

Problem Set 1: Linear Regression

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, 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. Here are some useful resources for [Markdown guide \(https://www.markdownguide.org/basic-syntax/\)](https://www.markdownguide.org/basic-syntax/) and [LaTeX tutorial \(https://www.overleaf.com/learn/latex/Learn_LaTeX_in_30_minutes\)](https://www.overleaf.com/learn/latex/Learn_LaTeX_in_30_minutes) if you are not familiar with the basic syntax.

If certain output is given for some cells, that means that you are expected to get similar results.

Only **PDF** files are accepted for ps1 submission. To print this notebook to a pdf file, you can go to "File" -> "Download as" -> "PDF via LaTeX(.pdf)" or simply use "print" in browser.

Total: 185 points.

1. Numpy Tutorial

1.1 [5pt] Modify the cell below to return a 5x5 matrix of ones. Put some code there and press `Ctrl+Enter` to execute contents of the cell. You should see something like the output above. [\[1\] \(https://docs.scipy.org/doc/numpy-1.13.0/user/basics.creation.html#arrays-creation\)](https://docs.scipy.org/doc/numpy-1.13.0/user/basics.creation.html#arrays-creation) [\[2\] \(https://docs.scipy.org/doc/numpy-1.13.0/reference/routines.array-creation.html#routines-array-creation\)](https://docs.scipy.org/doc/numpy-1.13.0/reference/routines.array-creation.html#routines-array-creation)

In [1]:

```
import numpy as np
from sklearn import preprocessing
import matplotlib.pyplot as plt

sample = np.ones((5,5), dtype = float)

print(sample)
```

```
[[1. 1. 1. 1. 1.]
 [1. 1. 1. 1. 1.]
 [1. 1. 1. 1. 1.]
 [1. 1. 1. 1. 1.]
 [1. 1. 1. 1. 1.]]
```

1.2 [5pt] Vectorizing your code is very important to get results in a reasonable time. Let A be a 10x10 matrix and x be a 10-element column vector. Your friend writes the following code. How would you vectorize this code to run without any for loops? Compare execution speed for different values of n with `%timeit` [\(http://ipython.readthedocs.io/en/stable/interactive/magics.html#magic-timeit\)](http://ipython.readthedocs.io/en/stable/interactive/magics.html#magic-timeit).

In [2]:

```
n = 10
def compute_something(A, x):
    v = np.zeros((n, 1))
    for i in range(n):
        for j in range(n):
            v[i] += A[i, j] * x[j]
    return v

A = np.random.rand(n, n)
x = np.random.rand(n, 1)
print(compute_something(A, x))
```

```
[[2.15631321]
 [3.25411493]
 [2.89031453]
 [2.82268099]
 [2.98111108]
 [2.67606369]
 [2.34242693]
 [2.96335072]
 [2.40192092]
 [2.1423447 ]]
```

In [3]:

```
def vectorized(A, x):
    return (A.dot(x))

print(vectorized(A, x))
assert np.max(abs(vectorized(A, x) - compute_something(A, x))) < 1e-3
```

```
[[2.15631321]
 [3.25411493]
 [2.89031453]
 [2.82268099]
 [2.98111108]
 [2.67606369]
 [2.34242693]
 [2.96335072]
 [2.40192092]
 [2.1423447 ]]
```

In [4]:

```

for n in [5, 10, 100, 500]:
    A = np.random.rand(n, n)
    x = np.random.rand(n, 1)
    %timeit -n 5 compute_something(A, x)
    %timeit -n 5 vectorized(A, x)
    print('---')

```

```

68.4 µs ± 3.55 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
The slowest run took 60.45 times longer than the fastest. This could mean
that an intermediate result is being cached.
10.3 µs ± 13.3 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
274 µs ± 27 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
The slowest run took 9.46 times longer than the fastest. This could mean
that an intermediate result is being cached.
1.41 µs ± 1.87 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
24.4 ms ± 327 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
The slowest run took 19.50 times longer than the fastest. This could mean
that an intermediate result is being cached.
11.8 µs ± 11.3 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
609 ms ± 1.69 ms per loop (mean ± std. dev. of 7 runs, 5 loops each)
The slowest run took 7.24 times longer than the fastest. This could mean
that an intermediate result is being cached.
86.3 µs ± 88.4 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---

```

2. Linear regression with one variable

In this part of this exercise, you will implement linear regression with one variable to predict profits for a food truck. Suppose you are the CEO of a restaurant franchise and are considering different cities for opening a new outlet. The chain already has trucks in various cities and you have data for profits and populations from the cities. You would like to use this data to help you select which city to expand to next. The file `ex1data.txt` contains the dataset for our linear regression problem. The first column is the population of a city and the second column is the profit of a food truck in that city. A negative value for profit indicates a loss.

2.1 [10pt] Get a plot similar to below : [1]

(https://matplotlib.org/devdocs/api/_as_gen/matplotlib.pyplot.scatter.html). [2]

(https://matplotlib.org/api/pyplot_api.html?highlight=xlim#matplotlib.pyplot.xlim). [3]

(https://matplotlib.org/api/pyplot_api.html?highlight=matplotlib%20pyplot%20xlabel#matplotlib.pyplot.xlabel).

Before starting on any task, it is often useful to understand the data by visualizing it. For this dataset, you can use a scatter plot to visualize the data, since it has only two properties to plot (profit and population). Many other problems that you will encounter in real life are multi-dimensional and can't be plotted on a 2-d plot.

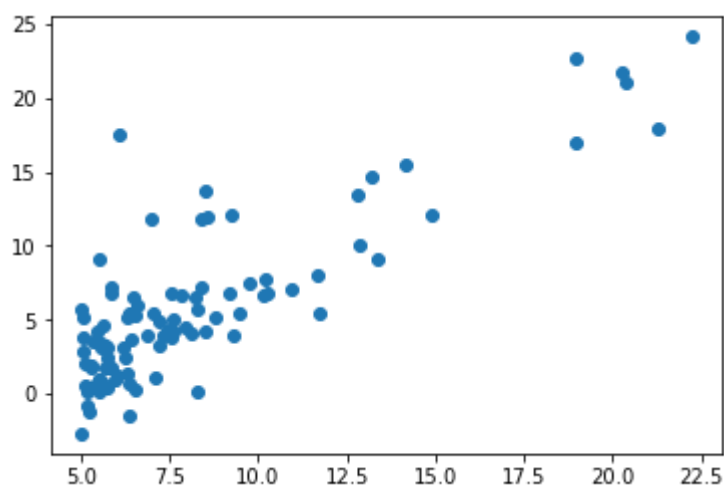
In [5]:

```
data = np.loadtxt('ex1data1.txt', delimiter=',')
X, y = data[:, 0, np.newaxis], data[:, 1, np.newaxis]
n = data.shape[0]
print(X.shape, y.shape, n)
print(X[:10], '\n', y[:10])

plt.scatter(X,y)
plt.show()
```

(97, 1) (97, 1) 97

```
[[ 6.1101]
 [ 5.5277]
 [ 8.5186]
 [ 7.0032]
 [ 5.8598]
 [ 8.3829]
 [ 7.4764]
 [ 8.5781]
 [ 6.4862]
 [ 5.0546]]
[[ 17.592 ]
 [  9.1302]
 [ 13.662 ]
 [ 11.854 ]
 [  6.8233]
 [ 11.886 ]
 [  4.3483]
 [ 12.     ]
 [  6.5987]
 [  3.8166]]
```



2.2 Gradient Descent

In this part, you will fit the linear regression parameter θ to our dataset using gradient descent.

The objective of linear regression is to minimize the cost function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}; \theta) - y^{(i)})^2$$

where the hypothesis $h(x; \theta)$ is given by the linear model (x' has an additional fake feature always equal to 1)

$$h(x; \theta) = \theta^T x' = \theta_0 + \theta_1 x$$

Recall that the parameters of your model are the θ_j values. These are the values you will adjust to minimize cost $J(\theta)$. One way to do this is to use the gradient descent algorithm. In batch gradient descent algorithm, each iteration performs the update.

$$\theta_j^{(k+1)} = \theta_j^{(k)} - \eta \frac{1}{m} \sum_i (h(x^{(i)}; \theta) - y^{(i)}) x_j^{(i)}$$

With each step of gradient descent, your parameter θ_j come closer to the optimal values that will achieve the lowest cost $J(\theta)$.

2.2.1 [5pt] Where does this update rule comes from?

Solution

In the gradient descent, the θ is updated iteratively to minimize the cost function until converge. Every time we update the θ according to the gradient:

$$\begin{aligned} \theta_j^{(k+1)} &= \theta_j^{(k)} - \eta \frac{\partial}{\partial \theta_j} J(\theta) \\ &= \theta_j^{(k)} - \eta \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}; \theta) - y^{(i)})^2 \\ &= \theta_j^{(k)} - \eta \frac{1}{m} \sum_i (h(x^{(i)}; \theta) - y^{(i)}) \frac{\partial}{\partial \theta_j} h(x^{(i)}; \theta) \\ &= \theta_j^{(k)} - \eta \frac{1}{m} \sum_i (h(x^{(i)}; \theta) - y^{(i)}) x_j^{(i)} \end{aligned}$$

2.2.2 [30pt] Cost Implementation

As you perform gradient descent to learn to minimize the cost function, it is helpful to monitor the convergence by computing the cost. In this section, you will implement a function to calculate $J(\theta)$ so you can check the convergence of your gradient descent implementation.

In the following lines, we add another dimension to our data to accommodate the intercept term and compute the prediction and the loss. As you are doing this, remember that the variables X and y are not scalar values, but matrices whose rows represent the examples from the training set. In order to get x' [add a column](https://docs.scipy.org/doc/numpy/reference/generated/numpy.insert.html) (<https://docs.scipy.org/doc/numpy/reference/generated/numpy.insert.html>) of ones to the data matrix X .

You should expect to see a cost of approximately 32.

In [6]:

```

# assertions below are true only for this
# specific case and are given to ease debugging!

def add_column(X):
    assert len(X.shape) == 2 and X.shape[1] == 1
    I = np.ones((X.shape[0],1), dtype = float)
    X = np.insert(X, [0], I, axis=1)
    return X

def predict(X, theta):
    """ Computes  $h(x; \theta)$  """
    assert len(X.shape) == 2 and X.shape[1] == 1
    assert theta.shape == (2, 1)

    X_prime = add_column(X)
    pred = X_prime.dot(theta)

    return pred

def loss(X, y, theta):
    assert X.shape == (n, 1)
    assert y.shape == (n, 1)
    assert theta.shape == (2, 1)

    X_prime = add_column(X)
    assert X_prime.shape == (n, 2)
    loss = (X_prime.dot(theta) - y).T.dot((X_prime.dot(theta) - y))/(2*n)
    return loss[0][0]

theta_init = np.zeros((2, 1))
print(loss(X, y, theta_init))

```

32.072733877455676

2.2.3 [40pt] GD Implementation

Next, you will implement gradient descent. The loop structure has been written for you, and you only need to supply the updates to θ within each iteration.

As you program, make sure you understand what you are trying to optimize and what is being updated. Keep in mind that the cost is parameterized by the vector θ not X and y . That is, we minimize the value of $J(\theta)$ by changing the values of the vector θ , not by changing X or y .

A good way to verify that gradient descent is working correctly is to look at the value of and check that it is decreasing with each step. Your value of $J(\theta)$ should never increase, and should converge to a steady value by the end of the algorithm. Another way of making sure your gradient estimate is correct is to check it against a [finite difference \(https://en.wikipedia.org/wiki/Finite_difference\)](https://en.wikipedia.org/wiki/Finite_difference) approximation.

We also initialize the initial parameters to 0 and the learning rate α to 0.01 .

In [7]:

```

import scipy.optimize
from functools import partial

def loss_gradient(X, y, theta):
    X_prime = add_column(X)
    loss_grad = (X_prime.T).dot((predict(X, theta) - y)) / X.shape[0]
    return loss_grad

assert loss_gradient(X, y, theta_init).shape == (2, 1)

def finite_diff_grad_check(f, grad, points, eps=1e-10):
    errs = []
    for point in points:
        point_errs = []
        grad_func_val = grad(point)
        for dim_i in range(point.shape[0]):
            diff_v = np.zeros_like(point)
            diff_v[dim_i] = eps
            dim_grad = (f(point+diff_v) - f(point-diff_v))/(2*eps)
            point_errs.append(abs(dim_grad - grad_func_val[dim_i]))
        errs.append(point_errs)
    return errs

test_points = [np.random.rand(2, 1) for _ in range(10)]
finite_diff_errs = finite_diff_grad_check(
    partial(loss, X, y), partial(loss_gradient, X, y), test_points
)

print('max grad comp error', np.max(finite_diff_errs))
assert np.max(finite_diff_errs) < 1e-3, "grad computation error is too large"

def run_gd(loss, loss_gradient, X, y, theta_init, lr=0.01, n_iter=1500):
    theta_current = theta_init.copy()
    loss_values = []
    theta_values = []

    for i in range(n_iter):
        loss_value = loss(X, y, theta_current)
        theta_current = theta_current - lr*loss_gradient(X, y, theta_current)
        loss_values.append(loss_value)
        theta_values.append(theta_current)

    return theta_current, loss_values, theta_values

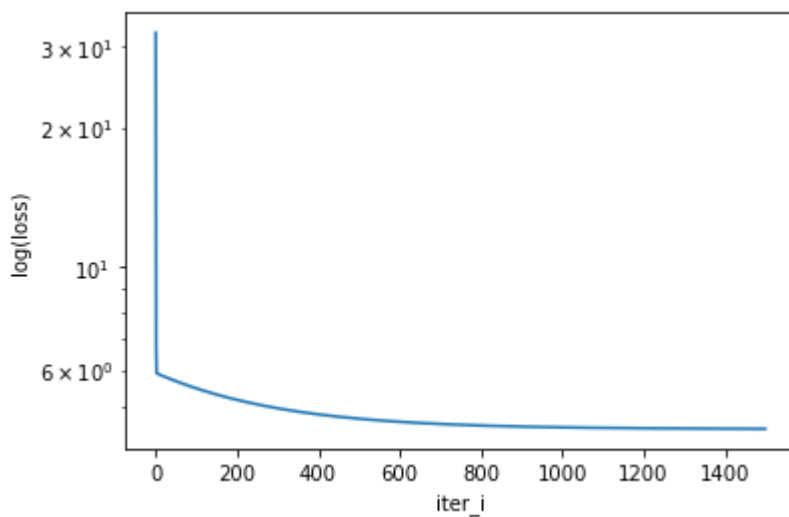
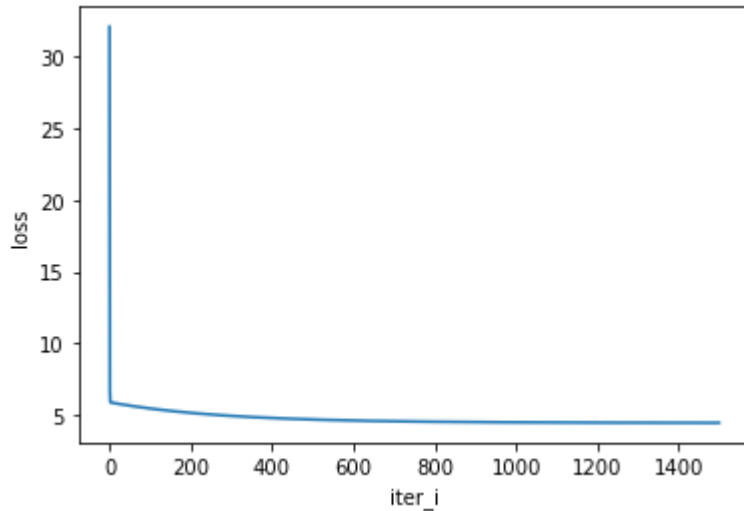
result = run_gd(loss, loss_gradient, X, y, theta_init)
theta_est, loss_values, theta_values = result

print('estimated theta value', theta_est.ravel())
print('resulting loss', loss(X, y, theta_est))
plt.ylabel('loss')
plt.xlabel('iter_i')
plt.plot(loss_values)
plt.show()

plt.ylabel('log(loss)')
plt.xlabel('iter_i')
plt.semilogy(loss_values)
plt.show()

```

```
max grad comp error 5.6789581272198575e-05  
estimated theta value [-3.63029144  1.16636235]  
resulting loss 4.483388256587726
```

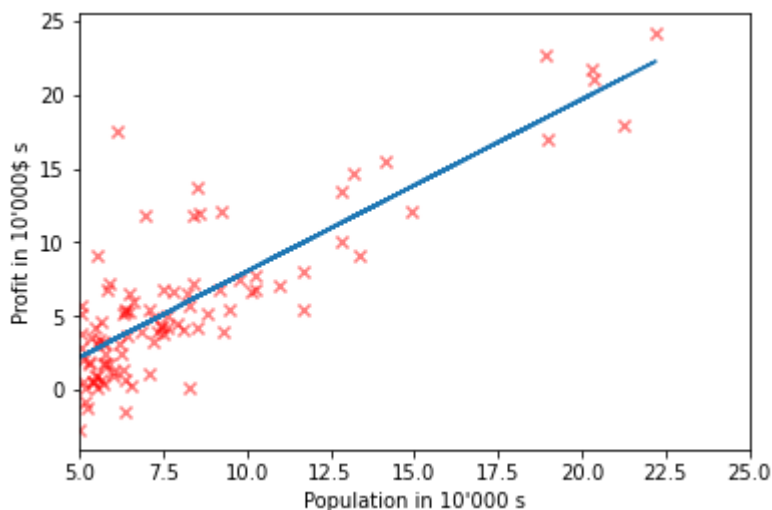


2.2.4 [10pt] After you are finished, use your final parameters to plot the linear fit. The result should look something like on the figure below. Use the `predict()` function.

In [8]:

```
plt.scatter(X, y, marker='x', color='r', alpha=0.5)
x_start, x_end = 5, 25
pred = predict(X, theta_est)

plt.xlabel('Population in 10\'000 s')
plt.ylabel('Profit in 10\'000$ s')
plt.plot(X, pred)
plt.xlim([x_start, x_end])
plt.show()
```



Now use your final values for θ and the `predict()` function to make predictions on profits in areas of 35,000 and 70,000 people.

In [9]:

```
x_1 = np.array([[35000]], dtype = float)
x_2 = np.array([[70000]], dtype = float)
result_1 = predict(x_1, theta_est)
result_2 = predict(x_2, theta_est)
print(result_1)
print(result_2)
```

```
[[40819.05197031]]
[[81641.73423205]]
```

To understand the cost function better, you will now plot the cost over a 2-dimensional grid of values. You will not need to code anything new for this part, but you should understand how the code you have written already is creating these images.

In [10]:

```

from mpl_toolkits.mplot3d import Axes3D
import matplotlib.cm as cm
limits = [(-10, 10), (-1, 4)]
space = [np.linspace(*limit, 100) for limit in limits]
theta_1_grid, theta_2_grid = np.meshgrid(*space)
theta_meshgrid = np.vstack([theta_1_grid.ravel(), theta_2_grid.ravel()])
loss_test_vals_flat = (((add_column(X) @ theta_meshgrid - y)**2).mean(axis=0)/2)
loss_test_vals_grid = loss_test_vals_flat.reshape(theta_1_grid.shape)
print(theta_1_grid.shape, theta_2_grid.shape, loss_test_vals_grid.shape)

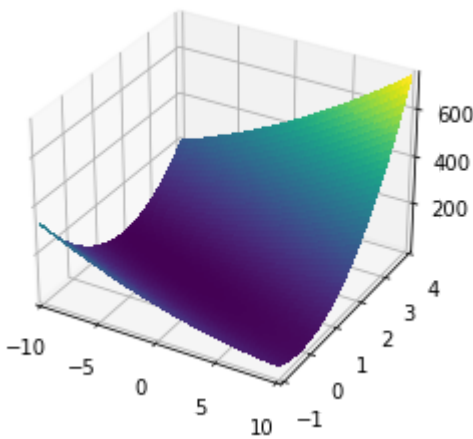
plt.gca(projection='3d').plot_surface(theta_1_grid, theta_2_grid,
                                      loss_test_vals_grid, cmap=cm.viridis,
                                      linewidth=0, antialiased=False)

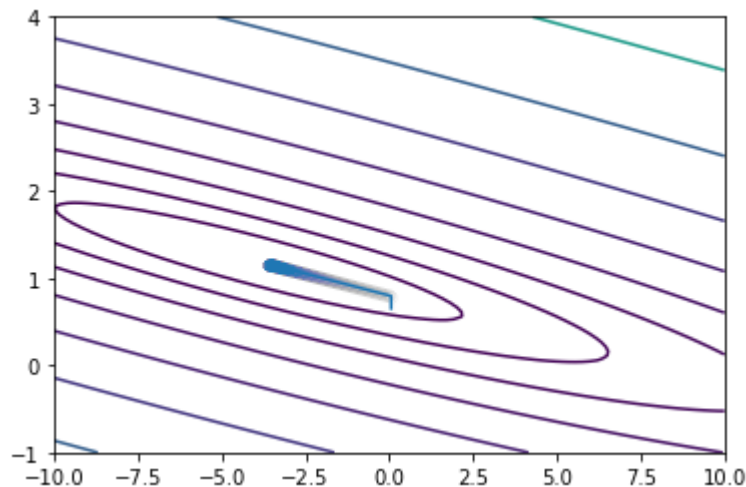
xs, ys = np.hstack(theta_values).tolist()
zs = np.array(loss_values)
plt.gca(projection='3d').plot(xs, ys, zs, c='r')
plt.xlim(*limits[0])
plt.ylim(*limits[1])
plt.show()

plt.contour(theta_1_grid, theta_2_grid, loss_test_vals_grid, levels=np.logspace(-2,
plt.plot(xs, ys)
plt.scatter(xs, ys, alpha=0.005)
plt.xlim(*limits[0])
plt.ylim(*limits[1])
plt.show()

```

(100, 100) (100, 100) (100, 100)





3. Linear regression with multiple input features

3.1 [20pt] Copy-paste your `add_column`, `predict`, `loss` and `loss grad` implementations from above and modify your code of linear regression with one variable to support any number of input features (vectorize your code.)

In [11]:

```
data = np.loadtxt('ex1data2.txt', delimiter=',')
X, y = data[:, :-1], data[:, -1, np.newaxis]
n = data.shape[0]
print(X.shape, y.shape, n)
print(X[:10], '\n', y[:10])
```

```
(47, 2) (47, 1) 47
[[2.104e+03 3.000e+00]
 [1.600e+03 3.000e+00]
 [2.400e+03 3.000e+00]
 [1.416e+03 2.000e+00]
 [3.000e+03 4.000e+00]
 [1.985e+03 4.000e+00]
 [1.534e+03 3.000e+00]
 [1.427e+03 3.000e+00]
 [1.380e+03 3.000e+00]
 [1.494e+03 3.000e+00]]
[[399900.]
 [329900.]
 [369000.]
 [232000.]
 [539900.]
 [299900.]
 [314900.]
 [198999.]
 [212000.]
 [242500.]]
```

In [12]:

```

#raise NotImplementedError("Implement new add_column(), predict(), loss(), loss_grad")
def add_column(X):
    assert len(X.shape) == 2
    I = np.ones((X.shape[0],1), dtype = float)
    X = np.insert(X, [0], I, axis=1)
    return X

def predict(X, theta):
    """ Computes h(x; theta) """
    X_prime = add_column(X)
    pred = X_prime.dot(theta)

    return pred

def loss_gradient(X, y, theta):
    X_prime = add_column(X)
    loss_grad = (X_prime.T).dot((predict(X, theta) - y)) / X.shape[0]
    return loss_grad

def loss(X, y, theta):

    X_prime = add_column(X)
    loss = (X_prime.dot(theta) - y).T.dot((X_prime.dot(theta) - y))/(2*n)
    return loss[0][0]

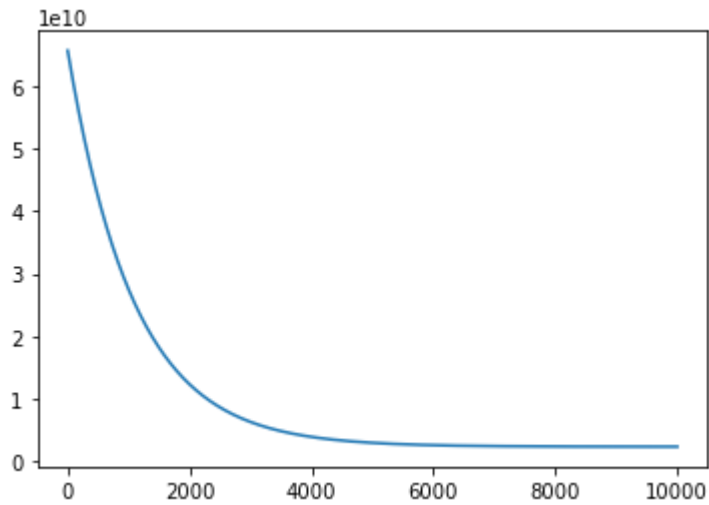
def run_gd(loss, loss_gradient, X, y, theta_init, lr=0.01, n_iter=1500):
    theta_current = theta_init.copy()
    loss_values = []
    theta_values = []

    for i in range(n_iter):
        loss_value = loss(X, y, theta_current)
        theta_current = theta_current - lr*loss_gradient(X, y, theta_current)
        loss_values.append(loss_value)
        theta_values.append(theta_current)

    return theta_current, loss_values, theta_values

theta_init = np.zeros((3, 1))
result = run_gd(loss, loss_gradient, X, y, theta_init, n_iter=10000, lr=1e-10)
theta_est, loss_values, theta_values = result
plt.plot(loss_values)
plt.show()

```

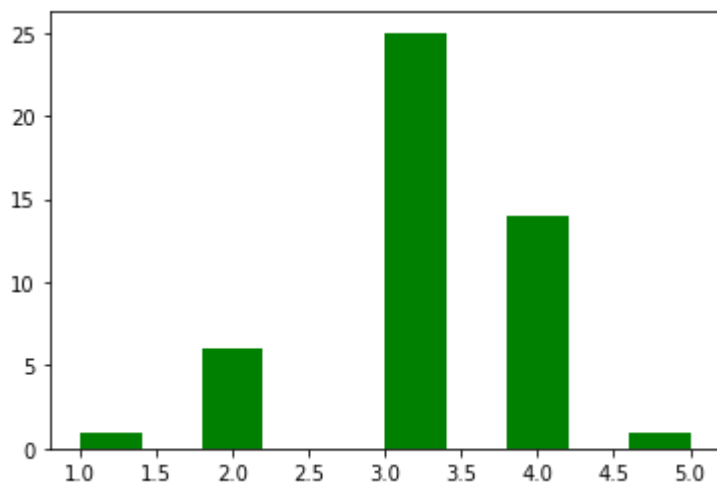
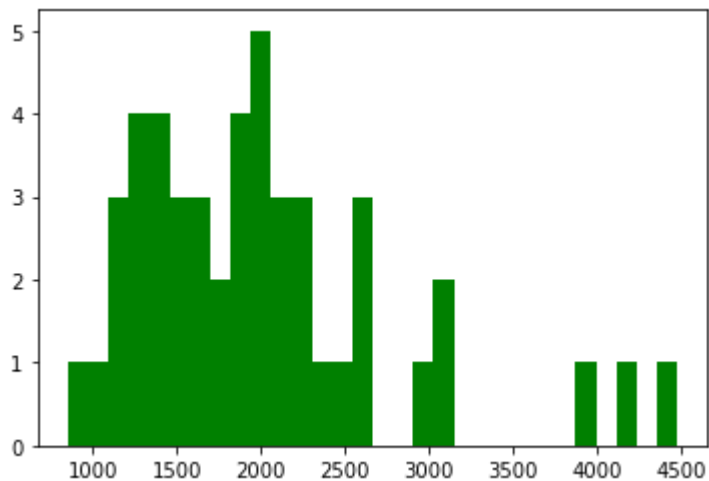


3.2 [20pt] Draw a histogram of values for the first and second feature. Why is feature normalization important? Normalize features and re-run the gradient decent. Compare loss plots that you get with and without feature normalization.

In [13]:

```
plt.hist(X[:,0],bins = 30,facecolor='green')
plt.show()

plt.hist(X[:,1],facecolor='green')
plt.show()
```



```
theta_init = np.zeros((3, 1)) X_normed = preprocessing.scale(X) result = run_gd(loss, loss_gradient, X_normed,
```

```
y, theta_init, n_iter=10000, lr=1e-3) theta_est, loss_values, theta_values = result plt.plot(loss_values) plt.show()
```

Feature normalization can reduce the distance between initial value and the minimum, so that the model will converge more quickly than the non-normalized case, as we can see in the loss value graph.

3.3 [10pt] How can we choose an appropriate learning rate? See what will happen if the learning rate is too small or too large for normalized and not normalized cases?

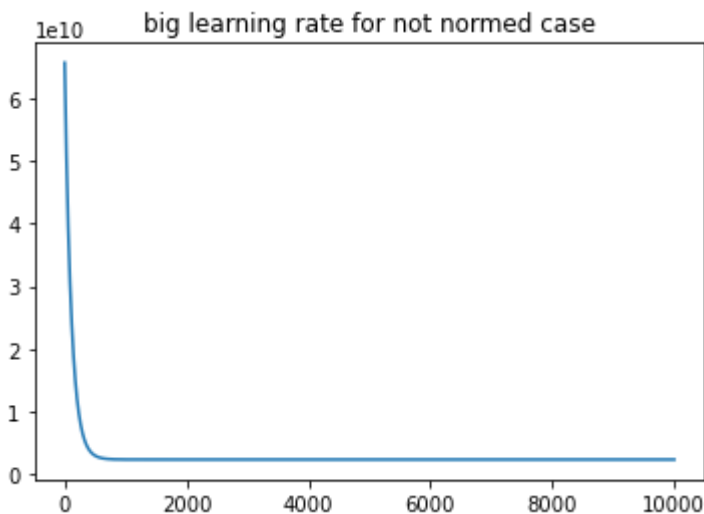
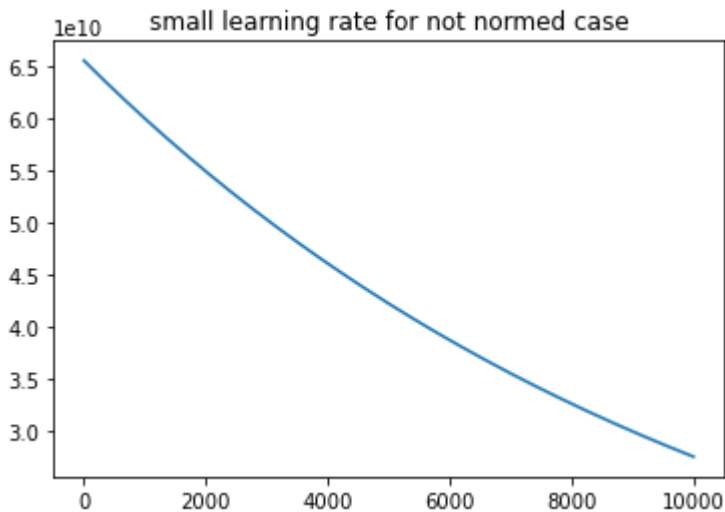
In [14]:

```

#Comparison between big learning rate and small learning rate in not normoalized case
big_result = run_gd(loss, loss_gradient, X, y, theta_init, n_iter=10000, lr=1e-09)
big_theta_est, big_loss_values, big_theta_values = big_result
small_result = run_gd(loss, loss_gradient, X, y, theta_init, n_iter=10000, lr=1e-11)
small_theta_est, small_loss_values, small_theta_values = small_result
plt.plot(small_loss_values)
plt.title('small learning rate for not normed case')
plt.show()
plt.plot(big_loss_values)
plt.title('big learning rate for not normed case')
plt.show()

#Comparison between big learning rate and small learning rate in normoalized case")
normed_big_result = run_gd(loss, loss_gradient, X_normed, y, theta_init, n_iter=10000, lr=1e-09)
normed_big_theta_est, normed_big_loss_values, normed_big_theta_values = normed_big_result
normed_small_result = run_gd(loss, loss_gradient, X_normed, y, theta_init, n_iter=10000, lr=1e-11)
normed_small_theta_est, normed_small_loss_values, normed_small_theta_values = normed_small_result
plt.plot(normed_small_loss_values)
plt.title('small learning rate for normed case')
plt.show()
plt.plot(normed_big_loss_values)
plt.title('big learning rate for normed case')
plt.show()

```

-----

NameError

Traceback (most recent call

```

last)
<ipython-input-14-53c92137666a> in <module>
    12
    13 #Comparison between big learning rate and small learning rate
in normoalized case")
--> 14 normed_big_result = run_gd(loss, loss_gradient, X_normed, y,
    theta_init, n_iter=10000, lr=1e-07)
    15 normed_big_theta_est, normed_big_loss_values, normed_big_theta_
a_values = normed_big_result
    16 normed_small_result = run_gd(loss, loss_gradient, X_normed, y
, theta_init, n_iter=10000, lr=1e-13)

NameError: name 'X_normed' is not defined

```

In non-normalized cases, the learning rate relates to the converge rate. Small learning rate will slower the converging process, while big learning rate will accelerate the converging process. In normalized cases, the effect on converge process when the learning rate is too big or too small is not obvious.

4. Written part

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

4.1 [10 pt] Maximum Likelihood Estimate for Coin Toss

The probability distribution of a single binary variable that takes value with probability is given by the Bernoulli distribution

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}$$

For example, we can use it to model the probability of seeing 'heads' ($x = 1$) or 'tails' ($x = 0$) after tossing a coin, with μ being the probability of seeing 'heads'. Suppose we have a dataset of independent coin flips $D = \{x^{(1)}, \dots, x^{(m)}\}$ and we would like to estimate μ using Maximum Likelihood. Recall that we can write down the likelihood function as

$$\mathcal{L}(x^{(i)}|\mu) = \mu^{x^{(i)}}(1 - \mu)^{1-x^{(i)}}$$

$$P(D|\mu) = \prod_i \mathcal{L}(x^{(i)}|\mu)$$

The log of the likelihood function is

$$\ln P(D|\mu) = \sum_i x^{(i)} \ln \mu + (1 - x^{(i)}) \ln(1 - \mu)$$

Show that the ML solution for μ is given by $\mu_{ML} = \frac{h}{m}$ where h is the total number of 'heads' in the dataset. Show all of your steps.

Solution

Given that h is the total number of 'heads' in the dataset, $(m-h)$ is the total number of 'tails' in the dataset.

$$\ln P(D|\mu) = \sum_i x^{(i)} \ln \mu + (1 - x^{(i)}) \ln(1 - \mu) = h \ln \mu + (m - h) \ln(1 - \mu)$$

Differentiate with μ :

$$\begin{aligned}\frac{\partial}{\partial \mu} \ln P(D|\mu) &= \frac{h}{\mu} - \frac{m-h}{1-\mu} \\ &= \frac{h - \mu m}{\mu(1-\mu)}\end{aligned}$$

Since $\mu \in (0, 1)$, when $\mu \in (0, \frac{h}{m})$, we have $\frac{\partial}{\partial \mu} \ln P(D|\mu) > 0$. When $\mu \in (\frac{h}{m}, 1)$, we have $\frac{\partial}{\partial \mu} \ln P(D|\mu) < 0$.

Therefore, when $\mu_{ML} = \frac{h}{m}$, the likelihood function is maximized, which means the ML solution for μ is given by $\mu_{ML} = \frac{h}{m}$.

4.2 [10 pt] Localized linear regression

Suppose we want to estimate localized linear regression by weighting the contribution of the data points by their distance to the query point x_q , i.e. using the cost

$$E(x_q) = \frac{1}{2} \sum_i^m \frac{(y^{(i)} - h(x^{(i)}|\theta))^2}{||x^{(i)} - x_q||^2}$$

where $\frac{1}{||x^{(i)} - x_q||^2} = w^{(i)}$ is the inverse Euclidean distance between the training point $x^{(i)}$ and query (test) point x_q .

Derive the modified normal equations for the above cost function $E(x_q)$. Hint: first, re-write the cost function in matrix/vector notation, using a diagonal matrix to represent the weights $w^{(i)}$.

your solution

Given that $w^{(i)} = \frac{1}{||x^{(i)} - x_q||^2}$, $E(x_q) = \frac{1}{2} \sum_{i=1}^m w^{(i)} (y^{(i)} - h(x^{(i)}|\theta))^2$.

We have:

$$E(x_q) = \begin{bmatrix} h(x^{(1)}|\theta) - y^{(1)} \\ h(x^{(2)}|\theta) - y^{(2)} \\ \vdots \\ h(x^{(m)}|\theta) - y^{(m)} \end{bmatrix}^T \begin{bmatrix} \frac{1}{2}w^{(1)} & 0 & \dots & 0 \\ 0 & \frac{1}{2}w^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{2}w^{(m)} \end{bmatrix} \begin{bmatrix} h(x^{(1)}|\theta) - y^{(1)} \\ h(x^{(2)}|\theta) - y^{(2)} \\ \vdots \\ h(x^{(m)}|\theta) - y^{(m)} \end{bmatrix}$$

And the weight can be represented as:

$$W = \begin{bmatrix} \frac{1}{2}w^{(1)} & 0 & \dots & 0 \\ 0 & \frac{1}{2}w^{(2)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{2}w^{(m)} \end{bmatrix}$$

Therefore, we can rewrite the cost function as:

$$\begin{aligned}E(X) &= (X\theta - y)^T (W)(X\theta - y) \\ &= W(\theta^T X^T X\theta - \theta^T X^T Y - (Y^T X^T \theta)^T + y^T y)\end{aligned}$$

Take derivative with respect to θ , setting it to 0 to get the minimum:

$$\frac{\partial}{\partial \theta} E(X) = X^T X\theta - X^T y = 0$$

Therefore, we have:

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

4.3 [10 pt] Betting on Trick Coins

A game is played with three coins in a jar: one is a normal coin, one has “heads” on both sides, one has “tails” on both sides. All coins are “fair”, i.e. have equal probability of landing on either side. Suppose one coin is picked randomly from the jar and tossed, and lands with “heads” on top. What is the probability that the bottom side is also “heads”? Show all your steps.

solution

Define event X as "The coin with 'heads' on both sides is picked".

Define event Y as "Pick randomly from the jar and toss, land with 'heads' on top".

According to the Bayesian rule:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

It can be derived that: $P(X|Y) = 1$, $P(X) = \frac{1}{3}$, $P(Y) = \frac{1}{3} * \frac{1}{2} + \frac{1}{3} * 1 = \frac{1}{2}$.

Therefore, we have:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)} = \frac{1 * \frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}$$