

In [1]:

```
1 import scipy.io as scio
2 import numpy as np
3 from functools import reduce
4 import operator
```

读数据集函数

In [2]:

```
1 """
2 xfile: x文件
3 yfile: y文件
4 rate: 训练集数量/测试集数量
5 """
6 def readData(xfile, yfile, rate):
7     # 加载文件, 此处是dict
8     x = scio.loadmat(xfile)
9     y = scio.loadmat(yfile)
10
11     x = x.get("mnist_train")
12     y = y.get("mnist_train_labels")
13
14     # 二值化
15     x = np.where(x > 0, 1, 0)
16
17     num = x.shape[0]
18     index = int(num * (rate / (rate + 1)))
19
20     x_train = x[0:index]
21     y_train = y[0:index]
22     x_test = x[index:num]
23     y_test = y[index:num]
24
25     return x_train, y_train, x_test, y_test
```

读取训练集和测试集

In [3]:

```
1 x_train, y_train, x_test, y_test = readData("dataset/mnist_train.mat", "dataset/mnist_train_lab
2 print("x_train.shape = {}".format(x_train.shape))
3 print("y_train.shape = {}".format(y_train.shape))
4 print("x_test.shape = {}".format(x_test.shape))
5 print("y_test.shape = {}".format(y_test.shape))
```

```
x_train.shape = (50000, 784)
y_train.shape = (50000, 1)
x_test.shape = (10000, 784)
y_test.shape = (10000, 1)
```

计算流程 $p(y|x) = p(y)*p(x|y)$

1. 计算 $p(y)$
2. 计算 $p(x|y)$
3. 计算 $p(y|x)$

1.计算 $p(y)$

In [4]:

```
1 # 按y把数据分成10组
2 # 初始化px_group
3 px_group = []
4 for i in range(10):
5     px_group.append(i)
6     px_group[i] = []
7
8 # px_group的第一维表示类别[0-9], 第二维表示相应的所有x样例
9 for i, y in enumerate(y_train):
10     px_group[int(y)].append(x_train[i])
11
12 py = []
13 for i in range(10):
14     py.append(len(px_group[i]) / x_train.shape[0])
15
16 py = np.array(py)
17 py = py.reshape(-1, 1)
18 print("py.shape = {}".format(py.shape))
19 print("py      = {}".format(py.squeeze()))
```

```
py.shape = (10, 1)
py      = [ 0.09864  0.11356  0.09936  0.10202  0.09718  0.09012  0.09902  0.1035
 0.09684  0.09976]
```

2.计算 $p(x|y)$

In [5]:

```
1 pxy = []
2 for i in range(10):
3     group = np.array(px_group[i])
4     # 按列求和
5     group = np.sum(group, axis=0) / len(px_group[i])
6     pxy.append(group)
7
8 pxy = np.array(pxy)
9 print("pxy.shape = {}".format(pxy.shape))
```

```
pxy.shape = (10, 784)
```

3.计算 $p(y|x)$

(1) 多项式朴素贝叶斯

In [6]:

```
1 def predict_mul(x):
2     result = []
3     for i in range(10):
4         # 把0换成1
5         temp = x * pxy[i]
6         resulti = []
7         for j in np.where(temp == 0, 1, temp):
8             py_pxy = py[i][0] * reduce(operator.mul, j)
9             resulti.append(py_pxy)
10        result.append(resulti)
11
12    result = np.argmax(result, axis=0)
13    return result
```

(2) 伯努力朴素贝叶斯

In [7]:

```
1 def predict_ber(x):
2
3     result = []
4
5     for i in range(10):
6         # 把0换成1
7         temp = x * pxy[i]
8         resulti = []
9         for j in np.where(temp == 0, 1-pxy[i], temp):
10            py_pxy = py[i][0] * reduce(operator.mul, j)
11            resulti.append(py_pxy)
12            result.append(resulti)
13
14    result = np.argmax(result, axis=0)
15    return result
```

In [8]:

```
1 # 多项式朴素贝叶斯
2 result_mul = predict_mul(x_test)
3 # 伯努力朴素贝叶斯
4 result_ber = predict_ber(x_test)
```

In [9]:

```
1 acc_mul = np.where(result_mul == np.squeeze(y_test), 1, 0)
2 print("多项式朴素贝叶斯 Acc_mul = {}".format(np.sum(acc_mul) / len(result_mul)))
3
4 acc_ber = np.where(result_ber == np.squeeze(y_test), 1, 0)
5 print("伯努力朴素贝叶斯 Acc_ber = {}".format(np.sum(acc_ber) / len(result_ber)))
```

多项式朴素贝叶斯 Acc_mul = 0.626

伯努力朴素贝叶斯 Acc_ber = 0.8469

