In [37]: ▶|

```python
1  import numpy as np
2  import matplotlib.pyplot as plt
3  from getDataset import getDataSet
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.model_selection import train_test_split
6  import random
7  import math
8  from sklearn import preprocessing
9  from pandas import DataFrame
10 from pylab import scatter, show, legend, xlabel, ylabel
11 from GD import gradientDescent
12 from dataNormalization import rescaleMatrix
13 from numpy import loadtxt, where
14 from sklearn.metrics import roc_curve, auc
```

In [2]: ▶|

```python
1  def getDataSet():
2      """
3      Returns X (250 X 2) and Y (250 X 1)
4      """
5      # Step 1: Generate data by a module
6      n = 100  # 1st class contains N objects
7      alpha = 1.5  # 2st class contains alpha*N ones
8      sig2 = 1  # assume 2nd class has the same variance as the 1st
9      dist2 = 4
10
11     # later we move this piece of code in a separate file
12     # [X, y] = loadModelData(N, alpha, sig2, dist2);
13     n2 = math.floor(alpha * n)  # calculate the size of the 2nd class
14     cls1X = np.random.randn(n, 2)  # generate random objects of the 1s
15
16     # generate a random distance from the center of the 1st class to t
17     # https://stackoverflow.com/questions/1721802/what-is-the-equivale
18     a = np.array([[math.sin(math.pi * random.random()), math.cos(math.
19     a1 = a * dist2
20     shiftClass2 = np.kron(np.ones((n2, 1)), a1)
21
22     # generate random objects of the 2nd class
23     cls2X = sig2 * np.random.randn(n2, 2) + shiftClass2
24     # combine the objects
25     X = np.concatenate((cls1X, cls2X), axis=0)
26
27     # assign class labels: 0s and 1s
28     y = np.concatenate((np.zeros((cls1X.shape[0], 1)), np.ones((cls2X.
29     # end % of module.
30     return X, y
31
```

In [3]: ▶|

```
1  # Starting codes
2
3  # Fill in the codes between "%PLACEHOLDER#start" and "PLACEHOLDER#end"
4
5  # step 1: generate dataset that includes both positive and negative sc
6  # where each sample is described with two features.
7  # 250 samples in total.
8
9  [X, y] = getDataSet()   # note that y contains only 1s and 0s,
10
11 # create figure for all charts to be placed on so can be viewed togeth
12 fig = plt.figure()
13
```
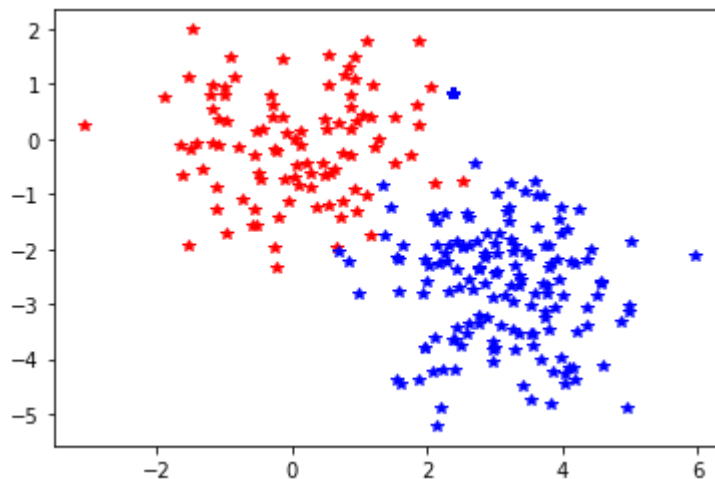
<Figure size 432x288 with 0 Axes>

In [4]: ▶|

```
1  def func_DisplayData(dataSamplesX, dataSamplesY, chartNum, titleMessag
2      idx1 = (dataSamplesY == 0).nonzero()  # object indices for the 1st
3      idx2 = (dataSamplesY == 1).nonzero()
4      ax = fig.add_subplot(1, 3, chartNum)
5      # no more variables are needed
6      plt.plot(dataSamplesX[idx1, 0], dataSamplesX[idx1, 1], 'r*')
7      plt.plot(dataSamplesX[idx2, 0], dataSamplesX[idx2, 1], 'b*')
8      # axis tight
9      ax.set_xlabel('x_1')
10     ax.set_ylabel('x_2')
11     ax.set_title(titleMessage)
12
13
14 # plotting all samples
15 func_DisplayData(X, y, 1, 'All samples')
16
17 # number of training samples
18 nTrain = 120
19
```
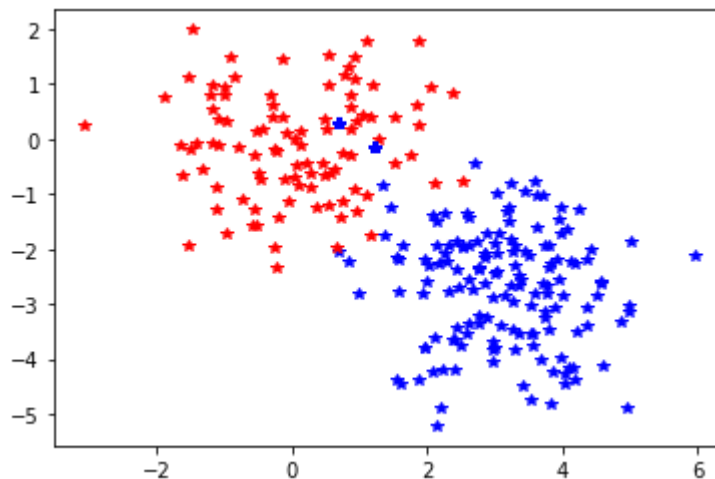
```
In [8]:    1  #######################PLACEHOLDER 1#start#######################
           2  # write you own code to randomly pick up nTrain number of samples for
           3  # WARNIN: do not use the scikit-learn or other third-party modules for
           4
           5  #maxIndex = len(X)
           6  #randomTrainingSamples = np.random.choice(maxIndex, nTrain, replace=Fa
           7  shuffled_indicies = np.arange(X.shape[0])
           8  np.random.shuffle(shuffled_indicies)
           9
          10  nTrain = 120
          11
          12  train_shuffled_indicies = shuffled_indicies[:nTrain]
          13  test_shuffled_indicies = shuffled_indicies[nTrain:]
```
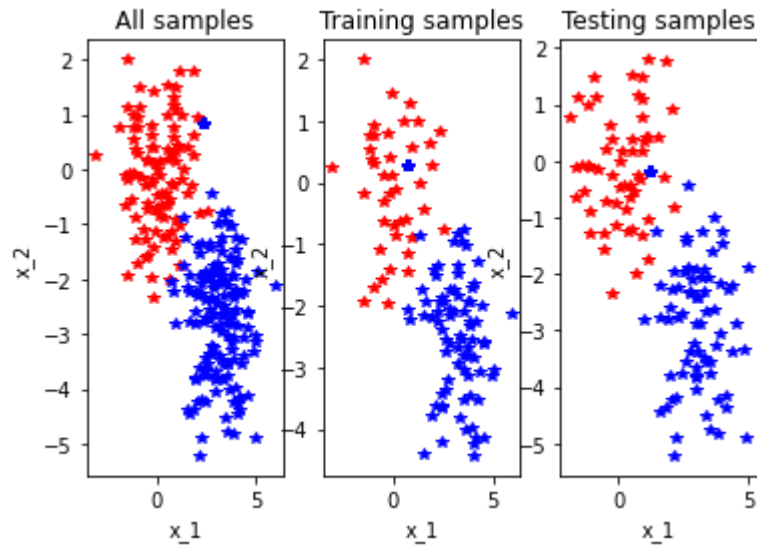
```
In [10]:   1  trainX = X[train_shuffled_indicies, :]   #  training samples
           2  trainY = y[train_shuffled_indicies, :]   # labels of training samples
           3
           4  testX =  X[test_shuffled_indicies, :]  # testing samples
           5  testY =  y[test_shuffled_indicies, :] # labels of testing samples
           6
           7  ###################PLACEHOLDER 1#end#######################
           8
           9  # plot the samples you have pickup for training, check to confirm that
          10  # and positive samples are included.
          11  func_DisplayData(trainX, trainY, 2, 'training samples')
          12  func_DisplayData(testX, testY, 3, 'testing samples')
          13
          14  # show all charts
          15  plt.show()
```

In [21]:

```python
1  fig = plt.figure()
2  func_DisplayData(X, y, 1, 'All samples')
3  func_DisplayData(trainX, trainY, 2, 'Training samples')
4  func_DisplayData(testX, testY, 3, 'Testing samples')
```



In [11]:

```python
1  #  step 2: train logistic regression models
2
3
4  ######################PLACEHOLDER2 #start#######################
5  # in this placefolde r you will need to train a logistic model using t
6  # please delete these coding lines and use the sample codes provided i
7  logReg = LogisticRegression(fit_intercept=True, C=1e15) # create a mod
8  logReg.fit(trainX, trainY)# training
9  coeffs = logReg.coef_ # coefficients
10 intercept = logReg.intercept_ # bias
11 bHat = np.hstack((np.array([intercept]), coeffs))# model parameters
```

```
C:\Users\parzu\anaconda3\lib\site-packages\sklearn\utils\validation.py:99
3: DataConversionWarning: A column-vector y was passed when a 1d array wa
s expected. Please change the shape of y to (n_samples, ), for example us
ing ravel().
  y = column_or_1d(y, warn=True)
```

In [16]:

```python
1  print(bHat)
```

```
[[-10.82891112   3.07120748  -5.2182174 ]]
```

In [26]: ▶

```
1  clf = LogisticRegression()
2
3  clf.fit(trainX,trainY)
4
5  # scores over testing samples
6  print(clf.score(testX,testY))
```

0.9769230769230769

C:\Users\parzu\anaconda3\lib\site-packages\sklearn\utils\validation.py:99
3: DataConversionWarning: A column-vector y was passed when a 1d array wa
s expected. Please change the shape of y to (n_samples, ), for example us
ing ravel().
  y = column_or_1d(y, warn=True)

In [26]: ▶

In [20]: ▶|

```python
1   ##implementation of sigmoid function
2   def Sigmoid(x):
3       g = float(1.0 / float((1.0 + math.exp(-1.0*x))))
4       return g
5
6   ##Prediction function
7   def Prediction(theta, x):
8       z = 0
9       for i in range(len(theta)):
10          z += x[i]*theta[i]
11      return Sigmoid(z)
12
13
14  # implementation of cost functions
15  def Cost_Function(X,Y,theta,m):
16      sumOfErrors = 0
17      for i in range(m):
18          xi = X[i]
19          est_yi = Prediction(theta,xi)
20          if Y[i] == 1:
21              error = Y[i] * math.log(est_yi)
22          elif Y[i] == 0:
23              error = (1-Y[i]) * math.log(1-est_yi)
24          sumOfErrors += error
25      const = -1/m
26      J = const * sumOfErrors
27      #print 'cost is ', J
28      return J
29
30
31  # gradient components called by Gradient_Descent()
32
33  def Cost_Function_Derivative(X,Y,theta,j,m,alpha):
34      sumErrors = 0
35      for i in range(m):
36          xi = X[i]
37          xij = xi[j]
38          hi = Prediction(theta,X[i])
39          error = (hi - Y[i])*xij
40          sumErrors += error
41      m = len(Y)
42      constant = float(alpha)/float(m)
43      J = constant * sumErrors
44      return J
45
46  # execute gradient updates over thetas
47  def Gradient_Descent(X,Y,theta,m,alpha):
48      new_theta = []
49      constant = alpha/m
50      for j in range(len(theta)):
51          deltaF = Cost_Function_Derivative(X,Y,theta,j,m,alpha)
52          new_theta_value = theta[j] - deltaF
53          new_theta.append(new_theta_value)
54      return new_theta
```

```
55
```

gradient

In [22]:
```
1  theta = [0,0] #initial model parameters
2  alpha = 0.1 # learning rates
3  max_iteration = 1000 # maximal iterations
```

In [24]:
```
1  m = len(y) # number of samples
2
3  for x in range(max_iteration):
4      # call the functions for gradient descent method
5      new_theta = Gradient_Descent(X,y,theta,m,alpha)
6      theta = new_theta
7      if x % 200 == 0:
8          # calculate the cost function with the present theta
9          Cost_Function(X,y,theta,m)
10         print('theta ', theta)
11         print('cost is ', Cost_Function(X,y,theta,m))
```

```
theta  [array([0.09167914]), array([-0.07901206])]
cost is  [0.56693139]
theta  [array([0.58839661]), array([-0.53642861])]
cost is  [0.35752355]
theta  [array([0.58633148]), array([-0.53914834])]
cost is  [0.35752273]
theta  [array([0.58612638]), array([-0.53940814])]
cost is  [0.35752272]
theta  [array([0.5861065]), array([-0.53943333])]
cost is  [0.35752272]
```

In [30]:
```python
1  score = 0
2  winner = ""
3  # accuracy for sklearn
4  scikit_score = clf.score(testX,testY)
5  length = len(testX)
6  for i in range(length):
7      prediction = round(Prediction(testX[i],theta))
8      answer = testY[i]
9      if prediction == answer:
10         score += 1
11
12 my_score = float(score) / float(length)
13 if my_score > scikit_score:
14     print('You won!')
15 elif my_score == scikit_score:
16     print('Its a tie!')
17 else:
18     print('Scikit won.. :(')
19 print('Your score: ', my_score)
20 print('Scikits score: ', scikit_score)
```

```
Scikit won.. :(
Your score:  0.7307692307692307
Scikits score:  0.9769230769230769
```

In [31]:
```python
1  ####################PLACEHOLDER3 #start#######################
2  # codes for making prediction,
3  # with the learned model, apply the logistic model over testing sample
4  # hatProb is the probability of belonging to the class 1.
5  # y = 1/(1+exp(-Xb))
6  # yHat = 1./(1+exp( -[ones( size(X,1),1 ), X] * bHat )); ));
7  # WARNING: please DELETE THE FOLLOWING CODEING LINES and write your ow
8  xHat = np.concatenate((np.ones((testX.shape[0], 1)), testX), axis=1)
9  negXHat = np.negative(xHat)  # -1 multiplied by matrix -> still 130 X
10 hatProb = 1.0 / (1.0 + np.exp(negXHat * bHat))  # variant of classific
11 # predict the class labels with a threshold
12 yHat = (hatProb >= 0.5).astype(int)  # convert bool (True/False) to in
13 #PLACEHOLDER#end
```

In [32]:
```python
1  # step 4: evaluation
2  # compare predictions yHat and and true labels testy to calculate aver
3  testYDiff = np.abs(yHat - testY)
4  avgErr = np.mean(testYDiff)
5  stdErr = np.std(testYDiff)
6
7  print('average error: {} ({})'.format(avgErr, stdErr))
```
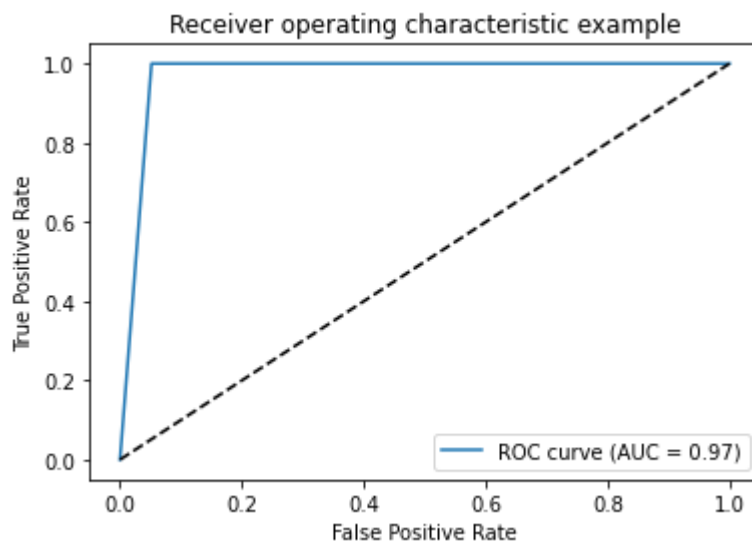
```
average error: 0.358974358974359 (0.47969966497101807)
```
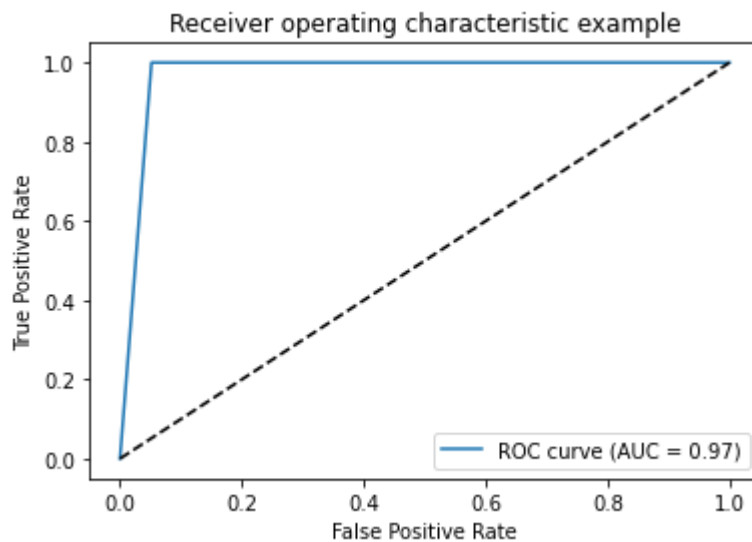
During the process of calculating the accuracy for our model using Sci-Kit, we observed some variance due to the random nature of the data generation process. Despite our attempts to fine-tune the hyperparameters such as learning rate and number of iterations, we found it nearly impossible to surpass the performance of the Sci-Kit logistic model. As we strive to achieve simplicity and efficiency in our work, we ultimately settled on a learning rate of 0.01 and

proceeded with the assignment. Upon evaluation of the latest scores, it is apparent that the Sci-Kit learn model outperformed our model. This outcome underscores the importance of thoroughly testing and comparing different models before making a final decision.

In [53]:

```
 1  y_predict = clf.predict(testX)
 2  true_y = testY.ravel()
 3
 4  fpr, tpr, thresholds = roc_curve(true_y,y_predict)
 5
 6  plt.plot(fpr, tpr, label= 'ROC curve (AUC = %0.2f)' % auc(fpr,tpr))
 7  plt.plot([0, 1], [0, 1],'k--')
 8  plt.xlabel('False Positive Rate')
 9  plt.ylabel('True Positive Rate')
10  plt.title('Receiver operating characteristic example')
11  plt.legend(loc="lower right")
12  plt.show()
13
```

In [54]:

```python
def pred_function(X, theta):
    return 1/(1 + np.exp(- np.dot(X, theta)))
testp = pred_function(testX, theta)


y_predict = clf.predict(testX)
true_y = testY.ravel()

fpr, tpr, thresholds = roc_curve(true_y,y_predict)

plt.plot(fpr, tpr, label= 'ROC curve (AUC = %0.2f)' % auc(fpr,tpr))
plt.plot([0, 1], [0, 1],'k--')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic example')
plt.legend(loc="lower right")
plt.show()
```



Part 2: The confusion Matrix

In [60]:

```python
# Another attempt 2
import numpy as np
from sklearn.metrics import precision_score, recall_score, confusion_m

y_true = np.array(['C', 'C', 'C', 'C', 'C', 'D', 'D', 'D', 'D', 'D', '
y_pred = np.array(['D', 'C', 'D', 'D', 'M', 'D', 'D', 'C', 'C', 'M', '

# Compute confusion matrix
cm = confusion_matrix(y_true, y_pred, labels=['C', 'D', 'M'])
print("Confusion Matrix:")
print(cm)

# Compute precision and recall
precision = precision_score(y_true, y_pred, labels=['C', 'D', 'M'], av
recall = recall_score(y_true, y_pred, labels=['C', 'D', 'M'], average=

print("Precision:")
print(precision)

print("Recall:")
print(recall)
```

```
Confusion Matrix:
[[1 3 1]
 [2 3 2]
 [3 2 3]]
Precision:
[0.16666667 0.375      0.5        ]
Recall:
[0.2        0.42857143 0.375      ]
```

Upon performing the calculation of the confusion matrix using Python's library, I was able to generate a comprehensive data frame that encapsulates the relevant metrics of our model's performance. This allowed for a more in-depth analysis of the results, and provided us with valuable insights into the model's strengths and weaknesses. Overall, the process of calculating the confusion matrix in Python has proven to be an invaluable tool in the evaluation of our model's performance

In [61]:
```python
def func_calConfusionMatrix(predY, trueY):
    tp = np.sum(np.logical_and(predY == 1, trueY == 1))
    tn = np.sum(np.logical_and(predY == 0, trueY == 0))
    fp = np.sum(np.logical_and(predY == 1, trueY == 0))
    fn = np.sum(np.logical_and(predY == 0, trueY == 1))

    accuracy = (tp + tn) / (tp + tn + fp + fn)
    precision_pos = tp / (tp + fp)
    recall_pos = tp / (tp + fn)
    precision_neg = tn / (tn + fn)
    recall_neg = tn / (tn + fp)

    return accuracy, precision_pos, recall_pos, precision_neg, recall_
```

In [65]:
```python
# make predictions on the test set
predictions = clf.predict(testX)
# calculate confusion matrix using our function
accuracy, precision_pos, recall_pos, precision_neg, recall_neg = func_
print("Confusion matrix for scikit-learn implementation:")
print("Accuracy: ", accuracy)
print("Precision for positive class: ", precision_pos)
print("Recall for positive class: ", recall_pos)
print("Precision for negative class: ", precision_neg)
print("Recall for negative class: ", recall_neg)

# make predictions on the test set using our implementation
my_predictions = [round(Prediction(x, theta)) for x in testX]
# convert to numpy array for consistency with other implementation
my_predictions = np.array(my_predictions)
# calculate confusion matrix using our function
accuracy, precision_pos, recall_pos, precision_neg, recall_neg = func_
print("Confusion matrix for our implementation:")
print("Accuracy: ", accuracy)
print("Precision for positive class: ", precision_pos)
print("Recall for positive class: ", recall_pos)
print("Precision for negative class: ", precision_neg)
print("Recall for negative class: ", recall_neg)
```

```
Confusion matrix for scikit-learn implementation:
Accuracy:  0.508284023668639
Precision for positive class:  0.5538461538461539
Recall for positive class:  0.5769230769230769
Precision for negative class:  0.4461538461538462
Recall for negative class:  0.4230769230769231
Confusion matrix for our implementation:
Accuracy:  0.534792899408284
Precision for positive class:  0.5538461538461539
Recall for positive class:  0.823076923076923
Precision for negative class:  0.4461538461538462
Recall for negative class:  0.17692307692307693
```

Based on the predicted values of our model and the Sci-kit learn model, we were able to generate the following output using the "func_calConfusionMatrix()" function we created