

哈尔滨工业大学计算机科学与技术学院

## 实验报告

课程名称： 机器学习

课程类型： 必修

实验题目： 逻辑回归实验

学号：

姓名：

## 一、实验目的

理解逻辑回归模型，掌握逻辑回归模型的参数估计算法。

## 二、实验要求及实验环境

### 2.1 要求：

实现两种损失函数的参数估计（1，无惩罚项；2. 加入对参数的惩罚），可以采用梯度下降、共轭梯度或者牛顿法等。

### 2.2 验证：

1. 可以手工生成两个分别类别数据（可以用高斯分布），验证你的算法。考察类条件分布不满足朴素贝叶斯假设，会得到什么样的结果。

2. 逻辑回归有广泛的用处，例如广告预测。可以到 UCI 网站上，找一实际数据加以测试。

### 2.3 实验环境：

PyCharm 2020.3.2 x64

Python 3.8 numpy1.20.3 matplotlib 3.4.2

## 三、设计思想（本程序中的用到的主要算法及数据结构）

### 3.1 理论推导：

在一般的假设条件下，类别 $C_1$ 的后验概率可以写为作用在特征向量 $\phi$ 的线性函数上的 logistics sigmoid 函数的形式，即

$$P(C_1|\phi) = y(\phi) = \sigma(w^T \phi) \quad (1)$$

且

$$P(C_2|\phi) = 1 - P(C_1|\phi) \quad (2)$$

由贝叶斯定理可得，X属于类 $C_1$ 的后验概率为：

$$P(C_2|X=x) = \frac{P(C_2)P(X=x|C_2)}{P(C_1)P(X=x|C_1) + P(C_2)P(X=x|C_2)} \quad (3)$$

$$= \frac{1}{1 + \frac{P(C_1)P(X=x|C_1)}{P(C_2)P(X=x|C_2)}} \quad (4)$$

$$= \frac{1}{1 + e^{\ln \frac{P(C_1)P(X=x|C_1)}{P(C_2)P(X=x|C_2)}}} \quad (5)$$

并定义

$$\sigma(a) = \frac{1}{1 + e^{-a}} \quad (6)$$

$$P(C_2|X=x) = \frac{1}{1+e^{-a}} \quad (7)$$

$$a = \ln \frac{P(C_2)P(X=x|C_2)}{P(C_1)P(X=x|C_1)} \quad (8)$$

其反函数为

$$a = \ln \left( \frac{\sigma}{1-\sigma} \right) \quad (9)$$

因此有

$$P(C_2|X) = \frac{1}{1 + e^{\ln \frac{P(X=x|C_1)}{P(X=x|C_2)} + \ln \frac{P(C_1)}{P(C_2)}}} \quad (10)$$

对于离散特征值，由于特征值相互独立，可得类条件分布，

$$P(X=x|C_k) = P(X^{(1)}=x^{(1)}, \dots, X^{(n)}=x^{(n)}|C_k) \quad (11)$$

$$= \prod_{i=1}^n P(X^{(i)}=x^{(i)}|C_k) \quad (12)$$

假设类条件概率密度是高斯分布，并假定所有的类别的协方差矩阵相同，此时类别 $C_k$ 的类条件概率为

$$P(X=x|C_k) = \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)} \quad (13)$$

其中 $x$ 表示维度为 $n$ 的向量， $\mu_k$ 是这些向量的平均值， $\Sigma$ 表示所有向量 $x$ 的协方差矩阵。

因此有

$$\ln \frac{P(X=x|C_1)}{P(X=x|C_2)} = -\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + \frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2) \quad (14)$$

$$= -\frac{1}{2}(x^T \Sigma^{-1}x - 2\mu_1^T \Sigma^{-1}x + \mu_1^T \Sigma^{-1}\mu_1) + \frac{1}{2}(x^T \Sigma^{-1}x - 2\mu_2^T \Sigma^{-1}x + \mu_2^T \Sigma^{-1}\mu_2) \quad (15)$$

$$= (\mu_1^T - \mu_2^T) \Sigma^{-1}x - \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1}\mu_2 \quad (16)$$

代入公式(8)中，可得

$$P(C_1|X) = \frac{1}{1 + e^{\ln \frac{P(X=x|C_1)}{P(X=x|C_2)} + \ln \frac{P(C_1)}{P(C_2)}}} \quad (17)$$

$$= \frac{1}{1 + e^{(\mu_1^T - \mu_2^T) \Sigma^{-1}x - \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1}\mu_2 + \ln \frac{P(C_1)}{P(C_2)}}} \quad (18)$$

假设

$$w = \Sigma^{-1}(\mu_1 - \mu_2) \quad (19)$$

$$w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1}\mu_2 + \ln \frac{P(C_1)}{P(C_2)} \quad (20)$$

因此有

$$\ln \frac{P(C_1)P(X=x|C_1)}{P(C_2)P(X=x|C_2)} = \ln \frac{P(X=x|C_1)}{P(X=x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \quad (21)$$

$$= w^T x + w_0 \quad (22)$$

代入公式(17)中，可得

$$P(C_1|X) = \frac{1}{1 + e^{w^T x + w_0}} \quad (23)$$

由于其归一化的特性，我们有

$$P(C_2|X) = 1 - P(C_1|X) \quad (24)$$

$$= \frac{e^{w^T x + w_0}}{1 + e^{w^T x + w_0}} \quad (25)$$

假设x各维度之间独立，则有各维度间的协方差为0，因此有如下定义

$$\Sigma = \begin{bmatrix} \sigma(x_1, x_1) & \dots & \sigma(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \sigma(x_n, x_1) & \dots & \sigma(x_n, x_n) \end{bmatrix} \quad (26)$$

$$= \begin{bmatrix} \sigma(x_1, x_1) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma(x_n, x_n) \end{bmatrix} \quad (27)$$

因此很容易得到Σ矩阵的逆为

$$\Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma(x_1, x_1)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\sigma(x_n, x_n)} \end{bmatrix} \quad (28)$$

根据对μ的定义，可以展开写为下式

$$\mu_i = \begin{bmatrix} \mu_{i1} \\ \vdots \\ \mu_{in} \end{bmatrix} \quad (29)$$

代入公式(20)，可得

$$\begin{aligned} w_0 &= -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(C_1)}{P(C_2)} \quad (20) \\ &= -\frac{1}{2} \begin{bmatrix} \mu_{11} \\ \vdots \\ \mu_{1n} \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma(x_1, x_1)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\sigma(x_n, x_n)} \end{bmatrix} \begin{bmatrix} \mu_{11} \\ \vdots \\ \mu_{1n} \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \mu_{21} \\ \vdots \\ \mu_{2n} \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma(x_1, x_1)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\sigma(x_n, x_n)} \end{bmatrix} \begin{bmatrix} \mu_{21} \\ \vdots \\ \mu_{2n} \end{bmatrix} + \ln \frac{P(C_1)}{P(C_2)} \\ &= \ln \frac{P(C_1)}{P(C_2)} + \sum_{i=1}^n \frac{\mu_{2i}^2 - \mu_{1i}^2}{2\sigma_i^2} \quad (31) \end{aligned}$$

代入公式(19)，可得

$$w = \Sigma^{-1}(\mu_1 - \mu_2) \quad (19)$$

$$= \begin{bmatrix} \frac{1}{\sigma(x_1, x_1)} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma(x_n, x_n)} \end{bmatrix} \begin{bmatrix} \mu_{11} - \mu_{21} \\ \vdots \\ \mu_{1n} - \mu_{2n} \end{bmatrix} \quad (32)$$

$$= \begin{bmatrix} \frac{\mu_{11} - \mu_{21}}{\sigma_1^2} \\ \vdots \\ \frac{\mu_{1n} - \mu_{2n}}{\sigma_n^2} \end{bmatrix} \quad (33)$$

此处对 $w$ 扩充，将其表示为

$$w = \begin{bmatrix} \frac{1}{\sigma_1^2} \\ \frac{\mu_{11} - \mu_{21}}{\sigma_1^2} \\ \vdots \\ \frac{\mu_{1n} - \mu_{2n}}{\sigma_n^2} \\ \frac{1}{\sigma_n^2} \end{bmatrix} \quad (34)$$

对  $x$  扩充，将其表示为

$$x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad (35)$$

由此，我们可以得到

$$P(C_1|X) = \frac{1}{1 + e^{w^T x}} \quad (36)$$

$$P(C_2|X) = \frac{e^{w^T x}}{1 + e^{w^T x}} \quad (37)$$

因此

$$\frac{P(C_2|X)}{P(C_1|X)} = e^{w^T x} \quad (38)$$

公式(6)的函数被称为sigmoid函数，它具有将参数映射到0,1之间某个值的性质，在 $x = 0$ 的左右， $y$ 分别以很快的速度向0,1逼近，从而得到近似0,1标签的离散值。

根据我们的定义，我们将 $x$ 分到后验概率最大的类中，因此若 $\frac{P(C_2|X)}{P(C_1|X)} > 1$ ，即

$e^{w^T x} > 1$ ，则有 $x$ 属于 $C_2$ 类，否则有 $x$ 属于 $C_1$ 类。

假设当前的数据为 $\{< X^1, Y^1 >, < X^2, Y^2 >, < X^3, Y^3 >, \dots, < X^l, Y^l >\}$ ，其中 $y^i$ 是 $x^i$ 对应的分类( $1 \leq i \leq l$ )。

要计算损失函数，使用极大似然估计(MLE)，计算 $P(< X, Y > |W)$ 是难的问题，因此可以转化为计算极大条件似然估计(MCLE)，只需计算 $P(X|Y, W)$ 。因此使用极大条件似然估计，对参数 $w$ 估计。

$$w_{MCLE} = \arg \max_w \prod_l P(Y^l | X^l, w) \quad (38)$$

于是有似然函数

$$P(c, x|w) = \prod_{i=1}^l P(Y^l | X^l, w) \quad (39)$$

对其取对数可得

$$l(w) = \ln \left( \prod_{i=1}^l P(Y^l | X^l, w) \right) \quad (40)$$

$$= \sum_{i=1}^l \left( \ln(P(Y^l | X^l, w)) \right) \quad (41)$$

$$= \sum_{i=1}^l \left( Y^l \ln(P(Y^l = 1 | X^l, w)) + (1 - Y^l) \ln(P(Y^l = 0 | X^l, w)) \right) \quad (42)$$

$$= \sum_{i=1}^l \left( Y^l \ln(P(Y^l = 1 | X^l, w)) - Y^l \ln(P(Y^l = 0 | X^l, w)) + \ln(P(Y^l = 0 | X^l, w)) \right) \quad (43)$$

$$= \sum_{i=1}^l \left( Y^l \ln \left( \frac{P(Y^l = 1 | X^l, w)}{P(Y^l = 0 | X^l, w)} \right) + \ln(P(Y^l = 0 | X^l, w)) \right) \quad (44)$$

$$= \sum_{i=1}^l \left( Y^l \ln(e^{w^T x}) + \ln(P(Y^l = 0 | X^l, w)) \right) \quad (45)$$

$$= \sum_{i=1}^l \left( Y^l w^T x + \ln \left( \frac{1}{1 + e^{w^T x}} \right) \right) \quad (46)$$

$$= \sum_{i=1}^l \left( Y^l w^T x - \ln(1 + e^{w^T x}) \right) \quad (47)$$

若要最大化公式(40)，只需取其相反数，将其转化为梯度下降，寻找极小值的问题。

我们令

$$L(w) = \sum_{i=1}^l (-Y^l w^T x + \ln(1 + e^{w^T x})) \quad (48)$$

为了避免过拟合，此处添加惩罚项

$$L(w) = \sum_{i=1}^l (-Y^l w^T x + \ln(1 + e^{w^T x})) + \frac{\lambda}{2} \|w\| \quad (50)$$

其中  $\|w\| = w^T w$ .

### 3.2 凸优化：

由于公式(50)是凸函数，因此可以找到其极小值。

(1) 梯度下降法：

类比于lab1，我们可以使用梯度下降法，寻找使 $L(w)$ 最小的 $w$ 的值。

由公式(50)我们可以推出

$$\frac{\partial}{\partial w} L(w) = \sum_{i=1}^l \left( -Y^l x + \frac{e^{w^T x}}{1 + e^{w^T x}} x \right) \quad (51)$$

实现梯度下降法：

$$w = w - \alpha \frac{\partial}{\partial w} L(w) \quad (52)$$

对于公式(52)，选择合适的 $\alpha$ ，使 $L(\theta_{k+1}) < L(\theta_k)$ 。

加上惩罚项得到梯度下降迭代公式

$$w = w - \alpha \left( \sum_{i=1}^l \left( -Y^l x + \frac{e^{w^T x}}{1 + e^{w^T x}} x \right) + \lambda w \right) \quad (53)$$

(2) 牛顿迭代法：

规定

$$f(w) = \sum_{i=1}^l \left( -Y^l x + \frac{e^{w^T x}}{1 + e^{w^T x}} x \right) + \lambda w \quad (54)$$

首先选择一个接近公式(51)零点的 $w_0$ ，计算其对应的 $f(w_0)$ 和切线斜率 $f'(w_0)$ ，然后计算穿过点 $(w_0, f(w_0))$ 并且斜率为 $f'(w_0)$ 的直线和 $x$ 轴的交点 $w$ 的坐标，即对如下方程求解：

$$0 = (w - w_0)f'(w_0) + f(w_0) \quad (55)$$

我们新求得的 $w$ 坐标命名为 $w_1$ ，通常 $w_1$ 会比 $w_0$ 更接近方程 $f(w) = 0$ 的解，因此可以用 $w_1$ 进行下一轮迭代。

迭代过程可以简化为如下过程：

$$w_{n+1} = w_n - \frac{f(w_n)}{f'(w_n)} \quad (56)$$

因此我们可以得到：

$$f'(w) = \sum_{i=1}^l \left( \frac{x e^{w^T x}}{1 + e^{w^T x}} x^T \right) + \lambda \quad (57)$$

### 3.3 溢出问题:

对于公式(50),可能会引起溢出问题,如图。

```
main
D:\Anaconda3\envs\MLExperiment\python.exe D:/Code/PycharmProjects/MachineLearningExperiment/Experiment2/main.py
D:\Code\PycharmProjects\MachineLearningExperiment\Experiment2\ComputeCost.py:2: RuntimeWarning: overflow encountered in exp
cost = cost+(-Y[i]* np.dot(theta,np.reshape(X[i],(-1,1))) + log(1+np.exp(np.dot(theta,np.reshape(X[i],(-1,1))))),e)[0]
D:\Code\PycharmProjects\MachineLearningExperiment\Experiment2\GradientDescent.py:23: RuntimeWarning: overflow encountered in exp
ExpWX = np.exp(np.matmul(theta, np.reshape(X[j], (-1, 1))))[0]
D:\Code\PycharmProjects\MachineLearningExperiment\Experiment2\GradientDescent.py:24: RuntimeWarning: invalid value encountered in double_scalars
iter = iter + (-X[j] * Y[j] + (ExpWX / (1 + ExpWX)) * X[j])
```

因此,我们需要解决以上问题。

$$w = w - \alpha \left( \sum_{i=1}^l \left( -Y^i x + \frac{e^{w^T x}}{1 + e^{w^T x}} x \right) + \lambda w \right) \quad (53)$$

我们有公式(53)的迭代公式,需要对公式(53)改进。

可能存在以下情况:

当 $w^T x > 0$ 且过大时, $\text{np.exp}(w^T x)$ 会向上溢出,因此此时可以将公式 53 修改为如下公式。

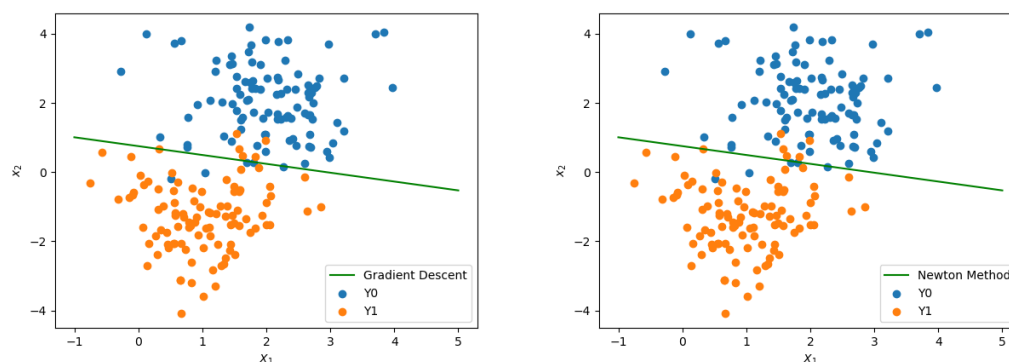
$$w = w - \alpha \left( \sum_{i=1}^l \left( -Y^i x + \frac{1}{1 + e^{-w^T x}} x \right) + \lambda w \right) \quad (58)$$

## 四、实验结果与分析

### 4.1 生成高斯分布的两类点,两类点中各特征相互独立:

设置训练集和测试集的比例为 3:7,生成特征相互独立的数据,每个类别 100 个数据。左侧使用梯度下降法,右侧使用牛顿迭代法。

绘图中样本点为训练集和测试集的合集。





```

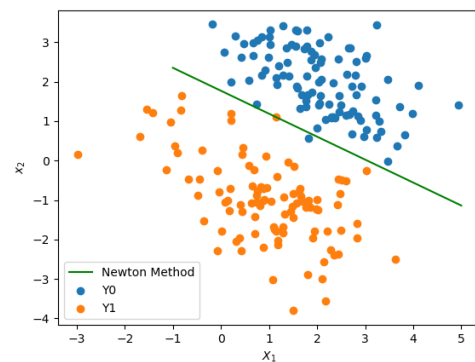
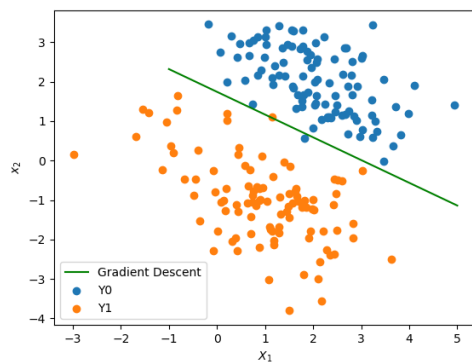
1004
迭代次数为1003
测试集总数60，预测正确个数59
正确率0.9833333333333333
3311
迭代次数为3310
测试集总数60，预测正确个数59
正确率0.9833333333333333

```

实验结果为：使用梯度下降法和使用牛顿迭代法得出曲线近乎一致，测试集共 60 个样本，预测正确 59 个。预测正确率为 98.3%。

#### 4.2 生成高斯分布的两类点，两类点中各特征不相互独立：

设置训练集和测试集的比例为 3:7，生成两类各特征不相互独立的点，每个类别 100 个数据。左侧使用梯度下降法，右侧使用牛顿迭代法。



```

81233
迭代次数为81232
测试集总数60，预测正确个数59
正确率0.9833333333333333
53743
迭代次数为53742
测试集总数60，预测正确个数59
正确率0.9833333333333333

```

实验结果为：使用梯度下降法和使用牛顿迭代法得出曲线近乎一致，测试集共 60 个样本，预测正确 59 个。预测正确率为 98.3%。

#### 4.3 使用uci上的数据集测试

##### (1) data\_banknote\_authentication数据集

```
main x
D:\Anaconda3\envs\MLExperiment\python.exe D:/C
./DataSet/data_banknote_authentication.csv
梯度下降法测试结果测试结果：
测试集总数412，预测正确个数406
正确率0.9854368932038835
./DataSet/data_banknote_authentication.csv
牛顿迭代法测试结果测试结果：
测试集总数412，预测正确个数405
正确率0.9830097087378641
```

实验数据：

- 第一列：图像经小波变换后的方差(variance) (连续值)；
- 第二列：图像经小波变换后的偏态(skewness) (连续值)；
- 第三列：图像经小波变换后的峰度(kurtosis) (连续值)；
- 第四列：图像的熵(entropy) (连续值)；
- 第五列：钞票所属的类别(整数，0 或 1)。

实验结果为：使用梯度下降法和使用牛顿迭代法测试结果几乎一致，测试集共 412 个样本，使用梯度下降法预测正确 406 个，预测正确率为 98.54%；使用牛顿迭代法预测正确 405 个，预测正确率为 98.3%

## (2) sonar数据集

```
main x
D:\Anaconda3\envs\MLExperiment\python.exe
./DataSet/sonar.csv
梯度下降法测试结果测试结果：
测试集总数63，预测正确个数50
正确率0.7936507936507936
./DataSet/sonar.csv
牛顿迭代法测试结果测试结果：
测试集总数63，预测正确个数50
正确率0.7936507936507936
```

实验结果为：使用梯度下降法和使用牛顿迭代法测试结果一致，测试集共 63 个样本，使用梯度下降法预测正确 50 个，预测正确率为 79.36%；使用牛顿迭代法预测正确 50 个，预测正确率为 79.36%

## 五、结论

1. 使用梯度下降法和牛顿迭代法所得分类结果近乎相同。
2. 选择不同的数据可能对分类结果有很不同的影响。

## 六、参考文献

- [1] Pattern Recognition and Machine Learning
- [2] 机器学习 周志华著 北京：清华大学出版社，2016 年 1 月.
- [2] <https://zhuanlan.zhihu.com/p/76639936>

## 七、附录：源代码（带注释）

main.py 主程序

```
1. from PredictResult import predictResult
2.
3. if __name__ == '__main__':
4.     # 使用生成的数据点求解
5.     # 符合朴素贝叶斯分布的点求解
6.     # predictResult(0)
7.     # 不符合朴素贝叶斯分布的点求解
8.     # predictResult(1)
9.     # 使用梯度下降法对保存的文件 求解
10.    predictResult(2, "./DataSet/data_banknote_authentication.csv")
11.    # predictResult(2, "./DataSet/sonar.csv")
12.    # predictResult(2, "./DataSet/heart.csv")
13.    # 使用牛顿迭代法对保存的文件 求解
14.    predictResult(3, "./DataSet/data_banknote_authentication.csv")
15.    # predictResult(3, "./DataSet/sonar.csv")
16.    # predictResult(3, "./DataSet/heart.csv")
```

ComputeCost.py 计算代价值

```
1. import numpy as np
2. from math import log
3. from math import e
4.
```

```

5. def computeCost(X, Y, theta, Lambda):
6.     cost = 0
7.     for i in range(0, X.shape[1]):
8.         cost = cost + (
9.             -Y[i] * np.dot(theta, np.reshape(X[i], (-
10. 1, 1))) + log(1 + np.exp(np.dot(theta, np.reshape(X[i], (-1, 1)))))
11.             e)[0]
12.     cost = cost + Lambda / 2 * np.dot(theta, np.reshape(theta, (-1, 1)))
13.     return cost

```

DataFromFile.py 从文件中读取数据

```

1. import pandas as pd
2. import numpy as np
3. from GenerateData import addUnitColumn
4.
5.
6. def readFromFile(DocumentName):
7.     df = pd.read_csv(DocumentName)
8.     df.rename(columns={df.columns.array[df.columns.shape[0] - 1]: 'Predict'},
9.               , inplace=True)
10.    Y = df['Predict']
11.    X = df.drop('Predict', axis=1)
12.    X = np.array(X.values)
13.    X = np.reshape(X, (1, X.shape[0], X.shape[1]))
14.    X = addUnitColumn(X)
15.    Y = Y.values
16.    return X, Y

```

Draw.py 绘画曲线 绘点

```

1. import matplotlib.pyplot as plt
2. import numpy as np
3.
4.
5. def drawTheta(Start, End, theta, method):
6.     # number
7.     number = 1000
8.     X = np.linspace(start=Start, stop=End, num=number)
9.     X = np.reshape(X, (-1, 1))
10.    UnitMatrix = np.reshape(np.ones(number), (-1, 1))
11.    Y = (theta[0] * UnitMatrix + theta[1] * X) * (-1) * (1 / theta[2])
12.    if method == 0:
13.        plt.plot(X, Y, 'g', label="Gradient Descent")

```

```

14.     elif method == 1:
15.         plt.plot(X, Y, 'g', label="Newton Method")
16.         plt.xlabel('$X_{1}$', fontsize=10)
17.         plt.ylabel('$x_{2}$', fontsize=10)
18.         plt.legend()
19.         plt.show()
20.
21.
22. def plotScatter(X1Y0, X2Y0, X1Y1, X2Y1):
23.     plt.scatter(X1Y0, X2Y0, label='Y{}'.format(0))
24.     plt.scatter(X1Y1, X2Y1, label='Y{}'.format(1))

```

#### GenerateData.py

```

1. import numpy as np
2. from sklearn.model_selection import train_test_split
3.
4.
5. def generate2DimensionalData(Mu1, Mu2, cov11, cov12, cov21, cov22, Num, noiseSigma1, noiseSigma2):
6.     """
7.     生成二维高斯分布数据
8.     :param Mu1: 维度 1 的平均值
9.     :param Mu2: 维度 2 的平均值
10.    :param cov11: 维度 1 的方差
11.    :param cov12: 维度 12 的协方差
12.    :param cov21: 维度 21 的协方差
13.    :param cov22: 维度 2 的方差
14.    :param Num: 生成样本点的个数
15.    :param Class: 类别
16.    :param noiseSigma1: 维度 1 噪声的方差
17.    :param noiseSigma2: 维度 2 噪声的方差
18.    :return: 分别返回两个维度的样本点
19.    """
20.    mean = np.array([Mu1, Mu2])
21.    cov = np.array([[cov11, cov12], [cov21, cov22]])
22.    Data = np.random.multivariate_normal(mean, cov, Num)
23.    X1 = Data[:, 0]
24.    X2 = Data[:, 1]
25.    # 添加高斯噪声
26.    GuassNoise1 = np.random.normal(0, scale=noiseSigma1, size=Num)
27.    GuassNoise2 = np.random.normal(0, scale=noiseSigma2, size=Num)
28.    X1 = X1 + GuassNoise1
29.    X2 = X2 + GuassNoise2

```

```

30.     return X1, X2
31.
32.
33. def generateData(method):
34.     """
35.     生成两类数据并满足朴素贝叶斯分布
36.     :param method: 分别生成符合朴素贝叶斯分布和不符合朴素贝叶斯分布的数据
37.     :return: 生成数据
38.     """
39.     if method == 1:
40.         # Y0Mu1 = 2
41.         # Y0Mu2 = 2
42.         # Y0cov11 = 1
43.         # Y0cov12 = -0.5
44.         # Y0cov21 = -0.5
45.         # Y0cov22 = 1
46.         # Y0Num = 200
47.         # Y0noiseSigma1 = 0.2
48.         # Y0noiseSigma2 = 0.4
49.         # X1Y, X2Y = generate2DimensionalData(Y0Mu1, Y0Mu2, Y0cov11, Y0cov12
        , Y0cov21, Y0cov22, Y0Num,
50.         #                                     Y0noiseSigma1, Y0noiseSigma2
        )
51.         # X1Y0 = X1Y[0:np.int64(Y0Num/2)]
52.         # X1Y1 = X1Y[np.int64(Y0Num/2):np.int64(Y0Num)]
53.         # X2Y0 = X2Y[0:np.int64(Y0Num/2)]
54.         # X2Y1 = X2Y[np.int64(Y0Num/2):np.int64(Y0Num)]
55.         # XY0 = np.dstack((X1Y0, X2Y0))
56.         # XY1 = np.dstack((X1Y1, X2Y1))
57.         # XY0 = addUnitColumn(XY0)
58.         # XY1 = addUnitColumn(XY1)
59.         # X = np.vstack((XY0, XY1))
60.         # Y_zero = np.zeros([XY0.shape[0]])
61.         # Y_ones = np.ones([XY1.shape[0]])
62.         # Y_zero = np.reshape(Y_zero, (-1, 1))
63.         # Y_ones = np.reshape(Y_ones, (-1, 1))
64.         # Y = np.vstack((Y_zero, Y_ones))
65.         # Y = np.reshape(Y, (-1))
66.         # return X, Y, X1Y0, X2Y0, X1Y1, X2Y1
67.
68.         # Y0 类生成数据的参数
69.         Y0Mu1 = 2
70.         Y0Mu2 = 2
71.         Y0cov11 = 1

```

```

72.         Y0cov12 = -0.5
73.         Y0cov21 = -0.5
74.         Y0cov22 = 0.7
75.         Y0Num = 100
76.         Y0noiseSigma1 = 0.2
77.         Y0noiseSigma2 = 0.4
78.         # Y1 类生成数据的参数
79.         Y1Mu1 = 1
80.         Y1Mu2 = -1
81.         Y1cov11 = 1
82.         Y1cov12 = -0.6
83.         Y1cov21 = -0.6
84.         Y1cov22 = 1.3
85.         Y1Num = 100
86.         Y1noiseSigma1 = 0.1
87.         Y1noiseSigma2 = 0.2
88.         X1Y0, X2Y0 = generate2DimensionalData(Y0Mu1, Y0Mu2, Y0cov11, Y0cov12
, Y0cov21, Y0cov22, Y0Num,
89.                                     Y0noiseSigma1, Y0noiseSigma2)

90.         X1Y1, X2Y1 = generate2DimensionalData(Y1Mu1, Y1Mu2, Y1cov11, Y1cov12
, Y1cov21, Y1cov22, Y1Num,
91.                                     Y1noiseSigma1, Y1noiseSigma2)

92.         XY0 = np.dstack((X1Y0, X2Y0))
93.         XY1 = np.dstack((X1Y1, X2Y1))
94.         XY0 = addUnitColumn(XY0)
95.         XY1 = addUnitColumn(XY1)
96.         X = np.vstack((XY0, XY1))
97.         Y_zero = np.zeros([XY0.shape[0]])
98.         Y_ones = np.ones([XY1.shape[0]])
99.         Y_zero = np.reshape(Y_zero, (-1, 1))
100.        Y_ones = np.reshape(Y_ones, (-1, 1))
101.        Y = np.vstack((Y_zero, Y_ones))
102.        Y = np.reshape(Y, (-1))
103.        return X, Y, X1Y0, X2Y0, X1Y1, X2Y1
104.
105.
106.
107.    elif method == 0:
108.        # Y0 类生成数据的参数
109.        Y0Mu1 = 2
110.        Y0Mu2 = 2
111.        Y0cov11 = 0.5

```

```

112.         Y0cov12 = 0
113.         Y0cov21 = 0
114.         Y0cov22 = 1
115.         Y0Num = 100
116.         Y0noiseSigma1 = 0.2
117.         Y0noiseSigma2 = 0.4
118.         # Y1 类生成数据的参数
119.         Y1Mu1 = 1
120.         Y1Mu2 = -1
121.         Y1cov11 = 0.5
122.         Y1cov12 = 0
123.         Y1cov21 = 0
124.         Y1cov22 = 1
125.         Y1Num = 100
126.         Y1noiseSigma1 = 0.1
127.         Y1noiseSigma2 = 0.2
128.         X1Y0, X2Y0 = generate2DimensionalData(Y0Mu1, Y0Mu2, Y0cov11, Y0cov1
129.         2, Y0cov21, Y0cov22, Y0Num,
130.         Y0noiseSigma1, Y0noiseSigma2)
131.
132.         X1Y1, X2Y1 = generate2DimensionalData(Y1Mu1, Y1Mu2, Y1cov11, Y1cov1
133.         2, Y1cov21, Y1cov22, Y1Num,
134.         Y1noiseSigma1, Y1noiseSigma2)
135.
136.         XY0 = np.dstack((X1Y0, X2Y0))
137.         XY1 = np.dstack((X1Y1, X2Y1))
138.         XY0 = addUnitColumn(XY0)
139.         XY1 = addUnitColumn(XY1)
140.         X = np.vstack((XY0, XY1))
141.         Y_zero = np.zeros([XY0.shape[0]])
142.         Y_ones = np.ones([XY1.shape[0]])
143.         Y_zero = np.reshape(Y_zero, (-1, 1))
144.         Y_ones = np.reshape(Y_ones, (-1, 1))
145.         Y = np.vstack((Y_zero, Y_ones))
146.         Y = np.reshape(Y, (-1))
147.         return X, Y, X1Y0, X2Y0, X1Y1, X2Y1
148.
149.
150. def addUnitColumn(X):
151.     """
152.     为 x 添加全为 1 的列
153.     :param X: 待添加的矩阵
154.     :return: 添加后的矩阵
155.     """

```



```

152.     UnitMatrix = np.ones(X.shape[1])
153.     # UnitMatrix = np.reshape(UnitMatrix, (X.shape[0], -1))
154.     X = np.dstack((UnitMatrix, X))
155.     X = np.reshape(X, (X.shape[1], X.shape[2]))
156.     return X
157.
158.
159. def TrainTestSplit(X, Y):
160.     """
161.     将数据集分为训练集和测试集
162.     :param X: 训练特征数据
163.     :param Y: 训练数据对应分类
164.     :return: 分开后的数据集
165.     """
166.     X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.
        7)
167.     return X_train, X_test, Y_train, Y_test

```

#### GradientDescent.py 梯度下降法

```

1. import numpy as np
2. from ComputeCost import computeCost
3.
4.
5. def gradientDescent(X, Y, alpha, Lambda, precision, iterNum):
6.     """
7.     使用梯度下降法进行迭代
8.     :param X: 数据特征
9.     :param Y: 数据分类
10.    :param alpha: 学习率
11.    :param Lambda: 惩罚项系数
12.    :param precision: 迭代精度
13.    :param iterNum: 迭代次数
14.    :return:
15.    """
16.    theta = np.zeros(X.shape[1])
17.    descentStore = np.zeros(iterNum + 1)
18.    iterStore = np.zeros(iterNum + 1)
19.    iterStore[0] = np.inf
20.    for i in range(1, iterNum):
21.        iter = 0
22.        for j in range(0, X.shape[0]):
23.
24.            # WX = np.matmul(theta, np.reshape(X[j], (-1, 1)))

```

```

25.         # ExpWX = np.exp(WX)[0]
26.         # iter = iter + (-X[j] * Y[j] + (ExpWX/ (1 + ExpWX)) * X[j])
27.
28.         WX = np.matmul(theta, np.reshape(X[j], (-1, 1)))
29.         ExpWX = np.exp(WX)[0]
30.
31.         if ExpWX == np.inf:
32.             iter = iter + (-X[j] * Y[j] + X[j])
33.         else:
34.             iter = iter + (-
X[j] * Y[j] + (ExpWX / (1 + ExpWX)) * X[j])
35.
36.         # if i % 1000 == 0:
37.         #     print("np.sum(iter)")
38.         #     print(abs(np.sum(iter)))
39.         iterStore[i + 1] = np.sum(iter)
40.         # if abs(iterStore[i+1]) > abs(iterStore[i]):
41.         #     alpha = alpha/2
42.         #     print("当前 alpha={}".format(alpha))
43.         # print("iterStore[i+1]")
44.         # print(iterStore[i+1])
45.         # print("iterStore[i+1]-iterStore[i]")
46.         # print(abs(iterStore[i+1]-iterStore[i]))
47.         # if abs(iterStore[i+1]-iterStore[i]) <= precision:
48.         #     # print(i)
49.         #     print("迭代次数为{}".format(i - 1))
50.         #     print(abs(descentStore[i]-descentStore[i-1]))
51.         #     break
52.         iter = iter + Lambda * theta
53.         theta = theta - alpha * iter
54.         descentStore[i] = computeCost(X, Y, theta, Lambda)
55.         # if abs(descentStore[i]-descentStore[i-1])<=precision:
56.         #     print(i)
57.         #     print("迭代次数为{}".format(i-1))
58.         #     # print("descentStore[i]-descentStore[i-1]")
59.         #     # print(abs(descentStore[i]-descentStore[i-1]))
60.         #     break
61.     return theta

```

NewtonMethod.py 牛顿迭代法

```

1. import numpy as np
2.
3.

```

```

4. def newtonMethod(X, Y, Lambda, precision, iterNum):
5.     """
6.     使用牛顿迭代法进行迭代
7.     :param X: 数据特征
8.     :param Y: 数据分类
9.     :param Lambda: 惩罚项系数
10.    :param precision: 预测精度
11.    :param iterNum: 迭代次数
12.    :return: theta
13.    """
14.    theta = np.zeros(X.shape[1])
15.    iterStore = np.zeros(iterNum + 1)
16.    iterStore[0] = np.inf
17.    for i in range(0, iterNum):
18.        iterComputeDifferential = 0
19.        iterCompute = 0
20.        for j in range(0, X.shape[0]):
21.            ExpWX = np.exp(np.dot(theta, np.reshape(X[j], (-1, 1))))
22.            iterComputeDifferential = iterComputeDifferential + (ExpWX * np.
                dot(X[j], np.reshape(X[j], (-1, 1)))) / (
23.                1 + ExpWX)
24.
25.            ExpWX = np.exp(np.matmul(theta, np.reshape(X[j], (-1, 1))))[0]
26.            iterCompute = iterCompute + (-
                X[j] * Y[j] + (ExpWX / (1 + ExpWX)) * X[j])
27.
28.        iterStore[i + 1] = np.sum(iterCompute)
29.        # if abs(iterStore[i+1]-iterStore[i]) <= precision:
30.        #     print(i)
31.        #     print("迭代次数为{}".format(i - 1))
32.        #     break
33.
34.        if abs(iterStore[i + 1] - iterStore[i]) <= precision:
35.            print(i)
36.            print("迭代次数为{}".format(i - 1))
37.            break
38.
39.        iterCompute = iterCompute + Lambda * theta
40.        iterComputeDifferential = iterComputeDifferential + Lambda
41.        theta = theta - iterCompute / iterComputeDifferential
42.    return theta

```

PredictResult.py 预测结果

```

1. from GradientDescent import gradientDescent
2. from GenerateData import TrainTestSplit
3. from DataFromFile import readFromFile
4. from VerificationPredict import predictPrecision
5. from NewtonMethod import newtonMethod
6. from GenerateData import generateData
7. from Draw import drawTheta
8. from Draw import plotScatter
9. import numpy as np
10. import matplotlib.pyplot as plt
11.
12.
13. def predictResult(method, FileName=""):
14.     """
15.     从文件中读取数据并根据不同的 method 使用不同方法求得极小值
16.     并通过调用自己写的精度函数验证预测准确率
17.     :param method: 使用的方法, 0: 对于符合朴素贝叶斯分布的数据 1: 对于不符合朴素贝
        叶斯分布的数据 2: 梯度下降法 3: 牛顿迭代法
18.     :param FileName: 文件名称
19.     """
20.
21.     if method == 0:
22.         X, Y, X1Y0, X2Y0, X1Y1, X2Y1 = generateData(0)
23.         X_train, X_test, Y_train, Y_test = TrainTestSplit(X, Y)
24.         theta = gradientDescent(X_train, Y_train, 0.01, 1e-7, 1e-7, 10000)
25.         plotScatter(X1Y0, X2Y0, X1Y1, X2Y1)
26.         drawTheta(-1, 5, theta, 0)
27.         predictPrecision(theta, X_test, Y_test)
28.         theta = newtonMethod(X_train, Y_train, 1e-7, 1e-7, 10000)
29.         plotScatter(X1Y0, X2Y0, X1Y1, X2Y1)
30.         drawTheta(-1, 5, theta, 1)
31.         predictPrecision(theta, X_test, Y_test)
32.
33.     elif method == 1:
34.         X, Y, X1Y0, X2Y0, X1Y1, X2Y1 = generateData(1)
35.         X_train, X_test, Y_train, Y_test = TrainTestSplit(X, Y)
36.         theta = gradientDescent(X_train, Y_train, 0.001, 1e-7, 1e-
            7, 10000)
37.         plotScatter(X1Y0, X2Y0, X1Y1, X2Y1)
38.         drawTheta(-1, 5, theta, 0)
39.         predictPrecision(theta, X_test, Y_test)
40.         theta = newtonMethod(X_train, Y_train, 1e-7, 1e-7, 10000)
41.         plotScatter(X1Y0, X2Y0, X1Y1, X2Y1)
42.         drawTheta(-1, 5, theta, 1)

```

```

43.         predictPrecision(theta, X_test, Y_test)
44.
45.
46.     elif method == 2:
47.         X, Y = readFromFile(fileName)
48.         X_train, X_test, Y_train, Y_test = TrainTestSplit(X, Y)
49.         theta = gradientDescent(X_train, Y_train, 0.0001, 1e-7, 1e-
            8, 50000)
50.         print(fileName)
51.         print("梯度下降法测试结果测试结果: ")
52.         predictPrecision(theta, X_test, Y_test)
53.     elif method == 3:
54.         X, Y = readFromFile(fileName)
55.         X_train, X_test, Y_train, Y_test = TrainTestSplit(X, Y)
56.         theta = newtonMethod(X_train, Y_train, 1e-7, 1e-8, 50000)
57.         print(fileName)
58.         print("牛顿迭代法测试结果测试结果: ")
59.         predictPrecision(theta, X_test, Y_test)

```

VerificationPredict.py 验证测试集结果

```

1. import numpy as np
2.
3.
4. def predictPrecision(theta, X_test, Y_test):
5.     """
6.     计算预测的准确率
7.     :param theta: 迭代出来的系数
8.     :param X_test: 待测试数据
9.     :param Y_test: 数据集所属真实分类
10.    """
11.
12.    predict = np.dot(X_test, theta)
13.    accuracy = 0
14.    for i in range(0, predict.shape[0]):
15.        truth = Y_test[i]
16.        if (predict[i] >= 0 and truth == 1) or (predict[i] <= 0 and truth ==
            0):
17.            accuracy += 1
18.    print("测试集总数{}, 预测正确个数{}".format(predict.shape[0], accuracy))
19.    print("正确率{}".format(accuracy / predict.shape[0]))

```