

Problem Set 2: Classification

Problem Set 2: Classification

To run and solve this assignment, one must have a working IPython Notebook installation. The easiest way to set it up for both Windows and Linux is to install [Anaconda](https://www.continuum.io/downloads) (<https://www.continuum.io/downloads>). Then save this file to your computer (use "Raw" link on gist\github), run Anaconda and choose this file in Anaconda's file explorer. Use `Python 3` version. Below statements assume that you have already followed these instructions. If you are new to Python or its scientific library, Numpy, there are some nice tutorials [here](https://www.learnpython.org/) (<https://www.learnpython.org/>) and [here](http://www.scipy-lectures.org/) (<http://www.scipy-lectures.org/>).

To run code in a cell or to render [Markdown](https://en.wikipedia.org/wiki/Markdown) (<https://en.wikipedia.org/wiki/Markdown>)+[LaTeX](https://en.wikipedia.org/wiki/LaTeX) (<https://en.wikipedia.org/wiki/LaTeX>) press `Ctrl+Enter` or `[>|]` (like "play") button above. To edit any code or text cell [double]click on its content. To change cell type, choose "Markdown" or "Code" in the drop-down menu above.

If a certain output is given for some cells, that means that you are expected to get similar results in order to receive full points (small deviations are fine). For some parts we have already written the code for you. You should read it closely and understand what it does.

Total: 100 points.

1. Logistic Regression

In this part of the exercise, you will build a logistic regression model to predict whether a student gets admitted into a university.

Suppose that you are the administrator of a university department and you want to determine each applicant's chance of admission based on their results on two exams. You have historical data from previous applicants in *ex2data1.txt* that you can use as a training set for logistic regression. For each training example, you have the applicant's scores on two exams and the admissions decision.

Your task is to build a classification model that estimates an applicant's probability of admission based on the scores from those two exams. This outline and code framework will guide you through the exercise.

1.1 Implementation

```
In [1]: import sys
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
print('Tested with:')
print('Python', sys.version)
print({x.__name__: x.__version__ for x in [np, matplotlib]})
```

Tested with:

```
Python 3.6.5 |Anaconda, Inc.| (default, Apr 26 2018, 08:42:37)
[GCC 4.2.1 Compatible Clang 4.0.1 (tags/RELEASE_401/final)]
{'numpy': '1.14.3', 'matplotlib': '2.2.2'}
```

1.1.1 Visualizing the data

Before starting to implement any learning algorithm, it is always good to visualize the data if possible. This first part of the code will load the data and display it on a 2-dimensional plot by calling the function `plotData`. The axes are the two exam scores, and the positive and negative examples are shown with different markers.

```
In [2]: #####
# Try to fit your code and comments into 80 charecters because
# - it is guaranteed to look as intened on any screen size
# - it encourages you to write "flater" logic that is easier to reason about
# - it encourages you to decompose logic into comprehansible blocks.
#
# Try to avoid reassinging/mutating variables because when you encounter an
# unexplainable bugs (you will) it is easier to have the whole history
# of values to reason about.
#
# Using %debug magic to run pdb might be useful for debugging.
# It is just like gdb but for Python.
```

```

In [3]: # it is good to isolate logical parts to avoid variables leaking into the
# global scope and messing up your logic later in weird ways

def read_classification_csv_data(fn, add_ones=False):
    # read comma separated data
    data = np.loadtxt(fn, delimiter=',')
    X_, y_ = data[:, :-1], data[:, -1, None] # a fast way to keep last dim

    # printing statistics of data before working with it might have saved
    # hundreds hours of of my time, do not repeat my errors :)
    print(X_.shape, X_.min(), X_.max(), X_.dtype)
    print(y_.shape, y_.min(), y_.max(), y_.dtype)
    # aha, y is float! this is not what we expected
    # what might go wrong with further y == 0 checks?
    # A: floating point equality comparison, that's what!

    # insert the column of 1's into the "X" matrix (for bias)
    X = np.insert(X_, X_.shape[1], 1, axis=1) if add_ones else X_
    y = y_.astype(np.int32)
    return X, y

X_data, y_data = read_classification_csv_data('ex2data1.txt', add_ones=True)
print(X_data.shape, X_data.min(), X_data.max(), X_data.dtype)
print(y_data.shape, y_data.min(), y_data.max(), y_data.dtype)

(100, 2) 30.05882244669796 99.82785779692128 float64
(100, 1) 0.0 1.0 float64
(100, 3) 1.0 99.82785779692128 float64
(100, 1) 0 1 int32

```

```

In [4]: # how does the *X[y.ravel()==1, :2].T trick work?
# https://docs.python.org/3/tutorial/controlflow.html#unpacking-argument-list
%matplotlib inline

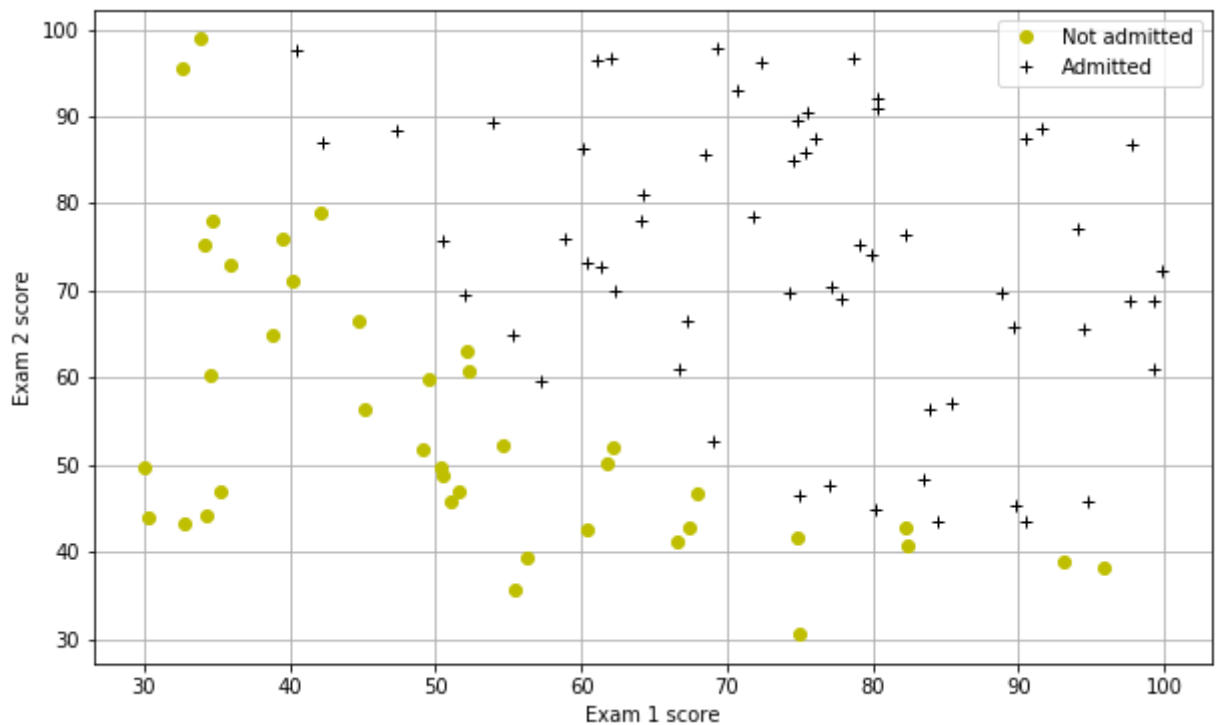
def plot_data(X, y, labels, markers, xlabel, ylabel, figsize=(10, 6), ax=None):
    if figsize is not None:
        plt.figure(figsize=figsize)

    ax = ax or plt.gca()
    for label_id, (label, marker) in enumerate(zip(labels, markers)):
        # ax.plot(*X_data[y_data.ravel()==label_id, :2].T, marker, label=label)
        ax.plot(*X[y.ravel()==label_id, :2].T, marker, label=label)
    ax.set_xlabel(xlabel)
    ax.set_ylabel(ylabel)
    plt.legend()
    ax.grid(True)

student_plotting_spec = {
    'X': X_data,
    'y': y_data,
    'xlabel': 'Exam 1 score',
    'ylabel': 'Exam 2 score',
    'labels': ['Not admitted', 'Admitted'],
    'markers': ['yo', 'k+'],
    'figsize': (10, 6)
}

plot_data(**student_plotting_spec)
plt.show()

```



1.1.2 [5pts] Sigmoid function

Before you start with the actual cost function, recall that the logistic regression hypothesis is defined as:

$$h_{\theta}(x) = g(\theta^T x)$$

where function g is the sigmoid function. The sigmoid function is defined as:

$$g(z) = \frac{1}{1 + e^{-z}}$$

Your first step is to implement/find a sigmoid function so it can be called by the rest of your program. Your code should also work with vectors and matrices. For a matrix, your function should perform the sigmoid function on every element.

When you are finished, (a) plot the sigmoid function, and (b) test the function with a scalar, a vector, and a matrix. For scalar large positive values of x , the sigmoid should be close to 1, while for scalar large negative values, the sigmoid should be close to 0. Evaluating $\text{sigmoid}(0)$ should give you exactly 0.5.

```

In [5]: # check out scipy.special for great variety of vectorized functions
# remember that sigmoid is the inverse of logit function
# maybe worth checking out scipy.special.logit first

from scipy.stats import logistic

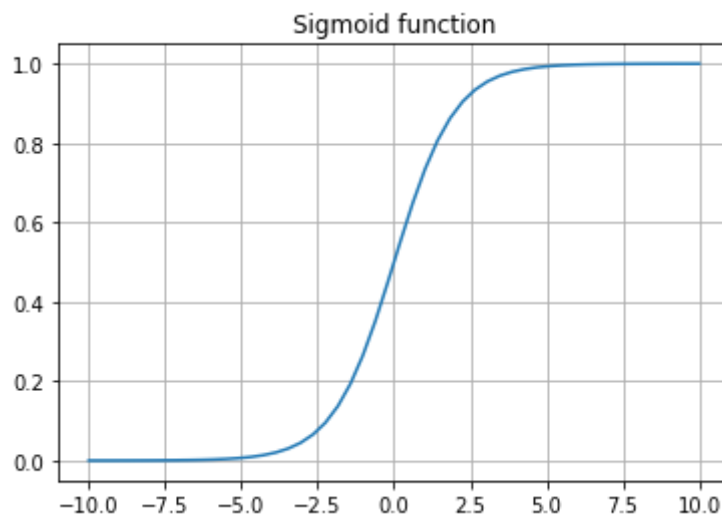
def sigmoid(z):
    return logistic.cdf(z)

def check_that_sigmoid_f(f):
    # don't use np.arange with float step because it works as
    # val_{i+1} = val_i + step while val_i < end
    # what might do wrong with float precision?
    x_test = np.linspace(-10, 10, 50)
    sigm_test = f(x_test)
    plt.plot(x_test, sigm_test)
    plt.title("Sigmoid function")
    plt.grid(True)
    plt.show()

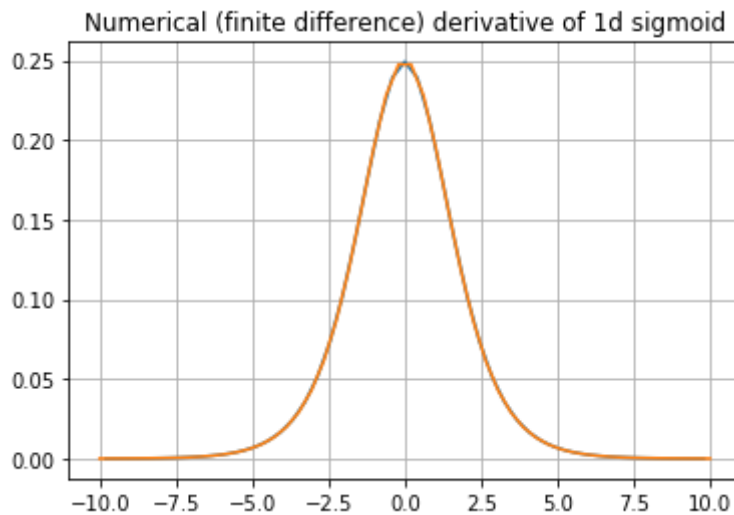
    # why should analytical_diff almost== finite_diff for sigmoid?
    analytical_diff = sigm_test*(1-sigm_test)
    finite_step = x_test[1]-x_test[0]
    finite_diff = np.diff(sigm_test) / finite_step
    print(x_test.shape, finite_diff.shape)
    plt.plot(x_test[:-1]+finite_step/2, finite_diff)
    plt.plot(x_test, analytical_diff)
    plt.title("Numerical (finite difference) derivative of 1d sigmoid")
    plt.grid(True)
    plt.show()

check_that_sigmoid_f(sigmoid)

```



(50,) (49,)



1.1.3 [15pts] Cost function and gradient

Now you will implement the cost function and gradient for logistic regression. Complete the code in the functions *hypothesis_function* and *binary_logistic_loss* below to return the value of the hypothesis function and the cost, respectively. Recall that the cost function in logistic regression is

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

and the gradient of the cost is a vector of the same length as θ where the j^{th} element (for $j = 0, 1, \dots, n$) is defined as follows:

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

where m is the number of points and n is the number of features. Note that while this gradient looks identical to the linear regression gradient, the formula is actually different because linear and logistic regression have different definitions of $h_{\theta}(x)$.

What should be the value of the loss for $\theta = \bar{0}$ regardless of input? Why? Make sure your code also outputs this value.

```

In [6]: # we are trying to fit a function that would return a
# "probability of "

# hyposesis_function describes parametric family of functions that we are
# going to pick our "best fitting function" from. It is parameterized by
# real-valued vector theta, i.e. we are going to pick
# h_best = argmin_{h \in H} logistic_loss_h(x, y, h)
# but because there exist a bijection between theta's and h's it is
# equivalent to choosing
# theta_best = argmin_{theta \in H} logistic_loss_theta(x, y, theta)

def hyposesis_function(x, theta):
    return sigmoid(np.dot(x, theta))

# negative log likelihood of observing sequence of integer
# y's given probabilities y_pred's of each Bernoulli trial
# recommendation: convert both variables to float's
# or weird sign stuff might happen like -1*y != -y for uint8
# use np.mean and broadcasting
def binary_logistic_loss(y, y_pred):
    assert y_pred.shape == y.shape
    m = y_pred.shape[0]
    return -1/m * np.sum( np.multiply(np.log(y_pred), y) + np.multiply(np.l

def logistic_loss_theta_grad(x, y, h, theta):
    y_pred = h(x, theta)
    point_wise_grads = (y_pred - y)*x
    grad = np.mean(point_wise_grads, axis=0)[: , None]
    assert grad.shape == theta.shape
    return grad

def logistic_loss_theta(x, y, h, theta):
    return binary_logistic_loss(y, h(x, theta))

```

```

In [7]: # Check that with theta as zeros, cost is about 0.693:
theta_init = np.zeros((X_data.shape[1], 1))
print(logistic_loss_theta(X_data, y_data, hyposesis_function, theta_init))
print(logistic_loss_theta_grad(X_data, y_data, hyposesis_function, theta_in

0.6931471805599453
[[-12.00921659]
 [-11.26284221]
 [-0.1         ] ]

```

Answer

When $\theta = \bar{0}$, the sigmoid function $g(z) = \frac{1}{1 + e^{-z}}$ when $z=0$ will always return 0.5 as the value for `hyposesis_function`. Thus the `hyposesis_function` will return a vector of 0.5 and pass to the cost function.

The cost function is defined below:

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

When $h_{\theta}(x) = [[0.5], [0.5], [0.5] \dots [0.5]]$, we will get

$$j(\theta) = \log(0.5) = 0.693$$

1.1.4 Learning parameters using *fmin*

In the previous assignment, you found the optimal parameters of a linear regression model by implementing gradient descent. You wrote a cost function and calculated its gradient, then took a gradient descent step accordingly. This time, instead of taking gradient descent steps, you will use a `scipy.optimize` built-in function called *fmin*.

The final θ value will then be used to plot the decision boundary on the training data, as seen in the figure below.

```
In [8]: import climin
        from functools import partial
```

```
In [9]: def optimize(theta_init, loss, loss_grad, max_iter=10000, print_every=1000,
                    theta = theta_init.copy()
                    opt_args = (theta, loss_grad)

                    if optimizer_fn is None:
                        optimizer_fn = partial(climin.GradientDescent, step_rate=1e-3, mome

                    optimizer = optimizer_fn(*opt_args)

                    loss_curve = []
                    for opt_info in optimizer:
                        n_iter = opt_info['n_iter']
                        f_value = loss(theta)
                        loss_curve.append(f_value)

                        if print_every != 0 and n_iter % print_every == 0:
                            print(n_iter, f_value)

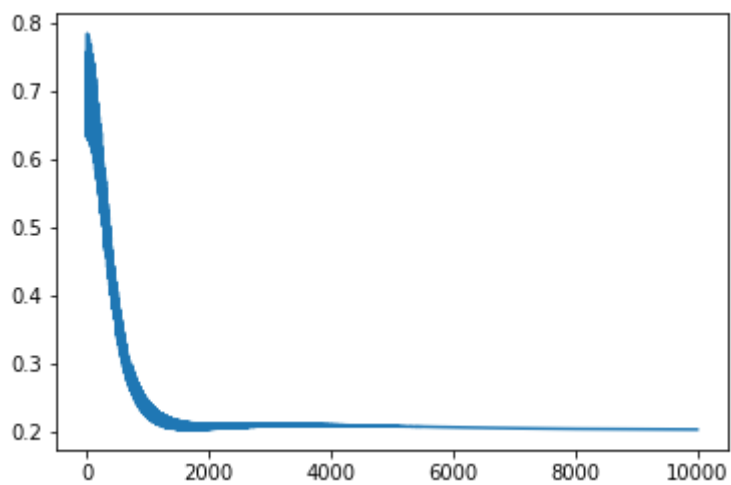
                        if n_iter == max_iter:
                            break

                    if show:
                        plt.plot(loss_curve)
                        plt.show()

                    return theta, f_value
```

```
In [10]: theta_init = np.zeros((3, 1))
loss = partial(logistic_loss_theta, X_data, y_data, hyposesis_function)
loss_grad = partial(logistic_loss_theta_grad, X_data, y_data, hyposesis_fun
theta, best_cost = optimize(theta_init, loss, loss_grad, show=True)
print(best_cost)
```

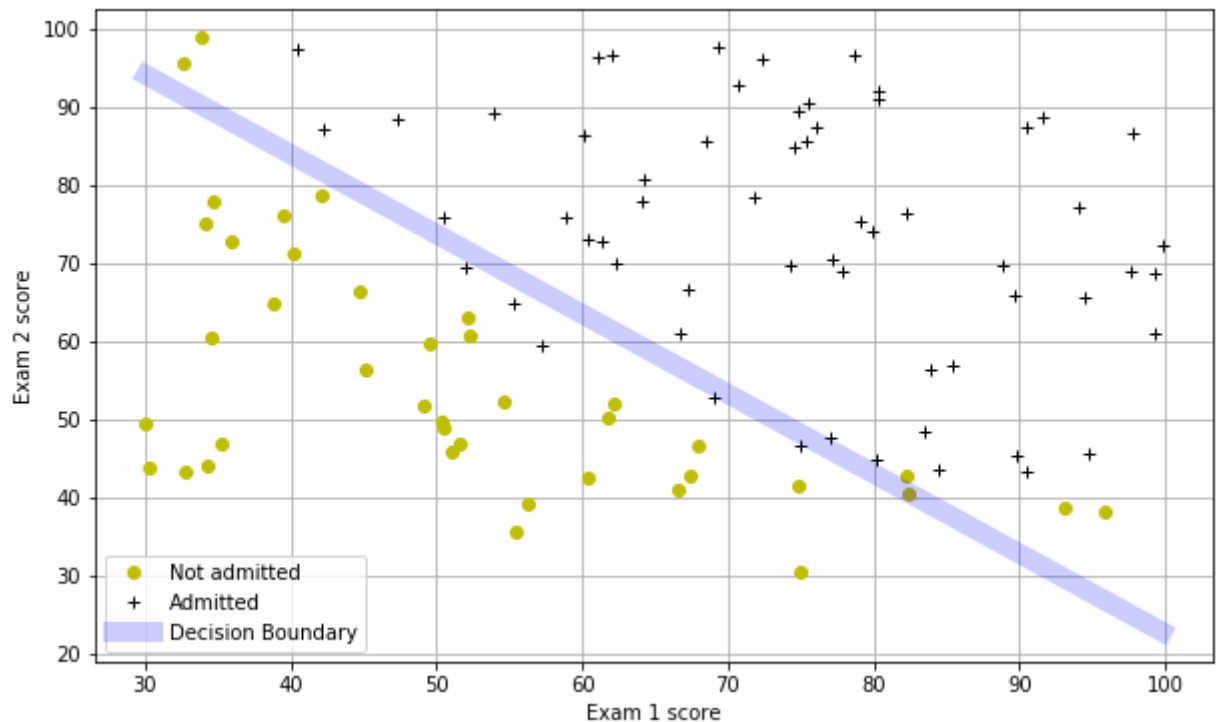
```
1000 0.2323549253449114
2000 0.20707436313671757
3000 0.20892645351895275
4000 0.20969077333394093
5000 0.2088788680198901
6000 0.20707752314457156
7000 0.2057098482536051
8000 0.2047598627173497
9000 0.2041540106268107
10000 0.20381418595046075
```



```
0.20381418595046075
```

```
In [11]: # Plotting the decision boundary: two points, draw a line between
# Decision boundary occurs when  $h = 0$ , or when
#  $\theta_0 x_1 + \theta_1 x_2 + \theta_2 = 0$ 
#  $y = mx + b$  is replaced by  $x_2 = (-1/\theta_1)(\theta_2 + \theta_0 x_1)$ 

line_xs = np.array([np.min(X_data[:,0]), np.max(X_data[:,0])])
line_ys = (-1./theta[1])*(theta[2] + theta[0]*line_xs)
plot_data(**student_plotting_spec)
plt.plot(line_xs, line_ys, 'b-', lw=10, alpha=0.2, label='Decision Boundary')
plt.legend()
plt.show()
```



1.1.5 [15pts] Evaluating logistic regression

After learning the parameters, you can use the model to predict whether a particular student will be admitted.

(a) [5 pts] Show that for a student with an Exam 1 score of 45 and an Exam 2 score of 85, you should expect to see an admission probability of 0.776.

Another way to evaluate the quality of the parameters we have found is to see how well the learned model predicts on our training set.

(b) [10 pts] In this part, your task is to complete the code in `makePrediction`. The predict function will produce “1” or “0” predictions given a dataset and a learned parameter vector θ . After you have completed the code, the script below will proceed to report the training accuracy of your classifier by computing the percentage of examples it got correct. You should also see a Training Accuracy of 89.0.

```
In [12]: # For a student with an Exam 1 score of 45 and an Exam 2 score of 85,
# you should expect to see an admission probability of 0.776.
check_data = np.array([[45., 85., 1]])
print(check_data.shape)
print(hypothesis_function(check_data, theta))

(1, 3)
[[0.78755263]]
```

```
In [13]: # use hypothesis function and broadcast compare operator
def predict(x, theta):
    y_pred = np.round(hypothesis_function(x, theta))
    return y_pred

def accuracy(x, y, theta):
    y_pred = predict(x, theta)
    return (y_pred == y).sum().astype(float) / len(y_pred)

print(accuracy(X_data, y_data, theta))

0.9
```

2. Regularized logistic regression

In this part of the exercise, you will implement regularized logistic regression to predict whether microchips from a fabrication plant pass quality assurance (QA). During QA, each microchip goes through various tests to ensure it is functioning correctly. Suppose you are the product manager of the factory and you have the test results for some microchips on two different tests. From these two tests, you would like to determine whether the microchips should be accepted or rejected. To help you make the decision, you have a dataset of test results on past microchips in *ex2data2.txt*, from which you can build a logistic regression model.

2.1 Visualizing the data

Similar to the previous parts of this exercise, `plotData` is used to generate the figure below, where the axes are the two test scores, and the positive ($y = 1$, accepted) and negative ($y = 0$, rejected) examples are shown with different markers.

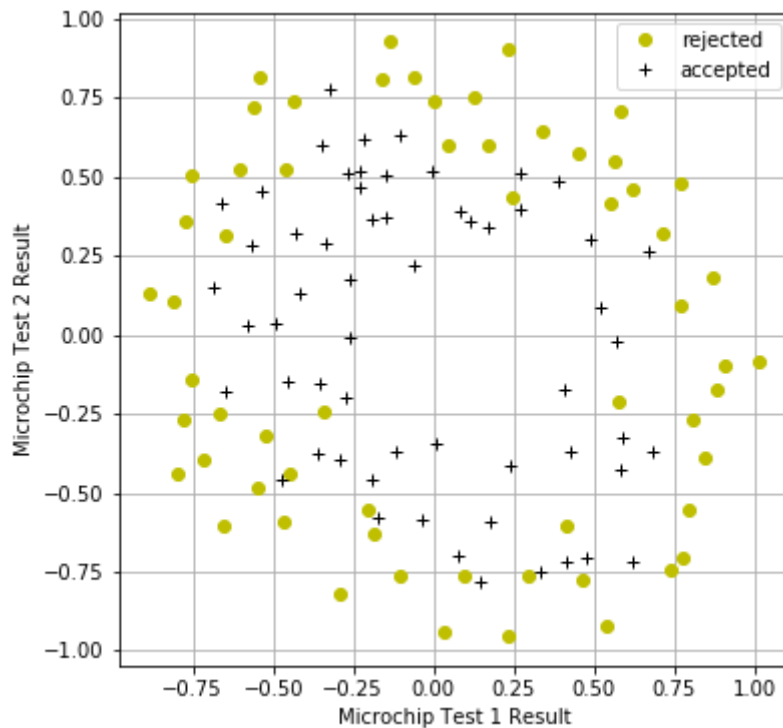
The figure below shows that our dataset cannot be separated into positive and negative examples by a straight line. Therefore, a straightforward application of logistic regression will not perform well on this dataset since logistic regression will only be able to find a linear decision boundary.

```
In [14]: X_data_, y_data = read_classification_csv_data('ex2data2.txt')
X_data = X_data_ - X_data_.mean(axis=0)[None, :]
print(X_data.shape, X_data.min(), X_data.max(), X_data.dtype)
print(y_data.shape, y_data.min(), y_data.max(), y_data.dtype)

(118, 2) -0.83007 1.1089 float64
(118, 1) 0.0 1.0 float64
(118, 2) -0.9528415593220338 1.0161210915254237 float64
(118, 1) 0 1 int32
```

```
In [15]: chip_plotting_spec = {
    'x': X_data,
    'y': y_data,
    'xlabel': 'Microchip Test 1 Result',
    'ylabel': 'Microchip Test 2 Result',
    'labels': ['rejected', 'accepted'],
    'markers': ['yo', 'k+'],
    'figsize': (6, 6)
}

plot_data(**chip_plotting_spec)
plt.show()
```



2.2 Nonlinear feature mapping

One way to fit the data better is to create more features from each data point. In *mapFeature* below, we will map the features into all polynomial terms of x_1 and x_2 up to the sixth power as follows:

$$\text{mapFeature}(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1 x_2 \\ x_2^2 \\ x_1^3 \\ \vdots \\ x_1 x_2^5 \\ x_2^6 \end{bmatrix}$$

As a result of this mapping, our vector of two features (the scores on two QA tests) has been transformed into a 28-dimensional vector. A logistic regression classifier trained on this higher-dimension feature vector will have a more complex decision boundary and will appear nonlinear when drawn in our 2-dimensional plot. While the feature mapping allows us to build a more expressive classifier, it is also more susceptible to overfitting. In the next parts of the exercise, you will implement regularized logistic regression to fit the data and also see for yourself how regularization can help combat the overfitting problem.

Either finite dimensional (or even infinite-dimensional, as you would see in the SVM lecture and the corresponding home assignment) feature mappings are usually denoted by Φ and therefore our hypothesis is now that the Bernoulli probability of chip malfunctioning might be described as

$$p_i = \sigma(\Phi(x_i)^T \theta)$$

```
In [16]: from itertools import combinations_with_replacement

def polynomial_feature_map(X_data, degree=20, show_me_ur_powers=False):
    assert len(X_data.shape) == 2
    group_size = X_data.shape[1]
    assert group_size == 2
    # hm.. how to get all ordered pairs (c, d) of non-negative ints
    # such that their sum is c + d <= degree?
    # it is equivalent to getting all groups of integers (a, b) such that
    # 0 <= a <= b <= degree and defining c = a, d = b - a
    # their sum is below degree, both are >= 0
    # then feature_i = (x_0 ^ c) * (x_1 ^ d)
    comb_iterator = combinations_with_replacement(range(degree+1), group_size)
    not_quite_powers = np.array(list(comb_iterator))
    powers_bad_order = not_quite_powers.copy()
    powers_bad_order[:, 1] -= not_quite_powers[:, 0]
    # let's reorder them so that lower power monomials come first
    rising_power_idx = np.argsort(powers_bad_order.sum(axis=1))
    powers = powers_bad_order[rising_power_idx]
    if show_me_ur_powers is True:
        print(powers.T)
        print('total power per monomial', powers.sum(axis=1))
    X_with_powers = np.power(X_data[:, :, None], powers.T[None])
    # tu tu power rangers (with replacement)
    X_poly = np.prod(X_with_powers, axis=1)
    return X_poly

X_pf = polynomial_feature_map(X_data, show_me_ur_powers=True)
print(X_pf.shape)
```

```
[[ 0  0  1  0  2  1  0  2  3  1  4  2  1  3  0  5  3  0  2  4  1  0  6  2
   5  4  3  1  4  0  1  2  6  7  3  5  2  3  4  6  8  5  7  0  1  4  7  3
   6  9  8  0  2  1  5  9  3  5  7  8  4 10  6  1  0  2  0  9  1  4  5  2
  10  7  8  6 11  3  7  0  1  9 12  4  6  2  8 11 10  3  5  6 13  1  5  7
   9  8  4 12  2 10  0 11  3  0 12 11 14 10  9  7  6  8 13  4  2  3  1  5
   5  4  9  0 15 14  7  1 10  8  3  2 11  6 12 13  8  3 11  2 14  9 12 13
  10 16  4  0  6  5  1 15  7 12  5  2 13 11  6 15  3 10  7  0  1 16 14  8
   4  9 17  5 12  2  4  8 14 13 16 11  9  6 15  0  1 10  3 18  7 17 15 13
  14 17 19 18 16  2 12  3  1  4  5  0  7  8  6  9 10 11 11 18  3  1 17 10
   4 16 13  5  0 19 15 12  6  9  7 14  2  8 20]]

[[ 0  1  0  2  0  1  3  1  0  2  0  2  3  1  4  0  2  5  3  1  4  6  0  4
   1  2  3  5  3  7  6  5  1  0  4  2  6  5  4  2  0  3  1  8  7  5  2  6
   3  0  1  9  7  8  4  1  7  5  3  2  6  0  4  9 10  8 11  2 10  7  6  9
   1  4  3  5  0  8  5 12 11  3  0  8  6 10  4  1  2  9  7  7  0 12  8  6
   4  5  9  1 11  3 13  2 10 14  2  3  0  4  5  7  8  6  1 10 12 11 13  9
  10 11  6 15  0  1  8 14  5  7 12 13  4  9  3  2  8 13  5 14  2  7  4  3
   6  0 12 16 10 11 15  1  9  5 12 15  4  6 11  2 14  7 10 17 16  1  3  9
  13  8  0 13  6 16 14 10  4  5  2  7  9 12  3 18 17  8 15  0 11  1  4  6
   5  2  0  1  3 17  7 16 18 15 14 19 12 11 13 10  9  8  9  2 17 19  3 10
  16  4  7 15 20  1  5  8 14 11 13  6 18 12  0]]

total power per monomial [ 0  1  1  2  2  2  3  3  3  3  4  4  4  4  4  5
   5  5  5  5  5  6  6  6
   6  6  6  6  7  7  7  7  7  7  7  7  8  8  8  8  8  8  8  8  8  9  9  9
   9  9  9  9  9  9  9 10 10 10 10 10 10 10 10 10 10 10 11 11 11 11 11 11
  11 11 11 11 11 11 12 12 12 12 12 12 12 12 12 12 12 12 12 12 13 13 13 13
  13 13 13 13 13 13 13 13 13 14 14 14 14 14 14 14 14 14 14 14 14 14 14 14
  15 15 15 15 15 15 15 15 15 15 15 15 15 15 15 15 16 16 16 16 16 16 16 16]
```

```

16 16 16 16 16 16 16 16 16 17 17 17 17 17 17 17 17 17 17 17 17 17 17 17
17 17 17 18 18 18 18 18 18 18 18 18 18 18 18 18 18 18 18 18 18 19 19
19 19 19 19 19 19 19 19 19 19 19 19 19 19 19 19 19 19 20 20 20 20 20
20 20 20 20 20 20 20 20 20 20 20 20 20 20 20 20]
(118, 231)

```

2.3 Cost function and gradient

Now you will implement code to compute the cost function and gradient for regularized logistic regression. Recall that the regularized cost function in logistic regression is:

$$J(\theta) = \left[\frac{1}{m} \sum_{i=1}^m \left[-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] \right] + \frac{\lambda}{2m} \sum_{j=2}^n \theta_j^2$$

Note that you should not regularize the parameter θ_0 (Why not? Think about why that would be a bad idea).

The gradient of the cost function is a vector where the j element is defined as follows (you should understand how to obtain this expression):

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \text{for } j = 0$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \quad \text{for } j \geq 1$$

2.3.1 [10pts] Implementing regularized logistic regression

Re-implement computeCost with regularization.

```

In [17]: # Cost function, default lambda (regularization) 0
def logistic_loss_theta_w_reg(x, y, h, theta, lambda_=0.0):
    no_lambda = logistic_loss_theta(x, y, h, theta)
    m = y.shape[0]
    w_lambda = no_lambda + lambda_ / (2 * m) * np.sum(np.square(theta))
    return w_lambda

def logistic_loss_theta_w_reg_grad(x, y, h, theta, lambda_=0.0):
    no_lambda = logistic_loss_theta_grad(x, y, h, theta)
    m = y.shape[0]
    w_lambda = (lambda_ / m) * theta
    w_lambda[0] = 0
    return no_lambda + w_lambda

```

Once you are done, you will call your cost function using the initial value of θ (initialized to all zeros). You should see that the cost is about 0.693.


```

In [18]: theta_init = np.zeros((X_pf.shape[1], 1))
print(logistic_loss_theta_w_reg(X_pf, y_data, hyposesis_function, theta_init))
print(logistic_loss_theta_w_reg_grad(X_pf, y_data, hyposesis_function, theta_init))

loss = partial(logistic_loss_theta_w_reg, X_pf, y_data, hyposesis_function)
loss_grad = partial(logistic_loss_theta_w_reg_grad, X_pf, y_data, hyposesis_function)
theta, best_cost = optimize(theta_init, loss, loss_grad, max_iter=10000, print_cost=True)
print('best loss', best_cost)
print('best acc', accuracy(X_pf, y_data, theta))

[-5.65355937e-05]
[ 5.97572413e-05]
[-2.28744025e-03]
[ 1.04308827e-03]
[-7.42326893e-04]
[-1.16041532e-04]
[-5.25177896e-05]
[-1.30965247e-04]
[ 5.49236357e-05]
[ 6.96814247e-03]
[-1.10125780e-04]
[ 6.33934609e-05]
[ 6.28673866e-04]
[ 1.86987655e-04]
[ 5.73885016e-05]
[ 9.92359449e-05]
[-5.45334072e-05]
[ 2.99854418e-04]
[-4.58889618e-05]
[ 4.58752628e-05]

```

2.3.2 [15pts] Learning parameters using *minimize*

You will use *optimize.minimize* to learn the optimal parameters θ . If you have completed the cost and gradient for regularized logistic regression correctly, you should be able to learn the parameters θ using *minimize*. Implement the function *optimizeRegularizedTheta* below.

2.4 Plotting the decision boundary

To help you visualize the model learned by this classifier, we have provided the function *plotBoundary* which plots the (non-linear) decision boundary that separates the positive and negative examples.

```
In [19]: def plot_boundary(theta, ax=None):
        """
        Function to plot the decision boundary for arbitrary theta, X, y, lambda
        Inside of this function is feature mapping, and the minimization routine
        It works by making a grid of x1 ("xvals") and x2 ("yvals") points,
        And for each, computing whether the hypothesis classifies that point as
        True or False. Then, a contour is drawn with a built-in pyplot function
        """
        ax = ax or plt.gca()
        x_range = np.linspace(-1, 1.5, 50)
        y_range = np.linspace(-1, 1.5, 50)
        xx, yy = np.meshgrid(x_range, y_range)
        X_fake = np.stack([xx, yy]).reshape(2, -1).T
        X_fake_fm = polynomial_feature_map(X_fake)
        y_pred_fake = hypothesis_function(X_fake_fm, theta)
        return ax.contour(x_range, y_range, y_pred_fake.reshape(50, 50).T, [0.
```

2.4.1 [10pts] Plot Decision Boundaries

- (a) [2 pts] Use *plotBoundary* to obtain four subplots of the decision boundary for the following values of the regularization parameter: $\lambda = 0, 1, 5, 10$
- (b) [2 pts] Comment on which plots are overfitting and which plots are underfitting.
- (c) [2 pts] Which is the model with the highest bias? The highest variance?
- (d) [2 pts] What is another way to detect overfitting?
- (e) [2 pts] Considering that later components of theta correspond to higher powers of monomials, plot values of theta and comment on effects of regularization

```

In [54]: # (a) Build a figure showing contours for various values of regularization

np.random.seed(2)
train_idx_mask = np.random.rand(X_pf.shape[0]) < 0.3
X_pf_train, y_train = X_pf[train_idx_mask], y_data[train_idx_mask]
X_pf_test, y_test = X_pf[-train_idx_mask], y_data[-train_idx_mask]
print([x.shape for x in (X_pf_train, y_train, X_pf_test, y_test)])

def silent_optimize_w_lambda(lambda_):
    theta_init = np.zeros((X_pf.shape[1], 1))
    data = (X_pf_train, y_train, hypothesis_function)
    loss = partial(logistic_loss_theta_w_reg, *data, lambda_=lambda_)
    loss_grad = partial(logistic_loss_theta_w_reg_grad, *data, lambda_=lambda_)
    optimizer_fn = partial(climin.GradientDescent, step_rate=1e-4, momentum=0.9)
    theta, final_loss = optimize(
        theta_init, loss, loss_grad, optimizer_fn=optimizer_fn,
        max_iter=1000, print_every=0, show=False
    )
    return theta, final_loss

thetas = []
plt.figure(figsize=(12,10))

# wow, I mutates an object used in the scope of another function (plot_data)
# don't do that! it is really hard to debug later

chip_plotting_spec['figsize'] = None

# you might find following lines useful:
#
# cnt_fmt = {0.5: 'Lambda = %d' % lambda_}
# ax.clabel(cnt, inline=1, fontsize=15, fmt=cnt_fmt)
#
# red dots indicate training samples

chip_plotting_train = {
    'X': X_train_plot,
    'y': y_train_plot,
    'xlabel': 'Microchip Test 1 Result',
    'ylabel': 'Microchip Test 2 Result',
    'labels': ['train rejected', 'train accepted'],
    'markers': ['ro', 'r+'],
    'figsize': None
}
X_train_plot, y_train_plot = X_data[train_idx_mask], y_data[train_idx_mask]

for id_, lambda_ in enumerate([0, 1, 5, 10]):
    ax = plt.subplot(2, 2, id_+1)
    theta, final_loss = silent_optimize_w_lambda(lambda_)
    thetas.append(theta)
    acc = accuracy(X_pf_test, y_test, theta)
    cnt_fmt = {0.5: 'Lambda = %d' % lambda_}
    cnt = plot_boundary(theta, ax=ax)
    ax.clabel(cnt, inline=1, fontsize=15, fmt=cnt_fmt)

```

```

ax.set_title('Decision Boundary, Accuracy=%.2f, Loss=%.2f' % (acc, final_loss))
plot_data(**chip_plotting_spec)
plot_data(**chip_plotting_train)

plt.show()

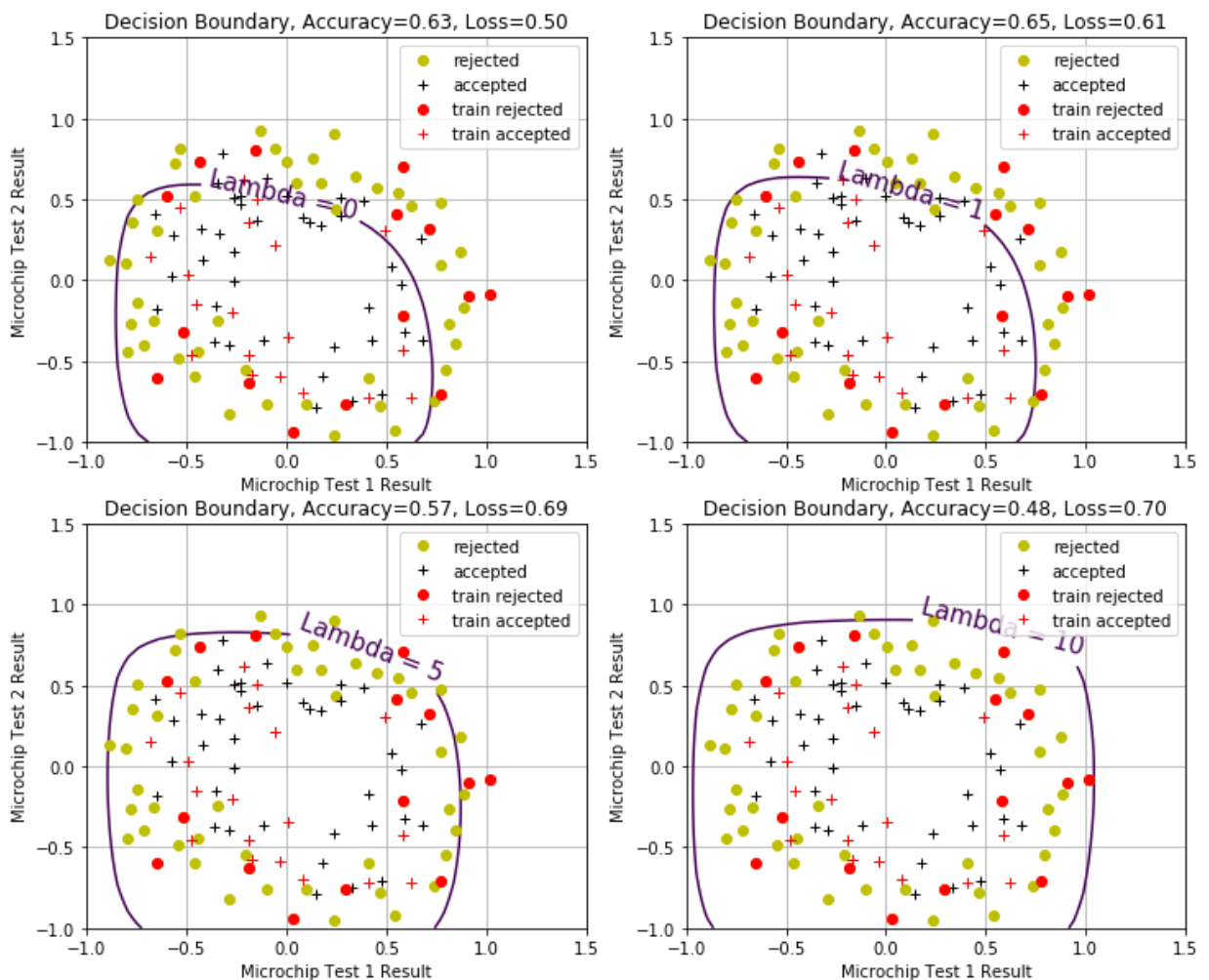
# (e) [2 pts] Considering that later components of theta correspond to high
# of monomials, plot values of theta and comment on effects of regularization

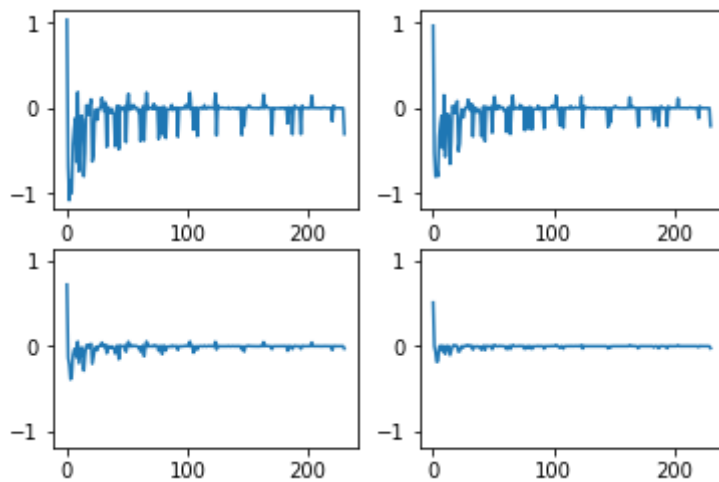
ax = None
for th_id, theta in enumerate(thetas):
    ax = plt.subplot(2, 2, th_id+1, sharey=ax)
    ax.plot(theta)

plt.show()

[(34, 231), (34, 1), (84, 231), (84, 1)]

```





- (a) [2 pts] Use `plotBoundary` to obtain four subplots of the decision boundary for the following values of the regularization parameter: $\lambda = 0, 1, 5, 10$
- (b) [2 pts] Comment on which plots are overfitting and which plots are underfitting.
- (c) [2 pts] Which is the model with the highest bias? The highest variance?
- (d) [2 pts] What is another way to detect overfitting?
- (e) [2 pts] Considering that later components of θ correspond to higher powers of monomials, plot values of θ and comment on effects of regularization

Answer

- (b) When $\lambda = 0$, the plot is overfitting, since the loss is low but accuracy is low as well. When $\lambda = 5, 10$, the plots are underfitting, where the loss is high and accuracy is low.
- (c) When $\lambda = 10$, the model has the highest bias. When $\lambda = 0$, the model has the highest variance.
- (d) Use cross validation to check if the validation error is high.
- (e) With higher λ , θ are less and more close to zero, which are penalized by the regularization. But if the λ is too high, most θ are close to zero and have no effect on the model, which will lead to underfitting.

3. Written part

These problems are extremely important preparation for the exam. Submit solutions to each problem by filling the markdown cells below.

3.1 [10pts] Maximum likelihood for Logistic Regression

Showing all steps, derive the LR cost function using maximum likelihood. Assume that the probability of y given x is described by:

$$P(y = 1 | x; \theta) = h_{\theta}(x)$$

$$P(y = 0 | x; \theta) = 1 - h_{\theta}(x)$$

Answer

Maximum likelihood can be written as

$$P(y | X; \theta) = \prod_{i=1}^m P(y_i | x_i; \theta) = \prod_{i=1}^m h_{\theta}(x)^{y_i} (1 - h_{\theta}(x))^{1-y_i}$$

After log both sides:

$$\log(P(y | X; \theta)) = \sum_{i=1}^m y_i \log(h_{\theta}(x)) + (1 - y_i) \log(1 - h_{\theta}(x))$$

Normalization and turn the negative value into positive, we will have the cost function:

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

3.2 [10pts] Logistic Regression Classification with Label Noise

Suppose you are building a logistic regression classifier for images of dogs, represented by a feature vector x , into one of two categories $y \in \{0, 1\}$, where 0 is “terrier” and 1 is “husky.” You decide to use the logistic regression model $p(y = 1 | x) = h_{\theta}(x) = \theta^T x$. You collected an image dataset $\mathbf{D} = \{x^{(i)}, t^{(i)}\}$, however, you were very tired and made some mistakes in assigning labels $t^{(i)}$. You estimate that you were correct in about τ fraction of all cases.

(a) Write down the equation for the posterior probability $p(t = 1 | x)$ of the label being 1 for some point x , in terms of the probability of the true class, $p(y = 1 | x)$.

(b) Derive the modified cost function in terms of θ , $x^{(i)}$, $t^{(i)}$ and τ .

Answer

(a)

$$\begin{aligned} p(t = 1 | x) &= p(y = 1 | x) \tau + p(y = 0 | x)(1 - \tau) \\ &= p(y = 1 | x) \tau + (1 - p(y = 1 | x))(1 - \tau) \end{aligned}$$

(b)

From Answer (a) we got :

$$p(t = 1 | x) = \tau (p(y = 1 | x) - p(y = 0 | x)) + p(y = 0 | x)$$

$$p(t = 0 | x) = \tau (p(y = 0 | x) - p(y = 1 | x)) + p(y = 1 | x)$$

Since

$$p(y = 1 | x) = h_{\theta}(x) = \theta^T x$$

$$p(y = 0 | x) = 1 - h_{\theta}(x)$$

We can rewrite the probability function below:

$$p(t = 1 | x) = h_{\theta}(x) (2\tau - 1) + (1 - \tau)$$

$$p(t = 0 | x) = h_{\theta}(x) (1 - 2\tau) + \tau$$

Maximum likelihood:

$$P(t | X; \theta) = \prod_{i=1}^m P(t_i | x_i) = \prod_{i=1}^m [h_{\theta}(x) (2\tau - 1) + (1 - \tau)]^{t_i} [h_{\theta}(x) (1 - 2\tau) + \tau]^{1-t_i}$$

Log both sides:

$$\log(P(t | X)) = \sum_{i=1}^m t_i \log[h_{\theta}(x) (2\tau - 1) + (1 - \tau)] + (1 - t_i) \log[h_{\theta}(x) (1 - 2\tau) + \tau]$$

Finally the loss function:

$$j(\theta) = -\frac{1}{m} \sum_{i=1}^m t_i \log[h_{\theta}(x) (2\tau - 1) + (1 - \tau)] + (1 - t_i) \log[h_{\theta}(x) (1 - 2\tau) + \tau]$$

3.3 [10pts] Cross-entropy loss for multiclass classification

This problem asks you to derive the cross-entropy loss for a multiclass classification problem using maximum likelihood. Consider the multiclass classification problem in which each input is assigned to one of K mutually exclusive classes. The binary target variables $y_k \in \{0, 1\}$ have a "one-hot" coding scheme, where the value is 1 for the indicated class and 0 for all others. Assume that we can interpret the network outputs as $h_k(x, \theta) = p(y_k = 1 | x)$, or the probability of the k th class.

Show that the maximum likelihood estimate of the parameters θ can be obtained by minimizing the multiclass *cross-entropy* loss function

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log(h_k(x_i, \theta))$$

where N is the number of examples $\{x_i, y_i\}$.

Answer

For one data point x :

$$p(y_k = 1 | x) = h_k(x, \theta) = \prod_{k=1}^K h_k(x, \theta)^{y_k}$$

For all data:

$$p(y_k = 1 | X) = \prod_{i=1}^N \prod_{k=1}^K h_k(x_i, \theta)^{y_{ik}}$$

Then we can get the maximum likelihood by taking log for both sides and normalization.

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log(h_k(x_i, \theta))$$

In []: