

E11 Naive Bayes (C++/Python)

20214966 Yangkai Lin 20214810 Suixin Ou

November 25, 2020

Contents

1	Datasets	2
2	Naive Bayes	3
3	Task	4
4	Codes and Results	5
4.1	Codes	5
4.1.1	连续值处理	5
4.1.2	缺失值处理	5
4.1.3	MAP	6
4.2	Results	7

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 , Trinidad&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 :

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

Using the naive conditional independence assumption that

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

, for all i , this relationship is simplified to

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

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

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

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

and we can use Maximum A Posteriori (MAP) estimation to estimate $P(y)$ and $P(x_i | y)$, 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 $P(x_i | y)$.

- When attribute values are discrete, $P(x_i | y)$ 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 μ_k be the mean of the values in x associated with class y_k , and let σ_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 , $P(x_i | y_k)$ can be computed by plugging x_i into the equation for a Normal distribution parameterized by μ_k and σ_k^2 . That is,

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

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

4.1 Codes

4.1.1 连续值处理

对每一种属性标签计算 μ , σ 。

```
mu = dict()
sigma = dict()
for attribute in attributes:
    name = attribute.name
    if attribute.domain == 'continuous':
        mu[name] = dict()
        sigma[name] = dict()
        for label in counting_table:
            counting = counting_table[label][name]
            n = len(counting)
            if label not in counting_table_total:
                counting_table_total[label] = n
            mu[name][label] = sum(counting) / n
            sigma[name][label] = (sum([(t - mu[name][label]) ** 2\
for t in counting]) / (n - 1)) ** 0.5
```

最后使用正态分布表示概率。

```
if attribute.domain == 'continuous':
    x = float(value)
    P *= math.exp(-(x - mu[name][label]) ** 2 / (2 * sigma[name][label] ** 2))
    / ((2 * math.pi) ** 0.5 * sigma[name][label])
```

4.1.2 缺失值处理

对于训练集，离散属性用众数，连续属性用均值。

```
mode = dict()
for attribute in attributes:
    name = attribute.name
    if attribute.domain == 'continuous':
```

```

        counting = counting_table_for_mode[name]
        mode[name] = sum(counting) / len(counting)
    else:
        mode[name] = max(counting_table_for_mode[name],\
            key = lambda k: counting_table_for_mode[name][k])
for label, name in missing_examples:
    if attribute.domain == 'continuous':
        counting_table[label][name].append(mode[name])
    else:
        counting_table[label][name][mode[name]] += 1

```

对于测试集，离散属性用训练集的众数，连续属性用训练集的均值。

```

if value == '?':
    value = mode[name]

```

4.1.3 MAP

$$\begin{aligned}
 h_{MAP} &= \arg \max_h \mathbf{Pr}(h|d) \\
 &= \arg \max_h \mathbf{Pr}(h) \mathbf{Pr}(d|h) \\
 &= \arg \max_h \mathbf{Count}(h) \prod_i \mathbf{Pr}(d_i|h)
 \end{aligned}$$

对于离散属性：

$$\mathbf{Pr}(d_i|h) = \frac{\mathbf{Count}(d_i|h)}{\mathbf{Count}(h)}$$

对于连续属性：

$$\mathbf{Pr}(d_i|h) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma}$$

代码如下：

```

P = counting_table_total[label]
for attribute in attributes:
    name = attribute.name
    value = example_attributes[name]
    if value == '?':
        value = mode[name]
    if attribute.domain == 'continuous':

```

```

x = float(value)
P *= math.exp(-(x - mu[name][label]) ** 2 / (2 * sigma[name][label]
        / ((2 * math.pi) ** 0.5 * sigma[name][label]))
else:
    P *= counting_table[label][name][value] / counting_table_total[label]

```

4.2 Results

完整代码和结果见[naive_bayes.html](#)。