

## ✓ Logistic Regression

### CS 4650 "Natural Language Processing" - Project 0

Georgia Tech, Spring 2025

(Instructor: Prof. Wei Xu; TAs: Yao Dou, Tarek Naous, Xiaofeng Wu, Jonathan Zheng)

In this assignment, we will walk you through the process of implementing logistic regression from scratch. You will also apply your implemented logistic regression model to a small dataset and predict whether a student will be admitted to a university. This dataset will allow you to visualize the data and debug more easily. You may find [this documentation](#) very helpful, though it is about how to implement logistic regression in Octave, other than Python.

This assignment also serves as a programming preparation test. We will use [Numpy](#) – a popular Python package for scientific computing and implementing machine learning algorithms. It provides very good support for matrix and vector operations. You need to feel comfortable working with matrices, vectors, and tensors in order to complete all the programming projects in CS 4650.

**IMPORTANT:** In this assignment, except Numpy and Matplotlib, no other external Python packages are allowed. Scipy package can be used in gradient checking, though, it is not allowed elsewhere.

## 0. Honor Code [1 points]

**Honor Code:** I hereby agree to abide the Georgia Tech's Academic Honor Code, promise that the submitted assignment is my own work, and understand that my code is subject to plagiarism test.

**Signature:** *(double click on this block and type your name here)*

## ✓ 1. Importing Numpy and Matplotlib [Code provided - do not change]

```
1 import sys
2
3 # Check what version of Python is running
4 print (sys.version)
```

3.11.11 (main, Dec 4 2024, 08:55:07) [GCC 11.4.0]

We will also import [Matplotlib](#), a Python package for data visualization.

```

1 # Run some setup code for this notebook. Don't modify anything in this cell.
2
3 import random
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 # This is a bit of magic to make matplotlib figures appear inline in the notebook
8 # rather than in a new window.
9 %matplotlib inline
10 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
11 plt.rcParams['image.interpolation'] = 'nearest'
12 plt.rcParams['image.cmap'] = 'gray'
13
14 # reload external python modules;
15 # http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
16 %load_ext autoreload
17 %autoreload 2

```

## ✓ 2. Visualizing the Data [Code provided - no need to change]

The provided dataset contains applicants' scores on two exams and the admission decisions for 100 students. This dataset will allow us to visualize in a 2D figure and showcase how the logistic regression algorithm works more intuitively.

```
1 !wget https://raw.githubusercontent.com/cocoxu/CS4650\_spring2025\_projects/refs/heads/main/p0\_data.txt
```

```

--2025-01-16 21:34:26-- https://raw.githubusercontent.com/cocoxu/CS4650\_spring2025\_projects/refs/heads/main/p0\_data.txt
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.1
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.1:
HTTP request sent, awaiting response... 200 OK
Length: 3775 (3.7K) [text/plain]
Saving to: 'p0_data.txt'

```

```
p0_data.txt          100%[=====>]    3.69K  --.-KB/s    in 0s
```

```
2025-01-16 21:34:26 (45.7 MB/s) - 'p0_data.txt' saved [3775/3775]
```

```


1 #load the dataset
2 data = np.loadtxt('p0_data.txt', delimiter=',')
3
4 train_X = data[:, 0:2]
5 train_y = data[:, 2]
6
7 # Get the number of training examples and the number of features
8 m_samples, n_features = train_X.shape
9 print("# of training examples = ", m_samples)
10 print("# of features = ", n_features)

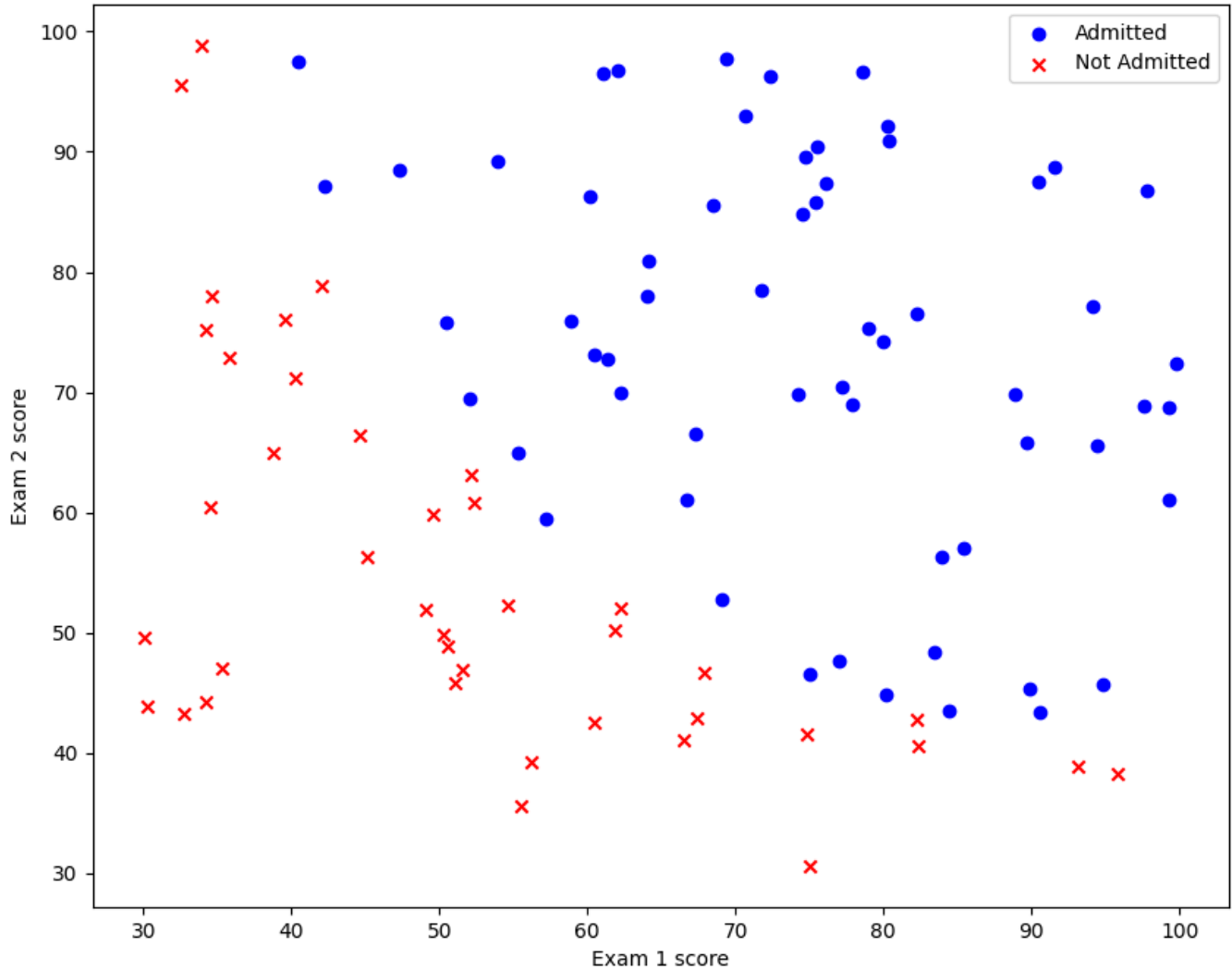
```

```

10 print ("# of features = ", n_features)
11
12 pos = np.where(train_y == 1)
13 neg = np.where(train_y == 0)
14 plt.scatter(train_X[pos, 0], train_X[pos, 1], marker='o', c='b')
15 plt.scatter(train_X[neg, 0], train_X[neg, 1], marker='x', c='r')
16 plt.xlabel('Exam 1 score')
17 plt.ylabel('Exam 2 score')
18 plt.legend(['Admitted', 'Not Admitted'])
19 plt.show()

```

 # of training examples = 100  
 # of features = 2



### ✓ 3. Cost Function [5 points]

You're going to first implement the sigmoid function, then the cost function for (binary) logistic regression.

The sigmoid function is defined as  $\text{sigmoid}(\mathbf{z}) = \frac{1}{1+e^{-z}}$ . It is important to handle potential underflow or overflow in the sigmoid implementation. It is also important to not take a log of 0.

Note that, you are asked to use the [Numpy](#) package for vector and matrix operations in order to ensure the **efficiency of the code**.

```

1 def sigmoid(z):
2     """ Sigmoid function """
3     #####
4     # Compute the sigmoid function for the input here.          #
5     #####
6
7     s = None
8     ### YOUR CODE HERE: be careful of the potential underflow or overflow here
9
10    if z.all() >= 0:
11        s = 1 / (1 + np.exp(-z))
12    else:
13        s = np.exp(z) / (1 + np.exp(z))
14
15    ### END YOUR CODE
16
17    return s
18
19 # Check your sigmoid implementation
20 z = np.array([[1, 2], [-1, -2]])
21 f = sigmoid(z)
22 print (f)

```

```

⇒ [[0.73105858 0.88079708]
   [0.26894142 0.11920292]]

```

```

1 def cost_function(theta, X, y):
2     """ The cost function for logistic regression """
3     #####
4     # Compute the cost given the current parameter theta on the training data set
5     #####
6
7     cost = None
8     ### YOUR CODE HERE
9
10    m = X.shape[0]
11    z = np.dot(X, theta)
12    h = sigmoid(z)
13    cost = (1.0 / m) * np.sum(-y * np.log(h) - (1 - y) * np.log(1 - h))
14
15    ### END YOUR CODE
16
17    return cost

```

```

18
19 # Check your cost function implementation
20
21 t_X = np.array([[1, 2], [-1, -2]])
22 t_y = np.array([0, 1])
23 t_theta1 = np.array([-10, 10])
24 t_theta2 = np.array([10, -10])
25 t_c1 = cost_function(t_theta1, t_X, t_y)
26 t_c2 = cost_function(t_theta2, t_X, t_y)
27 print (t_c1)
28 print (t_c2)

```

```

↩ 10.000045398899701
  4.539889921682063e-05

```

## ✓ 4. Gradient Computation [5 points]

Implement the gradient computations for logistic regression.

```

1 def gradient_update(theta, X, y):
2     """ The gradient update for logistic regression"""
3     #####
4     # Compute the gradient update #
5     #####
6
7     grad = None
8     ### YOUR CODE HERE
9     m = X.shape[0]
10    h = sigmoid(np.dot(X, theta))
11    grad = (1.0 / m) * np.dot(X.T, h - y)
12    ### END YOUR CODE
13
14    return grad
15
16 # Check your gradient computation implementation
17 t_X = np.array([[1, 2, 3], [-1, -2, -3]])
18 t_y = np.array([0, 1])
19 t_theta1 = np.array([-10, 10, 0])
20 t_theta2 = np.array([10, -10, 0])
21 t_g1 = gradient_update(t_theta1, t_X, t_y)
22 t_g2 = gradient_update(t_theta2, t_X, t_y)
23 print (t_g1)
24 print (t_g2)

```

```

↩ [0.9999546  1.9999092  2.99986381]
  [4.53978687e-05  9.07957374e-05  1.36193606e-04]

```

## ✓ 5. Gradient Checking [Code provided. Bonus 5 points if implemented from scratch]

You can use the code provided below to check the gradient of your logistic regression functions using [Scipy](#) package. Alternatively, you can implement the gradient checking from scratch by yourself (bonus 5 points). If you attempt the bonus, your implementation should replicate the behavior Scipy's implementation. Note: Copying Scipy's implementation does not count.

[Gradient checking](#) is an important technique for debugging the gradient computation. Logistic regression is a relatively simple algorithm where it is straightforward to derive and implement its cost function and gradient computation. For more complex models, the gradient computation can be notoriously difficult to debug and get right. Sometimes a subtly buggy implementation will manage to learn something that can look surprisingly reasonable, while performing less well than a correct implementation. Thus, even with a buggy implementation, it may not at all be apparent that anything is amiss.

```
1 # Check your gradient computation implementation
2 t_samples, t_features = 100, 10
3 t_X = np.random.randn(t_samples, t_features)
4 t_y = np.random.randint(2, size=t_samples)
5 t_theta = np.random.randn(t_features)
6
7 from scipy import optimize
8 print('Output of check_grad: %s' % optimize.check_grad(cost_function, gradient_upd,
```

↗ Output of check\_grad: 5.911868185637387e-07

## ✓ 6. Gradient Descent and Decision Boundary [10 points]

Implement the batch gradient descent algorithm for logistic regression. For every 'print\_iterations' number of iterations, also visualize the decision boundary and observe how it changes during the training. Please print the change between **10-20** times to fully demonstrate the learned decision boundary, along with the allowing enough iterations for the method to converge.

Please use the *x\_axis\_range* variable initialized at the end of the plotting code (on line 55 of the block below if line numbers are enabled) for the range of the x-axis when plotting.

Note that, you will need to carefully choose the learning rate and the total number of iterations (hint: without feature scaling, it will need a small learning rate and a large number of iterations), especially given that the starter code does not include feature scaling (e.g., scale each feature by its

maximum absolute value to convert feature value to  $[-1,1]$  range -- in order to make this homework simple and easier for you to write code to visualize.

```

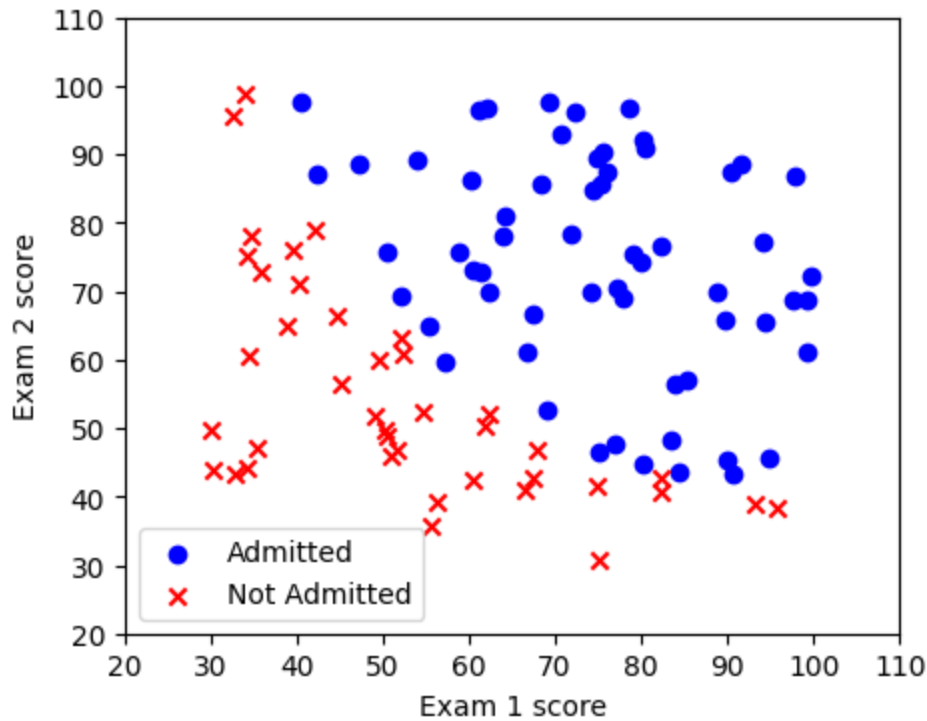
1 def gradient_descent(theta, X, y, alpha, max_iterations, print_iterations):
2     """ Batch gradient descent algorithm """
3     #####
4     # Update the parameter 'theta' iteratively to minimize the cost #
5     # 'alpha' is learning rate.                                     #
6     # Also visualize the decision boundary during learning          #
7     #####
8
9     #alpha *= m_samples
10    iteration = 0
11
12    ### YOUR CODE HERE: handle bias term, i.e., adding x0=1 into the X
13    ### If doing feature scaling (not required in this homework), this will be the
14    m_samples = X.shape[0]
15    X = np.hstack([ np.ones((m_samples, 1)), X ])
16
17    ### END YOUR CODE
18
19
20
21    while(iteration < max_iterations):
22        iteration += 1
23
24        ### YOUR CODE HERE: simultaneous update of partial gradients
25
26        grad = gradient_update(theta, X, y)
27        theta -= alpha * grad
28
29        ### END YOUR CODE
30
31
32        # For first iteration and every print_iterations
33        if iteration % print_iterations == 0 or iteration == 1:
34            cost = 0
35
36            ### YOUR CODE HERE: calculate the cost
37            ### IMPORTANT: The cost function is guaranteed to decrease after
38            ## every iteration of the gradient descent algorithm.
39
40            cost = cost_function(theta, X, y)
41
42            ### END YOUR CODE
43
44            print ("[ Iteration", iteration, "]", "cost =", cost)
45            plt.rcParams['figure.figsize'] = (5, 4)
46            plt.xlim([20,110])
47            plt.ylim([20,110])
48

```

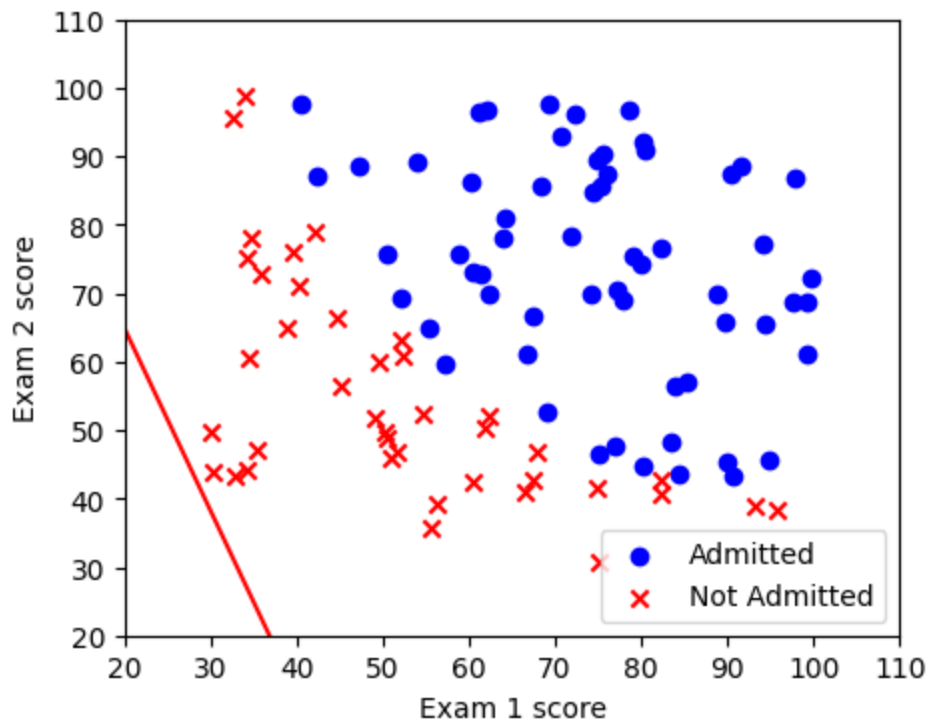
```
49     pos = np.where(y == 1)
50     neg = np.where(y == 0)
51
52     plt.scatter(X[pos, 1], X[pos, 2], marker='o', c='b')
53     plt.scatter(X[neg, 1], X[neg, 2], marker='x', c='r')
54     plt.xlabel('Exam 1 score')
55     plt.ylabel('Exam 2 score')
56     plt.legend(['Admitted', 'Not Admitted'])
57     x_axis_range = np.arange(10, 100, 0.1)
58
59
60     ### YOUR CODE HERE: plot the decision boundary using the variable 'x_a
61     ### try "plt.plot(x_axis_range, (x_axis_range * 3 - 10), '-r')" for an
62
63     boundary_y = -(theta[0] + theta[1] * x_axis_range) / theta[2]
64     plt.plot(x_axis_range, boundary_y, '-r')
65
66     ### END YOUR CODE
67
68     plt.show()
69
70     return theta
71
72
73 ### YOUR CODE HERE: initialize the parameters 'theta' to random values;
74 ### And set up learning rate, number of max iterations, number of iterations for p
75
76
77 initial_theta = np.zeros(3)
78 alpha_test = 0.001
79 max_iter = 100000
80 print_iter = 10000
81
82
83 ### END YOUR CODE
84
85
86 learned_theta = gradient_descent(initial_theta, train_X, train_y, alpha_test, max_
```



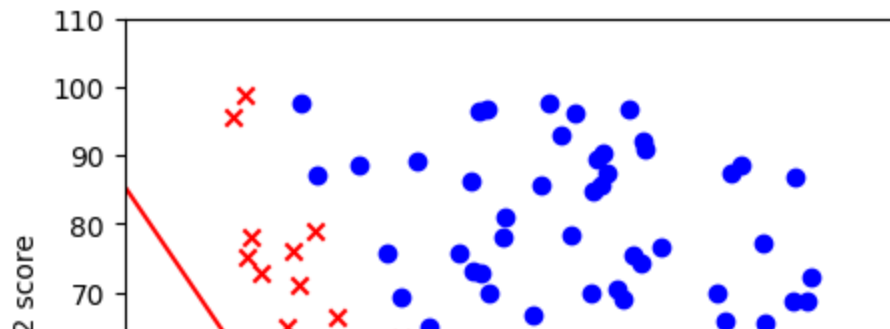
⇒ [ Iteration 1 ] cost = 0.6982906893667754

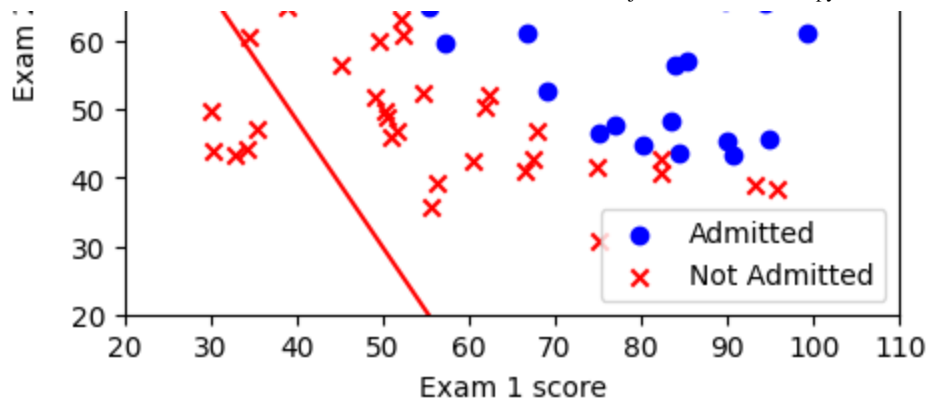


[ Iteration 10000 ] cost = 0.5850274988176747

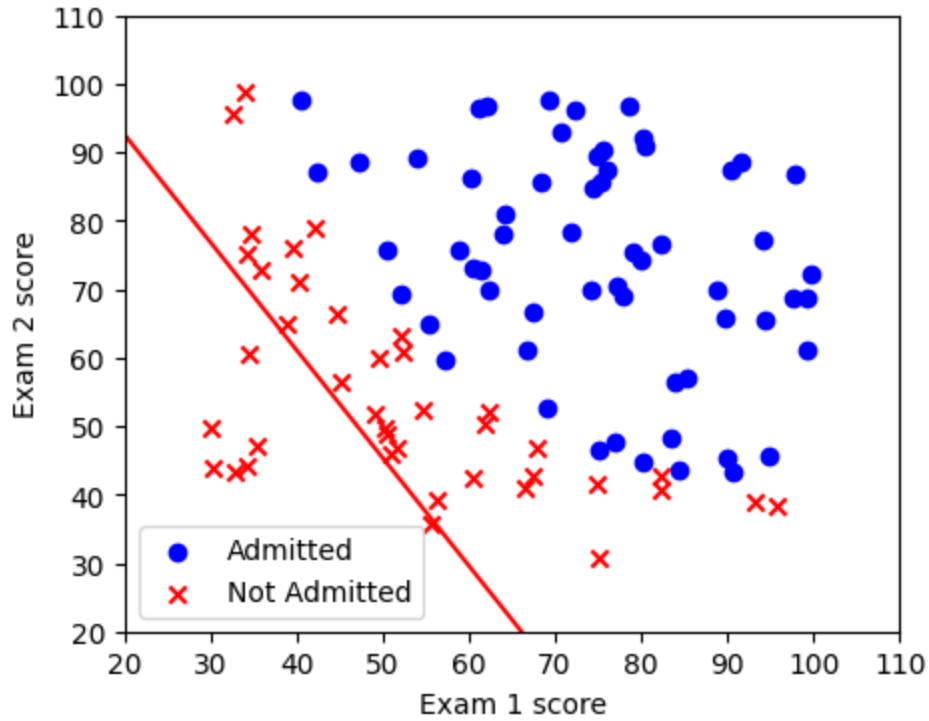


[ Iteration 20000 ] cost = 0.5472954636936678

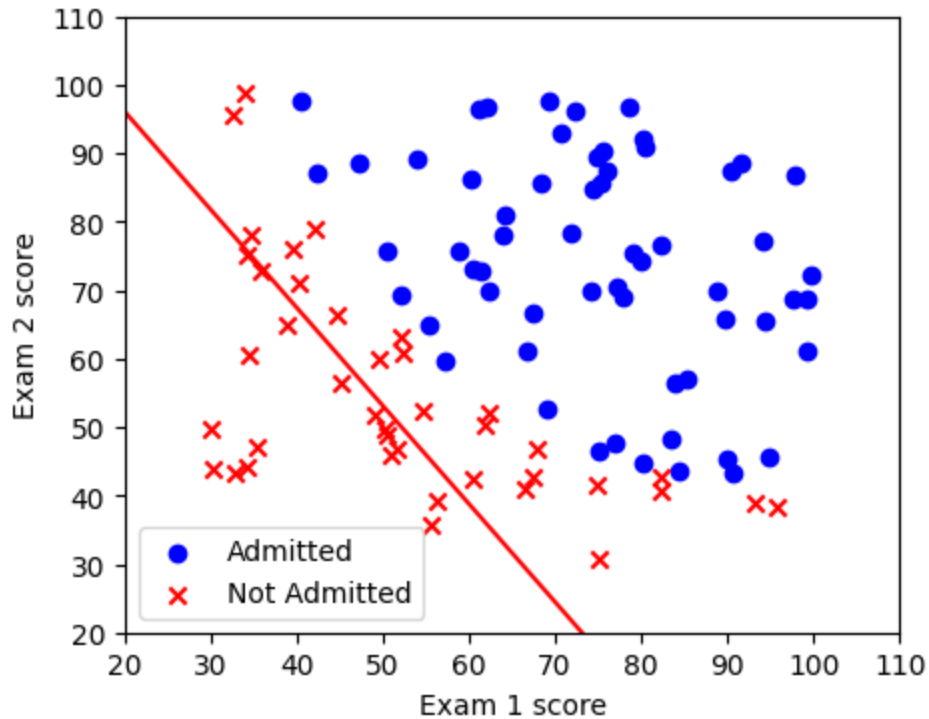




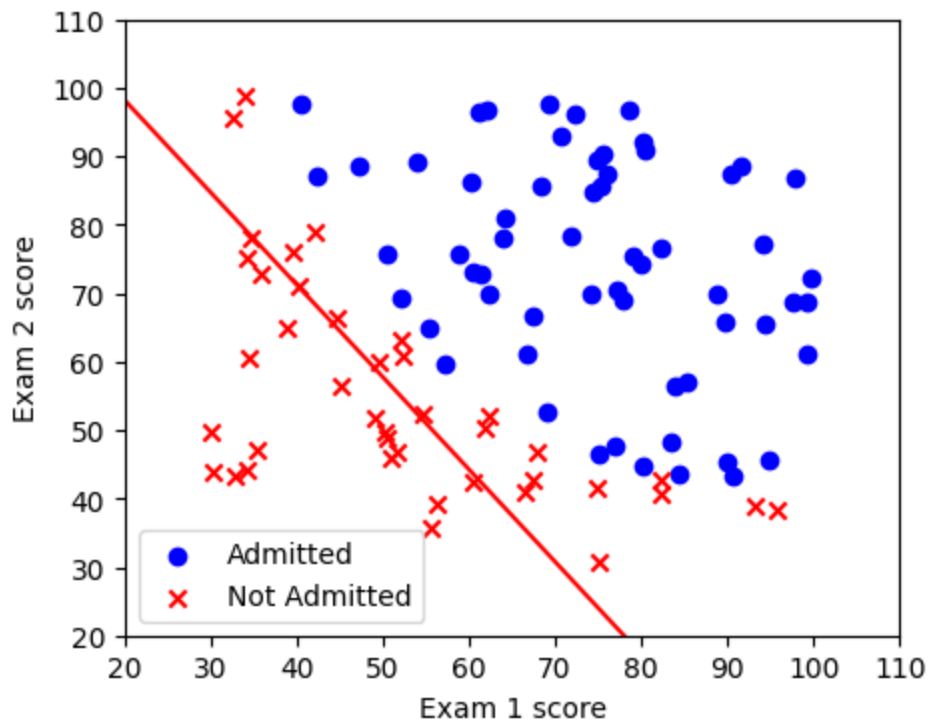
[ Iteration 30000 ] cost = 0.5154018233277201



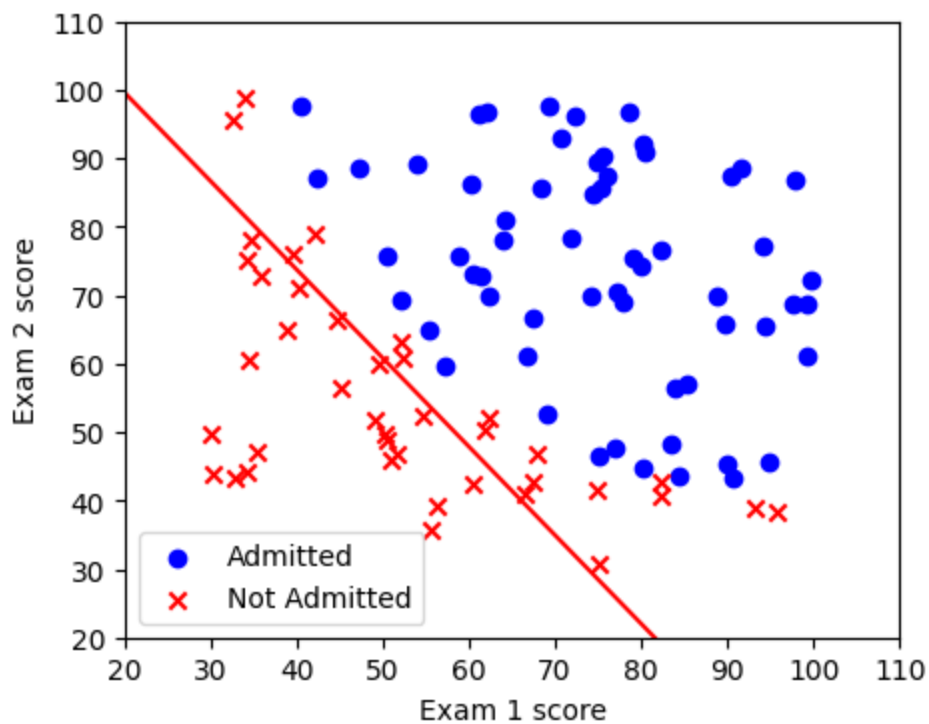
[ Iteration 40000 ] cost = 0.48829570918772613



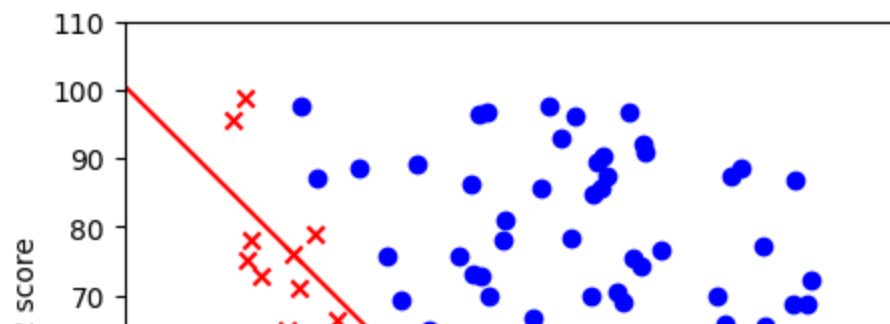
[ Iteration 50000 ] cost = 0.46510456912211945

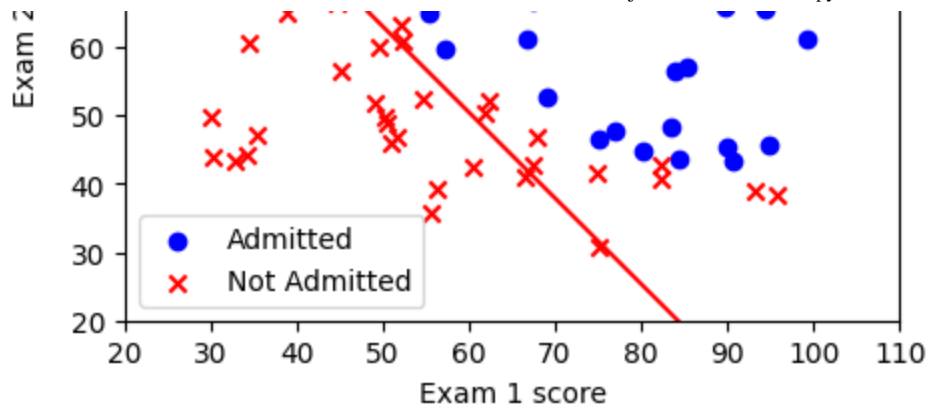


[ Iteration 60000 ] cost = 0.445118809496071

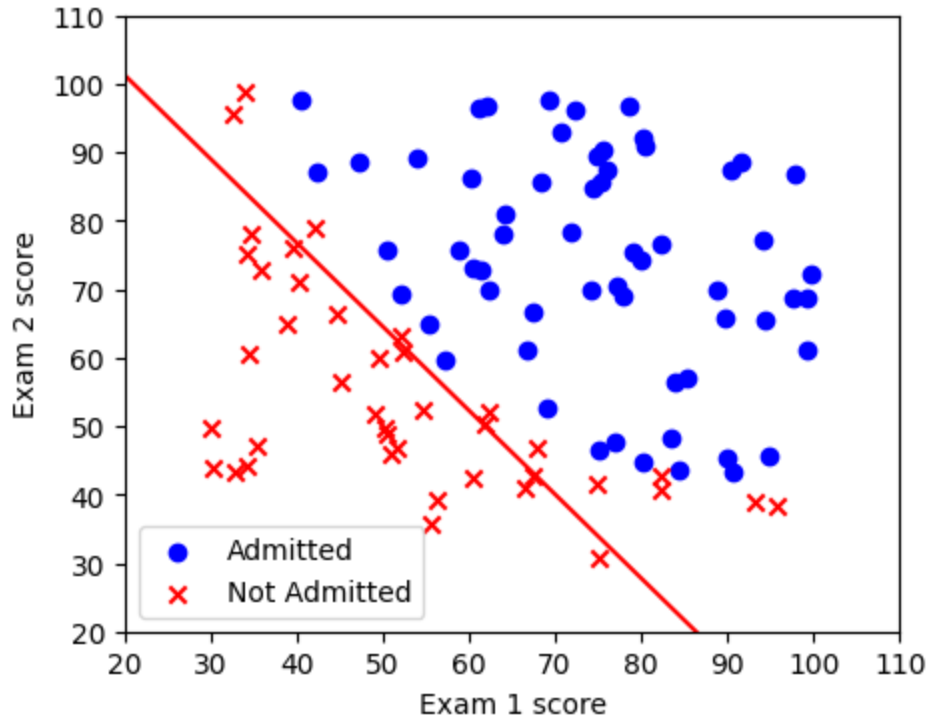


[ Iteration 70000 ] cost = 0.42776773553739156

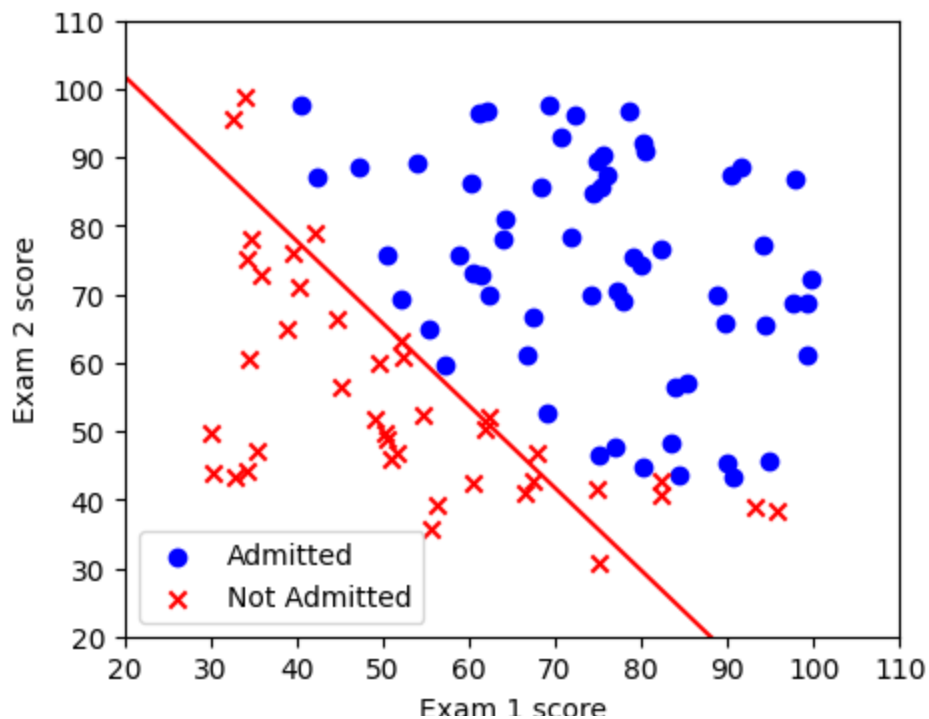




[ Iteration 80000 ] cost = 0.41259440170978545



[ Iteration 90000 ] cost = 0.39923293458511083



A scatter plot showing the relationship between Exam 1 score (X-axis) and Exam 2 score (Y-axis). The X-axis ranges from 20 to 110, and the Y-axis ranges from 20 to 110. Data points are categorized into two groups: 'Admitted' (blue circles) and 'Not Admitted' (red crosses). A red line represents the decision boundary, separating the two groups. The line starts at approximately (20, 102) and ends at (100, 20). Most 'Admitted' points are located above and to the right of the line, while most 'Not Admitted' points are below and to the left of the line. There is a small cluster of 'Not Admitted' points near the top-left corner, and a few 'Admitted' points near the bottom-right corner.

## ✓ 6. Predicting [5 points]

Now that you learned the parameters of the model, you can use the model to predict whether a particular student will be admitted.

```

1 def predict(theta, X):
2     """ Predict label using learned logistic regression parameters """
3     #####
4     # Predict a label of 0 or 1 using learned logistic regression parameters for
5     # Return the probabilities of X as a (1, N) array of floats
6     # Return the predicted_labels as a (1, N) array of 0 or 1 integers
7     #####
8
9     ### YOUR CODE HERE:
10
11     m_samples = X.shape[0]
12     X = np.hstack([np.ones((m_samples, 1)), X])
13     probabilities = sigmoid(np.dot(X, theta))
14     predicted_labels = (probabilities >= 0.5).astype(int)
15
16     ### END YOUR CODE
17
18
19     ## convert an array of booleans 'predicted_labels' into an array of 0 or 1 in
20     return probabilities, 1*predicted_labels
21
22 # Check your predication function implementation
23 t_X1 = np.array([[90, 90]])
24 t_X2 = np.array([[50, 60]])
25 t_X3 = np.array([[10, 50]])
26 print (predict(learned_theta, t_X1))
27 print (predict(learned_theta, t_X2))
28 print (predict(learned_theta, t_X3))
29
30 # Compute accuracy on the training dataset
31 t_prob, t_label = predict(learned_theta, train_X)
32 t_precision = t_label[np.where(t_label == train_y)].size / float(train_y.size) *
33 print('Accuracy on the training set: %s%%' % round(t_precision,2))

```

 (array([0.93706746]), array([1]))  
 (array([0.43627873]), array([0]))  
 (array([0.07948104]), array([0]))  
 Accuracy on the training set: 91.0%

## ✓ 7. Submit Your Homework

This is the end. Congratulations!