

大家好，这篇是有关Learning from data第三章习题的详解，这一章主要介绍了线性回归，Logistic回归以及特征转换。

我的github地址：

<https://github.com/Doraemonzzz>

个人主页：

<http://doraemonzzz.com/>

参考资料：

<https://blog.csdn.net/a1015553840/article/details/51085129>

<http://www.vynguyen.net/category/study/machine-learning/page/6/>

<http://book.caltech.edu/bookforum/index.php>

<http://beader.me/mlnotebook/>

Chapter 3 The Linear Model

Part 1:Exercise

Exercise 3.1 (Page 79)

Will PLA ever stop updating if the data is not linearly separable?

如果数据不是线性可分的，那么PLA不会终止，因为PLA每次挑选一个错分的数据做更新。

Exercise 3.2 (Page 80)

Take $d = 2$ and create a data set \mathcal{D} of size $N = 100$ that is not linearly separable. You can do so by first choosing a random line in the plane as your target function and the inputs x_n of the data set as random points in the plane. Then, evaluate the target function on each x_n to get the corresponding output y_n . Finally, flip the labels of $\frac{N}{10}$ randomly selected y_n 's and the data set will likely become non separable.

Now, try the pocket algorithm on your data set using $T = 1,000$ iterations. Repeat the experiment 20 times. Then, plot the average $E_{\text{in}}(w(t))$ and the average $E_{\text{in}}(\hat{w})$ (which is also a function of t) on the same figure and see how they behave when t increases. Similarly, use a test set of size 1,000 and plot a figure to show how $E_{\text{out}}(w(t))$ and $E_{\text{out}}(\hat{w})$ behave.

首先构造100个线性不可分的点，题目给出的方法是先随意取一条直线（这里我们选择的直线是 $y = x$ ），然后根据这条直线给出 $N = 100$ 个线性可分的点，再随机挑选其中 $\frac{N}{10} = 10$ 个数据，翻转他们的 y_n ，然后使用pocket PLA进行1000次迭代，画出 $E_{\text{in}}(w(t))$ ，以及平均值 $E_{\text{in}}(\hat{w})$ 随迭代次数的变化，在这个过程中同样计算1000个测试数据的 $E_{\text{out}}(w(t))$ 和平均值 $E_{\text{out}}(\hat{w})$ 并作图，这里的测试数据和之前的数据符合的规则应该一致，都是同样有10%的数据经过翻转。

```
# -*- coding: utf-8 -*-  
.....
```

Created on Sat Mar 2 10:59:02 2019

```

@author: qinzhen
"""

import numpy as np
import matplotlib.pyplot as plt
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号
from helper import Pocket_PLA
from helper import preprocess

#Step产生数据
def generatedata(n, flag=True):
    """
    产生n组数据，10%为噪声，这里直线选择为y=x，y>x则标签为1，否则为-1
    x，y的范围都介于(-2，2)，flag用于判断是否产生噪声
    """
    #产生x
    X = np.random.uniform(-2, 2, (n, 2))
    #计算标签
    y = np.ones(n)
    y[X[:, 0] <= X[:, 1]] = -1
    if flag:
        #让前10%数据误分
        y[: n//10] *= -1
    #合并数据
    Data = np.c_[X, y]
    #打乱数据
    np.random.shuffle(Data)

    return Data

```

产生100个以及1000的点，然后分别作图看看。

```

#Step2展示数据
#生成训练数据
D_train = generatedata(100)
X_train = D_train[:, :-1]
y_train = D_train[:, -1]
#后面两步是为了作图
#标签为+1的点
train_px = X_train[y_train > 0][:, 0]
train_py = X_train[y_train > 0][:, 1]
#标签为-1的点
train_nx = X_train[y_train < 0][:, 0]
train_ny = X_train[y_train < 0][:, 1]

#生成测试数据
D_test = generatedata(1000)
X_test = D_test[:, :-1]
y_test = D_test[:, -1]
#标签为+1的点
test_px = X_test[y_test > 0][:, 0]

```

```

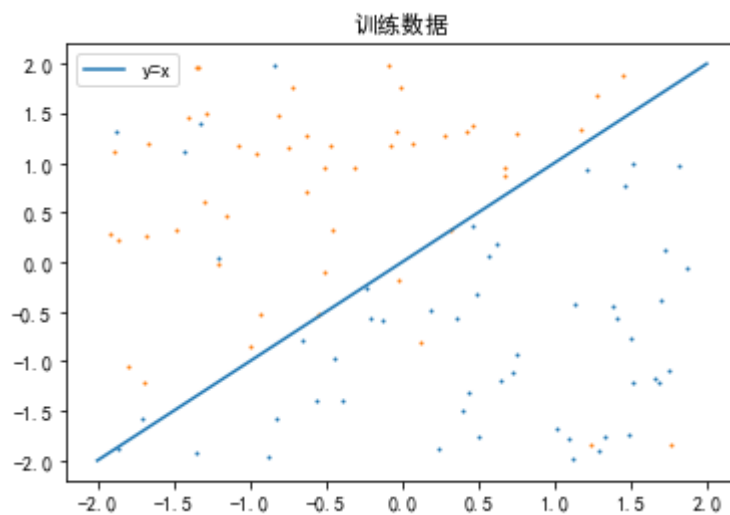
test_py = X_test[y_test > 0][:, 1]
#标签为-1的点
test_nx = X_test[y_test < 0][:, 0]
test_ny = X_test[y_test < 0][:, 1]

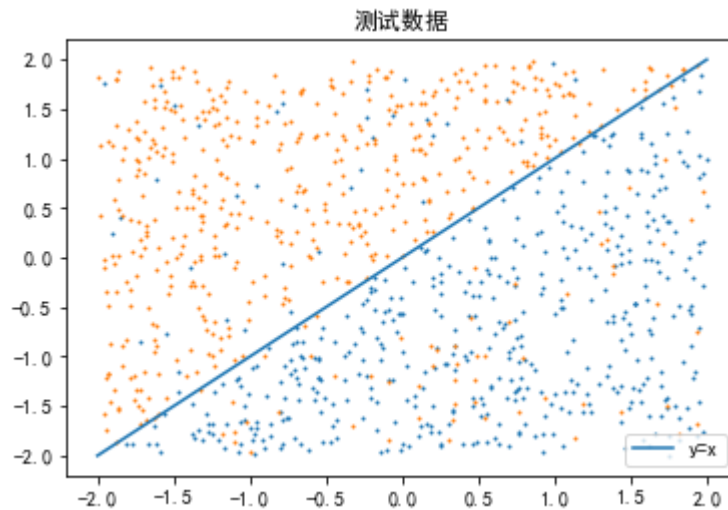
x = [-2, 2]
y = [-2, 2]

#训练数据
plt.scatter(train_px, train_py, s=1)
plt.scatter(train_nx, train_ny, s=1)
plt.title('训练数据')
plt.plot(x, y, label='y=x')
plt.legend()
plt.show()

#测试数据
plt.scatter(test_px, test_py, s=1)
plt.scatter(test_nx, test_ny, s=1)
plt.title('测试数据')
plt.plot(x, y, label='y=x')
plt.legend()
plt.show()

```





接着定义Pocket PLA，由于这该算法会复用，这里将其写在helper.py文件中，后续使用的时候只要导入

```
# -*- coding: utf-8 -*-
"""
Created on Wed Feb 13 01:31:19 2019

@author: qinzhen
"""

import numpy as np
import matplotlib.pyplot as plt
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号

def count(X, y, w):
    """
    统计错误数量
    """
    num = np.sum(X.dot(w) * y <= 0)
    return np.sum(num)

def preprocess(data):
    """
    数据预处理
    """
    #获取维度
    n, d = data.shape
    #分离x
    X = data[:, :-1]
    #添加偏置项1
    X = np.c_[np.ones(n), X]
    #分离y
    y = data[:, -1]

    return X, y

def Pocket_PLA(X, y, eta=1, max_step=np.inf):
    """
```

```

Pocket_PLA算法, X, y为输入数据, eta为步长, 默认为1, max_step为最多迭代次数, 默认为无穷
"""
#获得数据维度
n, d = X.shape
#初始化
w = np.zeros(d)
#记录最优向量
w0 = np.zeros(d)
#记录次数
t = 0
#记录最少错误数量
error = count(X, y, w0)
#记录元素的下标
i = 0
#记录每一步的w
w = []
#记录最优w
w_hat = []
while (error != 0 and t < max_step):
    if np.sign(X[i, :].dot(w) * y[i]) <= 0:
        w += eta * y[i] * X[i, :]
        #迭代次数增加
        t += 1
        #记录当前错误
        error_now = count(X, y, w)
        if error_now < error:
            error = error_now
            w0 = np.copy(w)

        #记录最优w
        w_hat.append(np.copy(w0))
        #记录w
        w.append(np.copy(w))
    #移动到下一个元素
    i += 1
    #如果达到n, 则重置为0
    if i == n:
        i = 0
return np.array(w), np.array(w_hat), w0, error

```

我们看下训练结果

```

#Step3训练数据
#n为迭代次数, k为步长, N为实验次数
n = 1000
k = 1
N = 20

def experiment(n, k):
    """
    模拟一次实验, n为迭代次数, k为步长
    """
    #训练数据

```

```

D_train = generatedata(100)
X_train, y_train = preprocess(D_train)

#测试数据
D_test = generatedata(1000)
X_test, y_test = preprocess(D_test)

#训练模型
W, W_hat, w_hat, error = Pocket_PLA(X_train, y_train, k, n)

#计算错误率
ein = np.mean(np.sign(W.dot(X_train.T)) != y_train, axis=1)
ein_hat = np.mean(np.sign(W_hat.dot(X_train.T)) != y_train, axis=1)
eout = np.mean(np.sign(W.dot(X_test.T)) != y_test, axis=1)
eout_hat = np.mean(np.sign(W_hat.dot(X_test.T)) != y_test, axis=1)
return ein, ein_hat, eout, eout_hat

#存储结果
Ein = np.zeros(n)
Ein_hat = np.zeros(n)
Eout = np.zeros(n)
Eout_hat = np.zeros(n)

for i in range(N):
    ein, ein_hat, eout, eout_hat = experiment(n, k)
    Ein += ein
    Ein_hat += ein_hat
    Eout += eout
    Eout_hat += eout_hat

#计算均值
Ein /= N
Ein_hat /= N
Eout /= N
Eout_hat /= N

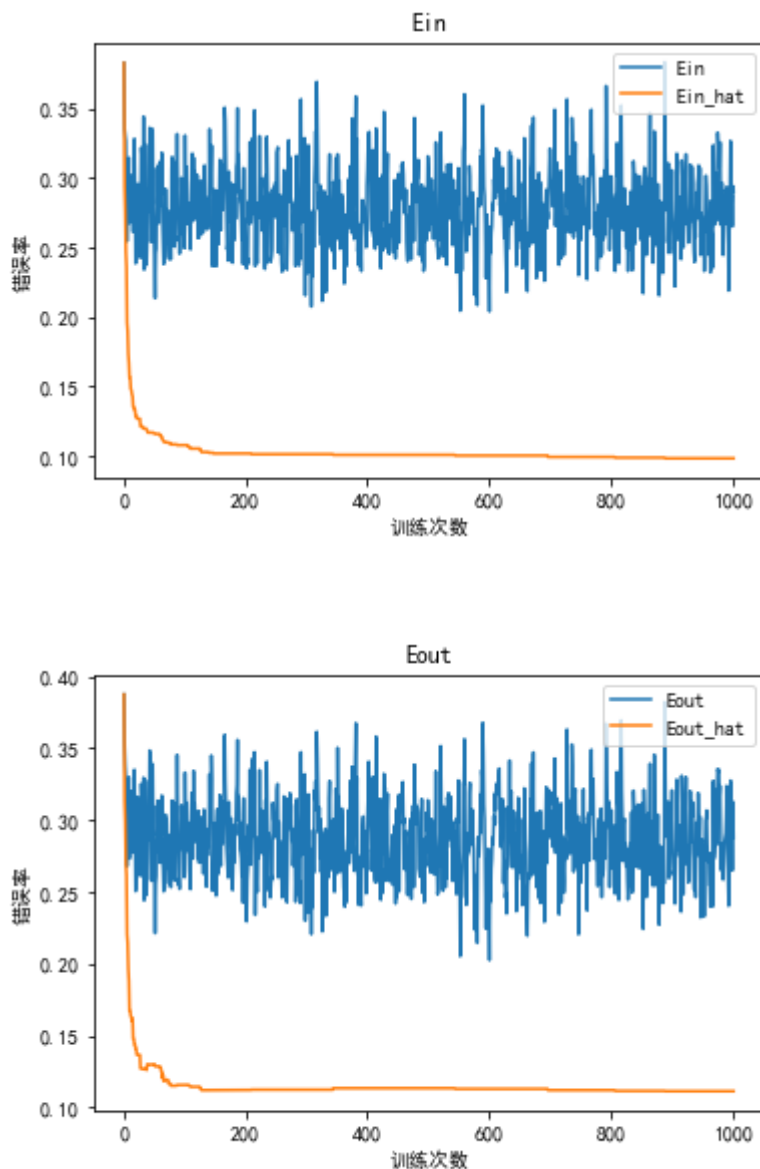
t = np.arange(1, n+1)

#Ein作图
plt.plot(t, Ein, label='Ein')
plt.plot(t, Ein_hat, label='Ein_hat')
plt.title('Ein')
plt.xlabel('训练次数')
plt.ylabel('错误率')
plt.legend()
plt.show()

#Eout作图
plt.plot(t, Eout, label='Eout')
plt.plot(t, Eout_hat, label='Eout_hat')
plt.title('Eout')
plt.xlabel('训练次数')
plt.ylabel('错误率')

```

```
plt.legend()
plt.show()
```



可以看到 $E_{\text{in}}(\hat{w})$, $E_{\text{out}}(\hat{w})$ 随着训练次数增加逐渐减少, 而 $E_{\text{in}}(w(t))$, $E_{\text{out}}(w(t))$ 随着训练次数增加则波动较大。

Exercise 3.3 (Page 87)

Consider the hat matrix $H = X(X^T X)^{-1} X^T$, where X is an N by $d + 1$ matrix, and $X^T X$ is invertible.

(a) Show that H is symmetric.

(b) Show that $H^K = H$ for any positive integer K .

(c) If I is the identity matrix of size N , show that $(I - H)^K = I - H$ for any positive integer K .

(d) Show that $\text{trace}(H) = d + 1$, where the trace is the sum of diagonal elements. [Hint: $\text{trace}(AB) = \text{trace}(BA)$]

(a)证明 H 是对称矩阵

$$\begin{aligned}
H^T &= (X(X^T X)^{-1} X^T)^T \\
&= X((X^T X)^{-1})^T X^T \\
&= X((X^T X)^T)^{-1} X^T \\
&= X(X^T X)^{-1} X^T \\
&= H
\end{aligned}$$

(b)证明 $H^K = H$, 直接验证即可, 先来看 $K = 2$ 的情形

$$\begin{aligned}
H^2 &= X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T \\
&= X(X^T X)^{-1} X^T \\
&= H
\end{aligned}$$

那么对于任意 K

$$\begin{aligned}
H^K &= H^2 H^{K-2} \\
&= H H^{K-2} \\
&= H^{K-1} \\
&= \dots \\
&= H
\end{aligned}$$

(c)利用二项式定理直接展开验证

$$\begin{aligned}
(I - H)^K &= \sum_{i=0}^{i=K} C_K^i I^{K-i} (-H)^i \\
&= \sum_{i=0}^{i=K} C_K^i (-1)^i H^i \\
&= I + H \sum_{i=1}^{i=K} C_K^i (-1)^i \text{ (注意 } H^i = H) \\
&= I + H[(1 - 1)^K - 1] \\
&= I - H
\end{aligned}$$

(d)利用迹(trace)的性质 $\text{trace}(AB) = \text{trace}(BA)$

$$\begin{aligned}
\text{trace}(H) &= \text{trace}(X(X^T X)^{-1} X^T) \\
&= \text{trace}(X^T X(X^T X)^{-1}) \\
&= \text{trace}(I_{d+1}) \text{ (注意 } H^T H \text{ 为 } (d+1) \times (d+1) \text{ 阶矩阵)} \\
&= d + 1
\end{aligned}$$

Exercise 3.4 (Page 88)

Consider a noisy target $y = w^{*T} x + \epsilon$ for generating the data, where ϵ is a noise term with zero mean and σ^2 variance, independently generated for every example (x, y) . The expected error of the best possible linear fit to this target is thus σ^2 .

For the data $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$, denote the noise in y_n as ϵ_n and let $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^T$; assume that $X^T X$ is invertible. By following the steps below, show that the expected in sample error of linear regression with respect to \mathcal{D} is given by .

$$\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(w_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right)$$

(a) Show that the in sample estimate of is given by $\hat{y} = Xw^* + H\epsilon$.

(b) Show that the in sample error vector $\hat{y} - y$ can be expressed by a matrix times ϵ . What is the matrix?

(c) Express $E_{\text{in}}(w_{\text{lin}})$ in terms of ϵ using (b), and simplify the expression using Exercise 3.3(c).

(d) Prove that $\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(w_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right)$ using (c) and the independence of $\epsilon_1, \epsilon_2, \dots, \epsilon_N$. [Hint: The sum of the diagonal elements of a matrix (the trace) will play a role. See Exercise 3.3]

For the expected out of sample error, we take a special case which is easy to analyze. Consider a test data set $\mathcal{D}_{\text{test}} = \{(x_1, y'_1), \dots, (x_N, y'_N)\}$ which shares the same input vectors x_n with \mathcal{D} but with a different realization of the noise terms. Denote the noise in y'_n as ϵ'_n and let $\epsilon' = [\epsilon'_1, \epsilon'_2, \dots, \epsilon'_N]^T$. Define $E_{\text{test}}(w_{\text{lin}})$ to be the average squared error on $\mathcal{D}_{\text{test}}$.

(e) Prove that $\mathbb{E}_{\mathcal{D}, \epsilon'}[E_{\text{in}}(w_{\text{lin}})] = \sigma^2 \left(1 + \frac{d+1}{N}\right)$.

The special test error E_{test} is a very restricted case of the general out-of sample error. Some detailed analysis shows that similar results can be obtained for the general case, as shown in Problem 3.1.

(a) 首先将 $y = w^{*T}x + \epsilon$ 改写为向量的形式, 记 $y = [y_1 \dots y_N]^T, X = [x_1 \dots x_N]^T$, 注意题目中给出 $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^T$

那么

$$y = Xw^* + \epsilon$$

我们知道 $w_{\text{lin}} = (X^T X)^{-1} X^T y$

那么

$$\begin{aligned} \hat{y} &= Xw_{\text{lin}} \\ &= X(X^T X)^{-1} X^T y \\ &= X(X^T X)^{-1} X^T (Xw^* + \epsilon) \\ &= X(X^T X)^{-1} X^T Xw^* + X(X^T X)^{-1} X^T \epsilon \\ &= Xw^* + H\epsilon \end{aligned}$$

其中 $H = X(X^T X)^{-1} X^T$ 的定义来自于 Exercise 3.3

(b) 直接计算即可

$$\begin{aligned} \hat{y} - y &= Xw^* + H\epsilon - (Xw^* + \epsilon) \\ &= (H - I)\epsilon \end{aligned}$$

(c) 直接计算即可, 注意要用到 Exercise 3.3 证明的性质

$$\begin{aligned}
E_{\text{in}}(w_{\text{lin}}) &= \frac{1}{N} \|\hat{y} - y\|^2 \\
&= \frac{1}{N} \|(H - I)\epsilon\|^2 \\
&= \frac{1}{N} ((H - I)\epsilon)^T ((H - I)\epsilon) \\
&= \frac{1}{N} \epsilon^T (H - I)(H - I)\epsilon \text{ (注意 } H \text{ 对称)} \\
&= \frac{1}{N} \epsilon^T (I - H)\epsilon \text{ (注意 } (I - H)^K = I - H) \\
&= \frac{1}{N} \epsilon^T (I - H)\epsilon
\end{aligned}$$

(d)这题也是直接计算，注意要用到trace的性质和上题结论

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(w_{\text{lin}})] &= \frac{1}{N} \mathbb{E}_{\mathcal{D}}(\epsilon^T (I - H)\epsilon) \\
&= \frac{1}{N} \mathbb{E}_{\mathcal{D}} \text{trace}(\epsilon^T (I - H)\epsilon) \text{ (这一步是由于 } \epsilon^T (I - H)\epsilon \text{ 是一个实数, 对于实数 } a, a = \text{trace}(a)) \\
&= \frac{1}{N} \mathbb{E}_{\mathcal{D}} \text{trace}(\epsilon^T \epsilon - \epsilon^T H \epsilon) \\
&= \frac{1}{N} [\mathbb{E}_{\mathcal{D}}(\sum_{i=1}^N \epsilon_i^2) - \mathbb{E}_{\mathcal{D}}(\sum_{i=1}^N \epsilon_i H_{ii} \epsilon_i)] \text{ (} H_{ii} \text{ 为 } H \text{ 第 } (i, i) \text{ 个元素)} \\
&= \frac{1}{N} [N\sigma^2 - (\sum_{i=1}^N H_{ii})\sigma^2] \\
&= \frac{1}{N} [N\sigma^2 - \text{trace}(H)\sigma^2] \text{ (注意上一题结论 } \text{trace}(H) = d + 1) \\
&= \frac{1}{N} [N\sigma^2 - (d + 1)\sigma^2] \\
&= \sigma^2(1 - \frac{d + 1}{N})
\end{aligned}$$

为了方便解决下一题，我们把上述计算中的两个结果单独列出

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}}(\epsilon^T \epsilon) &= N\sigma^2 \\
\mathbb{E}_{\mathcal{D}}(\epsilon^T H \epsilon) &= (d + 1)\sigma^2
\end{aligned}$$

(e)首先还是改写为向量的形式 $y' = [y'_1 \dots y'_N]^T$, $X = [x_1 \dots x_N]^T$, $\epsilon' = [\epsilon'_1, \epsilon'_2, \dots, \epsilon'_N]^T$, 那么

$$y' = Xw^* + \epsilon'$$

那么由(a)的结论知

$$\hat{y} = Xw_{\text{lin}} = Xw^* + H\epsilon$$

因此

$$\begin{aligned}
E_{\text{test}}(w_{\text{lin}}) &= \frac{1}{N} \|\hat{y} - y'\|^2 \\
&= \frac{1}{N} \|Xw^* + H\epsilon - (Xw^* + \epsilon')\|^2 \\
&= \frac{1}{N} \|H\epsilon - \epsilon'\|^2 \\
&= \frac{1}{N} (H\epsilon - \epsilon')^T (H\epsilon - \epsilon') \\
&= \frac{1}{N} (\epsilon'^T H - \epsilon'^T) (H\epsilon - \epsilon') \text{ (注意 } H \text{ 对称)} \\
&= \frac{1}{N} (\epsilon'^T H H \epsilon - 2\epsilon'^T H \epsilon + \epsilon'^T \epsilon') \\
&= \frac{1}{N} (\epsilon'^T H \epsilon - 2\epsilon'^T H \epsilon + \epsilon'^T \epsilon') \text{ (注意 } H^K = H)
\end{aligned}$$

接着我们计算 $\mathbb{E}_{\mathcal{D}, \epsilon'} [E_{\text{in}}(w_{\text{lin}})]$

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}, \epsilon'} [E_{\text{in}}(w_{\text{lin}})] &= \mathbb{E}_{\mathcal{D}, \epsilon'} \left[\frac{1}{N} (\epsilon'^T H \epsilon - 2\epsilon'^T H \epsilon + \epsilon'^T \epsilon') \right] \\
&= \frac{1}{N} [\mathbb{E}_{\mathcal{D}, \epsilon'} (\epsilon'^T H \epsilon) - 2\mathbb{E}_{\mathcal{D}, \epsilon'} (\epsilon'^T H \epsilon) + \mathbb{E}_{\mathcal{D}, \epsilon'} (\epsilon'^T \epsilon')]
\end{aligned}$$

回顾上一题单独列出的结论

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} (\epsilon^T \epsilon) &= N\sigma^2 \\
\mathbb{E}_{\mathcal{D}} (\epsilon^T H \epsilon) &= (d+1)\sigma^2
\end{aligned}$$

因此

$$\mathbb{E}_{\mathcal{D}, \epsilon'} [E_{\text{in}}(w_{\text{lin}})] = \sigma^2 \left(1 + \frac{d+1}{N} \right) - \frac{2}{N} \mathbb{E}_{\mathcal{D}, \epsilon'} (\epsilon'^T H \epsilon)$$

后面我们单独计算 $\mathbb{E}_{\mathcal{D}, \epsilon'} (\epsilon'^T H \epsilon)$, 注意 ϵ_i, ϵ'_i 独立且 $\mathbb{E}(\epsilon_i) = \mathbb{E}(\epsilon'_i) = 0$

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}, \epsilon'} (\epsilon'^T H \epsilon) &= \mathbb{E}_{\mathcal{D}, \epsilon'} (\text{trace}(\epsilon'^T H \epsilon)) \\
&= \mathbb{E}_{\mathcal{D}, \epsilon'} \left(\sum_{i=1}^N \epsilon'_i H_{ii} \epsilon_i \right) \\
&= \sum_{i=1}^N (\mathbb{E}(\epsilon'_i) H_{ii} \mathbb{E}(\epsilon_i)) \text{ (由独立性)} \\
&= 0 \text{ (数学期望为 0)}
\end{aligned}$$

因此

$$\mathbb{E}_{\mathcal{D}, \epsilon'} [E_{\text{in}}(w_{\text{lin}})] = \sigma^2 \left(1 + \frac{d+1}{N} \right)$$

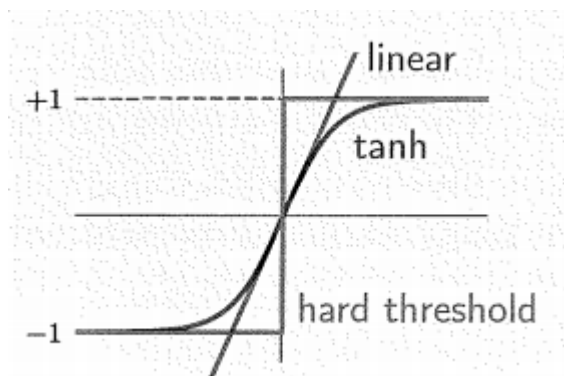
Exercise 3.5 (Page 90)

Another popular soft threshold is the hyperbolic tangent

$$\tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

(a) How is \tanh related to the logistic function θ ? [Hint: shift and scale]

(b) Show that $\tanh(s)$ converges to a hard threshold for large $|s|$, and converges to no threshold for small $|s|$ [Hint: Formalize the figure below.]



(a)我们先回顾 θ 的定义

$$\theta(s) = \frac{e^s}{1 + e^s}$$

那么

$$\begin{aligned} \tanh(s) &= \frac{e^s - e^{-s}}{e^s + e^{-s}} \\ &= \frac{e^{2s} - 1}{e^{2s} + 1} \\ &= \frac{2e^{2s} - (1 + e^{2s})}{e^{2s} + 1} \\ &= 2 \frac{e^{2s}}{1 + e^{2s}} - 1 \\ &= 2\theta(2s) - 1 \end{aligned}$$

所以 $\tanh(s)$ 相当于 $\theta(s)$ 先沿x轴方向压缩2倍，再y轴方向扩张2倍，最后沿y轴向下平移一个单位

(b)这题我不是很理解题目中所说的converges to no threshold for small $|s|$ ，我看了下论坛上老师的回复以及参考图片，猜测是threshold是水平渐近线的意思，所以我们计算 $\frac{\tanh(s)}{s}$ ，这里使用了泰勒展开。

$$\begin{aligned} \lim_{s \rightarrow \infty} \frac{\tanh(s)}{s} &= \lim_{s \rightarrow \infty} \frac{1 - e^{-2s}}{(1 + e^{-2s})s} \\ &= \lim_{s \rightarrow \infty} \frac{1 - (1 - 2s)}{(1 + 1 - 2s)s} \\ &= 0 \end{aligned}$$

$$\begin{aligned}
\lim_{s \rightarrow 0} \frac{\tanh(s)}{s} &= \lim_{s \rightarrow 0} \frac{e^s - e^{-s}}{e^s + e^{-s}} \\
&= \lim_{s \rightarrow 0} \frac{1 + s - (1 - s)}{(1 + s + (1 - s))s} \\
&= \lim_{s \rightarrow 0} \frac{2}{2} \\
&= 1
\end{aligned}$$

所以当 $|s| \rightarrow \infty$, $\tanh(s)$ 有水平渐近线, $|s| \rightarrow 0$, $\tanh(s)$ 没有水平渐近线

Exercise 3.6 (Page 92)

[Cross-entropy error measure]

(a) More generally, if we are learning from ± 1 data to predict a noisy target $P(y|x)$ with candidate hypothesis h , show that the maximum likelihood method reduces to the task of finding h that minimizes

$$E_{\text{in}}(w) = \sum_{n=1}^N \mathbb{I}[y_n = +1] \ln \frac{1}{h(x_n)} + \mathbb{I}[y_n = -1] \ln \frac{1}{1 - h(x_n)}$$

(b) For the case $h(x) = \theta(w^T x)$, argue that minimizing the in sample error in part (a) is equivalent to minimizing the one in (3.9).

For two probability distributions $p, 1 - p$ and $q, 1 - q$ with binary outcomes, the cross entropy (from information theory) is

$$p \log \frac{1}{q} + (1 - p) \log \frac{1}{1 - q}$$

The in sample error in part (a) corresponds to a cross entropy error measure on the data point (x_n, y_n) , with $p = \mathbb{I}[y_n = +1]$ and $q = h(x_n)$.

(a) 记最大似然函数为 L

$$\begin{aligned}
L &= \prod_{i=1}^N (h(x_n))^{[y_n=+1]} (1 - h(x_n))^{[y_n=-1]} \\
\ln L &= \sum_{n=1}^N \mathbb{I}[y_n = +1] \ln h(x_n) + \mathbb{I}[y_n = -1] \ln (1 - h(x_n))
\end{aligned}$$

因此要使得 L 最大, 只要使 $\sum_{n=1}^N \mathbb{I}[y_n = +1] \ln h(x_n) + \mathbb{I}[y_n = -1] \ln (1 - h(x_n))$ 最大即可, 也就是使得 $-(\sum_{n=1}^N \mathbb{I}[y_n = +1] \ln h(x_n) + \mathbb{I}[y_n = -1] \ln (1 - h(x_n)))$ 最小, 注意以下事实

$$\begin{aligned}
-(\sum_{n=1}^N \mathbb{I}[y_n = +1] \ln h(x_n) + \mathbb{I}[y_n = -1] \ln (1 - h(x_n))) &= \sum_{n=1}^N \mathbb{I}[y_n = +1] \ln \frac{1}{h(x_n)} + \mathbb{I}[y_n = -1] \ln \frac{1}{1 - h(x_n)} \\
&= E_{\text{in}}(w)
\end{aligned}$$

因此结论成立。

(b)回顾下3.9

$$E'_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$$

注意 $\theta(s) = \frac{e^s}{1+e^s}$, 那么

$$\theta(-s) = \frac{e^{-s}}{1+e^{-s}} = \frac{1}{1+e^s} = 1 - \theta(s)$$

$$h(x) = \theta(w^T x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

$$h(-x) = \theta(w^T(-x)) = 1 - \theta(w^T x) = 1 - h(x)$$

我们对这里的 $E_{\text{in}}(w)$ 进行一个处理

$$\begin{aligned} \mathbb{I}[y_n = +1] \ln \frac{1}{h(x_n)} &= \mathbb{I}[y_n = +1] \ln \frac{1}{h(y_n x_n)} \\ \mathbb{I}[y_n = -1] \ln \frac{1}{1 - h(x_n)} &= \mathbb{I}[y_n = -1] \ln \frac{1}{h(-x_n)} = \mathbb{I}[y_n = -1] \ln \frac{1}{h(y_n x_n)} \end{aligned}$$

代入我们这里的 $E_{\text{in}}(w)$

$$\begin{aligned} E_{\text{in}}(w) &= \sum_{n=1}^N \mathbb{I}[y_n = +1] \ln \frac{1}{h(x_n)} + \mathbb{I}[y_n = -1] \ln \frac{1}{1 - h(x_n)} \\ &= \sum_{n=1}^N (\mathbb{I}[y_n = +1] + \mathbb{I}[y_n = -1]) \ln \frac{1}{h(y_n x_n)} \\ &= \sum_{n=1}^N \ln \frac{1}{h(y_n x_n)} \\ &= \sum_{n=1}^N \ln \left(\frac{e^{y_n w^T x}}{1 + e^{y_n w^T x}} \right)^{-1} \\ &= \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n}) \\ &= N E'_{\text{in}}(w) \end{aligned}$$

因此最小化 $E_{\text{in}}(w)$ 等价于最小化 $E'_{\text{in}}(w)$ 。

题目最后的意思是这里的结论可以和信息论里的结论类比。

Exercise 3.7 (Page 92)

For logistic regression, show that

$$\begin{aligned}\nabla E_{\text{in}}(w) &= -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{y_n w^T x_n}} \\ &= \frac{1}{N} \sum_{n=1}^N -y_n x_n \theta(-y_n w^T x_n)\end{aligned}$$

Argue that a 'misclassified' example contributes more to the gradient than a correctly classified one.

这题就是对logistic的 $E_{\text{in}}(w)$ 求梯度，回顾下 $E_{\text{in}}(w)$

$$E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n})$$

注意 $\theta(s) = \frac{e^s}{1+e^s}$

$$\begin{aligned}\frac{\partial E_{\text{in}}(w)}{\partial w_i} &= \frac{1}{N} \sum_{n=1}^N \frac{\ln(1 + e^{-y_n w^T x_n})}{\partial w_i} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + e^{-y_n w^T x_n}} \frac{\partial e^{-y_n w^T x_n}}{\partial w_i} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + e^{-y_n w^T x_n}} (-e^{-y_n w^T x_n}) (y_n x_n^{(i)}) (x_n^{(i)} \text{表示 } x_n \text{ 的第 } i \text{ 个分量}) \\ &= \frac{1}{N} \sum_{n=1}^N -y_n x_n^{(i)} \frac{e^{-y_n w^T x_n}}{1 + e^{-y_n w^T x_n}} \\ &= \frac{1}{N} \sum_{n=1}^N -y_n x_n^{(i)} \frac{1}{1 + e^{y_n w^T x_n}} \\ &= \frac{1}{N} \sum_{n=1}^N -y_n x_n^{(i)} \theta(-y_n w^T x_n)\end{aligned}$$

因此

$$\begin{aligned}\nabla E_{\text{in}}(w) &= -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{y_n w^T x_n}} \\ &= \frac{1}{N} \sum_{n=1}^N -y_n x_n \theta(-y_n w^T x_n)\end{aligned}$$

当一个样例错分时, $y_n w^T x_n < 0$, 对应的 $\theta(-y_n w^T x_n) > \frac{1}{2}$; 当一个样例分类正确时, $y_n w^T x_n > 0$ 对应的 $\theta(-y_n w^T x_n) < \frac{1}{2}$, 因此错分的样例的权重比正确的样例要大。

Exercise 3.8 (Page 94)

The claim that \hat{v} is the direction which gives largest decrease in E_{in} only holds for small η . Why?

回顾下书上之前的叙述

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(w(0) + \eta \hat{v}) - E_{\text{in}}(w(0)) \\ &= \eta \nabla E_{\text{in}}(w(0))^T \hat{v} + O(\eta^2) \\ &\geq -\eta \|\nabla E_{\text{in}}(w(0))\| \\ \hat{v} &= -\frac{\nabla E_{\text{in}}(w(0))}{\|\nabla E_{\text{in}}(w(0))\|} \text{时等号成立}\end{aligned}$$

这个式子是泰勒展开，所以 $\eta \hat{v}$ 的模不能太大，如果太大，泰勒展开的偏差会很大。 \hat{v} 是单位向量，因此 $\|\eta \hat{v}\| = \eta$ ，所以上述式子只有当 η 较小的时候才有效。

Exercise 3.9 (Page 97)

Consider pointwise error measures $e_{\text{class}}(s, y) = \mathbb{I}[y \neq \text{sign}(s)]$, $e_{\text{sq}}(s, y) = (y - s)^2$, and $e_{\log}(s, y) = \ln(1 + \exp(-ys))$, where the signal $s = w^T x$

- (a) For $y = +1$, plot e_{class} , e_{sq} and $\frac{1}{\ln 2} e_{\log}$ versus s , on the same plot.
- (b) Show that $e_{\text{class}}(s, y) \leq e_{\text{sq}}(s, y)$, and hence that the classification error is upper bounded by the squared error.
- (c) Show that $e_{\text{class}}(s, y) \leq \frac{1}{\ln 2} e_{\log}(s, y)$, and, as in part (b), get an upper bound (up to a constant factor) using the logistic regression error. These bounds indicate that minimizing the squared or logistic regression error should also decrease the classification error, which justifies using the weights returned by linear or logistic regression as approximations for classification.

(a),(b),(c)三题作图即可

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 09:50:23 2019

@author: qinzheng
"""

#Step1 构造损失函数
import numpy as np
import matplotlib.pyplot as plt
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号

def Class(s,y):
    if s * y > 0:
        return 0
    else:
        return 1

def Sq(s, y):
    return (s - y) ** 2

def Log(s, y):
```



```

        return np.log(1 + np.exp(- y * s))

#Step2 构造点集并作图
#构造点
x = np.arange(-2,2,0.01)

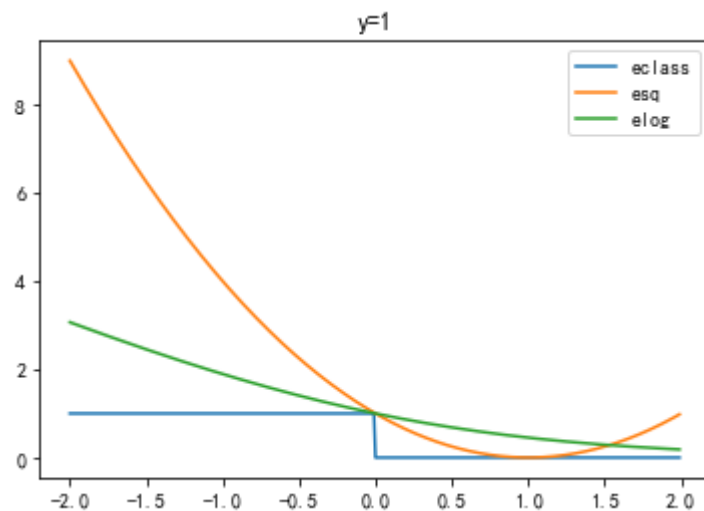
#y=1
eclass1 = [Class(i, 1) for i in x]
esq1 = [Sq(i, 1) for i in x]
elog1 = [Log(i, 1) / np.log(2) for i in x]

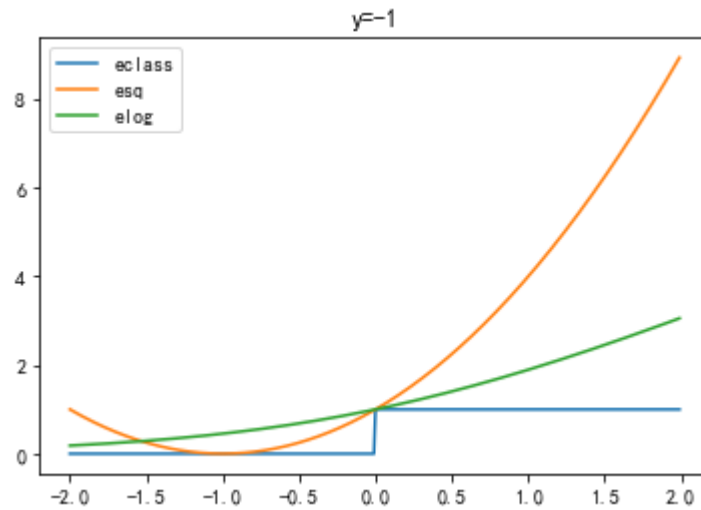
#y=-1
eclass2 = [Class(i, -1) for i in x]
esq2 = [Sq(i, -1) for i in x]
elog2 = [Log(i, -1) / np.log(2) for i in x]

plt.plot(x, eclass1, label='eclass')
plt.plot(x, esq1, label='esq')
plt.plot(x, elog1, label='elog')
plt.title('y=1')
plt.legend()
plt.show()

plt.plot(x, eclass2, label='eclass')
plt.plot(x, esq2, label='esq')
plt.plot(x, elog2, label='elog')
plt.title('y=-1')
plt.legend()
plt.show()

```





从图像中我们看出, e_{sq} 和 $\frac{1}{\ln 2} e_{log}$ 都是 e_{class} 的上界, 因此我们可以用线性回归或者logistic回归计算结果, 再用产生的结果喂给PLA。

Exercise 3.10 (Page 98)

(a) Define an error for a single data point (x_n, y_n) to be

$$e_n(w) = \max(0, -y_n w^T x_n)$$

Argue that PLA can be viewed as SGD on e_n with learning rate $\eta = 1$.

(b) For logistic regression with a very large w , argue that minimizing E_{in} using SGD is similar to PLA. This is another indication that the logistic regression weights can be used as a good approximation for classification.

(a) $e_n(w) = \max(0, -y_n w^T x_n)$ 的意思是对于分类正确的点 $e_n(w) = 0$, 对于分类不正确的点 $e_n(w) = -y_n w^T x_n$, 我们来求梯度

$$\frac{\partial(-y_n w^T x_n)}{\partial w_i} = -y_n x_n^{(i)} \quad (x_n^{(i)} \text{ 表示 } x_n \text{ 的第 } i \text{ 个分量})$$

$$\nabla(-y_n w^T x_n) = -y_n x_n$$

所以对于分类错误的点 (x_n, y_n) , 根据SGD, 更新规则为

$$w(t+1) = w(t) + \eta(-\nabla(-y_n w^T x_n)) = w(t) + \eta y_n x_n$$

所以PLA可以被看成 $e_n(w) = \max(0, -y_n w^T x_n)$ 的SGD且 $\eta = 1$ 的情形。

(b) 我们知道logistic的梯度公式有如下形式

$$\nabla E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n x_n}{1 + e^{y_n w^T x_n}}$$

那么对于SGD梯度即为

$$\nabla e_{in}(w) = \frac{-y_n x_n}{1 + e^{y_n w^T x_n}}$$

带入更新规则 $w \leftarrow w - \eta \nabla e_{\text{in}}(w)$

$$w(t+1) = w(t) + \frac{\eta y_n x_n}{1 + e^{y_n w^T x_n}}$$

我们注意更新的点一定是错误的, 即 $y_n w^T x_n < 0$, 那么当 w 非常大时, $e^{y_n w^T x_n} \approx 0$

因此对于非常大的 w , 上述更新规则可以近似为

$$w(t+1) = w(t) + \eta y_n x_n$$

和PLA一致, 这也从另一个角度说明了logistic是分类问题的一个近似。

Exercise 3.11 (Page 101)

Consider the feature transform ϕ in (3.12). What kind of boundary in \mathcal{X} does a hyperplane \hat{w} in \mathcal{Z} correspond to in the following cases?

Draw a picture that illustrates an example of each case.

(a) $\hat{w}_1 > 0, \hat{w}_2 < 0$

(b) $\hat{w}_1 > 0, \hat{w}_2 = 0$

(c) $\hat{w}_1 > 0, \hat{w}_2 > 0, \hat{w}_0 < 0$

(d) $\hat{w}_1 > 0, \hat{w}_2 > 0, \hat{w}_0 > 0$

回顾下3.12

$$\phi(x) = (1, x_1^2, x_2^2)$$

因此对应方程为

$$\hat{w}_0 + \hat{w}_1 x_1^2 + \hat{w}_2 x_2^2 = 0$$

后面的叙述实际上是高中解析几何的知识。

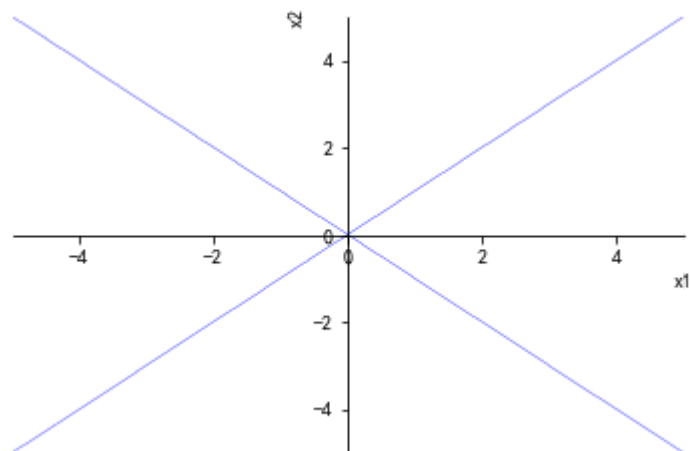
(a)有三种可能, 分别是 $\hat{w}_0 = 0, \hat{w}_0 > 0, \hat{w}_0 < 0$

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 09:51:21 2019

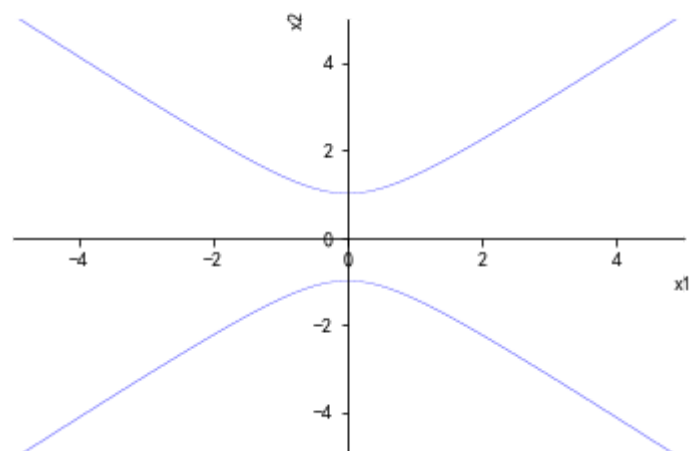
@author: qinzhen
"""

#隐函数作图库
from sympy.parsing.sympy_parser import parse_expr
from sympy import plot_implicit
ezplot = lambda exper: plot_implicit(parse_expr(exper))

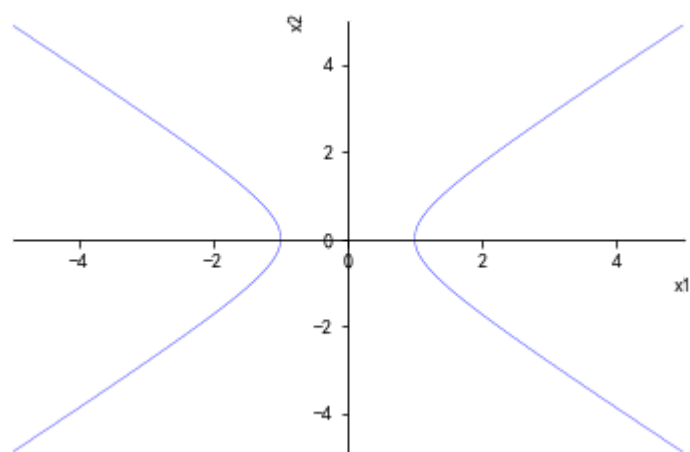
#a
#w0=0
ezplot('x1**2-x2**2')
```



```
#w0>0
ezplot('1+x1**2-x2**2')
```

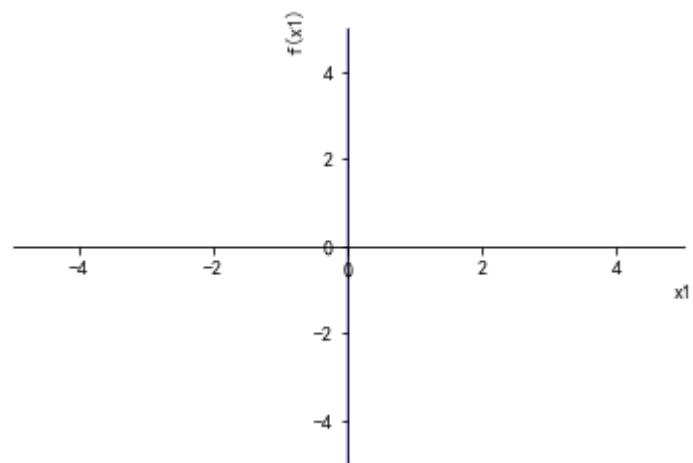


```
#w0<0
ezplot('-1+x1**2-x2**2')
```

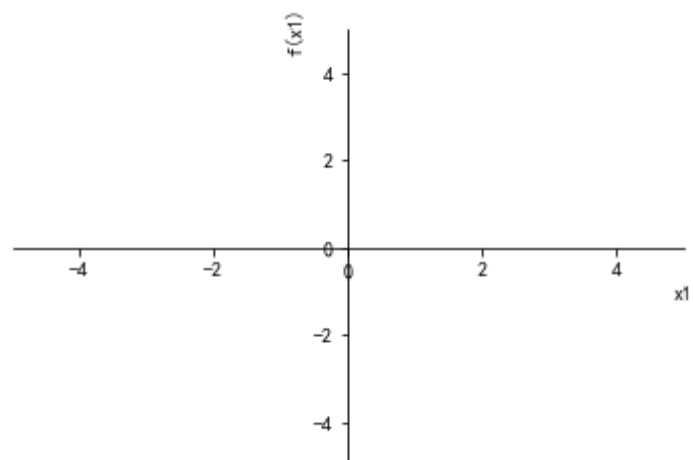


(b)有三种可能，分别是 $\hat{w}_0 = 0$, $\hat{w}_0 > 0$, $\hat{w}_0 < 0$

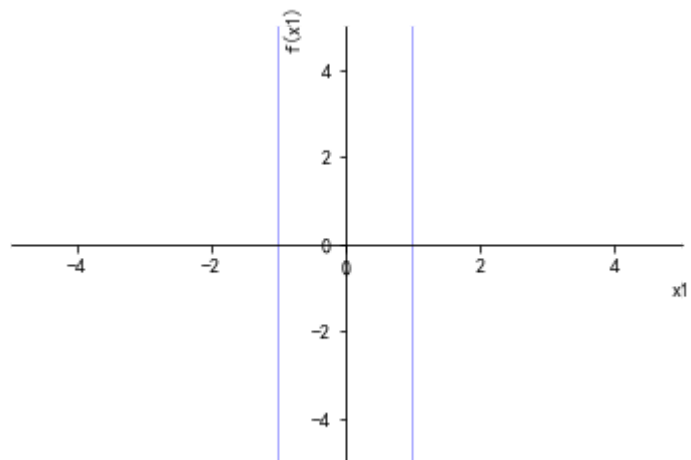
```
#b  
#w0=0  
ezplot('x1**2')
```



```
#w0>0  
ezplot('1+x1**2')
```

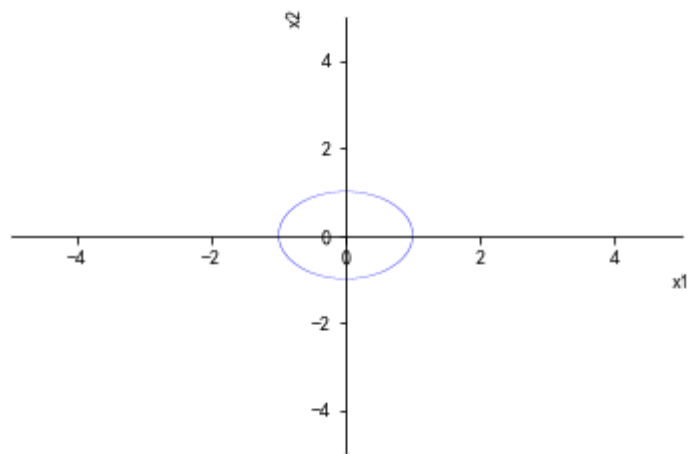


```
#w0<0  
ezplot('-1+x1**2')
```



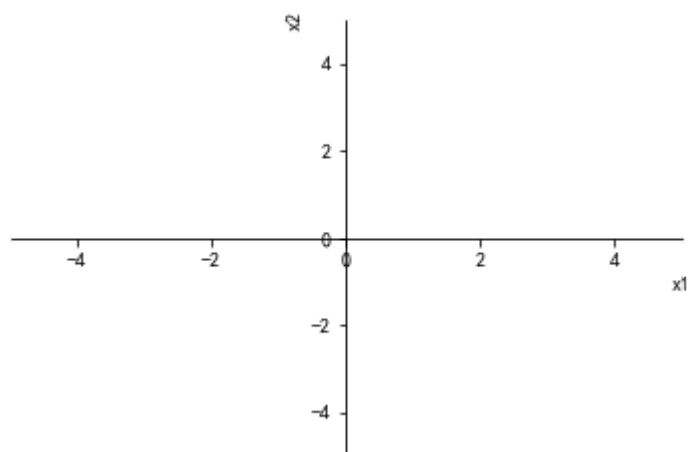
(c)条件已经定死

```
#c
ezplot('-1+x1**2+x2**2')
```



(d)条件已经定死

```
#d
ezplot('1+x1**2+x2**2')
```



Exercise 3.12 (Page 103)

We know that in the Euclidean plane, the perceptron model \mathcal{H} cannot implement all 16 dichotomies on 4 points. That is, $m_{\mathcal{H}}(4) < 16$. Take the feature transform Φ in (3.12).

(a) Show that $m_{\mathcal{H}_{\Phi}}(3) = 8$.

(b) Show that $m_{\mathcal{H}_{\Phi}}(4) < 16$.

(c) Show that $m_{\mathcal{H} \cup \mathcal{H}_{\Phi}}(4) = 16$.

That is, if you used lines, $d_{vc} = 3$; if you used ellipses, $d_{vc} = 3$; if you used lines and ellipses, $d_{vc} > 3$.

回顾下3.12

$$\Phi(x) = (1, x_1^2, x_2^2)$$

(a)作图即可，这里画图比较麻烦，略去了。

(b)注意 $\Phi(x)$ 不包括偏置项为2维，2维感知机最多shatter 3个点，所以

$$m_{\mathcal{H}_{\Phi}}(4) < 16$$

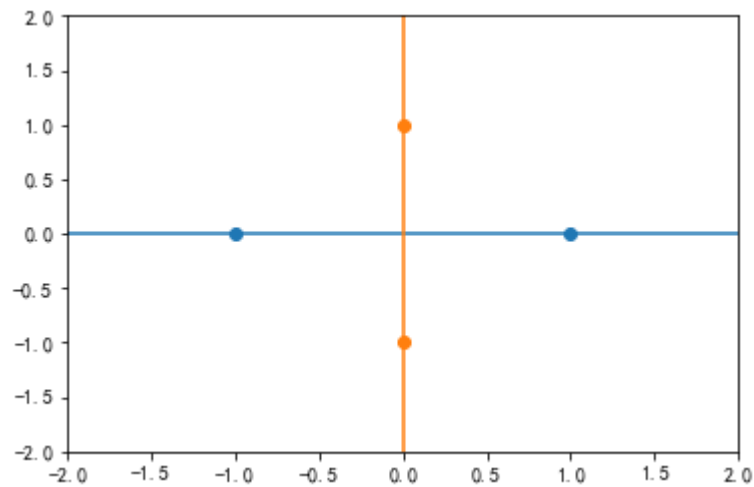
(c)这里只列一种我们之前感知机无法表示的情况

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 09:53:48 2019

@author: qinzhen
"""

import matplotlib.pyplot as plt

#(c)
plt.scatter([-1, 1],[0, 0])
plt.scatter([0, 0],[-1, 1])
plt.xlim(-2, 2)
plt.ylim(-2, 2)
plt.plot([-2, 2],[0, 0])
plt.plot([0, 0],[-2, 2])
plt.xticks()
plt.show()
```



这种形式可以用双曲线进行划分，其余的情形用直线可以轻松划分出来，这里不列出来了。

Exercise 3.13 (Page 104)

Consider the feature transform $z = \Phi_2(x)$ in (3.13). How can we use a hyperplane \hat{w} in \mathcal{Z} to represent the following boundaries in \mathcal{X}

- (a) parabola $(x_1 - 3)^2 + x_2 = 1$
- (b) The circle $(x_1 - 3)^2 + (x_2 - 4)^2 = 1$
- (c) The ellipse $2(x_1 - 3)^2 + (x_2 - 4)^2 = 1$
- (d) The hyperbola $(x_1 - 3)^2 - (x_2 - 4)^2 = 1$
- (e) The ellipse $2(x_1 + x_2 - 3)^2 + (x_1 - x_2 - 4)^2 = 1$
- (f) line $2x_1 + x_2 = 1$

回顾3.13

$$\phi_2(x) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$$

接下来分别打开上述6个式子即可

(a)

$$\begin{aligned} (x_1 - 3)^2 + x_2 &= 1 \\ x_1^2 - 6x_1 + 9 + x_2 - 1 &= 0 \\ 8 - 6x_1 + x_2 + x_1^2 &= 0 \\ \hat{w} &= (8, -6, 1, 1, 0, 0) \end{aligned}$$

(b)

$$\begin{aligned} (x_1 - 3)^2 + (x_2 - 4)^2 &= 1 \\ x_1^2 - 6x_1 + x_2^2 - 8x_2 + 24 &= 0 \\ \hat{w} &= (24, -6, -8, 1, 0, 1) \end{aligned}$$

(c)

$$\begin{aligned}2(x_1 - 3)^2 + (x_2 - 4)^2 &= 1 \\2(x_1^2 - 6x_1 + 9) + (x_2^2 - 8x_2 + 16) - 1 &= 0 \\2x_1^2 - 12x_1 + x_2^2 - 8x_2 + 33 &= 0 \\\hat{w} &= (33, -12, -8, 2, 0, 1)\end{aligned}$$

(d)

$$\begin{aligned}(x_1 - 3)^2 - (x_2 - 4)^2 &= 1 \\x_1^2 - 6x_1 + 9 - (x_2^2 - 8x_2 + 16) - 1 &= 0 \\x_1^2 - x_2^2 + 8x_2 - 6x_1 - 8 &= 0 \\\hat{w} &= (-8, -6, 8, 1, 0, -1)\end{aligned}$$

(e)

$$\begin{aligned}2(x_1 + x_2 - 3)^2 + (x_1 - x_2 - 4)^2 &= 1 \\2[(x_1 + x_2)^2 + 9 - 6(x_1 + x_2)] + (x_1 - x_2)^2 - 8(x_1 - x_2) + 16 - 1 &= 0 \\2(x_1^2 + x_2^2 + 2x_1x_2 + 9 - 6x_1 - 6x_2) + x_1^2 + x_2^2 - 2x_1x_2 - 8x_1 + 8x_2 + 15 &= 0 \\3x_1^2 + 3x_2^2 + 33 - 20x_1 - 4x_2 + 2x_1x_2 &= 0 \\\hat{w} &= (33, -20, -4, 3, 2, 3)\end{aligned}$$

(f)

$$\begin{aligned}2x_1 + x_2 &= 1 \\2x_1 + x_2 - 1 &= 0 \\\hat{w} &= (-1, 2, 1, 0, 0, 0)\end{aligned}$$

Exercise 3.14 (Page 105)

Consider the Q th order polynomial transform ϕ_Q for $\mathcal{X} = \mathbb{R}^d$. What is the dimensionality \tilde{d} of the feature space \mathcal{Z} (excluding the fixed coordinate $z_0 = 1$). Evaluate your result on $d \in \{2, 3, 5, 1\}$ and $Q \in \{2, 3, 5, 1\}$.

设 $x = (x_1 \dots x_d)$, 那么多项式转换的一般形式为 $\prod_{i=1}^d x_i^{n_i}$, 那么 Q 次多项式相当于对此加了一个条件 ($z_0 = 1$ 不算在内)

$$1 \leq \sum_{i=1}^d n_i \leq Q (n_i \geq 0)$$

我们记 $\sum_{i=1}^d n_i = q$ ($n_i \geq 0$) 的解的数量为 $f(q)$, 那么 $1 \leq \sum_{i=1}^d n_i \leq Q$ ($n_i \geq 0$) 的解的数量为

$$\sum_{q=1}^Q f(q)$$

我们接下来求解 $f(q)$, 对式子稍做变形

$$\sum_{i=1}^d n_i = q, n_i \geq 0$$

$$\sum_{i=1}^d (n_i + 1) = q + d, n_i \geq 0$$

令 $n_i + 1 = m_i$, 那么上式可化为

$$\sum_{i=1}^d m_i = q + d, m_i \geq 1$$

所以求正整数解即可, 隔板法可得

$$f(q) = C_{q+d-1}^{d-1}$$

因此

$$\hat{h} = \sum_{q=1}^Q f(q) = \sum_{q=1}^Q C_{q+d-1}^{d-1}$$

我们将 $d = 2$ 代入

$$\hat{h} = \sum_{q=1}^Q f(q) = \sum_{q=1}^Q C_{q+1}^1 = 2 + \dots + (Q + 1) = \frac{Q(Q + 3)}{2}$$

符合课本104页的叙述。

Exercise 3.15 (Page 106)

High-dimensional feature transforms are by no means the only transforms that we can use. We can take the tradeoff in the other direction, and use low dimensional feature transforms as well (to achieve an even lower generalization error bar). Consider the following feature transform, which maps a d -dimensional x to a one-dimensional z , keeping only the k th coordinate of x .

$$\phi_{(k)}(x) = (1, x_k)$$

Let \mathcal{H}_k be the set of perceptrons in the feature space.

(a) Prove that $d_{vc}(\mathcal{H}_k) = 2$.

(b) Prove that $d_{vc}(\bigcup_{k=1}^d \mathcal{H}_k) \leq 2(\log_2 d + 1)$.

\mathcal{H}_k is called the decision stump model on dimension k .

(a)这个比较简单, \mathcal{H}_k 是特征空间里的感知机, 并且特征空间的维度为1, 因此根据感知机的性质, 我们知道

$$d_{vc}(\mathcal{H}_k) = 1 + 1 = 2$$

(b)我们来看下 \mathcal{H}_k 的具体形式, 设参数为 (w_0, w_1) , 那么对应的边界平面为

$$\begin{aligned}w_0 + w_1 x_k &= 0 \\x_k &= -\frac{w_0}{w_1}\end{aligned}$$

因此 \mathcal{H}_k 划分方法可以理解为看第 k 个下标, 如果 x_k 大于阈值 $-\frac{w_0}{w_1}$, 标记为1, 反之标记为-1, 或者反过来(大于 $-\frac{w_0}{w_1}$ 标记为-1, 小于 $-\frac{w_0}{w_1}$ 标记为1)。

假设现在有 N 个点, 现在来计算 \mathcal{H}_k 能区分的数量, 先对这 N 个点的第 k 个坐标排序, **先不管全1或者全-1的两种情况**, 那么 \mathcal{H}_k 相当于在这 N 个 x_k 的 $N-1$ 个间隔挑选, 一共可以有 $N-1$ 种选择, 那么由于大于阈值可以为1, 也可以为-1, 所以一共可以区分 $2(N-1)$ 种情形, 因此

除去全1或者全-1的情况, 每个 \mathcal{H}_k 可以区分 $2N-2$ 种情形

那么 $\bigcup_{k=1}^d \mathcal{H}_k$ 一共可以表示 $f(N) = 2(N-1) \times d + 2$ 种情形, 注意我们这里为了更加准确, 全1或者全-1的情形合并在一起统计了。当 $N = 2(\log_2 d + 1)$ 时

$$\begin{aligned}f(2(\log_2 d + 1)) &= 2(2(\log_2 d + 1) - 1) \times d + 2 \\&= 2(2\log_2 d + 1) \times d + 2\end{aligned}$$

我们来证明

$$2(2\log_2 d + 1) \times d + 2 \leq 2^N = 2^{2(\log_2 d + 1)} = 4d^2$$

接着我们来进行一些处理

$$\begin{aligned}2(2\log_2 d + 1) \times d + 2 &\leq 4d^2 \Leftrightarrow \\2d(2\log_2 d + 1) &\leq 4d^2 - 2 \Leftrightarrow \\2\log_2 d + 1 &\leq 2d - \frac{1}{d} \Leftrightarrow \\2d - \frac{1}{d} - 2\log_2 d - 1 &\geq 0\end{aligned}$$

记 $g(d) = 2d - \frac{1}{d} - 2\log_2 d - 1$, 求导得

$$\begin{aligned}g'(d) &= 2 + \frac{1}{d^2} - \frac{2}{d \ln 2} \\g'(d) &= \left(\frac{1}{d} - \frac{1}{\ln 2}\right)^2 + 2 - \left(\frac{1}{\ln 2}\right)^2\end{aligned}$$

将 $\left(\frac{1}{d} - \frac{1}{\ln 2}\right)^2 + 2 - \left(\frac{1}{\ln 2}\right)^2$ 看成关于 $\frac{1}{d}$ 的二次函数, 由二次函数的性质可得

$$g'(d) \geq g'(1) = 3 - \frac{2}{\ln 2} > 0$$

所以 $g(d)$ 在 $[1, +\infty)$ 递增递增, $g(d) \geq g(1) = 0$

因此 $N = 2(\log_2 d + 1)$ 时,

$$f(N) \leq 2^N$$

从而可得

$$d_{vc}(\bigcup_{k=1}^d \mathcal{H}_k) \leq 2(\log_2 d + 1)$$

这是因为在 $N = 2(\log_2 d + 1)$ 表示的最大种类数量小于等于 2^N ，所以最多 shatter $2(\log_2 d + 1)$ 个点。

Exercise 3.16 (Page 106)

Write down the steps of the algorithm that combines ϕ_3 with linear regression. How about using ϕ_{10} instead? Where is the main computational bottleneck of the resulting algorithm?

这部分可以参考课本86页。我们直接对 ϕ_k 进行总结,记特征空间的维度为 \tilde{d} ，原始数据为 $(x_1 \dots x_N), (y_1 \dots y_N)$ 。

第一步进行特征转换，记得到的新的数据为 $(\tilde{x}_1 \dots \tilde{x}_N)$ ，构成的矩阵为 \tilde{X}

第二步计算 $(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T$

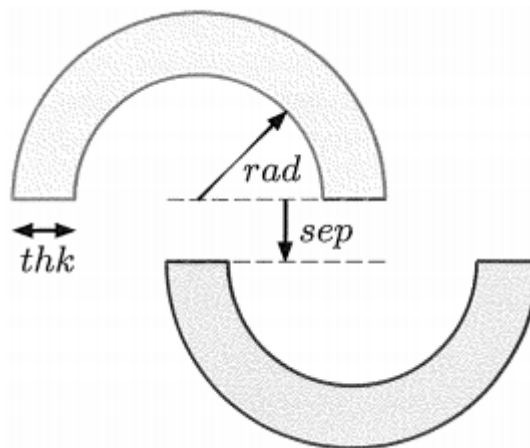
第三步计算 $(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$

主要计算时间应该是第二部求逆矩阵。

Part 2: Problems

Problem 3.1 (Page 109)

Consider the double semi-circle "toy" learning task below.



There are two semi circles of width thk with inner radius rad , separated by sep as shown (red is -1 and blue is +1). The center of the top semi circle is aligned with the middle of the edge of the bottom semi circle. This task is linearly separable when $sep \geq 0$, and not so for $sep < 0$. Set $rad = 10$, $thk = 5$ and $sep = 5$. Then, generate 2,000 examples uniformly, which means you will have approximately 1,000 examples for each class.

- Run the PLA starting from $w = 0$ until it converges. Plot the data and the final hypothesis.
- Repeat part (a) using the linear regression (for classification) to obtain w . Explain your observations.

这题看了非常久，楞是不明白什么意思，后来看了论坛里老师的回复才明白。题目的意思是有如图两个圆环，我们要在圆环内的区域随机生成点，然后用PLA和线性回归进行分类。我们先生成数据，这里我们将上半个圆环对应的圆心放在 (a, b) ，那么下半个圆环对应的圆心为 $(a + \text{rad} + \frac{\text{thk}}{2}, b - \text{sep})$ 。为方便叙述，这里记录上半个圆环的圆心为 (X_1, Y_1) ，下半个圆环的圆心为 (X_2, Y_2) 。

生成点的方式是用参数方程，首先从 $[0, 2\pi]$ 的均匀分布中生成角度 θ ，再从 $[\text{rad}, \text{rad} + \text{thk}]$ 生成距离圆心的距离 r 。如果 $\theta \in [0, \pi]$ ，那么属于上半圆弧，此时

$$x = X_1 + r \cos \theta, y = Y_1 + r \sin \theta$$

否则

$$x = X_2 + r \cos \theta, y = Y_2 + r \sin \theta$$

这里得到如下函数，依旧在helper.py文件中

```
def generatedata(rad, thk, sep, n, x1=0, y1=0):
    """
    产生课本109页的数据集，这里设置的参数为半径rad，圆环宽度thk，
    上下圆环间隔sep，n为数据集总数，x1, y1为上半圆环的圆心
    """
    #上半个圆的圆心
    x1 = x1
    y1 = y1

    #下半个圆的圆心
    x2 = x1 + rad + thk / 2
    y2 = y1 - sep

    #生成角度theta
    Theta = np.random.uniform(0, 2*np.pi, n)
    #生成距离r
    R = np.random.uniform(rad, rad+thk, n)

    #根据Theta生成标签
    y = 2 * (Theta < np.pi) - 1

    #生成点集合X，首先根据y的标签生成圆心
    X = np.zeros((n, 2))
    X[y > 0] = np.array([x1, y1])
    X[y < 0] = np.array([x2, y2])
    #其次用参数方程生成坐标
    X[:, 0] += np.cos(Theta) * R
    X[:, 1] += np.sin(Theta) * R

    return X, y
```

接着作图看下

```
import numpy as np
import matplotlib.pyplot as plt
from numpy.linalg import inv
from helper import generatedata
```

```

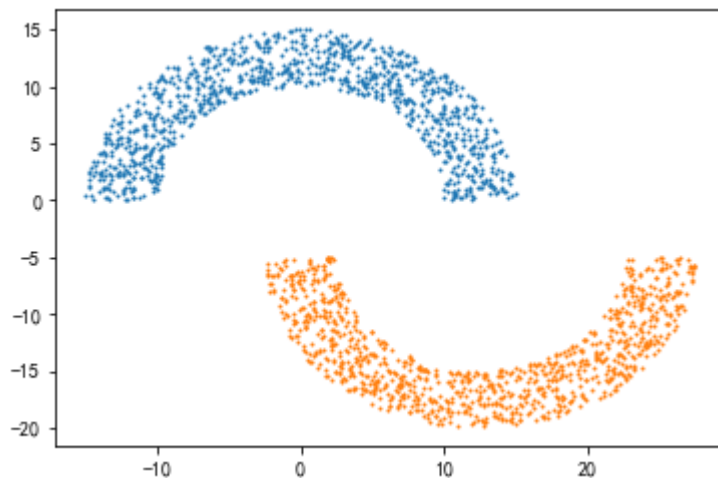
from helper import PLA

#Step1 产生数据
#参数
rad = 10
thk = 5
sep = 5
N = 2000

#产生数据
X, y = generatedata(rad, thk, sep, N)

#作图
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.show()

```



可以看到和圆环形状还是很接近的，接着实现PLA算法，如下代码存储在helper.py文件中：

```

def Judge(X, y, w):
    """
    判别函数，判断所有数据是否分类完成
    """
    n = X.shape[0]
    #判断是否分类完成
    num = np.sum(X.dot(w) * y > 0)
    return num == n

def PLA(X, y, eta=1, max_step=np.inf):
    """
    PLA算法，X, y为输入数据，eta为步长，默认为1，max_step为最多迭代次数，默认为无穷
    """
    #获取维度
    n, d = X.shape
    #初始化
    w = np.zeros(d)
    #记录迭代次数
    t = 0

```

```

#记录元素的下标
i = 0
#记录最后一个错误的下标
last = 0
while not(Judge(X, y, w)) and t < max_step:
    if np.sign(X[i, :].dot(w) * y[i]) <= 0:
        #迭代次数增加
        t += 1
        w += eta * y[i] * X[i, :]
        #更新最后一个错误
        last = i

    #移动到下一个元素
    i += 1
    #如果达到n, 则重置为0
    if i == n:
        i = 0

return t, last, w

```

接着我们处理(a)

```

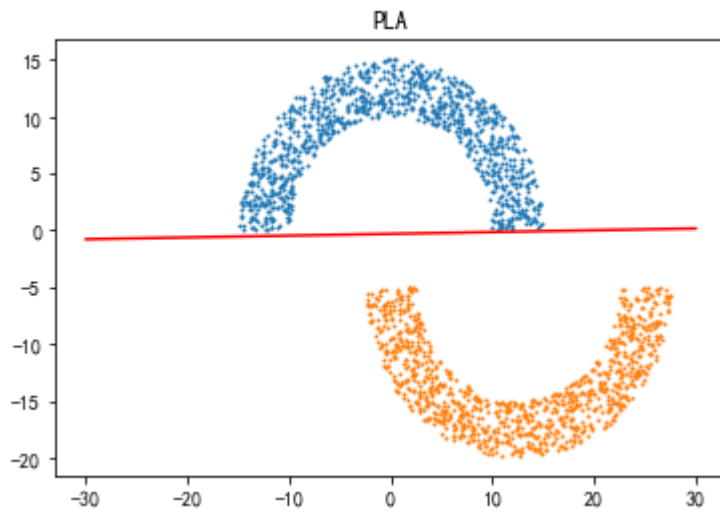
#Step2 训练数据
#(a)PLA
#对数据预处理, 加上偏置项1
X_treat = np.c_[np.ones(N), X]

#PLA
t, last, w = PLA(X_treat, y)

#作出直线
r = 2 * (rad + thk)
a1 = np.array([-r, r])
b1 = - (w[0] + w[1] * a1) / w[2]

plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.plot(a1, b1, c="red")
plt.title('PLA')
plt.show()

```

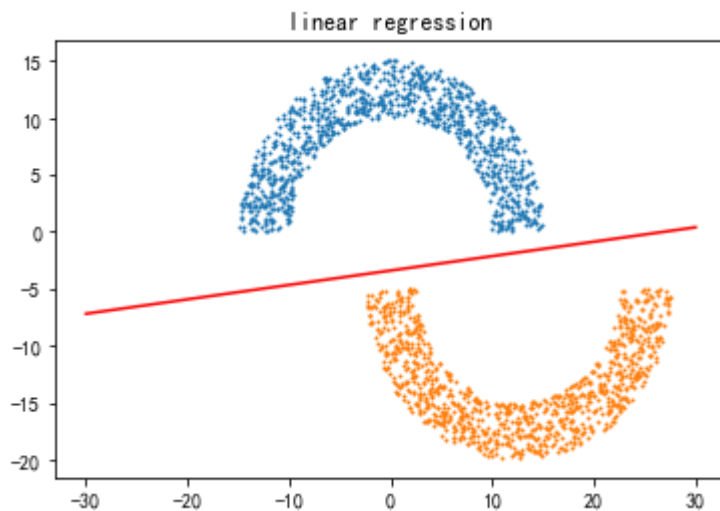


(b)直接带入线性回归的公式即可

```
#(b)linear regression
w1 = inv(X_treat.T.dot(X_treat)).dot(X_treat.T).dot(y)

#作图
a2 = np.array([-r,r])
b2 = - (w1[0] + w1[1] * a1) / w1[2]

plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.plot(a2, b2, c="red")
plt.title('linear regression')
plt.show()
```



可以看到，线性回归同样能解决这个分类问题，和课本里的描述一致，相比PLA，线性回归得到的直线距离点集的距离更大。

Problem 3.2 (Page 109)

For the double semi circle task in Problem 3.1, vary sep in the range $\{0.2, 0.4, \dots, 5\}$. Generate 2,000 examples and run the PLA starting with $w = 0$. Record the number of iterations PLA takes to converge.

Plot sep versus the number of iterations taken for PLA to converge. Explain your observations. [Hint: Problem 1.3.]

和上题一致，不过这里改变了 sep ，我们要观察 sep 和迭代次数的关系：

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 10:38:42 2019

@author: qinzhen
"""

import numpy as np
from helper import generatedata
from helper import PLA
import matplotlib.pyplot as plt
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号

#参数
rad = 10
thk = 5
Sep = np.arange(0.2,5.2,0.2)
N = 2000
#实验次数
n = 30

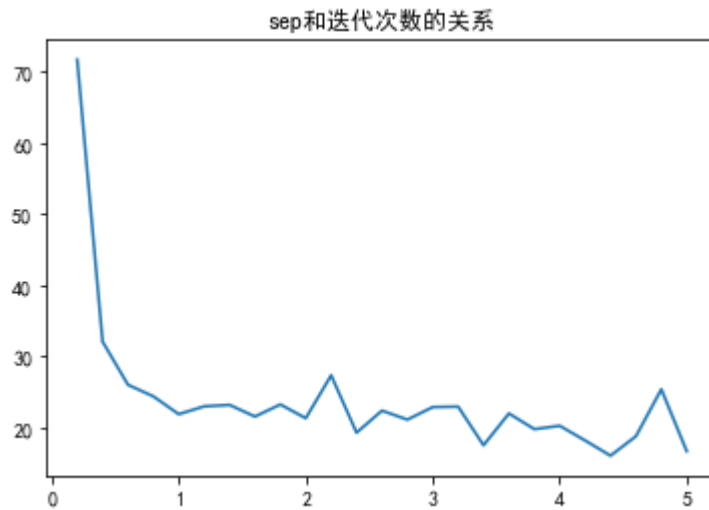
#记录迭代次数
T=np.array([])

for sep in Sep:
    t1 = 0
    for i in range(n):
        X, y = generatedata(rad, thk, sep, N)
        X_treat = np.c_[np.ones(N), X]

        t, last, w = PLA(X_treat, y)
        t1 += t

    T = np.append(T, t1 / n)

plt.plot(Sep, T)
plt.title('sep和迭代次数的关系')
plt.show()
```



可以看到 sep 越大，迭代次数总体来说在下降。现在来简单分析一下原因，回顾Problem 1.3(Page 33)

$$\begin{aligned}\rho &= \min_{1 \leq n \leq N} y_n (w^{*T} x_n) \\ R &= \max_{1 \leq n \leq N} \|x_n\| \\ t &\leq \frac{R^2 \|w^*\|^2}{\rho^2}\end{aligned}$$

这里圆环的位置固定，所以可以认为 R 是一个常数，我们来分析 $\frac{\|w^*\|}{\rho}$ 。

假设 $w = (w_0, w_1, \dots, w_m)^T$ ，那么由解析几何知识，我们知道

$$\frac{\|w^T x_n\|}{\sqrt{w_1^2 + \dots + w_m^2}}$$

等于 x_n 到平面 $w^T x = 0$ 的距离，所以 $\frac{\|w^T x_n\|}{\|w\|}$ 约等于 x_n 到平面 $w^T x = 0$ 的距离。而题目中的 $\frac{\|w^*\|}{\rho}$ 约等于点集到平面 $w^{*T} x = 0$ 的最小距离的倒数，因此如果 sep 越大，最小距离也越大，从而 $\frac{\|w^*\|}{\rho}$ 越小，迭代次数相应也会减少。

Problem 3.3 (Page 109)

For the double semi circle task in Problem 3.1, set $sep = -5$ and generate 2,000 examples.

- What will happen if you run PLA on those examples?
- Run the pocket algorithm for 100,000 iterations and plot E_{in} versus the iteration number t .
- Plot the data and the final hypothesis in part (b).
- Use the linear regression algorithm to obtain the weights w , and compare this result with the pocket algorithm in terms of computation time and quality of the solution.
- Repeat (b) - (d) with a 3rd order polynomial feature transform.

这题的 $sep < 0$ ，所以数据不可分。

(a) 因为数据不可分，所以如果运行PLA，那么算法不会停下来。

(b)回顾下Pocket PLA

The pocket algorithm:

- 1: Set the pocket weight vector $\hat{\mathbf{w}}$ to $\mathbf{w}(0)$ of PLA.
- 2: **for** $t = 0, \dots, T - 1$ **do**
- 3: Run PLA for one update to obtain $\mathbf{w}(t + 1)$.
- 4: Evaluate $E_{\text{in}}(\mathbf{w}(t + 1))$.
- 5: If $\mathbf{w}(t + 1)$ is better than $\hat{\mathbf{w}}$ in terms of E_{in} , set $\hat{\mathbf{w}}$ to $\mathbf{w}(t + 1)$.
- 6: **Return** $\hat{\mathbf{w}}$.

编程实现，题目要求迭代十万次，我发现十万次实在是太慢了，所以改为10000次，先作图看一下

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 10:44:29 2019

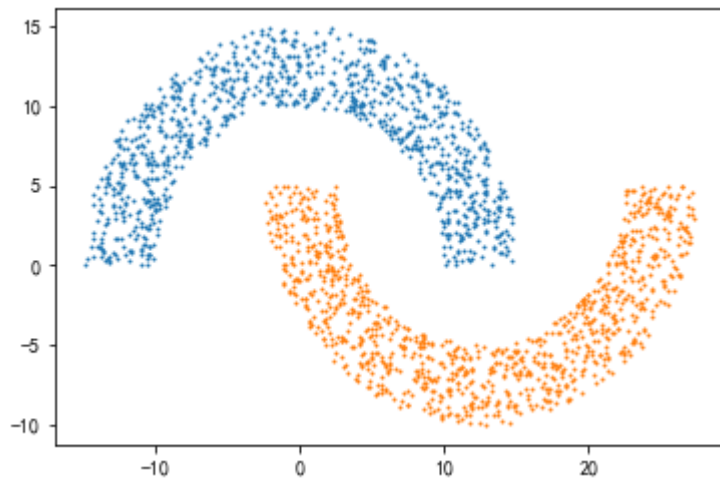
@author: qinzhen
"""

import numpy as np
from numpy.linalg import inv
from helper import generatedata
from helper import Pocket_PLA
from sklearn.preprocessing import PolynomialFeatures
import matplotlib.pyplot as plt
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号

#Step1 产生数据
#参数
rad = 10
thk = 5
sep = -5
N = 2000

#产生数据
X, y = generatedata(rad, thk, sep, N)

#作图
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.show()
```



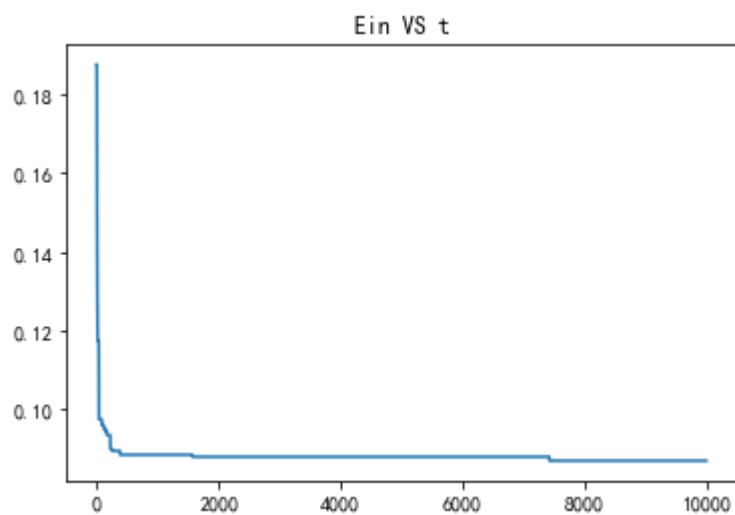
(b)运行算法并作图

```
#Step2 训练数据
#Pocket_PLA
#对数据预处理, 加上偏置项1
X_treat = np.c_[np.ones(N), X]

#迭代次数
max_step = 10000

#产生结果
w, w_hat, w, error = Pocket_PLA(X_treat, y, max_step=max_step)
ein = np.mean(np.sign(w_hat.dot(X_treat.T)) != y, axis=1)

#(b)
t = np.arange(max_step)
plt.plot(t, ein)
plt.title('Ein VS t')
plt.show()
```



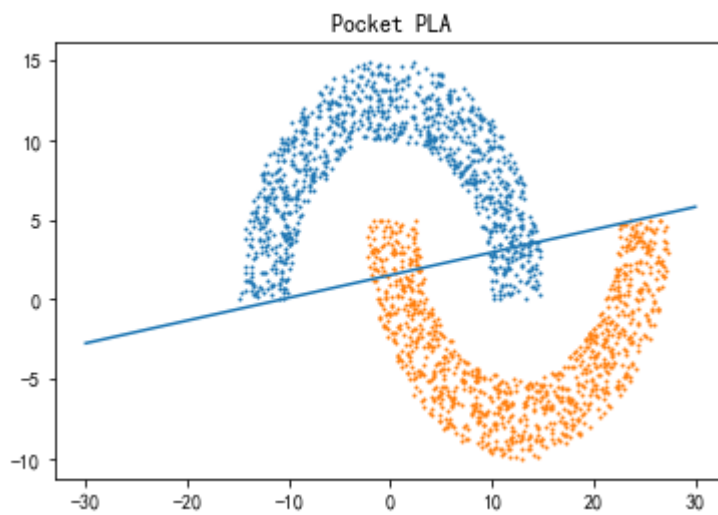
(c)

```

#(c)
r = 2 * (rad + thk)
a1 = np.array([-r,r])
b1 = - (w[0] + w[1] * a1) / w[2]

plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.plot(a1, b1)
plt.title('Pocket PLA')
plt.show()
print('Pocket PLA的错误率为' + str(error / N))

```



Pocket PLA的错误率为5e-05

(d)linear regression, 之前有处理过, 直接带公式即可

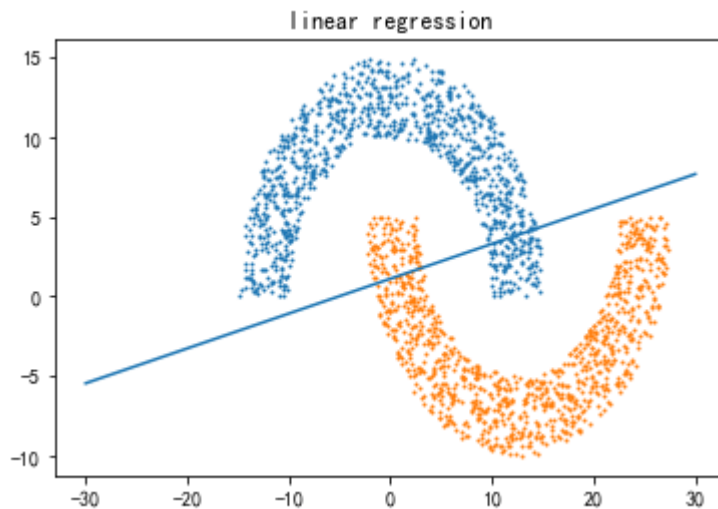
```

#(d)
#Linear regression
w_lr = inv(X_treat.T.dot(X_treat)).dot(X_treat.T).dot(y)

#作图
a2 = np.array([-r,r])
b2 = - (w_lr[0] + w_lr[1] * a1) / w_lr[2]

plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.plot(a2, b2)
plt.title('linear regression')
plt.show()
error = np.mean(np.sign(X_treat.dot(w_lr)) != y)
print('linear regression的错误率为' + str(error))

```



linear regression的错误率为0.1

(e)先做三次特征转换，再重复(b)到(d)，注意到三次特征转换为

$$\phi_3(x) = (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2, x_1^3, x_1^2 x_2, x_1 x_2^2, x_2^3)$$

这里特征转换使用了scikit-learn，后续作曲线的图用到了plt.contour函数，原本是用来绘制等高线的，这里可以用来画隐函数的图像

```
#(e)
#特征转换
poly = PolynomialFeatures(3)
X_poly = poly.fit_transform(X)

# 定义等高线高度函数
def f(x1, x2, w):
    #将网格拉直并拼接
    X = np.c_[x1.reshape(-1, 1), x2.reshape(-1, 1)]
    #多项式转换
    poly = PolynomialFeatures(3)
    X_poly = poly.fit_transform(X)

    #计算结果
    result = X_poly.dot(w)
    #恢复成网格形状
    result = np.reshape(result, np.shape(x1))
    return result

#数据数目
n = 2000
#定义a, b
a = np.linspace(-r, r, n)
b = np.linspace(-r, r, n)

#生成网格数据
A, B = np.meshgrid(a, b)
```

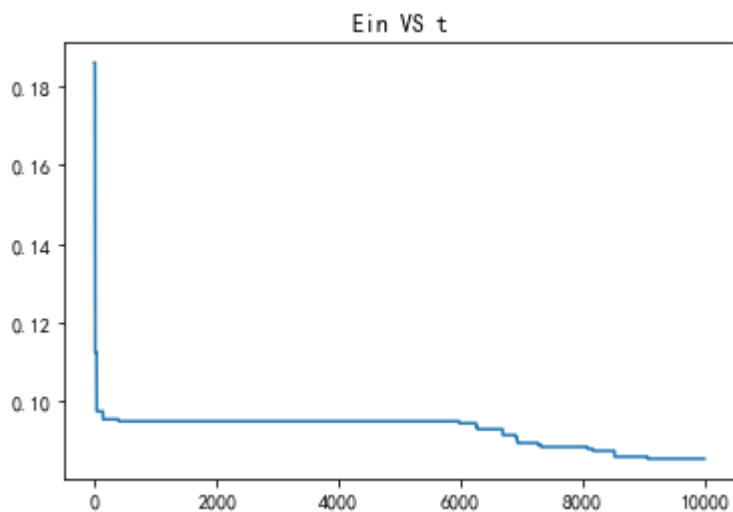
```

#迭代次数
max_step = 10000

#Pocket_PLA
w_poly, w_poly_hat, w_poly, error_poly = Pocket_PLA(X_poly, y, max_step=max_step)
ein_poly = np.mean(np.sign(w_poly_hat.dot(X_poly.T)) != y, axis=1)

#(b)
plt.plot(t, ein_poly)
plt.title('Ein VS t')
plt.show()

```

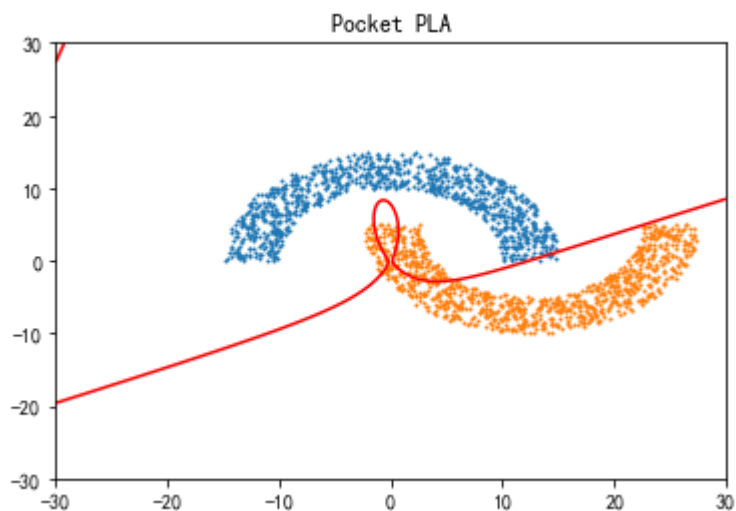


作图

```

#(c)
plt.contour(A, B, f(A, B, w_poly), 1, colors = 'red')
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.title('Pocket PLA')
plt.show()
print('特征转换后的Pocket PLA的错误率为' + str(error_poly / N))

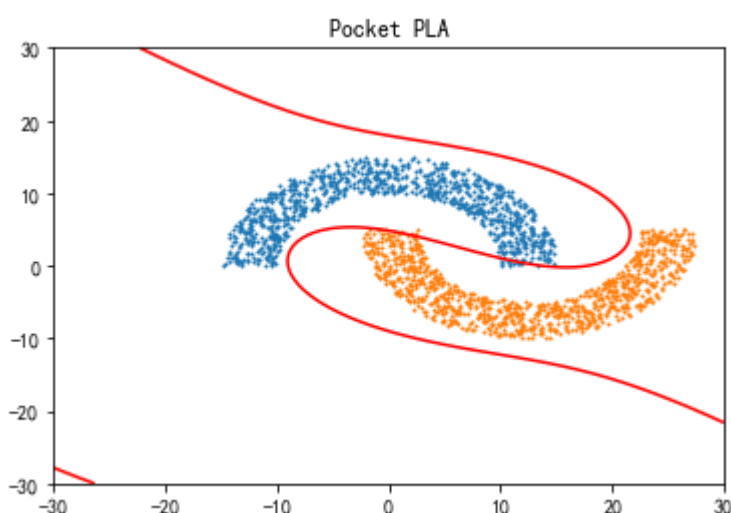
```



特征转换后的Pocket PLA的错误率为0.0855

将转换后的数据带入linear regression

```
#(d)Linear regression
w_poly_lr = inv(X_poly.T.dot(X_poly)).dot(X_poly.T).dot(y)
plt.contour(A, B, f(A, B, w_poly_lr), 1, colors = 'red')
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.title('Pocket PLA')
plt.show()
error = np.mean(np.sign(X_poly.dot(w_poly_lr)) != y)
print('特征转换后的linear regression的错误率为' + str(error))
```



特征转换后的linear regression的错误率为0.0085

Problem 3.4 (Page 110)

In Problem 1.5, we introduced the Adaptive Linear Neuron (Ada line) algorithm for classification. Here, we derive Ada line from an optimization perspective.

- (a) Consider $E_n(w) = (\max(0, 1 - y_n w^T x_n))^2$. Show that $E_n(w)$ is continuous and differentiable. Write down the gradient $\nabla E_n(w)$.
- (b) Show that $E_n(w)$ is an upper bound for $\mathbb{I}[\text{sign}(w^T x_n) \neq y_n]$. Hence, $\frac{1}{N} \sum_{n=1}^N E_n(w)$ is an upper bound for the in sample classification error $E_{\text{in}}(w)$.
- (c) Argue that the Adaline algorithm in Problem 1.5 performs stochastic gradient descent on $\frac{1}{N} \sum_{n=1}^N E_n(w)$.
- (a) $1 - y_n w^T x_n$ 关于 w 是连续的, $\max(a, x)$ 关于 x 是连续的, 所以 $\max(0, 1 - y_n w^T x_n)$ 关于 w 是连续的, 连续函数的平方也是连续的, 所以 $E_n(w) = (\max(0, 1 - y_n w^T x_n))^2$ 关于 w 连续。再来看可导性, 令 $s(w) = 1 - y_n w^T x_n$, $s(w)$ 显然关于 w 可导性,我们再看下 $f(s) = (\max(0, s))^2$

$$f(s) = \begin{cases} s^2 & s \geq 0 \\ 0 & s < 0 \end{cases}$$

显然这个函数也是可导的。所以 $E_n(w) = f(s(w))$ 也可导。接下来我们求梯度

$$\frac{\partial E_n(w)}{\partial w_k} = \begin{cases} \frac{\partial(1-y_n w^T x_n)^2}{\partial w_k} = (1 - y_n w^T x_n)(-y_n x_n^k) & y_n w^T x_n \leq 1 \\ 0 & y_n w^T x_n > 1 \end{cases}$$

因此

$$\nabla E_n(w) = \begin{cases} (1 - y_n w^T x_n)(-y_n x_n) & y_n w^T x_n \leq 1 \\ 0 & y_n w^T x_n > 1 \end{cases}$$

(b)作图，我们先对式子做点变形

$$\begin{aligned} \llbracket \text{sign}(w^T x_n) \neq y_n \rrbracket &\Leftrightarrow \\ \llbracket y_n \times \text{sign}(w^T x_n) \neq y_n \times y_n \rrbracket &\Leftrightarrow \\ \llbracket \text{sign}(y_n w^T x_n) \neq 1 \rrbracket \end{aligned}$$

令 $s = \text{sign}(y_n w^T x_n)$

说以上两个式子可以化为 $(\max(0, 1 - s))^2$ 以及 $\llbracket \text{sign}(s) \neq 1 \rrbracket$

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 12:02:33 2019

@author: qinzhen
"""

import matplotlib.pyplot as plt
import numpy as np
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号

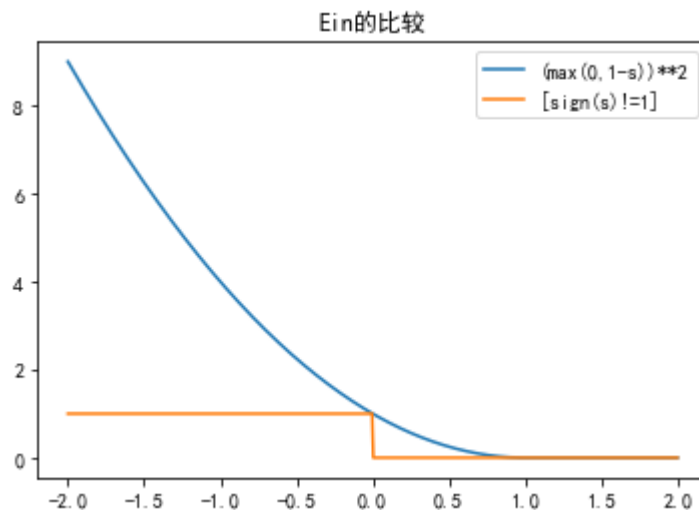
def f1(s):
    a = max(0, 1 - s)
    return a ** 2

def f2(s):
    if s > 0:
        return 0
    else:
        return 1

x = np.linspace(-2, 2, 500)
y1 = [f1(i) for i in x]
y2 = [f2(i) for i in x]

plt.plot(x, y1, label="(max(0,1-s))**2")
plt.plot(x, y2, label="[sign(s)!=1]")
```

```
plt.legend()
plt.title('Ein的比较')
plt.show()
```



(c)回顾Problem 1.5的更新规则

$$s(t) = w^T(t)x(t)$$

$$\text{当 } y(t) \cdot s(t) \leq 1 \text{ 时, } w(t+1) = w(t) + \eta(y(t) - s(t)) \cdot x(t)$$

再来看一下我们的梯度，稍微做下变形

$$(1 - y(n)w^T x(n))(-y(n)x(n)) = -(y(n) - w^T x(n))x(n)$$

$$\nabla E_n(w) = \begin{cases} -(y(n) - w^T x(n))x(n) & y(n)w^T x(n) \leq 1 \\ 0 & y(n)w^T x(n) > 1 \end{cases}$$

所以随机梯度下降法的更新规则为

$$w(t+1) = w(t) - \eta \nabla E_t(w) = \begin{cases} w(t) + \eta(y(t) - w(t)^T x(t))x(t) & y(t)w(t)^T x(t) \leq 1 \\ w(t) & y(t)w(t)^T x(t) > 1 \end{cases}$$

我们使用Problem 1.5一样的符号 $s(t) = w^T(t)x(t)$ ，那么随机梯度下降法的更新规则即为Problem 1.5的更新规则。

Problem 3.5 (Page 110)

(a) Consider

$$E_n(w) = \max(0, 1 - y_n w^T x_n)$$

Show that $E_n(w)$ is continuous and differentiable except when $y_n = w^T x_n$.

(b) Show that $E_n(w)$ is an upper bound for $\mathbb{I}[\text{sign}(w^T x_n) \neq y_n]$. Hence, $\frac{1}{N} \sum_{n=1}^N E_n(w)$ is an upper bound for the in sample classification error $E_{\text{in}}(w)$.

(c) Apply stochastic gradient descent on $\frac{1}{N} \sum_{n=1}^N E_n(w)$ (ignoring the singular case of $y_n = w^T x_n$) and derive a new perceptron learning algorithm.

(a)我们用上一题一样的思路, 令 $s = y_n w^T x_n$, $f(s) = \max(0, 1 - s)$ 关于 s 连续, s 关于 w 连续, 因此 $E_n(w) = f(s(w))$ 关于 w 连续。

s 关于 w 处处可导, 但 $f(s) = \max(0, 1 - s)$ 在 $s = 1$ 处不可导, 其余点均可导。我们来看下 $s = 1$ 的特点, 注意 $y_n \in \{1, -1\}$, 那么

$$\begin{aligned} s = 1 &\Leftrightarrow \\ y_n w^T x_n = 1 &\Leftrightarrow \\ y_n \times y_n w^T x_n = y_n &\Leftrightarrow \\ w^T x_n = y_n \end{aligned}$$

所以 $E_n(w) = f(s(w))$ 在 $s = 1$ 即 $y_n = w^T x_n$ 处不可导, 其余点均可导。

(b)同Problem 3.4方法, $s = y_n w^T x_n$

$$\begin{aligned} E_n(w) &= \max(0, 1 - y_n w^T x_n) = \max(0, 1 - s) \\ E_{\text{in}} &= \mathbb{I}[\text{sign}(s) \neq 1] \end{aligned}$$

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 12:02:33 2019

@author: qinzhen
"""

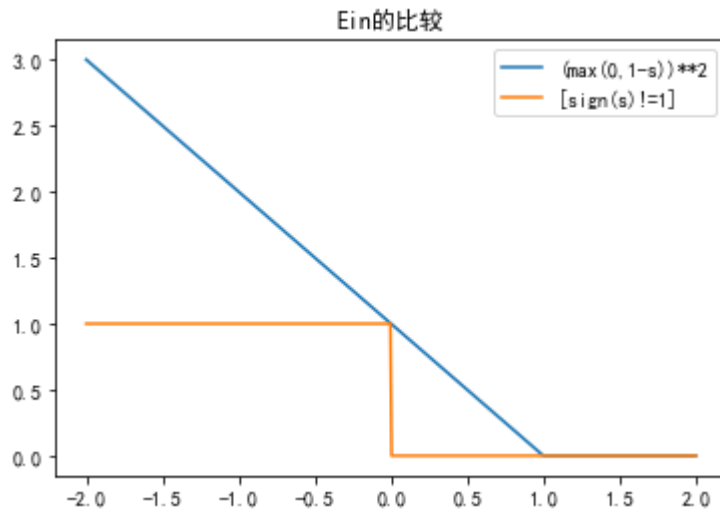
import matplotlib.pyplot as plt
import numpy as np
plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
plt.rcParams['axes.unicode_minus']=False #用来正常显示负号

def f1(s):
    a = max(0, 1 - s)
    return a

def f2(s):
    if s > 0:
        return 0
    else:
        return 1

x = np.linspace(-2, 2, 500)
y1 = [f1(i) for i in x]
y2 = [f2(i) for i in x]

plt.plot(x, y1, label="(max(0,1-s))**2")
plt.plot(x, y2, label="[sign(s)!=1]")
plt.legend()
plt.title('Ein的比较')
plt.show()
```



(c)先不管不可导点，我们来求梯度，只考虑 $y_n w^T x_n < 1$ 的情形，此时

$$E_n(w) = 1 - y_n w^T x_n$$

$$\frac{\partial E_n(w)}{\partial w_i} = \frac{\partial(1 - y_n w^T x_n)}{\partial w_i} = -y_n x_n^i$$

所以

$$\text{当 } y_n w^T x_n < 1 \text{ 时, } \nabla E_n(w) = -y_n x_n$$

$$\text{当 } y_n w^T x_n \geq 1 \text{ 时, } \nabla E_n(w) = 0$$

所以SGD（随机梯度下降法）的更新规则为

$$\text{当 } y(t)w(t)^T x(t) < 1 \text{ 时}$$

$$w(t+1) = w(t) - \eta \nabla E_t(w) = w(t) + \eta y(t)x(t)$$

$$\text{当 } y(t)w(t)^T x(t) \geq 1 \text{ 时不更新}$$

Problem 3.6 (Page 110)

Derive a linear programming algorithm to fit a linear model for classification using the following steps. A linear program is an optimization problem of the following form:

$$\min_z \quad c^T z$$

$$\text{subject to} \quad Az \leq b$$

A , b and c are parameters of the linear program and z is the optimization variable. This is such a well studied optimization problem that most mathematics software have canned optimization functions which solve linear programs.

(a) For linearly separable data, show that for some w , $y_n(w^T x_n) \geq 1$ for $n = 1, \dots, N$.

(b) Formulate the task of finding a separating w for separable data as a linear program. You need to specify what the parameters A , b , c are and what the optimization variable z is.

(c) If the data is not separable, the condition in (a) cannot hold for every n . Thus introduce the violation $\xi_n \geq 0$ to capture the amount of violation for example x_n . So, for $n = 1, \dots, N$,

$$\begin{aligned} y_n(w^T x_n) &\geq 1 - \xi_n \\ \xi_n &\geq 0 \end{aligned}$$

Naturally, we would like to minimize the amount of violation. One intuitive approach is to minimize $\sum_{n=1}^N \xi_n$, i.e., we want w that solves

$$\begin{aligned} \min_{w, \xi_n} \quad & \sum_{n=1}^N \xi_n \\ \text{subject to} \quad & y_n(w^T x_n) \geq 1 - \xi_n \\ & \xi_n \geq 0 \end{aligned}$$

where the inequalities must hold for $n = 1, \dots, N$. Formulate this problem as a linear program.

(d) Argue that the linear program you derived in (c) and the optimization problem in Problem 3.5 are equivalent.

这题的任务是要把我们之前讨论的线性分类问题转化为线性规划问题。

(a)由第一章的结论, 对于线性可分的数据, 存在 w_1 , 使得 $y_n w_1^T x_n > 0 (n = 1, \dots, N)$, 设 $\rho = \min_{1 \leq n \leq N} y_n w_1^T x_n$, 显然 $\rho > 0$, 现在取 $w = \frac{w_1}{\rho}$, 那么

$$y_n w^T x_n = y_n \left(\frac{w_1}{\rho} \right)^T x_n = \frac{y_n w_1^T x_n}{\rho} \geq 1$$

因此结论成立。

(b)这题的限制条件就是刚刚所说的 $y_n w^T x_n \geq 1$, 因此 $z = w$, 比较让人费解的是 c 应该取什么, 实际上思考下, 我们这里只要找到满足 $y_n w^T x_n \geq 1, n = 1, \dots, N$ 这个条件的 w 即可, 所以这里 c 可以取任意值。结合以上几点, 下面把 A, b, c 分别写出, 不妨设 $w, x_n \in \mathbb{R}^d, n = 1, \dots, N$ 。

$$\begin{aligned} z = w &= (w_1, \dots, w_d)^T \\ x_n &= (x_n^1, \dots, x_n^d)^T \\ A &= \begin{pmatrix} -y_1 x_1^T \\ \vdots \\ -y_N x_N^T \end{pmatrix} = \begin{pmatrix} -y_1 x_1^1 & \dots & -y_1 x_1^d \\ \vdots & \dots & \vdots \\ -y_N x_N^1 & \dots & -y_N x_N^d \end{pmatrix} \in \mathbb{R}^{N \times d} \\ b &= \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix} \in \mathbb{R}^N \\ c &\text{为 } \mathbb{R}^d \text{ 中任意向量} \end{aligned}$$

我们看下这里的 Az, b

$$Az = \begin{pmatrix} -y_1 x_1^T \\ \dots \\ -y_N x_N^T \end{pmatrix} w = \begin{pmatrix} -y_1 x_1^T w \\ \dots \\ -y_N x_N^T w \end{pmatrix} = \begin{pmatrix} -y_1 w^T x_1 \\ \dots \\ -y_N w^T x_N \end{pmatrix}$$

$$b = \begin{pmatrix} -1 \\ \dots \\ -1 \end{pmatrix}$$

因此 $Az \leq b$ 即为 $y_n w^T x_n \geq 1$ 。

(c) 依旧设 $w, x_n \in \mathbb{R}^d, n = 1, \dots, N$, 和上一题类似的思路

$$z = (w_1, \dots, w_d, \xi_1, \dots, \xi_N)^T \in \mathbb{R}^{N+d}$$

记 $A_1 = \begin{pmatrix} -y_1 x_1^T \\ \dots \\ -y_N x_N^T \end{pmatrix} \in \mathbb{R}^{N \times d}, I_{N \times N}$ 为 $N \times N$ 阶单位矩阵

$$A_2 = \left(A_1 \mid -I_{N \times N} \right) \in \mathbb{R}^{N \times (N+d)}$$

$$A_3 = \left(0 \mid -I_{N \times N} \right) \in \mathbb{R}^{N \times (N+d)} \text{ (0 为 } N \times d \text{ 阶 0 矩阵)}$$

$$A = \begin{pmatrix} A_2 \\ A_3 \end{pmatrix} \in \mathbb{R}^{(2N) \times (N+d)}$$

$$b = (-1 \dots -1, 0 \dots 0)^T \in \mathbb{R}^{2N}, \text{ 其中前 } N \text{ 个分量为 } -1, \text{ 其余为 } 0$$

$$c = (0, \dots, 0, 1, \dots, 1)^T \in \mathbb{R}^{N+d}, \text{ 其中 } c \text{ 的前 } d \text{ 个分量为 } 0, \text{ 后 } N \text{ 的分量为 } 1$$

同上一题的验证方法可以知此问题即为原来的问题。

(d) 回顾下 3.5

$$E_n(w) = \max(0, 1 - y_n w^T x_n)$$

我们的目标是最小化 $\frac{1}{N} \sum_{n=1}^N E_n(w)$

这里我们令 $\xi_n = E_n(w) = \max(0, 1 - y_n w^T x_n)$, 那么

$$\begin{aligned} 1 - y_n w^T x_n &\leq \xi_n \\ 0 &\leq \xi_n \\ \frac{1}{N} \sum_{n=1}^N E_n(w) &= \frac{1}{N} \sum_{n=1}^N \xi_n \end{aligned}$$

注意 N 为常数, 所以 3.5 即为

$$\text{在条件 } 1 - y_n w^T x_n \leq \xi_n \text{ 和 } 0 \leq \xi_n \text{ 下最小化 } \sum_{n=1}^N \xi_n$$

这就是我们刚刚考虑的问题。

Problem 3.7 (Page 111)

Use the linear programming algorithm from Problem 3.6 on the learning task in Problem 3.1 for the separable ($sep = 5$) and the non separable ($sep = -5$) cases. Compare your results to the linear regression approach with and without the 3rd order polynomial feature transform.

这题的意思是使用线性规划方法对3.1再实践一遍，这里遇到个坑，一开始使用scipy的linprog函数做优化，发现一直显示无解，后来网上一查发现有人说这个函数如果数据量大于100就会求不出解，然后我将我的数据量由1000改为50，果然一下有解了，最后的解决方法是cvxopt包，这个包的效果非常好。后面记第一个算法为算法1，第二个算法为算法2,分别给出结果。

首先给出辅助函数

```
# -*- coding: utf-8 -*-
"""
Created on Sun Mar  3 12:27:30 2019

@author: qinzhen
"""

import numpy as np
import matplotlib.pyplot as plt
from cvxopt import matrix, solvers
from helper import generatedata
from sklearn.preprocessing import PolynomialFeatures

def algorithm1(X, y):
    """
    算法1
    """
    N, d = X.shape
    c = np.array(np.ones(d))
    A = X * y.reshape(-1, 1)
    b = np.ones(N) * (-1.0)

    #转化为cvxopt中的数据结构
    c = matrix(c)
    A = matrix(A)
    b = matrix(b)

    sol = solvers.lp(c, A, b)

    w = np.array((sol['x']))

    return w

def algorithm2(X, y):
    """
    算法2
    """
    N, d = X.shape
```

```

A1 = x * y.reshape(-1, 1)
A2 = np.c_[A1, (-1) * np.eye(N)]
A3 = np.c_[np.zeros((N, d)), (-1) * np.eye(N)]

A = np.r_[A2, A3]
c = np.array([0.0] * d + [1.0] * N)
b = np.array([-1.0] * N + [0.0] * N)

#带入算法求解
c = matrix(c)
A = matrix(A)
b = matrix(b)

sol = solvers.lp(c, A, b)

#返回向量
w = np.array((sol['x']))[:d]

return w

def draw(w, x, y, r, text, num):
    """
    作图
    """
    #作出直线
    a1 = np.array([-r, r])
    b1 = - (w[0] + w[1] * a1) / w[2]

    plt.scatter(x[y>0][:, 0], x[y>0][:, 1], s=1)
    plt.scatter(x[y<0][:, 0], x[y<0][:, 1], s=1)
    plt.plot(a1, b1, c="red")
    plt.title('sep={},algorithm{}'.format(text, num))
    plt.show()

    print(w)

# 定义等高线高度函数
def f(x1, x2, w):
    #将网格拉直并拼接
    x = np.c_[x1.reshape(-1, 1), x2.reshape(-1, 1)]
    #多项式转换
    poly = PolynomialFeatures(3)
    x_poly = poly.fit_transform(x)

    #计算结果
    result = x_poly.dot(w)
    #恢复成网格形状
    result = np.reshape(result, np.shape(x1))
    return result

```

作图


```

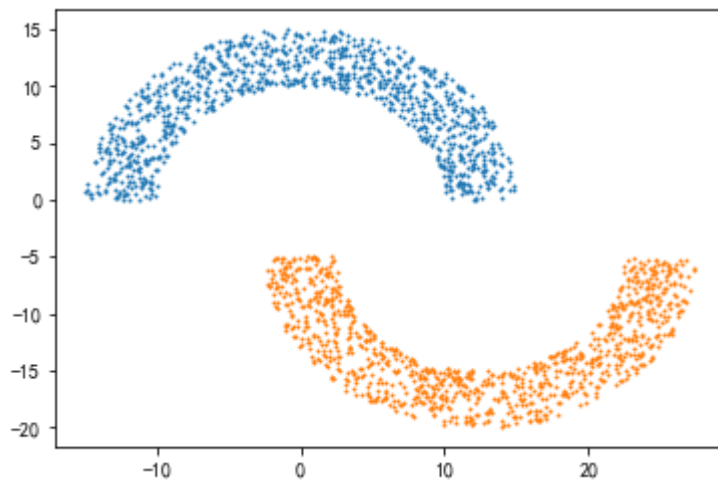
#参数
rad = 10
thk = 5
sep = 5
N = 2000
r = 2 * (rad + thk)

# =====
# 特征转换之前
# =====

#产生数据
X, y = generatedata(rad, thk, sep, N)

#作图
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.show()

```



(1)特征转换之前, 算法1, sep=5

```

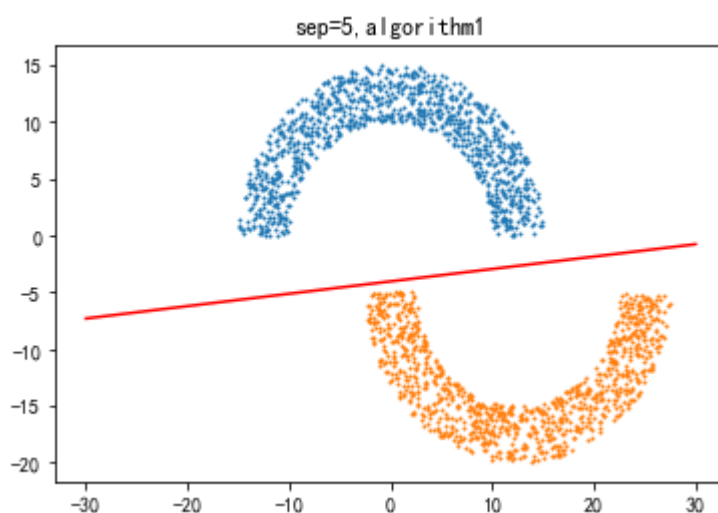
#特征转换之前, 算法1, sep=5
#对数据预处理, 加上偏置项1
X_treat = np.c_[np.ones(N), X]
w = algorithm1(X_treat, y)

#作图
draw(w, X, y, r, sep, 1)

```

	pcost	dcost	gap	pres	dres	k/t
0:	-3.1759e-01	2.0016e+03	3e+03	2e+00	1e+04	1e+00
1:	-3.0300e+01	9.8235e+04	2e+07	9e+01	7e+05	2e+02
2:	-5.8143e-01	8.9037e+02	1e+03	7e-01	5e+03	1e+02
3:	-9.6129e-01	1.3741e+03	4e+03	1e+00	8e+03	1e+02
4:	-2.5363e+00	1.9692e+02	5e+02	9e-02	7e+02	1e+02
5:	-2.7228e+00	5.0243e+00	2e+01	3e-03	2e+01	5e+00
6:	-9.9669e+00	2.5819e+00	4e+01	2e-03	1e+01	1e+01
7:	-7.7044e+02	5.0447e+00	6e+03	4e-03	3e+01	8e+02
8:	-7.7050e+04	5.0452e+00	6e+05	4e-03	3e+01	8e+04
9:	-7.7051e+06	5.0452e+00	6e+07	4e-03	3e+01	8e+06
10:	-7.7051e+08	5.0452e+00	6e+09	4e-03	3e+01	8e+08

Certificate of dual infeasibility found.



```
[[ -0.81975795]
 [ 0.02205008]
 [-0.20229213]]
```

可以看到线性规划产生的直线比感知机的结果更好一些，因为离两组数据更远。

(2)特征转换之前，算法2，sep=5

```
#特征转换之前，算法2，sep=5
w = algorithm2(X_treat, y)

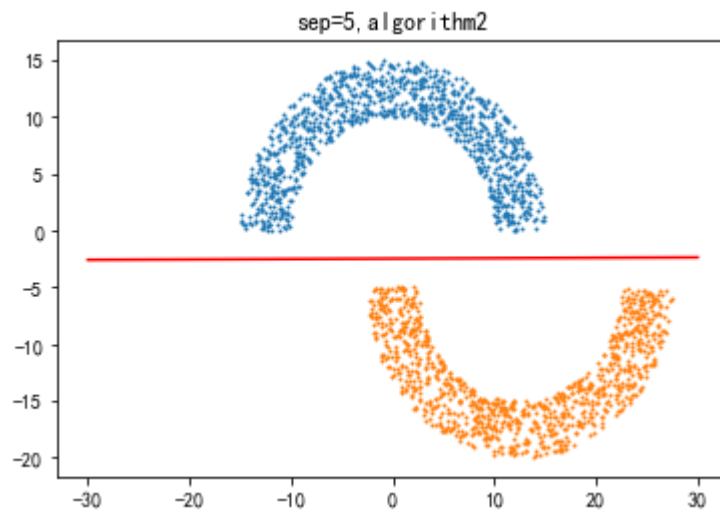
#作图
draw(w, X, y, r, sep, 2)
```

	pcost	dcost	gap	pres	dres	k/t
0:	1.2392e+02	2.5777e+03	1e+04	2e+00	7e+02	1e+00
1:	2.0224e+02	5.3602e+02	8e+02	3e-01	9e+01	1e+00
2:	1.3304e+02	3.0029e+02	4e+02	1e-01	4e+01	7e-01
3:	9.4426e+01	2.0200e+02	3e+02	9e-02	3e+01	4e-01
4:	6.9714e+01	1.4714e+02	3e+02	6e-02	2e+01	3e-01
5:	4.9776e+01	1.0343e+02	2e+02	4e-02	1e+01	2e-01
6:	3.8943e+01	8.0358e+01	2e+02	3e-02	1e+01	2e-01

```

7: 3.0024e+01 6.1635e+01 1e+02 3e-02 8e+00 1e-01
8: 2.5337e+01 5.2056e+01 1e+02 2e-02 7e+00 9e-02
9: 1.6925e+01 3.5077e+01 8e+01 1e-02 5e+00 5e-02
10: 1.0277e+01 2.1343e+01 5e+01 9e-03 3e+00 2e-02
11: 7.4621e+00 1.5521e+01 3e+01 7e-03 2e+00 1e-02
12: 2.6330e+00 5.4722e+00 1e+01 2e-03 8e-01 5e-03
13: 1.3329e-01 2.7832e-01 6e-01 1e-04 4e-02 2e-04
14: 1.3452e-03 2.8092e-03 6e-03 1e-06 4e-04 2e-06
15: 1.3452e-05 2.8092e-05 6e-05 1e-08 4e-06 2e-08
16: 1.3452e-07 2.8092e-07 6e-07 1e-10 4e-08 2e-10
17: 1.3452e-09 2.8092e-09 6e-09 1e-12 4e-10 2e-12
Optimal solution found.

```



```

[[-1.00401277]
 [ 0.00145698]
 [-0.40313334]]

```

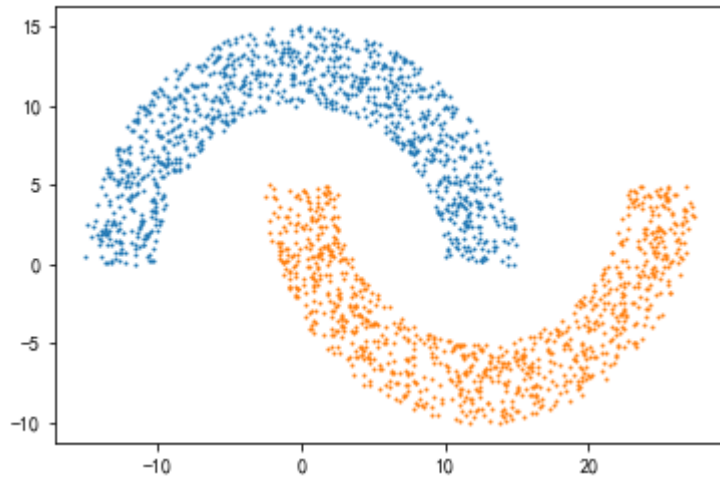
(3)作图

```

#特征转换之前, 算法1, sep=-5
#产生数据
sep = -5
X, y = generatedata(rad, thk, sep, N)

#作图
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.show()

```



特征转换之前, 算法1, sep=-5

```
#对数据预处理, 加上偏置项1
X_treat = np.c_[np.ones(N), X]
#特征转换之前, 算法1, sep=-5
w = algorithm1(X_treat, y)
print(w)
```

	pcost	dcost	gap	pres	dres	k/t
0:	3.8892e-02	2.0027e+03	4e+03	2e+00	9e+03	1e+00
1:	3.7512e+00	1.2454e+05	3e+07	2e+02	6e+05	3e+02
2:	1.0853e-01	1.8549e+03	4e+03	2e+00	7e+03	2e+02
3:	9.3465e-02	2.2731e+03	6e+03	2e+00	9e+03	4e+02
4:	4.1975e-01	4.2439e+03	2e+04	3e+00	1e+04	1e+03
5:	8.1723e-01	7.0654e+03	2e+04	3e+00	1e+04	5e+03
6:	7.4479e-01	1.0870e+04	3e+04	3e+00	1e+04	8e+03
7:	2.6590e+00	1.9385e+05	1e+06	6e+00	2e+04	2e+05
8:	2.7433e+00	1.9794e+07	1e+08	7e+00	2e+04	2e+07
9:	2.7432e+00	1.9800e+09	1e+10	7e+00	2e+04	2e+09
10:	2.7432e+00	1.9800e+11	1e+12	7e+00	2e+04	2e+11
11:	2.7432e+00	1.9800e+13	1e+14	7e+00	2e+04	2e+13

Certificate of primal infeasibility found.
None

此时无解, 因为不可分。

(4)特征转换之前, 算法2, sep=-5

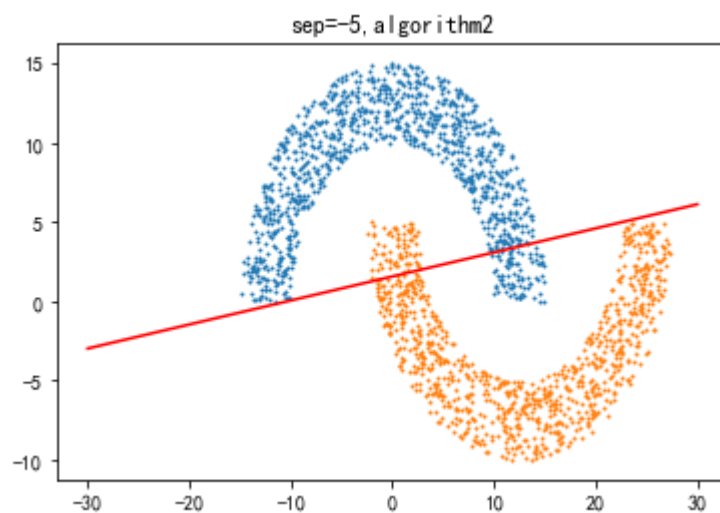
```
#特征转换之前, 算法2, sep=-5
w = algorithm2(X_treat, y)
draw(w, X, y, r, sep, 2)
```

	pcost	dcost	gap	pres	dres	k/t
0:	2.9880e+02	2.7548e+03	1e+04	2e+00	4e+02	1e+00
1:	5.2678e+02	8.1395e+02	8e+02	3e-01	5e+01	1e+00
2:	4.7265e+02	5.8905e+02	3e+02	1e-01	2e+01	5e-01
3:	4.4976e+02	5.2151e+02	2e+02	7e-02	1e+01	3e-01

```

4: 4.3785e+02 4.9056e+02 2e+02 5e-02 9e+00 2e-01
5: 4.2829e+02 4.6615e+02 1e+02 4e-02 7e+00 2e-01
6: 4.2304e+02 4.5307e+02 1e+02 3e-02 5e+00 1e-01
7: 4.1840e+02 4.4158e+02 8e+01 2e-02 4e+00 9e-02
8: 4.1427e+02 4.3077e+02 6e+01 2e-02 3e+00 6e-02
9: 4.1062e+02 4.2090e+02 4e+01 1e-02 2e+00 3e-02
10: 4.0860e+02 4.1522e+02 2e+01 7e-03 1e+00 2e-02
11: 4.0707e+02 4.1083e+02 1e+01 4e-03 7e-01 9e-03
12: 4.0569e+02 4.0676e+02 4e+00 1e-03 2e-01 1e-03
13: 4.0529e+02 4.0546e+02 7e-01 2e-04 3e-02 2e-04
14: 4.0521e+02 4.0523e+02 7e-02 2e-05 3e-03 2e-05
15: 4.0520e+02 4.0521e+02 1e-03 4e-07 7e-05 3e-07
16: 4.0520e+02 4.0520e+02 1e-05 4e-09 7e-07 3e-09
17: 4.0520e+02 4.0520e+02 1e-07 4e-11 7e-09 3e-11
Optimal solution found.

```



```

[[ 0.71176934]
 [ 0.06899415]
 [-0.45587043]]

```

第二种算法没有要求可分，所以是有解的。

(5)特征转换后，算法1，sep=-5

```

# =====
# 特征转换后
# =====
#特征转换器
poly = PolynomialFeatures(3)

#特征转换后，算法1，sep=-5
#特征转换
sep = 5
X, y = generatedata(rad, thk, sep, N)
X_poly = poly.fit_transform(X)
w_poly = algorithm1(X_poly, y)

```

```

#数据数目
n = 2000

#定义a, b
a = np.linspace(-r, r, n)
b = np.linspace(-r, r, n)

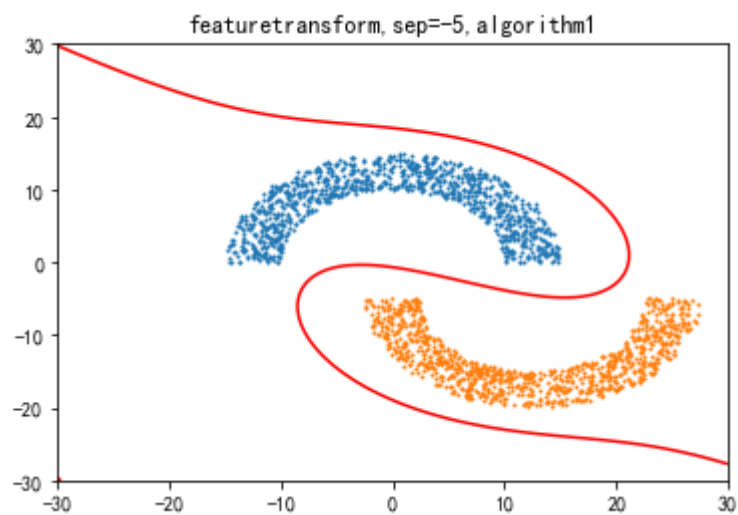
#生成网格数据
A, B = np.meshgrid(a, b)

plt.contour(A, B, f(A, B, w_poly), 1, colors = 'red')
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.title('featuretransform,sep=-5,algorithm1')
plt.show()
print(w_poly)

```

	pcost	dcost	gap	pres	dres	k/t
0:	-2.6260e-01	2.0090e+03	3e+03	2e+00	2e+06	1e+00
1:	-2.6438e+01	6.6850e+04	1e+07	5e+01	7e+07	2e+02
2:	-2.8418e+00	1.5497e+03	2e+04	1e+00	2e+06	3e+01
3:	-4.3127e+00	3.0736e+01	5e+02	2e-02	3e+04	5e+00
4:	-4.3566e+02	4.6856e+01	7e+04	4e-02	5e+04	4e+02
5:	-4.3566e+04	4.6861e+01	7e+06	4e-02	5e+04	4e+04
6:	-4.3566e+06	4.6861e+01	7e+08	4e-02	5e+04	4e+06
7:	-4.3566e+08	4.6861e+01	7e+10	4e-02	5e+04	4e+08

Certificate of dual infeasibility found.



```
[[ -0.35594275]
 [ -0.11587162]
 [ -0.50877516]
 [ -0.01480921]
 [ -0.01189737]
 [  0.0017957 ]
 [  0.00100017]
 [  0.00154855]
 [  0.00149536]
 [  0.00145634]]
```

效果还是相当不错的，最后使用算法2。

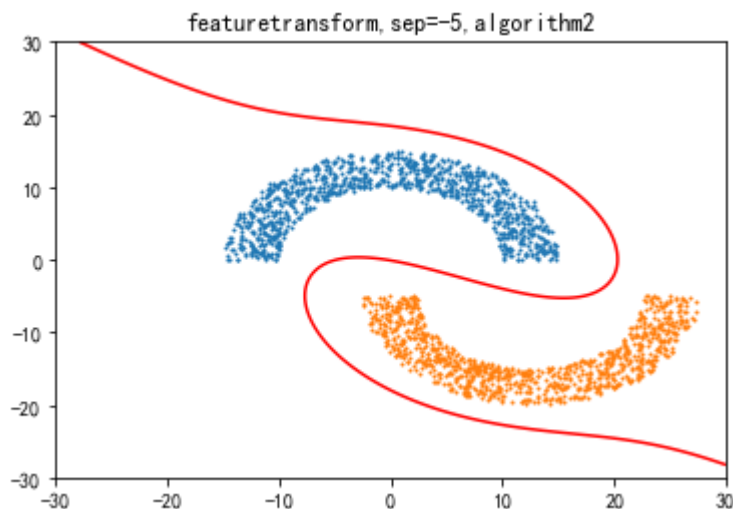
(6)特征转换后，算法2，sep=-5

```
#特征转换后，算法2，sep=-5
#根据之前所述构造矩阵
w_poly = algorithm2(X_poly, y)

plt.contour(A, B, f(A, B, w_poly), 1, colors = 'red')
plt.scatter(X[y>0][:, 0], X[y>0][:, 1], s=1)
plt.scatter(X[y<0][:, 0], X[y<0][:, 1], s=1)
plt.title('featuretransform,sep=-5,algorithm2')
plt.show()
print(w_poly)
```

	pcost	dcost	gap	pres	dres	k/t
0:	2.2028e+01	2.4218e+03	9e+03	2e+00	2e+05	1e+00
1:	2.8606e+01	4.2654e+02	8e+02	3e-01	3e+04	1e+00
2:	8.9041e+00	1.2097e+02	2e+02	8e-02	9e+03	3e-01
3:	2.3262e+00	3.3645e+01	6e+01	2e-02	2e+03	8e-02
4:	1.7501e-01	2.8593e+00	5e+00	2e-03	2e+02	6e-03
5:	1.7943e-03	2.9351e-02	5e-02	2e-05	2e+00	6e-05
6:	1.7943e-05	2.9351e-04	5e-04	2e-07	2e-02	6e-07
7:	1.7943e-07	2.9351e-06	5e-06	2e-09	2e-04	6e-09
8:	1.7943e-09	2.9351e-08	5e-08	2e-11	2e-06	6e-11
9:	1.7943e-11	2.9351e-10	5e-10	2e-13	2e-08	6e-13

Optimal solution found.



```

[[-2.82682691e-02]
 [-6.96701940e-02]
 [-2.59551250e-01]
 [-9.96959284e-03]
 [-8.34923122e-03]
 [-2.50263962e-04]
 [ 6.63553681e-04]
 [ 1.03560253e-03]
 [ 9.59105454e-04]
 [ 7.83260924e-04]]

```

算法2的效果更加的好，因为它限制了距离。

Problem 3.8 (Page 111)

For linear regression, the out of sample error is

$$E_{\text{out}}(h) = \mathbb{E}[(h(x) - y)^2]$$

Show that among all hypotheses, the one that minimizes $E_{\text{out}}(h)$ is given by

$$h^*(x) = \mathbb{E}[y|x]$$

The function h^* can be treated as a deterministic target function, in which case we can write $y = h^*(x) + \epsilon(x)$ where $\epsilon(x)$ is an (input dependent) noise variable. Show that $\epsilon(x)$ has expected value zero.

这题其实是统计学习里一个比较常见的结论。

$$\begin{aligned}
 E_{\text{out}}(h) &= \mathbb{E}[(h(x) - y)^2] \\
 &= \mathbb{E}[(h(x) - \mathbb{E}[y|x] + \mathbb{E}[y|x] - y)^2] \\
 &= \mathbb{E}[(h(x) - \mathbb{E}[y|x])^2 + (\mathbb{E}[y|x] - y)^2 + 2(h(x) - \mathbb{E}[y|x])(\mathbb{E}[y|x] - y)] \\
 &= \mathbb{E}[(h(x) - \mathbb{E}[y|x])^2] + \mathbb{E}[(\mathbb{E}[y|x] - y)^2] + 2\mathbb{E}[(h(x) - \mathbb{E}[y|x])(\mathbb{E}[y|x] - y)] \\
 &= \mathbb{E}[(h(x) - h^*(x))^2] + \mathbb{E}[(h^*(x) - y)^2] + 2\mathbb{E}[(h(x) - h^*(x))(h^*(x) - y)]
 \end{aligned}$$

下面分析 $\mathbb{E}[(h(x) - h^*(x))(h^*(x) - y)]$, 注意 $\mathbb{E}(\mathbb{E}(y|x)) = \mathbb{E}(y)$, 因此

$$\begin{aligned}\mathbb{E}[(h(x) - h^*(x))(h^*(x) - y)] &= \mathbb{E}[\mathbb{E}[(h(x) - h^*(x))(h^*(x) - y)|x]] \\ &= \mathbb{E}[(h(x) - h^*(x))\mathbb{E}[(h^*(x) - y)|x]]\end{aligned}$$

接着分析 $\mathbb{E}[(h^*(x) - y)|x]$, 注意到

$$\begin{aligned}\mathbb{E}[(h^*(x) - y)|x] &= \mathbb{E}(h^*(x)|x) - \mathbb{E}(y|x) \\ &= h^*(x) - h^*(x) \\ &= 0\end{aligned}$$

所以

$$\mathbb{E}[(h(x) - h^*(x))(h^*(x) - y)] = \mathbb{E}[(h(x) - h^*(x))\mathbb{E}[(h^*(x) - y)|x]] = \mathbb{E}[(h(x) - h^*(x)) \times 0] = 0$$

综上

$$\begin{aligned}E_{\text{out}}(h) &= \mathbb{E}[(h(x) - h^*(x))^2] + \mathbb{E}[(h^*(x) - y)^2] \geq \mathbb{E}[(h^*(x) - y)^2] \\ &\text{当且仅当 } h(x) = h^*(x) \text{ 时等号成立}\end{aligned}$$

接着证明另一个结论。首先 $y = y - h^*(x) + h^*(x) = h^*(x) + \epsilon(x)$, 所以只需计算 $\epsilon(x)$ 的数学期望, 注意 $\mathbb{E}(\mathbb{E}(y|x)) = \mathbb{E}(y)$, 因此

$$\begin{aligned}\mathbb{E}(\epsilon(x)) &= \mathbb{E}[\mathbb{E}(\epsilon(x)|x)] \\ &= \mathbb{E}[\mathbb{E}[(y - h^*(x))|x]] \\ &= \mathbb{E}[\mathbb{E}(y|x) - \mathbb{E}(h^*(x)|x)] \\ &= \mathbb{E}[h^*(x) - h^*(x)] \\ &= 0\end{aligned}$$

所以 $\epsilon(x)$ 满足条件, 因此结论成立。

Problem 3.9 (Page 112)

Assuming that $X^T X$ is invertible, show by direct comparison with Equation (3.4) that $E_{\text{in}}(w)$ can be written as

$$E_{\text{in}}(w) = (w - (X^T X)^{-1} X^T y)^T (X^T X) (w - (X^T X)^{-1} X^T y) + y^T (I - X(X^T X)^{-1} X^T) y$$

Use this expression for E_{in} to obtain w_{lin} . What is the in sample error? [Hint: The matrix $X^T X$ is positive definite.]

回顾等式3.4, 这里把 $\frac{1}{N}$ 这个常数略去

$$E_{\text{in}}(w) = w^T X^T X w - 2w^T X^T y + y^T y$$

令 $u = (X^T X)^{-1} X^T y, v = w - u$, 那么 $w = v + u$, 所以

$$\begin{aligned}E_{\text{in}}(w) &= w^T X^T X w - 2w^T X^T y + y^T y \\ &= (v + u)^T X^T X (v + u) - 2(v + u)^T X^T y + y^T y \\ &= v^T X^T X v + u^T X^T X v + v^T X^T X u + u^T X^T X u - 2v^T X^T y - 2u^T X^T y + y^T y \\ &= v^T X^T X v + 2v^T X^T X u - 2v^T X^T y + u^T X^T X u - 2u^T X^T y + y^T y\end{aligned}$$

先看下 $2v^T X^T Xu - 2v^T X^T y$, 将 $u = (X^T X)^{-1} X^T y$ 带入

$$\begin{aligned} 2v^T X^T Xu - 2v^T X^T y &= 2v^T (X^T Xu - X^T y) \\ &= 2v^T (X^T X (X^T X)^{-1} X^T y - X^T y) \\ &= 2v^T (X^T y - X^T y) \\ &= 0 \end{aligned}$$

再看下 $u^T X^T Xu - 2u^T X^T y$, 将 $u = (X^T X)^{-1} X^T y$ 带入

$$\begin{aligned} u^T X^T Xu - 2u^T X^T y &= y^T X (X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T y - 2y^T X (X^T X)^{-1} X^T y \\ &= y^T X (X^T X)^{-1} X^T y - 2y^T X (X^T X)^{-1} X^T y \\ &= -y^T X (X^T X)^{-1} X^T y \end{aligned}$$

所以

$$\begin{aligned} E_{\text{in}}(w) &= v^T X^T X v - y^T X (X^T X)^{-1} X^T y + y^T y \\ &= (w - (X^T X)^{-1} X^T y)^T (X^T X) (w - (X^T X)^{-1} X^T y) + y^T (I - X (X^T X)^{-1} X^T) y \end{aligned}$$

接着从这个式子来分析最优解, 因为 $X^T X$ 半正定, 所以

$$\begin{aligned} E_{\text{in}}(w) &\geq y^T (I - X (X^T X)^{-1} X^T) y \\ \text{当且仅当 } w - (X^T X)^{-1} X^T y &= 0 \text{ 时等号成立} \\ \text{即 } w_{\text{lin}} &= (X^T X)^{-1} X^T y \end{aligned}$$

可以看到和之前求导的结果一样, 显然求导要快很多。

Problem 3.10 (Page 112)

Exercise 3.3 studied some properties of the hat matrix $H = X(X^T X)^{-1} X^T$, where X is a N by $d + 1$ matrix, and $X^T X$ is invertible. Show the following additional properties.

(a) Every eigenvalue of H is either 0 or 1. [Hint: Exercise 3.3(b).]

(b) Show that the trace of a symmetric matrix equals the sum of its eigenvalues. [Hint: Use the spectral theorem and the cyclic property of the trace. Note that the same result holds for non-symmetric matrices, but is a little harder to prove.]

(c) How many eigenvalues of H are 1? What is the rank of H ? [Hint: Exercise 3.3(d).]

(a)由3.3(b)我们知道 $H^K = H$, 所以对于 H 的任意特征值 λ

$$\begin{aligned} \lambda^K &= \lambda \\ \lambda &= 0 \text{ 或 } 1 \end{aligned}$$

(b)直接对一般的矩阵证明结论, 利用标准型Jordan标准型的结论即可

任意方阵 A 可以相似于 J , $A = PJP^{-1}$, 其中 J 可以表示为如下形式

$$J = \text{diag}(J_1, J_2, \dots, J_k)$$

$$J_i = \begin{bmatrix} \lambda_i & 1 & 0 & \dots & 0 \\ 0 & \lambda_i & 1 & \dots & 0 \\ 0 & 0 & \lambda_i & \dots & 0 \\ \dots & \dots & \dots & \dots & 1 \\ 0 & 0 & 0 & \dots & \lambda_i \end{bmatrix} \in \mathbb{R}^{n_i \times n_i}$$

更具体的部分可以参考[维基百科](#)

下面仅从 trace 的角度利用这个结论

$$\begin{aligned} \text{trace}(A) &= \text{trace}(PJP^{-1}) \\ &= \text{trace}(P^{-1}PJ) \\ &= \text{trace}(J) \\ &= \text{trace}(\text{diag}(J_1, J_2, \dots, J_k)) \\ &= \sum_{i=1}^k \text{trace}(J_i) \end{aligned}$$

由特征值的性质我们知道 λ_i 为 A 的特征值, 而 $\sum_{i=1}^k \text{trace}(J_i)$ 即为特征值之和, 所以矩阵的 trace 等于特征值之和

(c) 由 3.3(d), 特征值之和 $\text{trace}(H) = d + 1 =$ 特征值之和, 因为特征值只能取 0 或 1, 所以特征值中一共有 $d + 1$ 个 1。

Problem 3.11 (Page 112)

Consider the linear regression problem setup in Exercise 3.4, where the data comes from a genuine linear relationship with added noise. The noise for the different data points is assumed to be iid with zero mean and variance σ^2 . Assume that the 2nd moment matrix $\Sigma = \mathbb{E}_x[xx^T]$ is non-singular. Follow the steps below to show that, with high probability, the out-of-sample error on average is

$$E_{\text{out}}(w_{\text{lin}}) = \sigma^2 \left(1 + \frac{d+1}{N} + O\left(\frac{1}{N}\right) \right)$$

(a) For a test point x , show that the error $y - g(x)$ is

$$\epsilon' - x^T (X^T X)^{-1} X^T \epsilon$$

where ϵ' is the noise realization for the test point and ϵ is the vector of noise realizations on the data.

(b) Take the expectation with respect to the test point, i.e., x and ϵ' , to obtain an expression for E_{out} . Show that

$$E_{\text{out}} = \sigma^2 + \text{trace}(\Sigma(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})$$

[Hints: $a = \text{trace}(a)$ for any scalar a ; $\text{trace}(AB) = \text{trace}(BA)$; expectation and trace commute.]

(c) What is $\mathbb{E}_\epsilon[\epsilon \epsilon^T]$?

(d) Take the expectation with respect to ϵ to show that, on average,

$$E_{\text{out}} = \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\Sigma(\frac{1}{N} X^T X)^{-1})$$

Note that $\frac{1}{N} X^T X = \frac{1}{N} \sum_{n=1}^N x_n x_n^T$ is an N sample estimate of Σ . So $\frac{1}{N} X^T X \approx \Sigma$. If $\frac{1}{N} X^T X = \Sigma$, then what is E_{out} on average?

(e) Show that (after taking the expectation over the data noise) with high probability,

$$E_{\text{out}} = \sigma^2(1 + \frac{d+1}{N} + O(\frac{1}{N}))$$

[Hint: By the law of large numbers: $\frac{1}{N} X^T X$ converges in probability to Σ , and so by continuity of the inverse at Σ , $(\frac{1}{N} X^T X)^{-1}$ converges in probability to Σ^{-1} .]

这题实际上是对Exercise 3.4的推广，注意这里 X 为训练数据， x 为测试数据，为了方便区分，我将 x 的噪音记录为 ϵ' ，原题为 ϵ 。

(a)同Exercise 3.4, 记 $y = [y_1 \dots y_N]^T$, $X = [x_1 \dots x_N]^T$, 注意题目中给出 $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^T$

那么

$$y = Xw^* + \epsilon$$

我们知道 $w_{\text{lin}} = (X^T X)^{-1} X^T y$

那么

$$\begin{aligned} g(x) &= x^T w_{\text{lin}} \\ &= x^T (X^T X)^{-1} X^T y \\ &= x^T (X^T X)^{-1} X^T (Xw^* + \epsilon) \\ &= x^T (X^T X)^{-1} X^T Xw^* + x^T (X^T X)^{-1} X^T \epsilon \\ &= x^T w^* + x^T (X^T X)^{-1} X^T \epsilon \end{aligned}$$

从而

$$\begin{aligned} y - g(x) &= x^T w^* + \epsilon' - (x^T w^* + x^T (X^T X)^{-1} X^T \epsilon) \\ &= \epsilon' - x^T (X^T X)^{-1} X^T \epsilon \end{aligned}$$

(b)利用定义计算即可，注意这题是关于 ϵ' , x 求期望

$$\begin{aligned}
E_{\text{out}} &= \mathbb{E}(\|y - g(x)\|^2) \\
&= \mathbb{E}(\|\epsilon' - x^T (X^T X)^{-1} X^T \epsilon\|^2) \\
&= \mathbb{E}((\epsilon' - x^T (X^T X)^{-1} X^T \epsilon)^T (\epsilon' - x^T (X^T X)^{-1} X^T \epsilon)) \\
&= \mathbb{E}[(\epsilon'^T X (X^T X)^{-1} x + \epsilon'^T) (\epsilon' - x^T (X^T X)^{-1} X^T \epsilon)] \\
&= \mathbb{E}[-\epsilon'^T X (X^T X)^{-1} x \epsilon' + \epsilon'^T \epsilon' + \epsilon'^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon - \epsilon'^T x^T (X^T X)^{-1} X^T \epsilon] \text{ (注意 } \epsilon' \sim \mathcal{N}(0, \sigma^2)) \\
&= -2\mathbb{E}[\epsilon'^T X (X^T X)^{-1} x \epsilon'] + \mathbb{E}((\epsilon')^2) + \text{trace}(\mathbb{E}(\epsilon'^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon)) \\
&= -2\mathbb{E}[\epsilon'^T X (X^T X)^{-1} x] \mathbb{E}(\epsilon') + \sigma^2 + \mathbb{E}[\text{trace}(\epsilon'^T X (X^T X)^{-1} x x^T (X^T X)^{-1} X^T \epsilon)] \text{ (注意 } \text{trace}(AB) = \text{trace}(BA), \mathbb{E}(\epsilon') = 0) \\
&= \sigma^2 + \mathbb{E}[\text{trace}(x x^T (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})] \\
&= \sigma^2 + \text{trace}(\mathbb{E}(x x^T (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})) \\
&= \sigma^2 + \text{trace}(\mathbb{E}(x x^T) (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}) \\
&= \sigma^2 + \text{trace}(\Sigma (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}) (\Sigma = \mathbb{E}_x[x x^T])
\end{aligned}$$

(c)直接计算即可，注意到

$$\mathbb{E}(\epsilon_i \epsilon_j) = \begin{cases} \sigma^2 & (i = j) \\ 0 & (i \neq j) \end{cases}$$

$$\begin{aligned}
\mathbb{E}_\epsilon[\epsilon \epsilon^T] &= \mathbb{E}_\epsilon[(\epsilon_1, \epsilon_2, \dots, \epsilon_N)(\epsilon_1, \epsilon_2, \dots, \epsilon_N)^T] \\
&= (\mathbb{E}[\epsilon_i \epsilon_j])_{ij} \\
&= \sigma^2 I
\end{aligned}$$

(d)利用c，对b计算的 E_{out} 关于 ϵ 取数学期望可得

$$\begin{aligned}
E'_{\text{out}} &= \mathbb{E}_\epsilon(E_{\text{out}}) \\
&= \mathbb{E}_\epsilon[\sigma^2 + \text{trace}(\Sigma (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})] \\
&= \sigma^2 + \text{trace}[\mathbb{E}_\epsilon(\Sigma (X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})] \\
&= \sigma^2 + \text{trace}[\Sigma (X^T X)^{-1} X^T \mathbb{E}_\epsilon(\epsilon \epsilon^T) X (X^T X)^{-1}] \\
&= \sigma^2 + \sigma^2 \text{trace}[\Sigma (X^T X)^{-1} X^T X (X^T X)^{-1}] \\
&= \sigma^2 + \sigma^2 \text{trace}[\Sigma (X^T X)^{-1}] \\
&= \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\Sigma (\frac{1}{N} X^T X)^{-1})
\end{aligned}$$

由计算我们知道 $\frac{1}{N} X^T X = \frac{1}{N} \sum_{n=1}^N x_n x_n^T$ 是 Σ 的极大似然估计，因此 $\frac{1}{N} X^T X \approx \Sigma$ 。如果 $\frac{1}{N} X^T X = \Sigma$

$$E'_{\text{out}} = \sigma^2 + \frac{\sigma^2}{N} \text{trace}(I_{d+1}) = \sigma^2(1 + \frac{d+1}{N})$$

(e)由大数定律我们知道 $\frac{1}{N} X^T X$ 依概率收敛于 Σ ，所以 $(\frac{1}{N} X^T X)^{-1}$ 依概率收敛于 Σ^{-1} ，从而有很高的概率

$$\Sigma(\frac{1}{N} X^T X)^{-1} = \Sigma(\Sigma^{-1} + S) = I_{d+1} + \Sigma S$$

其中 S 为一个矩阵，那么有很高的概率

$$\text{trace}(\Sigma(\frac{1}{N} X^T X)^{-1}) = \sigma^2(d+1 + O(1))$$

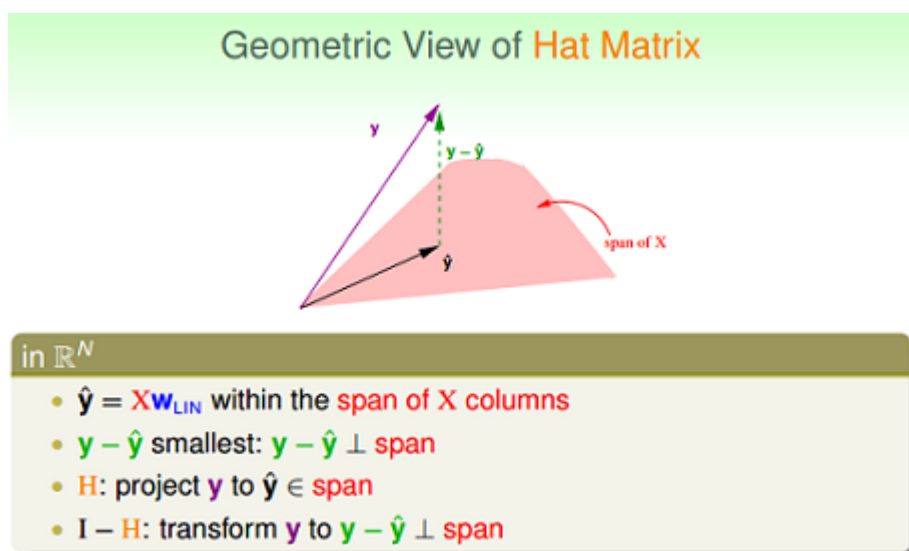
从而有很高的概率

$$E_{\text{out}} = \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\Sigma(\frac{1}{N} X^T X)^{-1}) = \sigma^2(1 + \frac{d+1}{N} + O(\frac{1}{N}))$$

Problem 3.12 (Page 113)

In linear regression, the in sample predictions are given by $\hat{y} = Hy$, where $H = X(X^T X)^{-1} X^T$. Show that H is a projection matrix, i.e. $H^2 = H$. So \hat{y} is the projection of y onto some space. What is this space?

回忆Exercies 3.3可知, $H^K = H$ 。关于投影可以参考林老师的课件



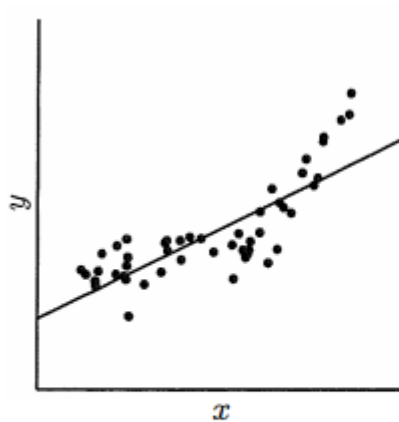
这里简单说明这几个结论。 $\hat{y} = Hy = X(X^T X)^{-1} X^T y = Xw_{\text{lin}}$, 由线性代数知识我们知道 \hat{y} 属于 X 的列张成的子空间。接着考虑 $y - \hat{y}$, 注意 H 为对称矩阵

$$\hat{y}^T (y - \hat{y}) = y^T H^T (I - H) y = y^T (H^T - H^T H) y = y^T (H - H^2) y = 0$$

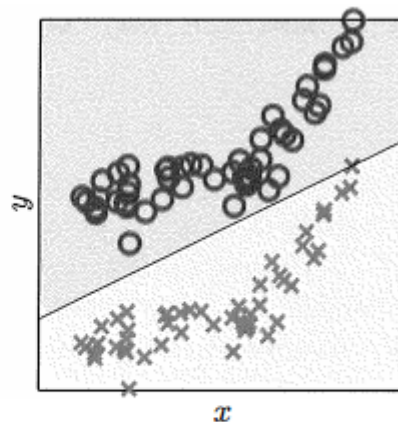
所以 $y - \hat{y}$ 垂直于 \hat{y} 。

Problem 3.13 (Page 113)

This problem creates a linear regression algorithm from a good algorithm for linear classification. As illustrated, the idea is to take the original data and shift it in one direction to get the +1 data points; then, shift it in the opposite direction to get the -1 data points.



Original data for the one dimensional regression problem



Shifted data viewed as a two dimensional classification problem

More generally, The data (x_n, y_n) can be viewed as data points in \mathbb{R}^{d+1} by treating the y value as the $(d + 1)$ th coordinate.

Now, construct positive and negative points

$$\begin{aligned}\mathcal{D}_+ &= (x_1, y_1) + a, \dots, (x_N, y_N) + a \\ \mathcal{D}_- &= (x_1, y_1) - a, \dots, (x_N, y_N) - a\end{aligned}$$

where a is a perturbation parameter. You can now use the linear programming algorithm in Problem 3.6 to separate \mathcal{D}_+ from \mathcal{D}_- . The resulting separating hyperplane can be used as the regression 'fit' to the original data.

(a) How many weights are learned in the classification problem? How many weights are needed for the linear fit in the regression problem?

(b) The linear fit requires weights w , where $h(x) = w^T x$. Suppose the weights returned by solving the classification problem are w_{class} . Derive an expression for w as a function of w_{class} .

(c) Generate a data set $y_n = x_n^2 + \sigma \epsilon_n$ with $N = 50$, where x_n is uniform on $[0, 1]$ and ϵ_n is zero mean Gaussian noise; set $\sigma = 0.1$. Plot \mathcal{D}_+ and \mathcal{D}_- for $a = [0, 0.1]^T$.

(d) Give comparisons of the resulting fits from running the classification approach and the analytic pseudo-inverse algorithm for linear regression.

题目的意思是对于回归问题的点 (x_n, y_n) , 有个一个偏移量 a , 构造两个点集,

$$\begin{aligned}\mathcal{D}_+ &= (x_1, y_1) + a, \dots, (x_N, y_N) + a \\ \mathcal{D}_- &= (x_1, y_1) - a, \dots, (x_N, y_N) - a\end{aligned}$$

我们对于这两个点集作分类问题, 利用分类问题得到的参数来做回归。

(a)由题设知 $x_n \in \mathbb{R}^d$, 所以 $(x_n, y_n) = (x_n^1 \dots x_n^d, y_n) \in \mathbb{R}^{d+1}$, 注意学习的时候还要加一个1分量, 数据变为 $(1, x_n, y_n) = (1, x_n^1 \dots x_n^d, y_n) \in \mathbb{R}^{d+2}$, 从而对于分类问题我们需要学习 $d + 2$ 个权重 $w = (w_0, \dots, w_d, w_{d+1})$ 。

计算完 w 之后, 我们要回到原来的回归问题, 注意此时分类边界为

$$w_0 + w_1 x_n^1 + \dots + w_d x_n^d + w_{d+1} y_n = 0$$

$$y_n = -\frac{w_0}{w_{d+1}} - \frac{w_1}{w_{d+1}} x_n^1 - \dots - \frac{w_d}{w_{d+1}} x_n^d$$

所以我们的回归直线为

$$w' = -\left(\frac{w_0}{w_{d+1}}, \dots, \frac{w_d}{w_{d+1}}\right) \in \mathbb{R}^{d+1}$$

(b)由(a)我们知道

$$w = -\left(\frac{w_{class}^0}{w_{class}^{d+1}}, \dots, \frac{w_{class}^d}{w_{class}^{d+1}}\right)$$

(c)编程处理，首先作图

```
# -*- coding: utf-8 -*-
"""
Created on Wed Mar  6 13:57:11 2019

@author: qinzhen
"""

import numpy as np
import matplotlib.pyplot as plt
from numpy.linalg import inv
from cvxopt import matrix, solvers

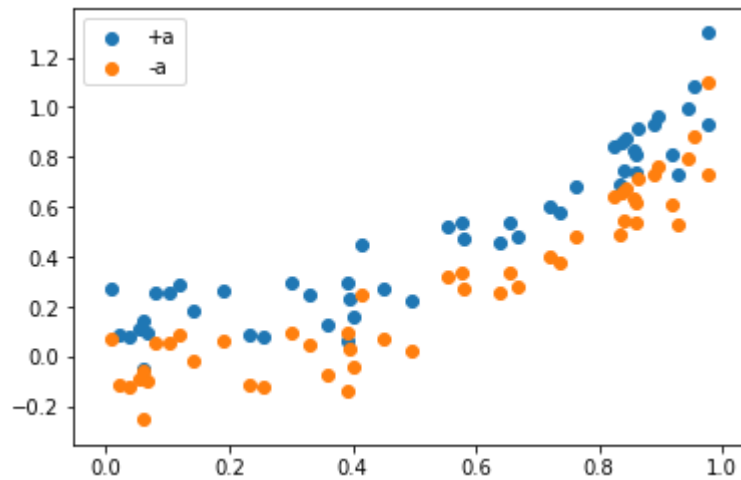
#(c)
def generate(n, delta):
    x = np.random.uniform(size=n)
    epsilon = np.random.normal(size=n)
    y = x * x + delta * epsilon
    data = np.c_[X, y]
    return X, y, data

#参数
n = 50
delta = 0.1

#生成数据
X, y, data = generate(n, delta)

#构造D1,D2
a = np.array([0, 0.1])
D1 = data + a
D2 = data - a

plt.scatter(D1[:, 0], D1[:, 1], label='+a')
plt.scatter(D2[:, 0], D2[:, 1], label='-a')
plt.legend()
plt.show()
```

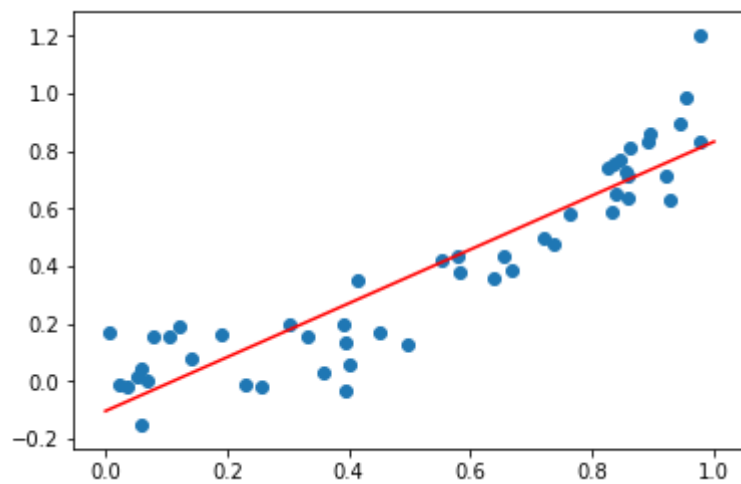



(d)分别使用线性回归的公式以及这题介绍的方法求解并进行比较。先看下线性回归的效果。

```
#(d)
#线性回归
#添加偏置项
X_treat = np.c_[np.ones(n), x]
w = inv(X_treat.T.dot(X_treat)).dot(X_treat.T).dot(y)

a1 = np.array([0, 1])
b1 = w[0] + w[1] * a1

plt.scatter(x, y)
plt.plot(a1, b1, 'r')
plt.show()
```



接着使用此题以及Problem 3.6介绍的优化方法求解该分类问题。

```
#使用此题介绍的分类方法
#Problem 3.5的算法2
def algorithm2(x, y):
    """
    算法2
    """
```

```

N, d = X.shape
A1 = X * y.reshape(-1, 1)
A2 = np.c_[A1, (-1) * np.eye(N)]
A3 = np.c_[np.zeros((N, d)), (-1) * np.eye(N)]

```

```

A = np.r_[A2, A3]
c = np.array([0.0] * d + [1.0] * N)
b = np.array([-1.0] * N + [0.0] * N)

```

```

#带入算法求解
c = matrix(c)
A = matrix(A)
b = matrix(b)

```

```
sol = solvers.lp(c, A, b)
```

```

#返回向量
w = np.array((sol['x']))[:d]

```

```
return w
```

```
#构造数据
```

```

X1 = np.r_[D1, D2]
X1 = np.c_[np.ones(2 * n), X1]
y1 = np.r_[np.ones(n), -1 * np.ones(n)]

```

```
#算法2
```

```
w2 = algorithm2(X1, y1)
```

```
#处理后的w
```

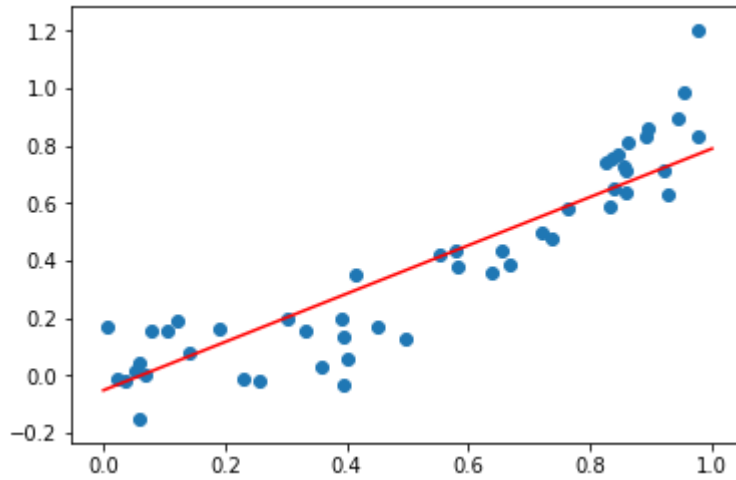
```

w_2 = - w2[:-1] / w2[-1]
a2 = np.array([0, 1])
b2 = w_2[0] + w_2[1] * a2
plt.scatter(X, y)
plt.plot(a2, b2, 'r')
plt.show()

```

	pcost	dcost	gap	pres	dres	k/t
0:	3.2035e+01	1.7006e+02	8e+02	3e+00	3e+00	1e+00
1:	5.1358e+01	7.5676e+01	7e+01	5e-01	5e-01	2e+00
2:	5.2417e+01	5.6617e+01	9e+00	8e-02	9e-02	3e-01
3:	5.2569e+01	5.4479e+01	4e+00	4e-02	4e-02	1e-01
4:	5.2600e+01	5.3116e+01	1e+00	1e-02	1e-02	3e-02
5:	5.2656e+01	5.2781e+01	3e-01	2e-03	3e-03	6e-03
6:	5.2676e+01	5.2692e+01	4e-02	3e-04	3e-04	7e-04
7:	5.2679e+01	5.2680e+01	2e-03	2e-05	2e-05	1e-05
8:	5.2679e+01	5.2679e+01	2e-05	2e-07	2e-07	2e-07
9:	5.2679e+01	5.2679e+01	2e-07	2e-09	2e-09	2e-09

Optimal solution found.



Problem 3.14 (Page 114)

In a regression setting, assume the target function is linear, so $f(x) = x^T w_f$, and $y = Xw_f + \epsilon$, where the entries in ϵ are zero mean, iid with variance σ^2 . In this problem derive the bias and variance as follows.

(a) Show that the average function is $\bar{g}(x) = f(x)$, no matter what the size of the data set, as long as $X^T X$ is invertible. What is the bias?

(b) What is the variance? [Hint: Problem 3.11]

这里 $\bar{g}(x)$ 的定义要参考课本63页，对于这题来说，实际上就是求 $\mathbb{E}(w_f)$ ，回顾之前的结论可知 $w_f = (X^T X)^{-1} X^T y$ 。

(a)注意 $(X^T X)^{-1}$ 可逆， ϵ 的数学期望为0

$$\begin{aligned}\mathbb{E}(w_f) &= \mathbb{E}((X^T X)^{-1} X^T y) \\ &= \mathbb{E}((X^T X)^{-1} X^T (Xw_f + \epsilon)) \\ &= \mathbb{E}((X^T X)^{-1} X^T Xw_f) + (X^T X)^{-1} X^T \mathbb{E}(\epsilon) \\ &= w_f\end{aligned}$$

所以 $\bar{g}(x) = x^T \mathbb{E}(w_f) = x^T w_f = f(x)$

$$\begin{aligned}\text{bias}(x) &= (\bar{g}(x) - x^T w_f - \epsilon)^2 = \epsilon^2 \\ \text{bias} &= \mathbb{E}[\text{bias}(x)] = \mathbb{E}[\epsilon^2] = \sigma^2\end{aligned}$$

(b)注意

$$E_{\text{out}} = \text{bias} + \text{var}$$

由Problem 3.11可得

$$E_{\text{out}} = \sigma^2 + \text{trace}(\Sigma(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})$$

因此

$$\text{var} = E_{\text{out}} - \text{bias} = \text{trace}(\Sigma(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1})$$

Problem 3.15 (Page 114)

In the text we derived that the linear regression solution weights must satisfy $X^T X w = X^T y$. If $X^T X$ is not invertible, the solution $w_{\text{lin}} = (X^T X)^{-1} X^T y$ won't work. In this event, there will be many solutions for w that minimize E_{in} . Here, you will derive one such solution. Let ρ be the rank of X . Assume that the singular value decomposition (SVD) of X is $X = U \Gamma V^T$ where $U \in \mathbb{R}^{N \times \rho}$ satisfies $U^T U = I_\rho$, $V \in \mathbb{R}^{(d+1) \times \rho}$ satisfies $V^T V = I_\rho$, and $\Gamma \in \mathbb{R}^{\rho \times \rho}$ is a positive diagonal matrix.

(a) Show that $\rho < d + 1$.

(b) Show that $w_{\text{lin}} = V \Gamma^{-1} U^T y$ satisfies $X^T X w_{\text{lin}} = X^T y$, and hence is a solution.

(c) Show that for any other solution that satisfies $X^T X w = X^T y$, $\|w_{\text{lin}}\| < \|w\|$. That is, the solution we have constructed is the minimum norm set of weights that minimizes E_{in} .

(a)由题设我们知道 $X^T X$ 不可逆, 又因为 $X^T X \in \mathbb{R}^{(d+1) \times (d+1)}$, 所以

$$r(X^T X) < d + 1$$

注意线性代数中的恒等式 $r(X) = r(X^T X)$, 因此

$$\rho = r(X) = r(X^T X) < d + 1$$

(b)注意 Γ 为对角阵, 所以 $\Gamma^T = \Gamma$

$$\begin{aligned} X^T X w_{\text{lin}} &= V \Gamma^T U^T U \Gamma V^T V \Gamma^{-1} U^T y \\ &= V \Gamma U^T y \\ &= X^T y \end{aligned}$$

(c)先对 $X^T X w = X^T y$ 进行处理, 将 $X = U \Gamma V^T$ 带入

$$\begin{aligned} X^T X w &= V \Gamma^T U^T U \Gamma V^T w = V \Gamma^2 V^T w \\ X^T y &= V \Gamma^T U^T y = V \Gamma U^T y \end{aligned}$$

从而 $X^T X w = X^T y$ 可以化为

$$V \Gamma^2 V^T w = V \Gamma U^T y$$

两边同乘 V^T , 注意 $V^T V = I_\rho$ 以及 Γ 为正定阵, 从而可逆, 所以

$$\begin{aligned} V^T V \Gamma^2 V^T w &= V^T V \Gamma U^T y \\ \Gamma^2 V^T w &= \Gamma U^T y \\ V^T w &= \Gamma^{-1} U^T y \end{aligned}$$

接着我们将 w , 分解为两项

$$w = w - w_{\text{lin}} + w_{\text{lin}}$$

下面证明 $w_{\text{lin}}^T (w - w_{\text{lin}}) = 0$

$$\begin{aligned} w_{\text{lin}}^T (w - w_{\text{lin}}) &= (V\Gamma^{-1}U^T y)^T (w - w_{\text{lin}}) \\ &= y^T U\Gamma^{-1} (V^T w - V^T w_{\text{lin}}) \end{aligned}$$

由我们刚刚推出的等式可知，对于每个满足 $X^T X w = X^T y$ 的 w 均满足如下等式

$$V^T w = \Gamma^{-1} U^T y$$

所以

$$w_{\text{lin}}^T (w - w_{\text{lin}}) = y^T U\Gamma^{-1} (V^T w - V^T w_{\text{lin}}) = y^T U\Gamma^{-1} (\Gamma^{-1} U^T y - \Gamma^{-1} U^T y) = 0$$

因此

$$\begin{aligned} \|w\|^2 &= \|w - w_{\text{lin}} + w_{\text{lin}}\|^2 \\ &= \|w - w_{\text{lin}}\|^2 + \|w_{\text{lin}}\|^2 + 2w_{\text{lin}}^T (w - w_{\text{lin}}) \\ &= \|w - w_{\text{lin}}\|^2 + \|w_{\text{lin}}\|^2 \\ &\geq \|w_{\text{lin}}\|^2 \end{aligned}$$

当且仅当 $w = w_{\text{lin}}$ 时等号成立

Problem 3.16 (Page 115)

In Example 3.4, it is mentioned that the output of the final hypothesis $g(x)$ learned using logistic regression can be thresholded to get a 'hard' (± 1) classification. This problem shows how to use the risk matrix introduced in Example 1.1 to obtain such a threshold.

Consider fingerprint verification, as in Example 1.1. After learning from the data using logistic regression, you produce the final hypothesis

$$g(x) = \mathbb{P}[y = +1|x]$$

which is your estimate of the probability that $y = +1$. Suppose that the cost matrix is given by

		True classification	
		+1 (correct person)	-1 (intruder)
you say	+1	0	c_a
	-1	c_r	0

For a new person with fingerprint x , you compute $g(x)$ and you now need to decide whether to accept or reject the person (i.e., you need a hard classification). So, you will accept if $g(x) \geq \kappa$, where κ is the threshold.

(a) Define the cost(accept) as your expected cost if you accept the person. Similarly define cost(reject). Show that

$$\begin{aligned} \text{cost(accept)} &= (1 - g(x))c_a \\ \text{cost(reject)} &= g(x)c_r \end{aligned}$$

(b) Use part (a) to derive a condition on $g(x)$ for accepting the person and hence show that

$$\kappa = \frac{c_a}{c_a + c_r}$$

(c) Use the cost matrices for the Supermarket and CIA applications in Example 1.1 to compute the threshold κ for each of these two cases. Give some intuition for the thresholds you get.

(a) 如果accept, 那么

$$\text{cost}(\text{accept}) = \mathbb{P}[y = +1|x] \times 0 + \mathbb{P}[y = -1|x] \times c_a = (1 - g(x))c_a$$

如果reject, 那么

$$\text{cost}(\text{reject}) = \mathbb{P}[y = +1|x] \times c_r + \mathbb{P}[y = -1|x] \times 0 = g(x)c_r$$

(b) 令只有当 $\text{cost}(\text{accept}) \leq \text{cost}(\text{reject})$ 时才应该接受, 解这个不等式可得

$$\begin{aligned} (1 - g(x))c_a &\leq g(x)c_r \\ c_a &\leq g(x)(c_a + c_r) \\ \frac{c_a}{c_a + c_r} &\leq g(x) \end{aligned}$$

所以当 $\frac{c_a}{c_a + c_r} \leq g(x)$ 时接受, $\frac{c_a}{c_a + c_r} > g(x)$ 时拒绝, 对照题目可知

$$\kappa = \frac{c_a}{c_a + c_r}$$

(c) 回顾课本上有关超市和CIA的图

		f	
		+1	-1
h	+1	0	1
	-1	10	0

Supermarket

		f	
		+1	-1
h	+1	0	1000
	-1	1	0

CIA

那么对于超市, $\kappa = \frac{1}{1+10} = \frac{1}{11}$, 对于CIA, $\kappa = \frac{1000}{1000+1} = \frac{1000}{1001}$, 所以只有当 $g(x)$ 很大时才能被CIA接受, 符合之前的讨论。

Problem 3.17 (Page 115)

Consider a function

$$E(u, v) = e^u + e^{2v} + e^{uv} + u^2 - 3uv + 4v^2 - 3u - 5v$$

(a) Approximate $E(u + \Delta u, v + \Delta v)$ by $\hat{E}_1(\Delta u, \Delta v)$, where \hat{E}_1 is the first-order Taylor's expansion of E around $(u, v) = (0, 0)$. Suppose $\hat{E}_1(\Delta u, \Delta v) = a_u \Delta u + a_v \Delta v + a$. What are the values of a_u, a_v , and a ?

(b) Minimize \hat{E}_1 over all possible $(\Delta u, \Delta v)$ such that $\|(\Delta u, \Delta v)\| = 0.5$. In this chapter, we proved that the optimal column vector $\begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix}$ is parallel to the column vector $-\nabla E(u, v)$, which is called the negative gradient direction. Compute the optimal $(\Delta u, \Delta v)$ and the resulting $E(u + \Delta u, v + \Delta v)$.

(c) Approximate $E(u + \Delta u, v + \Delta v)$ by $\hat{E}_2(\Delta u, \Delta v)$, where E_2 is the second order Taylor's expansion of E around $(u, v) = (0, 0)$. Suppose

$$\hat{E}_2(\Delta u, \Delta v) = b_{uu}(\Delta u)^2 + b_{vv}(\Delta v)^2 + b_{uv}(\Delta u)(\Delta v) + b_u \Delta u + b_v \Delta v + b$$

What are the values of $b_{uu}, b_{vv}, b_{uv}, b_u, b_v$, and b ?

(d) Minimize \hat{E}_2 over all possible $(\Delta u, \Delta v)$ (regardless of length). Use the fact that $\nabla^2 E(u, v)|_{(0,0)}$ (the Hessian matrix at $(0, 0)$) is positive definite to prove that the optimal column vector

$$\begin{bmatrix} \Delta u^* \\ \Delta v^* \end{bmatrix} = -(\nabla^2 E(u, v))^{-1} \nabla E(u, v)$$

which is called the Newton direction.

(e) Numerically compute the following values:

(i) the vector $(\Delta u, \Delta v)$ of length 0.5 along the Newton direction, and the resulting $E(u + \Delta u, v + \Delta v)$.

(ii) the vector $(\Delta u, \Delta v)$ of length 0.5 that minimizes $E(u + \Delta u, v + \Delta v)$, and the resulting $E(u + \Delta u, v + \Delta v)$. (Hint: Let $\Delta u = 0.5 \sin \theta$.)

Compare the values of $E(u + \Delta u, v + \Delta v)$ in (b), (e i), and (e ii). Briefly state your findings.

The negative gradient direction and the Newton direction are quite fundamental for designing optimization algorithms. It is important to understand these directions and put them in your toolbox for designing learning algorithms.

(a) 进行一阶泰勒展开，由公式可知

$$a_u = \frac{\partial E(u, v)}{\partial u} \Big|_{(u,v)=(0,0)} = e^u + v e^{uv} + 2u - 3v - 3 \Big|_{(u,v)=(0,0)} = -2$$

$$a_v = \frac{\partial E(u, v)}{\partial v} \Big|_{(u,v)=(0,0)} = 2e^v + u e^{uv} - 3u + 8v - 5 \Big|_{(u,v)=(0,0)} = -3$$

$$a = E(u, v) \Big|_{(u,v)=(0,0)} = 3$$

(b) 因为和负梯度同向时 \hat{E}_1 最小，所以

$$\begin{aligned}
\begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix} &= -k \nabla E(u, v)|_{(u,v)=(0,0)} \\
&= -k \begin{bmatrix} \frac{\partial E(u,v)}{\partial u}|_{(u,v)=(0,0)} \\ \frac{\partial E(u,v)}{\partial v}|_{(u,v)=(0,0)} \end{bmatrix} \\
&= -k \begin{bmatrix} -2 \\ -3 \end{bmatrix} \\
&= k \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad (k > 0)
\end{aligned}$$

因为 $\|(\Delta u, \Delta v)\| = 0.5$, 所以

$$\begin{aligned}
\Delta u &= \frac{2}{\sqrt{13}} \times 0.5, \Delta v = \frac{3}{\sqrt{13}} \times 0.5 \\
E(u + \Delta u, v + \Delta v) &= E\left(\frac{1}{\sqrt{13}}, \frac{3}{2\sqrt{13}}\right) = 2.25085973499
\end{aligned}$$

代码如下

```
# -*- coding: utf-8 -*-
"""
Created on Wed Mar  6 14:31:30 2019

@author: qinzhen
"""

import numpy as np
from numpy.linalg import inv

#(b)
def E(u,v):
    return np.exp(u) + np.exp(2 * v) + np.exp(u * v) + u * u - 3 * u * v + 4 * v * v - 3 * u
    - 5 * v

u = 1 / np.sqrt(13)
v = 3 / (2 * np.sqrt(13))
print(E(u,v))
```

2.25085973499

(c)二阶泰勒公式

$$\begin{aligned}
b_{uu} &= \frac{1}{2} \frac{\partial^2 E(u, v)}{\partial u^2} \Big|_{(u,v)=(0,0)} = \frac{1}{2} \frac{\partial}{\partial u} \frac{\partial E(u, v)}{\partial u} \Big|_{(u,v)=(0,0)} = \frac{3}{2} \\
b_{vv} &= \frac{1}{2} \frac{\partial^2 E(u, v)}{\partial v^2} \Big|_{(u,v)=(0,0)} = \frac{1}{2} \frac{\partial}{\partial v} \frac{\partial E(u, v)}{\partial v} \Big|_{(u,v)=(0,0)} = 5 \\
b_{uv} &= \frac{\partial^2 E(u, v)}{\partial u \partial v} \Big|_{(u,v)=(0,0)} = \frac{\partial}{\partial v} \frac{\partial E(u, v)}{\partial u} \Big|_{(u,v)=(0,0)} = -2 \\
b_u &= a_u = -2 \\
b_v &= a_v = -3 \\
b &= a = 3
\end{aligned}$$

(d)先判断正定性

$$\nabla^2 E(u, v) = \begin{bmatrix} \frac{\partial^2 E(u, v)}{\partial u^2} & \frac{\partial^2 E(u, v)}{\partial u \partial v} \\ \frac{\partial^2 E(u, v)}{\partial u \partial v} & \frac{\partial^2 E(u, v)}{\partial v^2} \end{bmatrix} = \begin{bmatrix} 3 & -2 \\ -2 & 10 \end{bmatrix}$$

显然 $\nabla^2 E(u, v)$ 正定, 接着带入牛顿公式求解

$$\begin{bmatrix} \Delta u^* \\ \Delta v^* \end{bmatrix} = -(\nabla^2 E(u, v))^{-1} \nabla E(u, v) = - \begin{bmatrix} b_{uu} & b_{uv} \\ b_{uv} & b_{vv} \end{bmatrix}^{-1} \begin{bmatrix} \frac{\partial E(u, v)}{\partial u} \\ \frac{\partial E(u, v)}{\partial v} \end{bmatrix} \Big|_{(u,v)=(0,0)} = - \begin{bmatrix} 3 & -2 \\ -2 & 10 \end{bmatrix}^{-1} \begin{bmatrix} -2 \\ -3 \end{bmatrix}$$

带入公式求解,

$$\Delta u = 0.4472136, \Delta v = 0.2236068, E(u + \Delta u, v + \Delta v) = 1.87339277$$

代码如下

```
#(c)
m1 = np.array([[3, -2], [-2, 10]])
m2 = np.array([[2], [3]])
d = inv(m1).dot(m2)

l = np.sqrt((d * d).sum())
d1 = 0.5 * d / l

u = d1[0]
v = d1[1]

print(E(u, v))
print(d1)
```

```
[ 1.87339277]
[[ 0.4472136]
 [ 0.2236068]]
```

(e)将之前计算出来的系数带入 $\hat{E}_2(\Delta u, \Delta v) = b_{uu}(\Delta u)^2 + b_{vv}(\Delta v)^2 + b_{uv}(\Delta u)(\Delta v) + b_u \Delta u + b_v \Delta v + b$, 令 $\Delta u = t, \Delta v = s$

$$\begin{aligned}
\hat{E}_2(\Delta u, \Delta v) &= b_{uu}(\Delta u)^2 + b_{vv}(\Delta v)^2 + b_{uv}(\Delta u)(\Delta v) + b_u \Delta u + b_v \Delta v + b \\
&= 1.5(t^2) + 5(s^2) - 2st - 2t - 3s - 2 \\
&= 1.5(t - a)^2 + 5(s - b)^2 - 2(t - a)(s - b) + C
\end{aligned}$$

其中 a, b, C 均为常数, 后续会求解出来, 令 $t_1 = t - a, s_1 = s - b$

$$\begin{aligned}
\hat{E}_2(\Delta u, \Delta v) &= 1.5(t - a)^2 + 5(s - b)^2 - 2(t - a)(s - b) + C \\
&= 1.5t_1^2 + 5s_1^2 - 2t_1s_1 + C \\
&= \frac{3}{2}(t_1 - \frac{2}{3}s_1)^2 + (5 - \frac{2}{3})s_1^2 + C
\end{aligned}$$

题目中要使得 $E_2(u, v) = E_2(0, 0) + \hat{E}_2(\Delta u, \Delta v)$ 最小, 所以求 $\hat{E}_2(\Delta u, \Delta v)$ 的最小值即可。

由上式, 当 $s_1 = 0, t_1 - \frac{2}{3}s_1 = 0$ 时, 即 $s_1 = t_1 = 0$ 时 $\hat{E}_2(\Delta u, \Delta v)$ 最小, 注意 t_1, s_1 的定义可得此时

$$t = a, s = b$$

而 $\Delta u = t, \Delta v = s$, 所以等号成立的条件为

$$\Delta u = a, \Delta v = b$$

接下来求解 a, b

$$\begin{aligned}
1.5(t - a)^2 + 5(s - b)^2 - 2(t - a)(s - b) + C &= 1.5(t^2 - 2at + a^2) + 5(s^2 - 2sb + b^2) - 2(ts - as - bt + ab) + C \\
&= 1.5t^2 + 5s^2 - 2st - (3a - 2b)t - (10b - 2a)s + 1.5a^2 + 5b^2 - 2ab + C \\
&= 1.5t^2 + 5s^2 - 2st - 2t - 3s - 2
\end{aligned}$$

那么

$$\begin{aligned}
&\begin{cases} 3a - 2b = 2 \\ -2a + 10b = 3 \end{cases} \\
&\begin{bmatrix} 3 & -2 \\ -2 & 10 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \\
&\begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 3 & -2 \\ -2 & 10 \end{bmatrix}^{-1} \begin{bmatrix} 2 \\ 3 \end{bmatrix}
\end{aligned}$$

所以这种方法可以得出牛顿方法同样的解, 这也验证了牛顿方法的正确性, 计算结果同(d)。

Problem 3.18 (Page 116)

Take the feature transform ϕ_2 in Equation (3.13) as ϕ .

(a) Show that $d_{vc}(\mathcal{H}_\phi) \leq 6$

(b) Show that $d_{vc}(\mathcal{H}_\phi) > 4$. [Hint: Exercise 3.12]

(c) Give an upper bound on $d_{vc}(\mathcal{H}_{\phi_k})$ for $\mathcal{X} = \mathbb{R}^d$.

(d) Define

$$\tilde{\phi}_2 : x \rightarrow (1, x_1, x_2, x_1 + x_2, x_1 - x_2, x_1^2, x_1 x_2, x_2 x_1, x_2^2) \text{ for } x \in \mathbb{R}^2$$

Argue that $d_{vc}(\mathcal{H}_{\tilde{\phi}_2}) = d_{vc}(\mathcal{H}_{\phi_2})$. In other words, while $\tilde{\phi}_2(\mathcal{X}) \in \mathbb{R}^9$, $d_{vc}(\mathcal{H}_{\tilde{\phi}_2}) \leq 6 < 9$. Thus, the dimension of $\phi(\mathcal{X})$ only gives an upper bound of $d_{vc}(\mathcal{H}_{\phi})$, and the exact value of $d_{vc}(\mathcal{H}_{\phi})$ can depend on the components of the transform.

回顾103页的 ϕ_2

$$\phi_2(x) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$$

(a)可以将 $\phi_2(x) \in \mathbb{R}^5$ 看成 \mathbb{R}^5 上的感知机，所以

$$d_{vc}(\mathcal{H}_{\phi_2}) \leq 6$$

(b)对于每种分类，实际上我们是在找到满足如下条件的 w ,

$$\text{sign}(w^T x^{(i)}) = y^{(i)} (i = 1, 2 \dots N)$$

其中 $(x^{(1)}, \dots, x^{(N)})$ 为输入数据, $(y^{(1)}, \dots, y^{(N)})$ 为对应的分类($y^{(i)} \in \{1, -1\}$), 对于此题, 我们取 $N = 5$, 且求解一个更强的条件

$$w^T x^{(i)} = y^{(i)} (i = 1, 2 \dots 5)$$

结合 $\phi_2(x) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$ 可得

$$w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)} + w_3 (x_1^{(i)})^2 + w_4 x_1^{(i)} x_2^{(i)} + w_5 (x_2^{(i)})^2 = y^{(i)} (i = 1, 2 \dots 5)$$

这是关于 $w_j (j = 0, 1, \dots, 5)$ 的六元一次方程组, 且方程组有五个, 所以必然有解。从而对于5个点, 任意一种分类均可以表示出来, 所以

$$d_{vc}(\mathcal{H}_{\phi}) \geq 5 > 4$$

(c)可以将 \mathcal{H}_{ϕ_k} 看成 \mathbb{R}^d 上的感知机, 所以

$$d_{vc}(\mathcal{H}_{\phi_k}) \leq d + 1$$

(d)我们每一种分类实际上对应了一个 w , 使得

$$\text{sign}(w^T x^{(i)}) = y^{(i)} (i = 1, 2 \dots N)$$

其中 $(x^{(1)}, \dots, x^{(N)})$ 为输入数据, $(y^{(1)}, \dots, y^{(N)})$ 为对应的分类($y^{(i)} \in \{1, -1\}$)。所以如果我们能证明 ϕ_2 对应的 w 与 $\tilde{\phi}_2$ 对应的 \tilde{w} 可以形成一一对应关系, 那么即可证明结论, 因为两种特征转化下的分类可以一一对应。

先证明 ϕ_2 对应的 w 与 $\tilde{\phi}_2$ 对应的 \tilde{w} 可以形成一一对应关系

$$\begin{aligned} \tilde{w}^T \tilde{\phi}_2 &= \tilde{w}_0 + \tilde{w}_1 x_1 + \tilde{w}_2 x_2 + \tilde{w}_3 (x_1 + x_2) + \tilde{w}_4 (x_1 - x_2) + \tilde{w}_5 (x_1^2) + \tilde{w}_6 (x_1 x_2) + \tilde{w}_7 (x_2 x_1) + \tilde{w}_8 (x_2^2) \\ &= \tilde{w}_0 + (\tilde{w}_1 + \tilde{w}_3 + \tilde{w}_4) x_1 + (\tilde{w}_2 + \tilde{w}_3 - \tilde{w}_4) x_2 + \tilde{w}_5 (x_1^2) + (\tilde{w}_6 + \tilde{w}_7) (x_1 x_2) + \tilde{w}_8 (x_2^2) \end{aligned}$$

那么 \tilde{w} 可以对应为 $(\tilde{w}_0, \tilde{w}_1 + \tilde{w}_3 + \tilde{w}_4, \tilde{w}_2 + \tilde{w}_3 - \tilde{w}_4, \tilde{w}_5, \tilde{w}_6 + \tilde{w}_7, \tilde{w}_8)$

接着证明 $\tilde{\phi}_2$ 对应的 \tilde{w} 与 ϕ_2 对应的 w 可以形成一一对应关系

$$\begin{aligned} w^T \phi_2 &= w_0 + w_1 x_1 + w_2 x_2 + w_3 (x_1^2) + w_4 (x_1 x_2) + w_5 (x_2^2) \\ &= w_0 + w_1 x_1 + w_2 x_2 + 0(x_1 + x_2) + 0(x_1 - x_2) + w_3 (x_1^2) + w_4 (x_1 x_2) + 0(x_2 x_1) + w_5 (x_2^2) \end{aligned}$$

所以 \tilde{w} 可以对应为 $(w_0, w_1, w_2, 0, 0, w_3, w_4, 0, w_5)$ 。

因此两种特征变化可以一一对应，那么

$$d_{vc}(\mathcal{H}_{\tilde{\phi}_2}) = d_{vc}(\mathcal{H}_{\phi_2}) \leq 6 < 9$$

Problem 3.19 (Page 117)

A Transformer thinks the following procedures would work well in learning from two-dimensional data sets of any size. Please point out if there are any potential problems in the procedures:

(a) Use the feature transform

$$\phi(x) = \begin{cases} (\underbrace{0, \dots, 0}_{n-1}, 1, 0, \dots) & \text{if } x = x_n \\ (0, \dots, 0) & \text{otherwise} \end{cases}$$

before running PLA.

(b) Use the feature transform ϕ with

$$\phi_n(x) = \exp\left(-\frac{\|x - x_n\|^2}{2\gamma^2}\right)$$

using some very small γ .

(c) Use the feature transform ϕ that consists of all

$$\phi_{i,j}(x) = \exp\left(-\frac{\|x - (i,j)\|^2}{2\gamma^2}\right)$$

before running PLA, with $i \in \{0, \frac{1}{100}, \dots, 1\}$ and $j \in \{0, \frac{1}{100}, \dots, 1\}$.

这题有点没有完全理解透，所以仅供参考。

(a)这题和台大的作业3第12题很像，我们来看三个点 x_1, x_2, x_3 的情形

$$\phi(x_1) = (1, 0, 0)$$

$$\phi(x_2) = (0, 1, 0)$$

$$\phi(x_3) = (0, 0, 1)$$

可以看到这里将3个点映射到了 \mathbb{R}^3 ，同理可知 N 个点可以映射到 \mathbb{R}^N ， \mathbb{R}^N 的感知机 $d_{vc} = N + 1$ ，所以 N 个点一定能被shatter，从而这种特征转换是好的。

(b)(c)暂时没有思路，略过。