# MA333 Project 1

11711416 MingjiHan

## Overview

In the project, we are going to implement the linear regression algorithm and test it on dataset.

## Formula Derivation

### 1.The log-likelihood function

$$
\begin{aligned}
l(\mathbf{w}) &= \sum_{i=1}^{n} log(p_1(\mathbf{x}_i; \mathbf{w})^{y_i} p_0(\mathbf{x}_i; \mathbf{w})^{1-y_i}) \\
&= \sum_{i=1}^{n} y_i log(p_1(\mathbf{x}_i; \mathbf{w})) + (1-y_i)log(p_0(\mathbf{x}_i; w)) \\
&= \sum_{i=1}^{n} y_i log(p(\mathbf{x}_i; \mathbf{w})) + (1-y_i)log(1-p(\mathbf{x}_i; w)) \\
&= \sum_{i=1}^{n} y_i log(\frac{e^{\mathbf{w}^T \mathbf{x}_i}}{1+e^{\mathbf{w}^T \mathbf{x}_i}}) + (1-y_i)log(\frac{1}{1+e^{\mathbf{w}^T \mathbf{x}_i}}) \\
&= \sum_{i=1}^{n} y_i log(e^{\mathbf{w}^T \mathbf{x}_i}) - y_i log(1+e^{\mathbf{w}^T \mathbf{x}_i}) + (y_i-1)log(1+e^{\mathbf{w}^T \mathbf{x}_i}) \\
&= \sum_{i=1}^{n} y_i \mathbf{w}^T \mathbf{x}_i - log(1+e^{\mathbf{w}^T \mathbf{x}_i})
\end{aligned}
$$

### 2. Score Equation

Based on the formula derived from part 1,we can continue expand it

$$
\begin{aligned}
l(\mathbf{w}) &= \sum_{i=1}^{n} y_i \mathbf{w}^T \mathbf{x}_i - log(1+e^{\mathbf{w}^T \mathbf{x}_i}) \\
&= \sum_{i=1}^{n} y_i \sum_{j=1}^{d} x_{ij} w_j - log(1+e^{\mathbf{w}^T \mathbf{x}_i})
\end{aligned}
$$

Then we find the partial derivative of $ w\_j (j = 0,1,2,3 .... d) $

$$
\begin{aligned}
\frac{\partial L(\mathbf{w})}{\partial w_j} &= \sum_{i=1}^{n} y_i x_{ij} - \frac{x_{ij} e^{\mathbf{w}^T \mathbf{x}_i}}{1+e^{\mathbf{w}^T \mathbf{x}_i}} \\
&= \sum_{i=1}^{n} x_{ij}(y_i - p(\mathbf{x}_i; \mathbf{w}))
\end{aligned}
$$

$$
\frac{\partial L(\mathbf{w})}{\partial w_j} = 0
$$

$$
\sum_{i=1}^{n} x_{ij}(y_i - p(\mathbf{x}_i; \mathbf{w})) = 0
$$

# 3.Hessian matrix

Consider the j th element of k th row in the Hessian Matrix. We want to find:

$$\frac{\partial L(\mathbf{w})}{\partial w_k \partial w_j}$$

Fron Problem 2 we know

$$\frac{\partial L(\mathbf{w})}{\partial w_j} = \sum_{i=1}^{n} x_{ij}(y_i - p(\mathbf{x}_i; \mathbf{w}))$$

So

$$\begin{aligned}
\frac{\partial L(\mathbf{w})}{\partial w_k \partial w_j} &= -\sum_{i=1}^{n} x_{ij} \frac{x_{ik} e^{\mathbf{w}^T \mathbf{x}_i}(1 + e^{\mathbf{w}^T \mathbf{x}_i}) - e^{\mathbf{w}^T \mathbf{x}_i} e^{\mathbf{w}^T \mathbf{x}_i} x_{ik}}{(1 + e^{\mathbf{w}^T \mathbf{x}_i})^2} \\
&= -\sum_{i=1}^{n} x_{ik} x_{ij} \frac{e^{\mathbf{w}^T \mathbf{x}_i}}{(1 + e^{\mathbf{w}^T \mathbf{x}_i})^2} \\
&= -\sum_{i=1}^{n} x_{ik} x_{ij} \frac{e^{\mathbf{w}^T \mathbf{x}_i}}{1 + e^{\mathbf{w}^T \mathbf{x}_i}} \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}_i}} \\
&= -\sum_{i=1}^{n} x_{ik} x_{ij} p(\mathbf{x}_i; \mathbf{w})(1 - p(\mathbf{x}_i; \mathbf{w}))
\end{aligned}$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = -\sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^T p(\mathbf{x}_i; \mathbf{w})(1 - p(\mathbf{x}_i; \mathbf{w}))$$

# 4.Matrix Notation of IRLS

Fristly we convert the derivative into matrix form

$$\begin{aligned}
\frac{\partial L(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} &= -\sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^T p(\mathbf{x}_i; \mathbf{w})(1 - p(\mathbf{x}_i; \mathbf{w})) \\
&= -(X^T D X)^{-1}
\end{aligned}$$

$$\frac{\partial L(\mathbf{w})}{\partial w_j} = \sum_{i=1}^{n} x_{ij}(y_i - p(\mathbf{x}_i; \mathbf{w}))$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = X^{-1}(\mathbf{y} - \mathbf{p})$$

Then we get the matrix form of $ \textbf{w}^{(k)} $

$$\begin{aligned}
\mathbf{w}^{(k)} &= \mathbf{w}^{(k-1)} - \left(\frac{\partial^2 l(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T}\right)^{-1}\Big|_{\mathbf{w}^{(k-1)}} \left(\frac{\partial l(\mathbf{w})}{\partial \mathbf{w}}\right)^{-1}\Big|_{\mathbf{w}^{(k-1)}} \\
&= \mathbf{w}^{(k-1)} + (X^T D X)^{-1} X^T (\mathbf{y} - \mathbf{p}) \\
&= (X^T D X)^{-1}((X^T D X)\mathbf{w}^{(k-1)} + X^T(\mathbf{y} - \mathbf{p})) \\
&= (X^T D X)^{-1} X^T D(X\mathbf{w}^{(k-1)} + D^{-1}(\mathbf{y} - \mathbf{p})) \\
&= (X^T D X)^{-1} X^T D \mathbf{z} \\
&= \underset{\mathbf{w}}{argmin}(\mathbf{z} - X\mathbf{w})^T D(\mathbf{z} - X\mathbf{w})
\end{aligned}$$

# IRLS Algorithm

```python
import numpy as np
from matplotlib import pyplot as plt
from sklearn.preprocessing import normalize
def exp(x):
    return np.exp(x)


'''IRLS Alogrithm for Logistic Regression

    Args:
        X:Training data, a n *(d+1) matrix
        y:Class of training data, n*1 column vector
        error: Error bound of w
    Return:
        w: Weights, a (d+1) * 1 column vector
    Notice: In order to avoid finding the inverse singular matrix which will crash the program,
            we use np.linalg.pinv() to get the "inverse" of martix
'''
def IRIS(X,y,print_loss=True,error=1e-3,step=10):
    print("Algorithm starts.")

    w = np.random.random((X.shape[1],1))
    w_pre = w
    loss = []
    iter = 0

    while True:
        inner_product = X.dot(w)

        p = exp(inner_product) / (1.0 + exp(inner_product))

        D = np.diag(p.T[0])

        w = w + np.dot(np.dot(np.linalg.pinv(np.dot(np.dot(X.T,D),X)),X.T),y - p)

        diff = np.sum(np.abs(w - w_pre))

        if diff < error:
            break
        else:
            w_pre = w

        if (iter+1) % 2 == 0:
            if print_loss:
                print(diff)
            loss.append(diff)
        iter += 1

    print("Algorithm finished.")
    return w,loss
```

```python
X = np.random.random((10,6))#np.array([[0.1,0.7,2.4],[4.1,5.8,6.2],[7.7,8.5,9.3]])
y = np.array([[1.0],[0.0],[0.0],[1.0],[0.0],[1.0],[0.0],[0.0],[1.0],[0.0]])
w,loss = IRIS(X,y,False,1e-3)
```
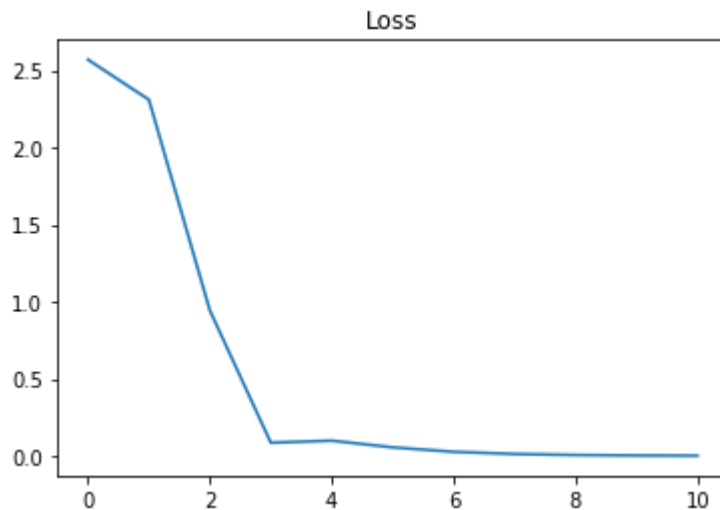
```
Algorithm starts.
Algorithm finished.
```

```
1  plt.title("Loss")
2  plt.plot(loss)
```

```
1  [<matplotlib.lines.Line2D at 0x1df383464a8>]
```



# Test algorithm on South African Hearth dataset

1.Define function for reading the dataset

```
1   import os
2   from sklearn import metrics
3   def read_dataset(file_name):
4       f = open(file_name)
5       data_x = []
6       data_y = []
7       for line in f.readlines():
8           if line[0] !='@':
9               line = line.strip('\n')
10              items = line.split(', ')
11
12              data_row = [1] # for w0
13              for i in range(len(items)):
14                  if i == 4:
15                      if items[i] == 'Absent':
16                          data_row.append(1.0)
17                          data_row.append(0.0)
18                      else:
19                          data_row.append(0.0)
20                          data_row.append(1.0)
21                  elif i == 9:
22                      data_y.append([int(items[i])])
23                  else:
24                      data_row.append(float(items[i]))
25              data_x.append(data_row)
26          # print(items)
27      data_x = np.array(data_x)
28      data_y = np.array(data_y)
```

```
29        return (data_x,data_y)
```
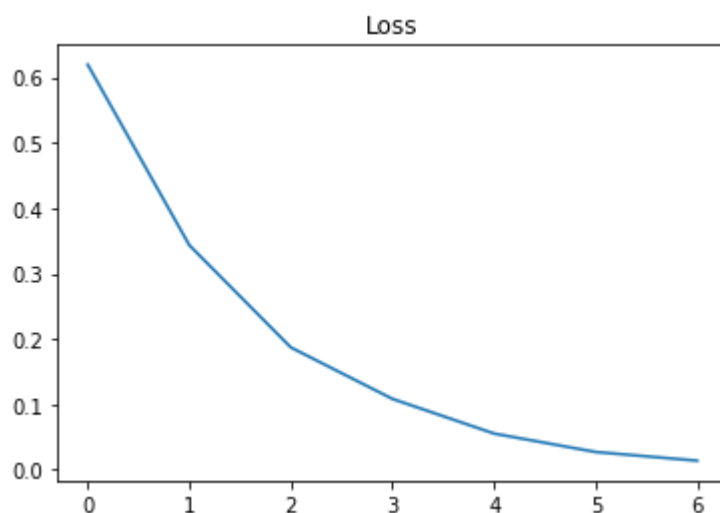
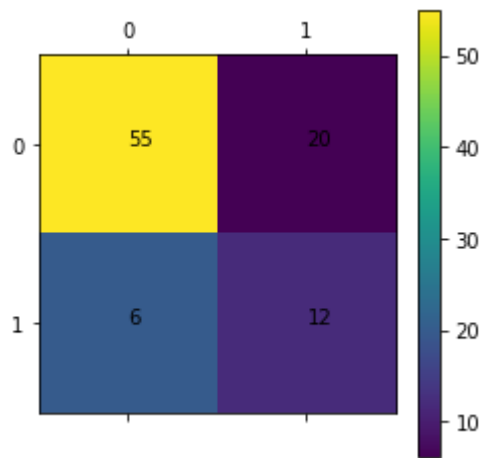2.Define the function for traininng and testing

```
1   from sklearn import preprocessing
2   def train_and_test(num,threashold=0.75):
3       (data_x,data_y) = read_dataset('./saheart-5-fold/saheart-5-'+str(num)+'tra.dat')
4       data_x = preprocessing.scale(data_x)
5       w1,loss1 = IRIS(data_x,data_y,False,1e-2)
6       plt.title("Loss")
7       plt.plot(loss1)
8       (test_x,test_y) = read_dataset('./saheart-5-fold/saheart-5-'+str(num)+'tst.dat')
9       test_x = preprocessing.scale(test_x)
10      y = test_x.dot(w1)
11      y  = exp(y) / (1.0 + exp(y))
12      y_pred = y.T > 0.75
13      y_pred = 1 * y_pred
14
15      score = metrics.accuracy_score(test_y.T[0],y_pred[0])
16      print("Accuracy: " + str(score))
17      print("Precision score:"+str(precision_score(y_pred[0],test_y)))
18      mat = confusion_matrix(test_y, y_pred[0])
19      plt.matshow(confusion_matrix(test_y, y_pred[0]))
20      plt.text(0,0,mat[0][0])
21      plt.text(0,1,mat[0][1])
22      plt.text(1,0,mat[1][0])
23      plt.text(1,1,mat[1][1])
24      plt.colorbar()
```
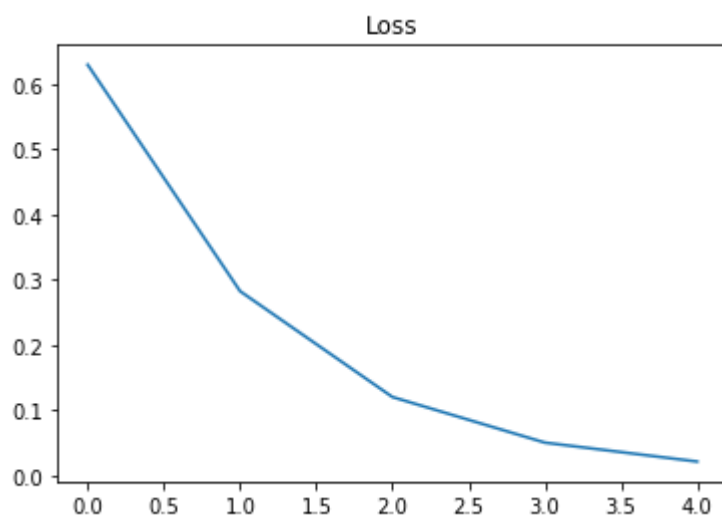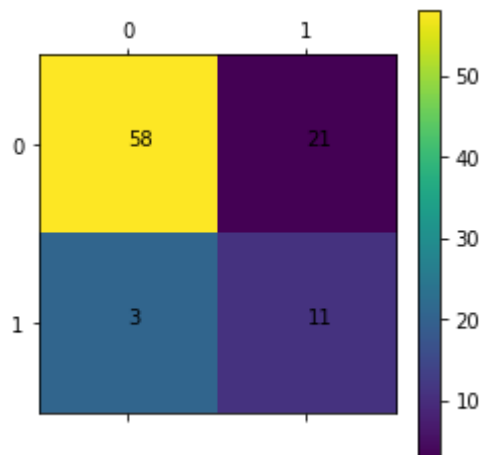
3.Train and test

```
1   train_and_test(1)
```

```
1   Algorithm starts.
2   Algorithm finished.
3   Accuracy: 0.7204301075268817
4   Precision score:0.375
```
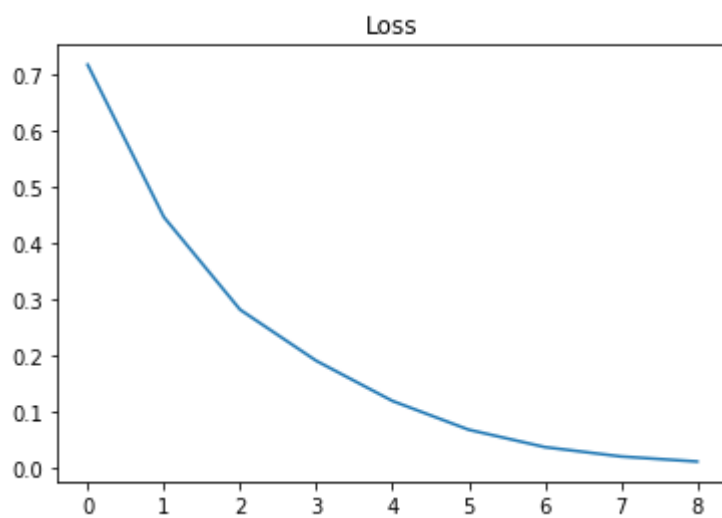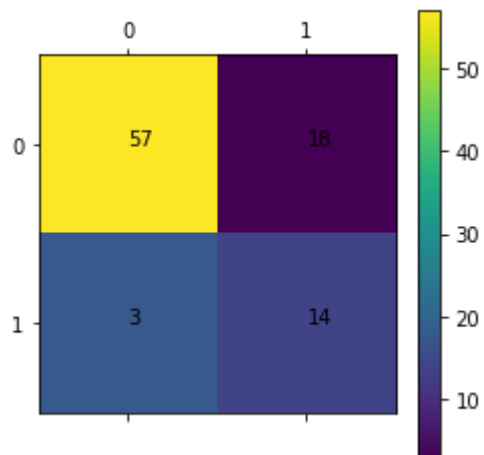
```
1  train_and_test(2)
```

```
1  Algorithm starts.
2  Algorithm finished.
3  Accuracy: 0.7419354838709677
4  Precision score:0.34375
```
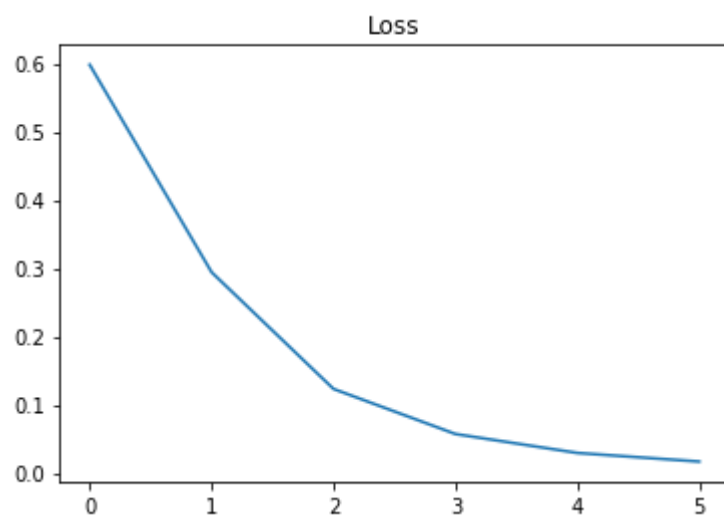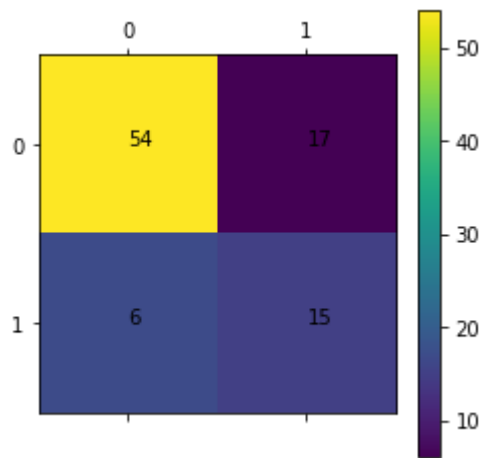
```
1  train_and_test(3)
```

```
1  Algorithm starts.
2  Algorithm finished.
3  Accuracy: 0.7717391304347826
4  Precision score:0.4375
```
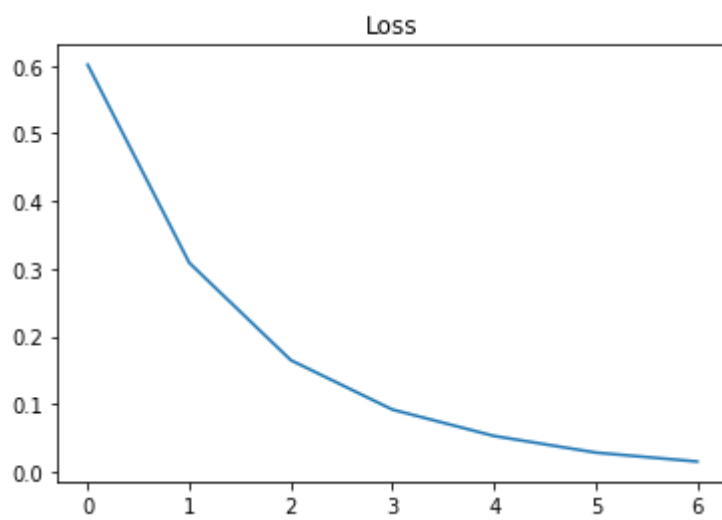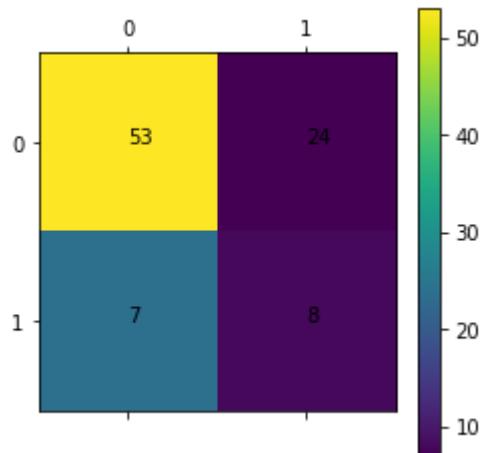
```
1  train_and_test(4)
```

```
1  Algorithm starts.
2  Algorithm finished.
3  Accuracy: 0.75
4  Precision score:0.46875
```



Loss

```
1  train_and_test(5)
```

```
1  Algorithm starts.
2  Algorithm finished.
3  Accuracy: 0.6630434782608695
4  Precision score:0.25
```

4.Train and test on sklearn

```
1  from sklearn.linear_model import LogisticRegression
2  from sklearn.metrics import confusion_matrix
3  from sklearn.metrics import precision_score
```
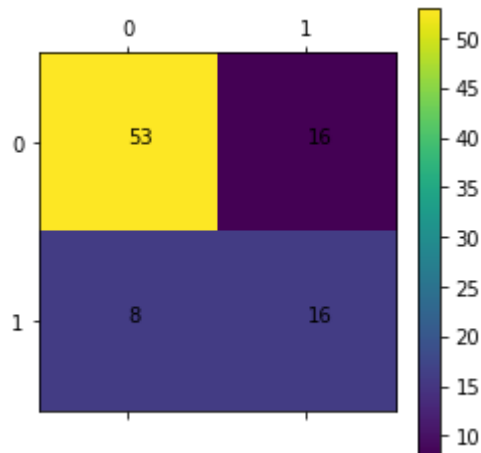
```
1   def sk_train_and_test(num,threashold=0.75):
2       (data_x,data_y) = read_dataset('./saheart-5-fold/saheart-5-'+str(num)+'tra.dat')
3       data_x = preprocessing.scale(data_x)
4       clf = LogisticRegression(random_state=0, solver='lbfgs',
5                                multi_class='ovr').fit(data_x, data_y)
6       (test_x,test_y) = read_dataset('./saheart-5-fold/saheart-5-'+str(num)+'tst.dat')
7       test_x = preprocessing.scale(test_x)
8       y_pred = clf.predict(test_x)
9       score = metrics.accuracy_score(test_y.T[0],y_pred)
10      print("Accuracy: " + str(score))
11      print("Precision score:"+str(precision_score(y_pred,test_y)))
12      mat = confusion_matrix(test_y, y_pred)
13      plt.matshow(confusion_matrix(test_y, y_pred))
14      plt.text(0,0,mat[0][0])
15      plt.text(0,1,mat[0][1])
16      plt.text(1,0,mat[1][0])
17      plt.text(1,1,mat[1][1])
18      plt.colorbar()
```

```
1  sk_train_and_test(1)
```

```
1  Accuracy: 0.7419354838709677
2  Precision score:0.5
```

```
1  c:\users\mingji han\appdata\local\programs\python\python37\lib\site-
   packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was
   passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
   example using ravel().
2    y = column_or_1d(y, warn=True)
```
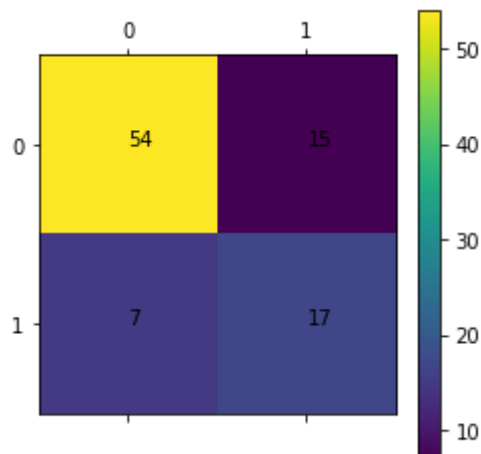
```
1  sk_train_and_test(2)
```

```
1  Accuracy: 0.7634408602150538
2  Precision score:0.53125
```

```
1  c:\users\mingji han\appdata\local\programs\python\python37\lib\site-
   packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was
   passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
   example using ravel().
2    y = column_or_1d(y, warn=True)
```
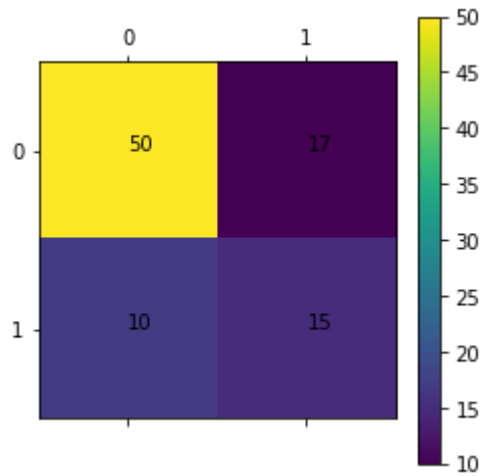


```
1  sk_train_and_test(3)
```

```
1  Accuracy: 0.7065217391304348
2  Precision score:0.46875
```

```
1  c:\users\mingji han\appdata\local\programs\python\python37\lib\site-
   packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was
   passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
   example using ravel().
2    y = column_or_1d(y, warn=True)
```
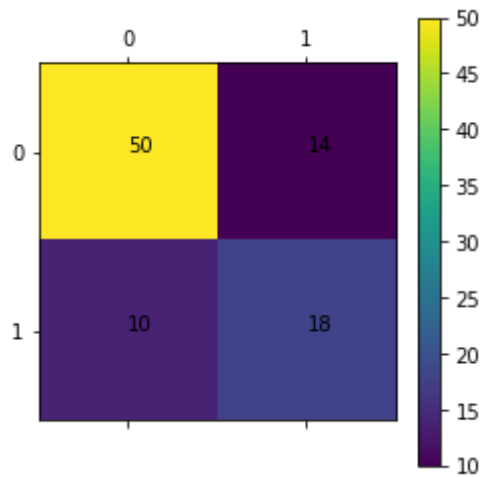
```
1  sk_train_and_test(4)
```

```
1  Accuracy: 0.7391304347826086
2  Precision score:0.5625
```

```
1  c:\users\mingji han\appdata\local\programs\python\python37\lib\site-
   packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was
   passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
   example using ravel().
2    y = column_or_1d(y, warn=True)
```
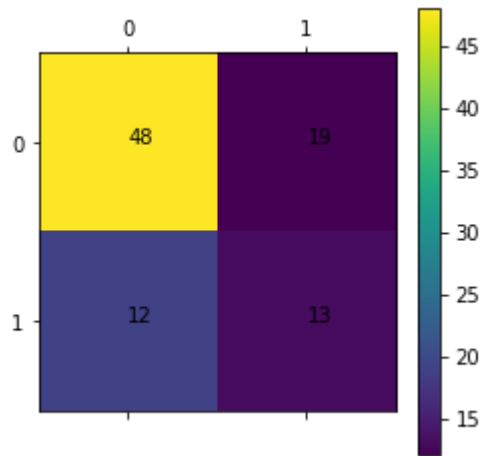


```
1  sk_train_and_test(5)
```

```
1  Accuracy: 0.6630434782608695
2  Precision score:0.40625
```

```
1  c:\users\mingji han\appdata\local\programs\python\python37\lib\site-
   packages\sklearn\utils\validation.py:761: DataConversionWarning: A column-vector y was
   passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
   example using ravel().
2    y = column_or_1d(y, warn=True)
```

# Conclusion

The algorithm we design is as good as the logistic regression in scikit-learn.

Our alogrithm can be improved through ensemble method.

```
1
```