
                                                        (2*math.pi)*sig)
101                                         # else:
102                                         #     sig = self.con_list[(
                                                    ind, label)][1]
103                                         #     ph[label] *= np.exp(0)
                                                    / (math.sqrt(2 * math.pi
                                                    ) * sig)
104                         rst.append(max(ph.items(), key=lambda x:x[1])[0])
105                 return rst
106
107     def main():
108             train = load('dataSet/adult.data')
109             labels, con_list, dis_map = load_attributes('dataSet/adult.names')
110
111             nb = NB(labels, con_list, dis_map, train)
112
113             test = load('dataSet/adult.test')
114             test_label = nb.test([i[:-1] for i in test])
115             length = len(test)
116             count = 0
117             for i in zip(test_label, [i[-1] for i in test]):
118                     if i[0]==i[1]:
119                             count += 1
120             print("Accuracy:", count/length)
121
122
123     if __name__ == "__main__":
124             main()
```

- 结果如图

- 朴素贝叶斯分类器的效果比决策树分类器的效果略差，原因是该问题模型和决策树模型更相似。

- 若不使用连续数据，结果约 0.77，感觉贝叶斯分类器比决策树分类器更适合处理连续数据。