



شناسایی اماری الگو

تمرین 1

بهاره غلامی

40033626

مهسا ملایم

40032718

1401/9/13

قسمت اول :

در این قسمت باید رگرسیون خطی را با استفاده از روش close form solution و gradient descent به دست بیاوریم و سپس خطای مدل که به دست آوردیم را با استفاده از mean square error محاسبه کنیم .

: Closed Form Linear Regression

در رگرسیون خطی ما به دنبال خطی هستیم به فرم زیر :

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

که تنها در فرمول بالا نماینگر پارامترهای رگرسیون هستند و در صورتی که ورودی های ما x باشد و به سبب ورودی را با n نمایش دهیم و y نماینگر متغیروابسته ما باشد تنها با استفاده از فرمول زیر به دست می آید :

Normal equation

$$\Theta = (X^T X)^{-1} X^T y$$

پیاده سازی :

برای خواندن دیتا از دستور زیر استفاده کرده ایم :

load dataset

```
In [17]: df = pd.read_csv('data.csv')
```

و با استفاده از دستور زیر دیتاهای آموزش و تست را از هم جدا کردیم که هفتاد در صد دیتا ها را آموزش و سی درصد را برای تست در نظر گرفتیم .

split_data

```
In [18]: train_x, test_x, train_y, test_y = ms.train_test_split(df['x'], df['y'], train_size=0.7)
```

نرمال سازی دیتا :

از روش minmax استفاده کردیم با استفاده از تابع زیر متغیر مستقل را به بازه 0 تا 1 میبریم .

```
] : def normalize(x):  
    x = (x - x.min() ) / (x.max() - x.min())  
    return x
```

:addbias

برای اضافه کردن ستون x0 به دیتای آموزش که برای محاسبه بایاس ضروری است از تابع زیر استفاده می کنیم :

```
] : def addbias(X):  
    matrix = X.reshape(X.shape[0], 1)  
    ones = np.ones((matrix.shape[0],1))  
    concat=np.concatenate((ones,matrix), axis=1)  
    return concat
```

خروجی این تابع به این صورت هست که یک ردیف یک به x_train یا x_test اضافه میکند .

: Leastsqared

از تابع زیر برای پیاده سازی فرمول بالا و به دست آوردن پارامتر های تتا استفاده میکنیم ورودی تابع x-train و y-train است و خروجی تابع مجموعه تتا :

```
def leastsquared(x,y):  
    y=y.reshape(y.shape[0],1)  
    thetaset = inv(x.T.dot(x)).dot(x.T).dot(y)  
    return thetaset
```

پیش بینی :

برای محاسبه پیش بینی دیتا های آموزش و تست و به دست آوردن y_{hat} از دستورات زیر استفاده کردیم

```
yHat_trn = leastSquared_theta[0] + leastSquared_theta[1] * train_x  
yHat_tst = leastSquared_theta[0] + leastSquared_theta[1] * test_x
```

Mse :

برای محاسبه mean square error از تابع زیر استفاده میکنیم :

```
def Mse(y,y_hat):  
    diff=y-y_hat  
    mse_pow=np.power(diff, 2,dtype='float64')  
    mse = np.mean(mse_pow)  
    return mse
```

نتایج :

```
print('MSE Train: ', mseTrain)  
print('MSE Test: ', mseTest)  
|
```

```
MSE Train: 17.019203409440927  
MSE Test: 13.990372372323291
```

پارامترهای به دست آمده :

```
print('theta0: ',leastSquared_theta[0])  
print('theta1: ',leastSquared_theta[1])
```

```
theta0: [79.75387007]  
theta1: [-49.62369159]
```

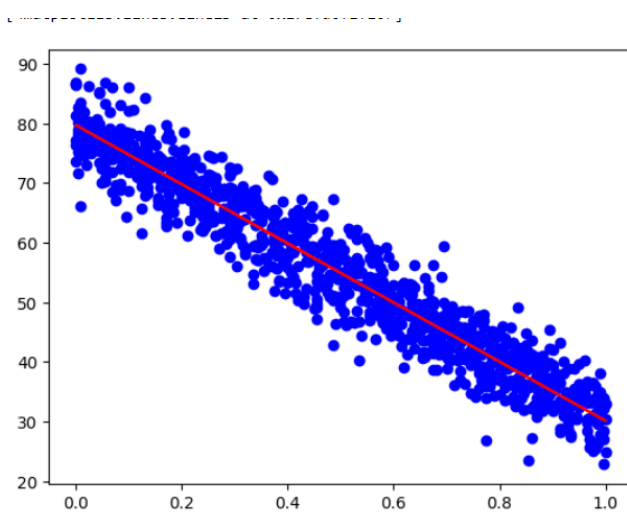
فرمول خط :

```
theta1: [-49.62369159]
```

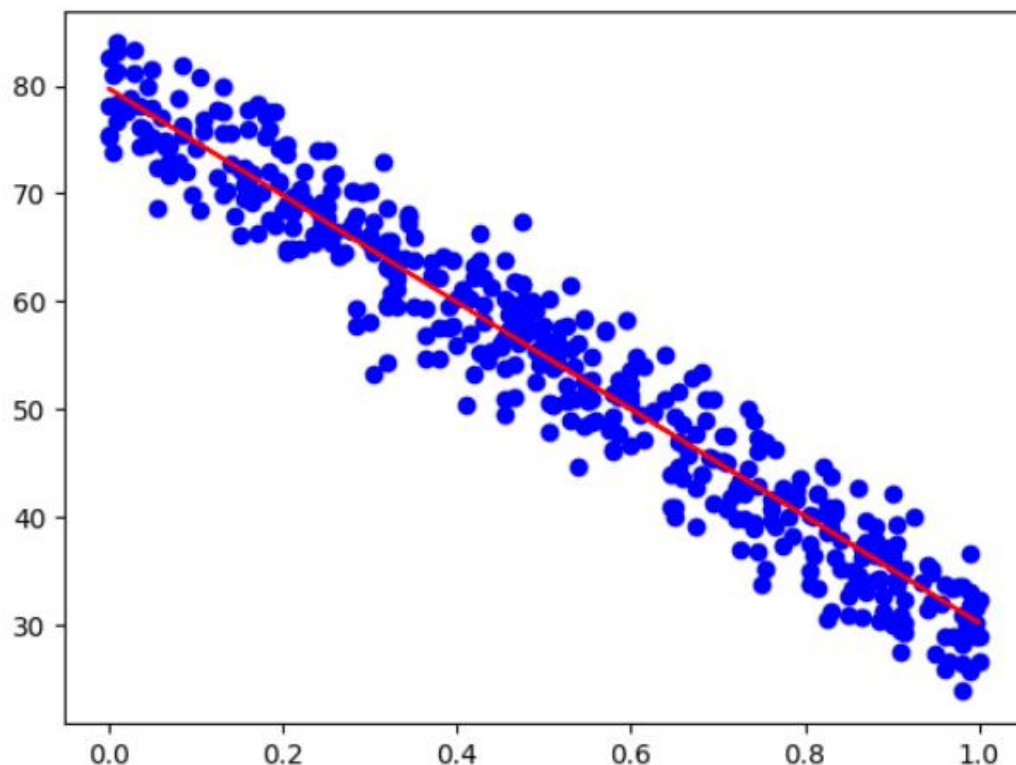
```
print('decision boundary formula:')  
print( 'y =' +str(leastSquared_theta[0]) + ' + ' + str(leastSquared_theta[1]) + '*x1')
```

```
decision boundary formula:  
y =[79.75387007] + [-49.62369159]*x1
```

نمودار دیتای آموزش :



نمودار دیتای تست :



Gradient Descent Linear Regression

ممکن هست که در بعضی از مسائل رگرسیون غیرخطی راه حل closed form وجود نداشته باشد یا حتی برای رگرسیون خطی هم در بعضی مواقع راه حل closed form مناسب نیست به دلیل اینکه ممکن ورودی ها خیلی بزرگ باشند و انجام محاسبات برای این راه حل پرهزینه میتواند باشد .

دونمونه gradient descent داریم :

: Stochastic Gradient Descent

این روش برای داده های حجیم مناسب است که آوردن همه آن ها بصورت یکجا در رم و انجام محاسبات منطقی نیست. برای همین در هر epoch ، محاسبات بر روی یک نمونه از داده انجام شده و پارامتر ها آپدیت میشوند.

: Batch Gradient Descent

در این روش در هر epoch از کل دیتا ها بصورت یکجا برای محاسبه پارامتر ها استفاده میشود.

در این قسمت ما از batch برای پیاده سازی استفاده کردیم .

از تابع هزینه mse استفاده کردیم :

$$\text{Cost} = J(w) = \frac{1}{2m} \sum_{i=1}^m (h(w)^{(i)} - y^{(i)})^2$$

و گرادیان مشتق این تابع بر حسب w است . (w همان تتا است)

$$\text{Gradient} = \frac{\partial J(w)}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (h(w)^{(i)} - y^{(i)}) \cdot X_j^{(i)}$$

برای آپدیت کردن وزن ها از فرمول زیر استفاده میکنیم :

$$w_j = w_j - lr \cdot \left(\frac{1}{m} \sum_{i=1}^m (h(w)^{(i)} - y^{(i)}) \cdot X_j^{(i)} \right)$$

پیاده سازی :

خواندن دیتا :

load dataset

```
In [17]: df = pd.read_csv('data.csv')
```

جداسازی دیتای تست از آموزش :

split_data

```
In [18]: train_x, test_x, train_y, test_y = ms.train_test_split(df['x'], df['y'], train_size=0.7)
```

نرمال کردن دیتا :

```
] : def normalize(x):  
    x = (x - x.min() ) / (x.max() - x.min())  
    return x
```

: Mse

```
def Mse(y,y_hat):  
    diff=y-y_hat  
    mse_pow=np.power(diff, 2,dtype='float64')  
    mse = np.mean(mse_pow)  
    return mse
```

: Addbias


```
] : def addbias(X):  
    matrix = X.reshape(X.shape[0], 1)  
    ones = np.ones((matrix.shape[0],1))  
    concat=np.concatenate((ones,matrix), axis=1)  
    return concat
```

Cost

از این تابع برای محاسبه هزینه استفاده میکنیم :

```
: def cost(m,a):  
    p=np.power(a,2,dtype='float64')  
    err = np.sum(p)/(2*m)  
    return err
```

ورودی این تابع a است که برابر

$$a=(y_hat-y)$$

و m برابر سائز دیتای آموزش هست .

: gradientDescent

از این تابع استفاده میکنیم برای پیاده سازی فرمول ها

```
def gradientDescent(x, y):
    alpha=0.01
    m,n=x.shape
    theta = np.random.rand(n).reshape(n, 1)
    errors=[]
    i=0
    for i in range(10000):
        x1=x[:,[1]]
        y_hat = theta[1]*x1 +theta[0]
        a=(y_hat-y.reshape(y.shape[0],1))
        error= cost(m,a)
        errors.append(error)
        #a=(y_hat-y.reshape(y.shape[0],1))
        gra_theta1= np.sum(x1*a)/m
        gra_theta0 = np.sum(a)/m
        theta[1] = theta[1] - alpha* gra_theta1
        theta[0] = theta[0] - alpha*gra_theta0

    errors = np.array(errors)
    return theta,errors
```

مقدار lr را 0.01 و تعداد ایپاک را 10000 در نظر گرفتیم .

آموزش :

```
: |theta,errors = gradientDescent( x_trnbias,train_y)
```

برای پیش بینی از کد زیر استفاده :

```
y_hat_Train =theta[0] + theta[1]*train_x
y_hat_Test = theta[0] +theta[1]*test_x
```

نتایج :

پارامترهای به دست آمده :

```
# Print Theta
print('theta0: ',theta[0])
print('theta1: ',theta[1])
```

```
theta0: [79.73857454]
theta1: [-49.54207482]
```

خط به دست آمده :

```
: print('decision boundary formula:')
print( 'y =' + str(theta[0]) + ' + ' + str(theta[1]) + '*x1')
```

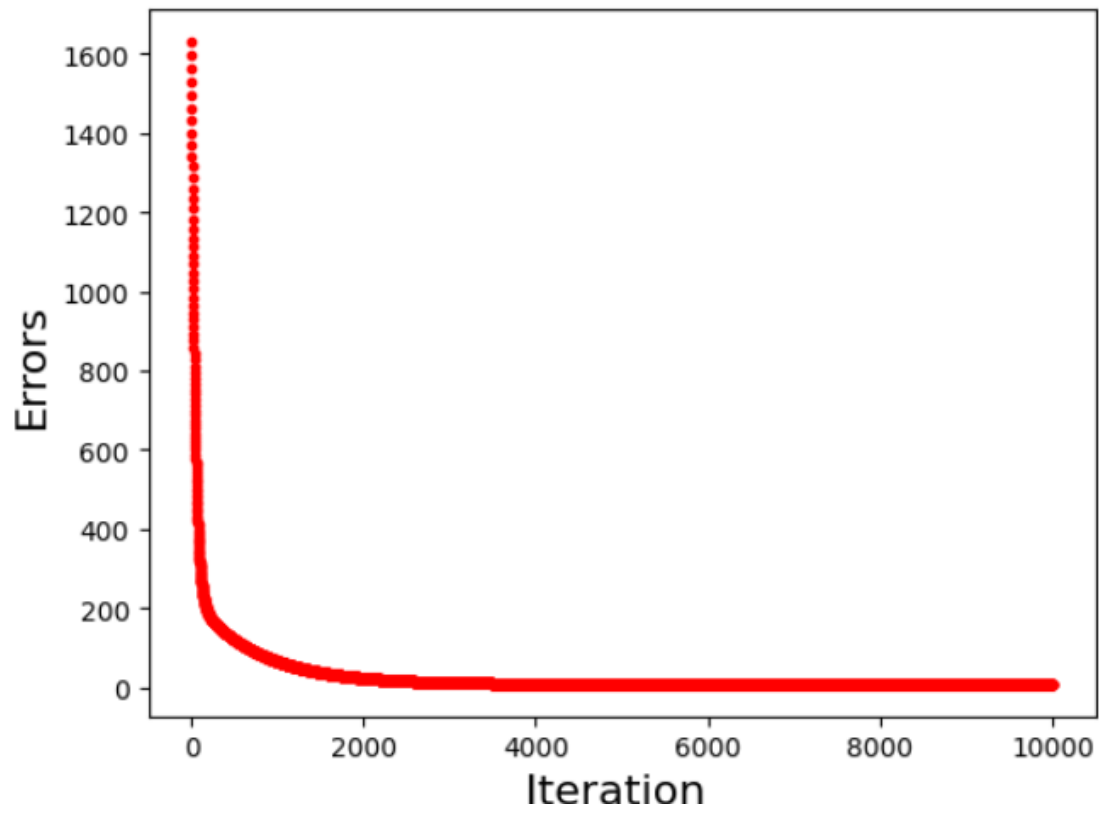
```
decision boundary formula:
y =[79.73857454] + [-49.54207482]*x1
```

: Mse

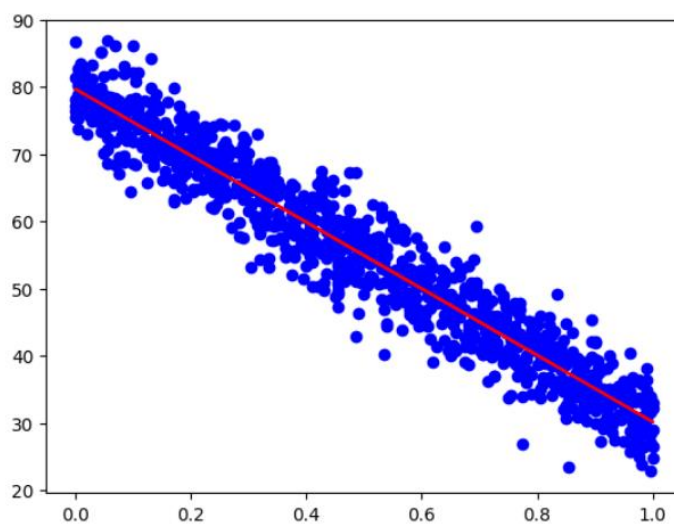
```
mseTrain = Mse(train_y,y_hat_Train)
mseTest = Mse(test_y,y_hat_Test)
print('MSE Train : ', mseTrain)
print('MSE Test : ', mseTest)
```

```
MSE Train : 16.028842593548845
MSE Test : 16.32632733058033
```

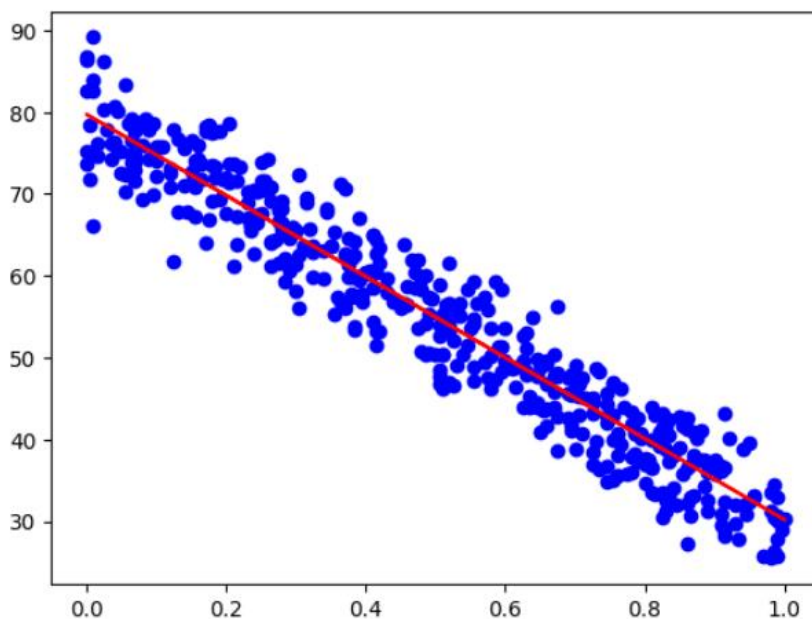
نمودار هزینه :



نمودار دیتای آموزش :



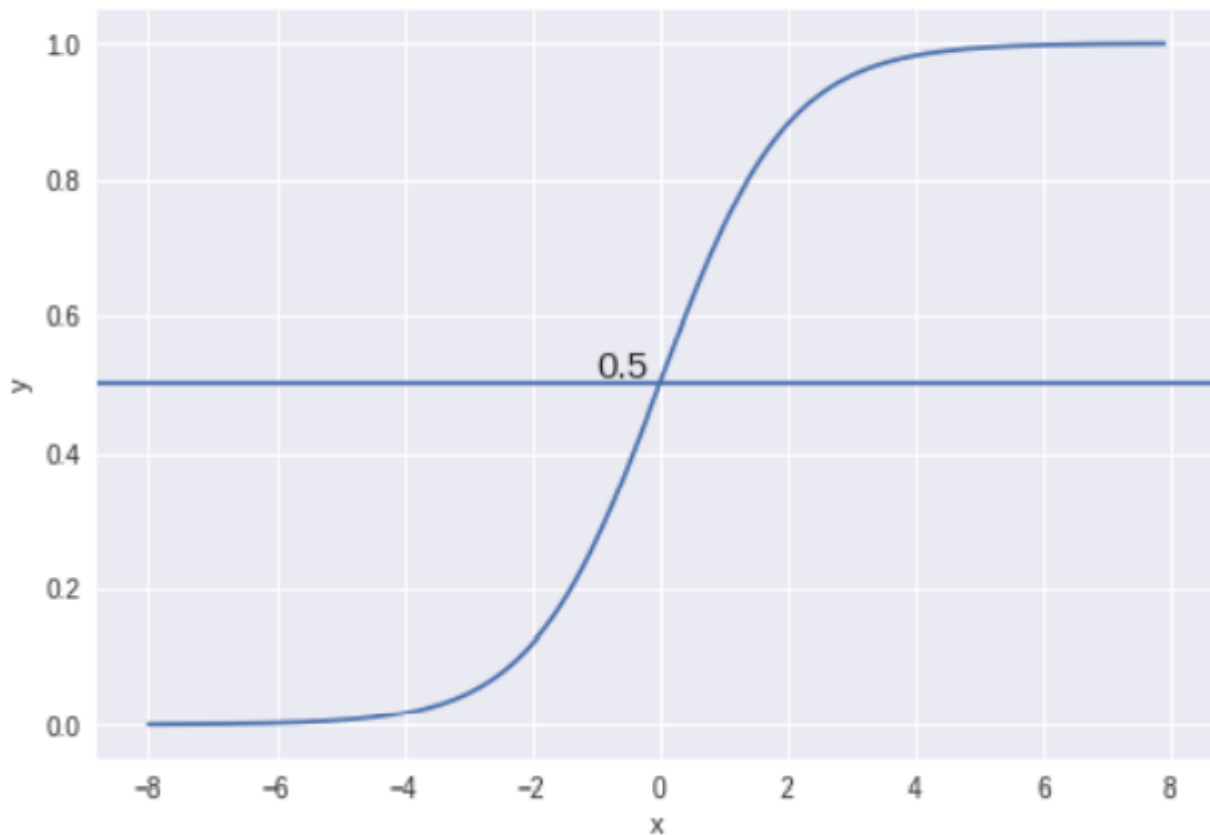
نمودار دیتای تست :



قسمت دوم :

Logistic Regression

از تابع سیگموئید در این روش استفاده میکنیم که نمودار آن به شکل زیر است :



این تابع همیشه مقداری بین صفر و یک دارد و از این تابع استفاده میکنیم در صورتی که مقدار احتمالی که به دست می آوریم بیشتر از 0.5 بود مربوط به کلاس یک و در صورتی که کوچکتر از 0.5 بود مربوط به کلاس 0 هست .

فرمول تابع سیگموئید به صورت زیر است :

$$h_{\theta}(x) = \underset{\uparrow}{g}(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad 0 \leq h_{\theta}(x) \leq 1$$



مانند رگرسیون خطی باید یک تابع هزینه مشخص کنیم و آن را کمینه کنیم از تابع زیر استفاده میکنیم :

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

برای کمینه کردن معادله از gradient descent استفاده میکنیم .

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad j=0, 1, \dots, n$$

}

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

پیاده سازی :

: load_data

از این تابع استفاده میکنیم و تنها فیچر ردیف سوم و هفتم را جدا کرده و دیتاهایی که مربوط به کلاس سه هست را حذف میکنیم و سپس دیتای آموزش و تست را از هم جدا میکنیم:

```
def load_data():
    df=pd.read_excel("seed.xlsx")
    df=df[['x3','x7','class']]
    cdf=df[df['class']<3]
    cdf_1=df[df['class'] == 1]
    cdf_2=df[df['class'] == 2]
    cdf_2=cdf_2.replace(to_replace = 2, value =0)
    g=cdf_1.append(cdf_2, ignore_index=True)
    train_x, test_x, train_y, test_y = ms.train_test_split(g[['x3','x7']],g['class'], train_size=0.7)
    return train_x, test_x, train_y, test_y
```

: Normalize

```
] def normalize(x):
    x = (x - x.min() ) / (x.max() - x.min())
    return x
```

: Mse

```
def Mse(y,y_hat):
    diff=y-y_hat
    mse_pow=np.power(diff, 2,dtype='float64')
    mse = np.mean(mse_pow)
    return mse
```

: Sigmoid

```
def sigmoid(x,theta,theta0):
    temp=np.dot(x,theta )+theta0
    y = 1 / (1 + np.exp(-temp))
    return y
```

ورودی این تابع متغیر های مستقل آموزش و مجموعه تتا هست .

: Cost


```
1]: def cost(y_hat,y,m):  
    d=-1/m * np.sum(y * np.log(y_hat) + (1 - y) * np.log(1-y_hat))  
    return d
```

ورودی تابع مقدار y و y_hat و سائز دیتاهای آموزش است و مقدار هزینه را مشخص میکند در هر تکرار .

: LogisticRegression

```
def LogisticRegression(x, y,alpha,itr):  
    m, n = x.shape  
    theta=np.random.rand(n)  
    theta0=0.001  
    errors= []  
    for i in range(itr):  
        y_hat = sigmoid(x,theta,theta0)  
        error= cost(y_hat,y,m)  
        errors.append(error)  
        theta= theta - alpha*(1/m * np.dot(x.T, (y_hat - y)))  
        theta0=theta0 - alpha* (1/m * np.sum(y_hat - y))  
    errors = np.array(errors)  
    return theta,theta0,errors
```

بر اساس فرمول که توضیح دادیم مدل آموزش میبیند .

آموزش مدل :

```
: theta,theta0,errors = LogisticRegression(train_x,train_y, 0.05,70000)
```

برای آموزش مدل تعداد ایپاک را 70000 و الفا را 0.05 گرفتیم .

: decision_boundary

برای نمایش خط جداکننده دو کلاس از این تابع استفاده میکنیم .

```
def decision_boundary(x,theta,theta0):  
    return - (theta0 + np.dot(theta[0], x)) / theta[1]
```

```
print('Decision boundary formula:')  
print('-(' + str(theta0) + ' + ' + str(  
    theta[0]) + '*x1) / (' + str(theta[1]) + '*x2)')
```

Decision boundary formula:

$-(22.55000292540777 + 0.10106886184755208*x1) / (-27.298907690551186*x2)$

نتایج :

پارامترهای به دست آمده :

```
: print('theta0: ',theta0)  
print('theta1: ',theta[0])  
print('theta2: ',theta[1])
```

theta0: 22.55000292540777
theta1: 0.10106886184755208
theta2: -27.298907690551186

: Mse

```
: print("MSE Train = ", Mse(train_y,sigmoid(train_x,theta,theta0)))  
print("MSE Test = ",Mse( test_y,sigmoid(test_x,theta,theta0)))
```

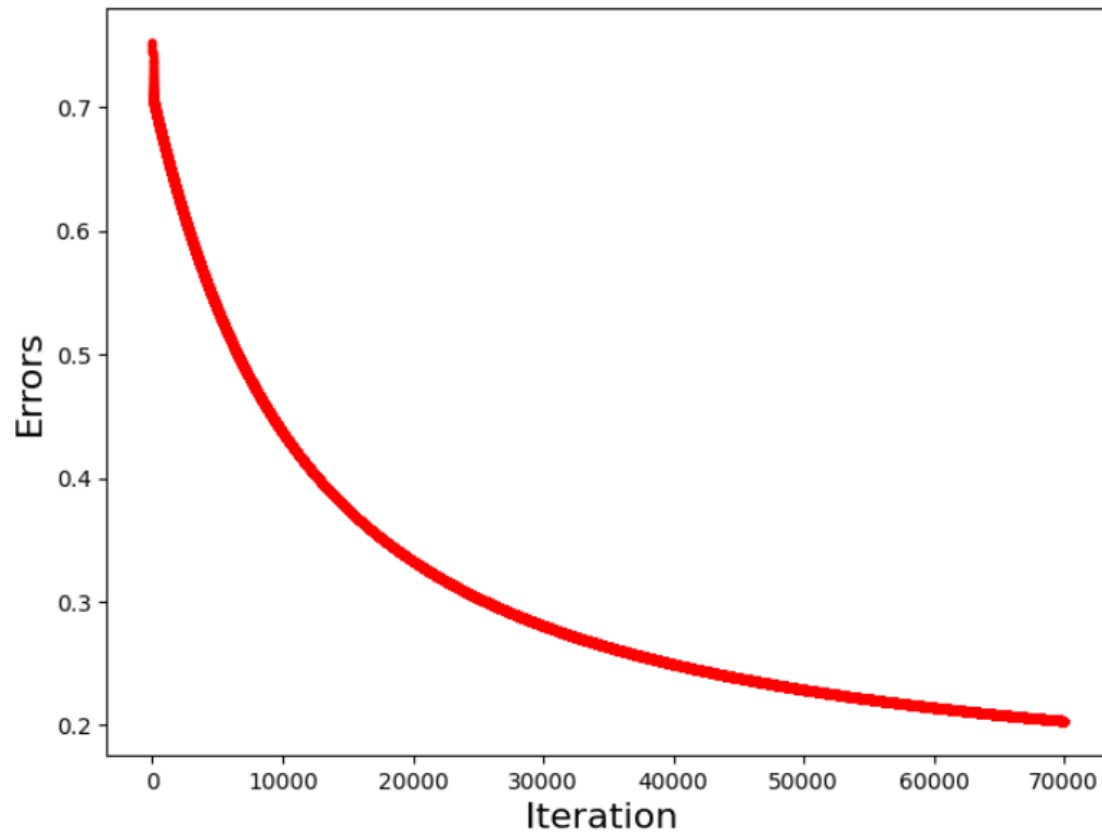
MSE Train = 0.04957891020411598
MSE Test = 0.022715285122802936

فرمول خط :

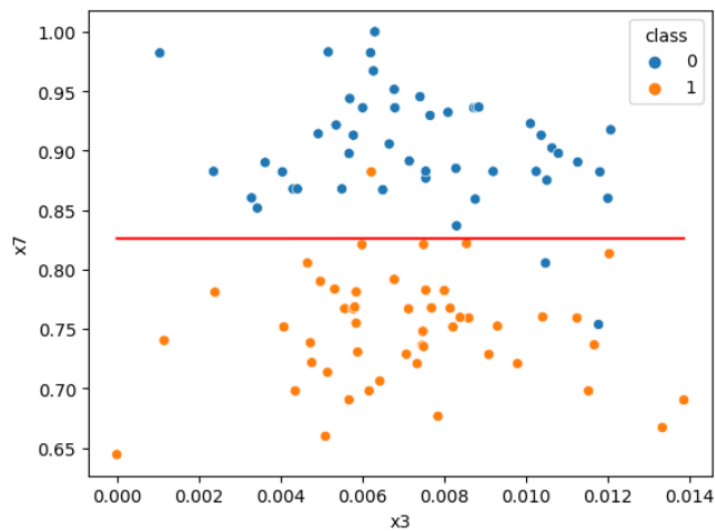
```
]: print('decision boundary formula:')  
print('y =' + str(theta0) + ' + ' + str(theta[0]) + '*x1' + ' + ' + str(theta[1]) + '*x2')
```

decision boundary formula:
 $y = 22.55000292540777 + 0.10106886184755208*x1 + -27.298907690551186*x2$

نمودار هزینه :

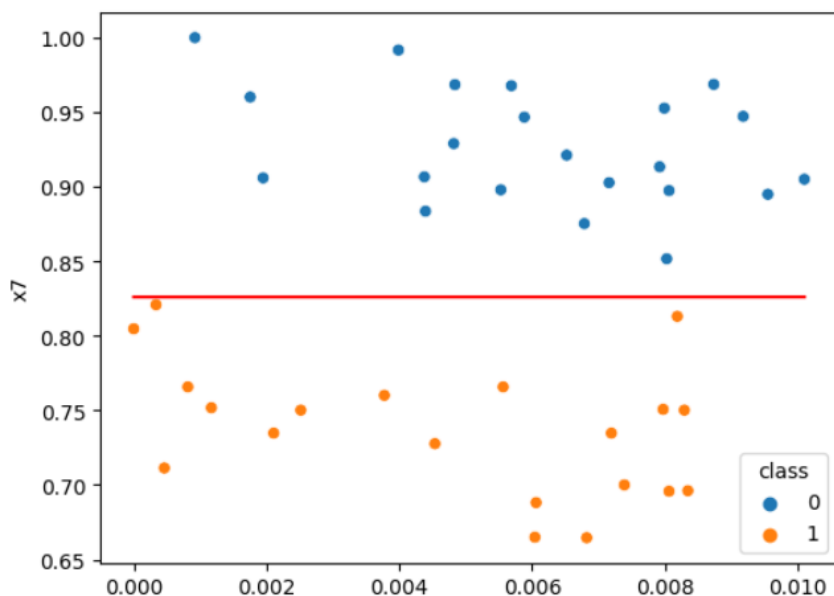


نمودار دیتای آموزش :



```
[0.35] x_values = pd.DataFrame(data=test_x, columns=["x3", "x7"])
```

نمودار دیتای تست :



دقت دیتای آموزش و تست :

که از فرمول گفته شده سر کلاس استفاده کردیم .

```
y_hat_train = sigmoid(train_x, theta, theta0)
y_hat_train = [0 if y <= 0.5 else 1 for y in y_hat_train]
accuracy_train = np.sum(train_y == y_hat_train) / (train_y.shape[0])
print('Train accuracy:', accuracy_train)
```

Train accuracy: 0.9693877551020408

```
y_hat_test = sigmoid(test_x, theta, theta0)
y_hat_test = [0 if y <= 0.5 else 1 for y in y_hat_test]
accuracy_test = np.sum(test_y == y_hat_test) / (test_y.shape[0])
print('Test accuracy:', accuracy_test)
```

Test accuracy: 1.0