

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 (use "Raw" link on [gist\github](https://gist.github.com)), run Anaconda and choose this file in Anaconda's file explorer. Use the `Python 3` version. Everything that follows assumes 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 certain output is given for some cells, that means that you are expected to get similar results.

Total: 155 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 below. [1] (<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>)

```
In [384]: import numpy as np
import matplotlib.pyplot as plt

print(np.array([[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 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>).

```
In [385]: import timeit

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

def mycompute_something(A, x):
    return np.matmul(A, x)

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

```
In [386]: def vectorized(A, x):
            return np.matmul(A, x)

loopstart = timeit.default_timer()
compute_something(A, x)
loopend = timeit.default_timer()

looptime = loopend - loopstart

vecstart = timeit.default_timer()
vectorized(A, x)
vecend = timeit.default_timer()

vectime = vecend - vecstart

difftime = looptime - vectime

print(difftime)

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

0.0007082560041453689
```

```
In [240]: 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('---')
```

```
143 µs ± 9.21 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
The slowest run took 4.77 times longer than the fastest. This could mean
that an intermediate result is being cached.
4.29 µs ± 2.47 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
622 µs ± 72.1 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
2.79 µs ± 1.83 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
41.7 ms ± 8.24 ms per loop (mean ± std. dev. of 7 runs, 5 loops each)
3.25 µs ± 1.38 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
999 ms ± 60 ms per loop (mean ± std. dev. of 7 runs, 5 loops each)
28.8 µs ± 11.8 µs per loop (mean ± std. dev. of 7 runs, 5 loops each)
---
```

2. Linear regression with one variable

In this part of the 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] Generate a plot similar to the one 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.

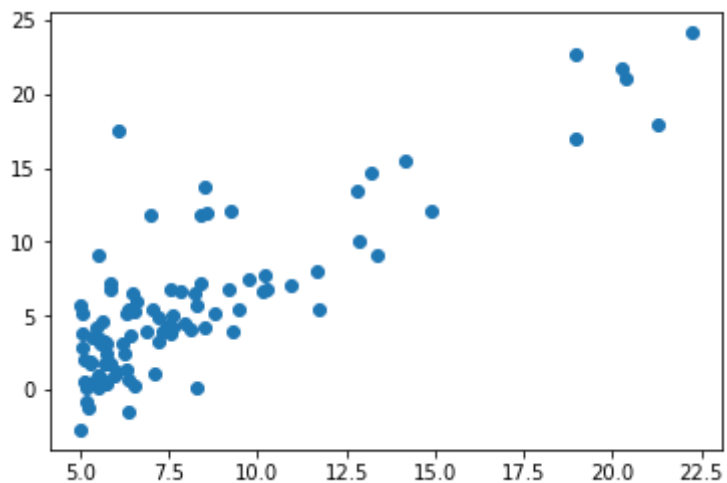
```
In [387]: data = np.loadtxt('ex1data1.txt', delimiter=',')
X, y = data[:, 0, np.newaxis], data[:, 1, np.newaxis]
n = data.shape[0]
print(X.shape, y.shape, n)
print(X[:10], '\n', y[:10])

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 the cost $J(\theta)$. One way to do this is to use the gradient descent. In batch gradient descent algorithm, value of θ is updated iteratively using the gradient of $J(\theta)$.

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

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

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

This update rule comes from the understanding that the moving along the gradient of the loss function will result in a theta with parameters that minimize the loss function, resulting in a stronger fitting linear regression. Mathematically, we arrive at the formula by taking the partial derivative with respect to theta j of $J(\theta)$, allowing us to see how the function changes with respect to each theta. Then, we decide how much we want to change theta by in one iteration by manipulating the learning rate, and subtract the learning rate times the derivative from the current theta. A derivation of the partial derivative with respect to theta j of $J(\theta)$ is shown on the next page.

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.

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} (h_{\theta}(x^i) - y^i)^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{2m} \sum_{i=1}^m 2(h_{\theta}(x^i) - y^i) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x^i) - y^i)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{2m} \sum_{i=1}^m 2(h_{\theta}(x^i) - y^i) \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^m \theta_i x_i - y \right)$$

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

```

In [388]: # 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
    a = np.insert(X, 0, 1, 1)
    return a

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 = np.matmul(X_prime, 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)

    predictedvals = predict(X, theta)
    diffsquared = (y - predictedvals) ** 2
    summation = np.sum(diffsquared)
    loss = summation / (2 * len(predictedvals))

    return loss

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 write your code, 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 $J(\theta)$ 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 alpha to 0.01 .

```

In [389]: import scipy.optimize
from functools import partial

def loss_gradient(X, y, theta):
    X_prime = add_column(X)
    Xsub1 = X_prime[:, 0]
    Xsub2 = X_prime[:, 1]
    mat = np.array([Xsub1, Xsub2])

    predictedvals = predict(X, theta)
    diff = predictedvals - y
    diffmat = diff.transpose()
    diffmat = np.array([diffmat[0], diffmat[0]])

    difftimesx = np.multiply(mat, diffmat)

    # print(difftimesx)

    summation = difftimesx.sum(axis=1)

    summation = np.array([summation]).T

    loss_grad = summation / len(y)

    return loss_grad

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)

    lossgrad = loss_gradient(X, y, theta_current)

    subterm = lossgrad * lr

    theta_current = theta_current - subterm

    loss_values.append(loss_value)
    theta_values.append(theta_current)

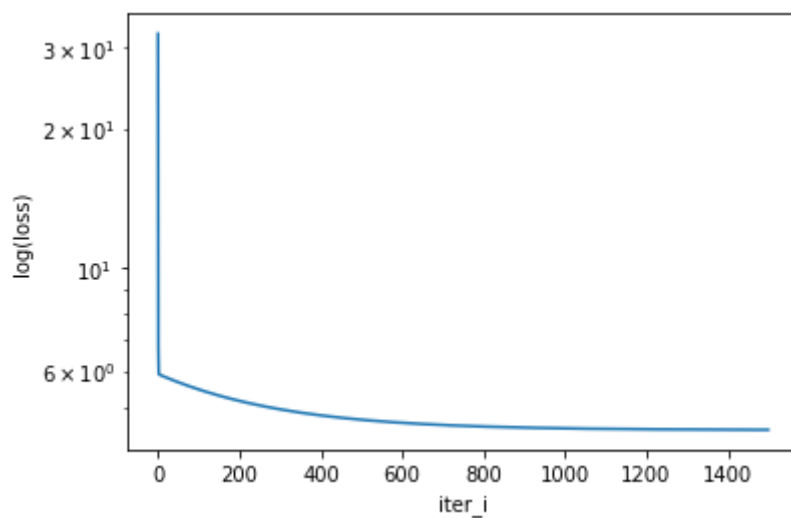
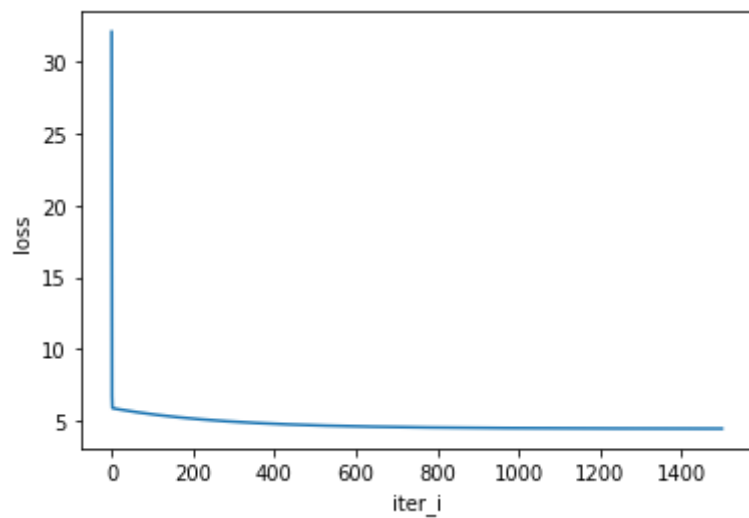
    return theta_current, loss_values, theta_values

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 3.389840206580175e-05  
estimated theta value [-3.63029144  1.16636235]  
resulting loss 4.483388256587726
```



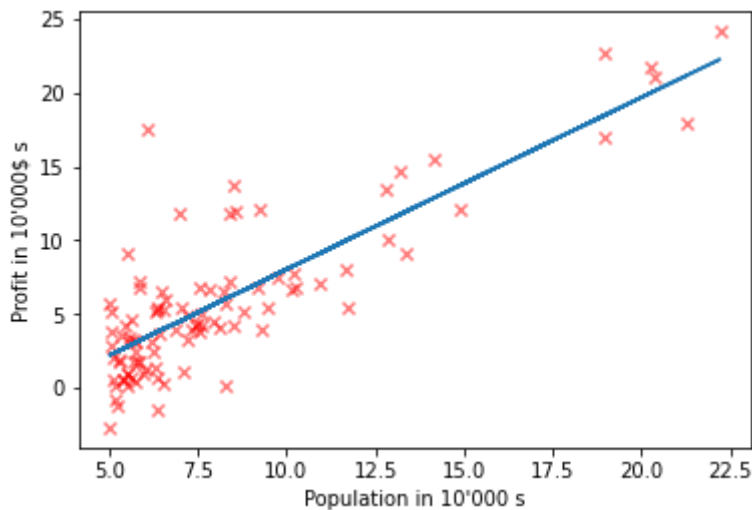
2.2.4 [10pt] After you are finished, use your final parameters to plot the linear fit. The result should look something like on the figure below. Use the `predict()` function.

```
In [390]: plt.scatter(X, y, marker='x', color='r', alpha=0.5)
           x_start, x_end = 5, 25

           prediction = predict(X, np.array([[-3.63029144, 1.16636235]]).T)

           plt.plot(X, prediction)

           plt.xlabel('Population in 10\'000 s')
           plt.ylabel('Profit in 10\'000$ s')
           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 [391]: Xval = np.array([[35000, 70000]]).T

           thet = np.array([[-3.63029144, 1.16636235]]).T

           prediction = predict(Xval, thet)

           print(prediction)

[[40819.05195856]
 [81641.73420856]]
```

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.

```

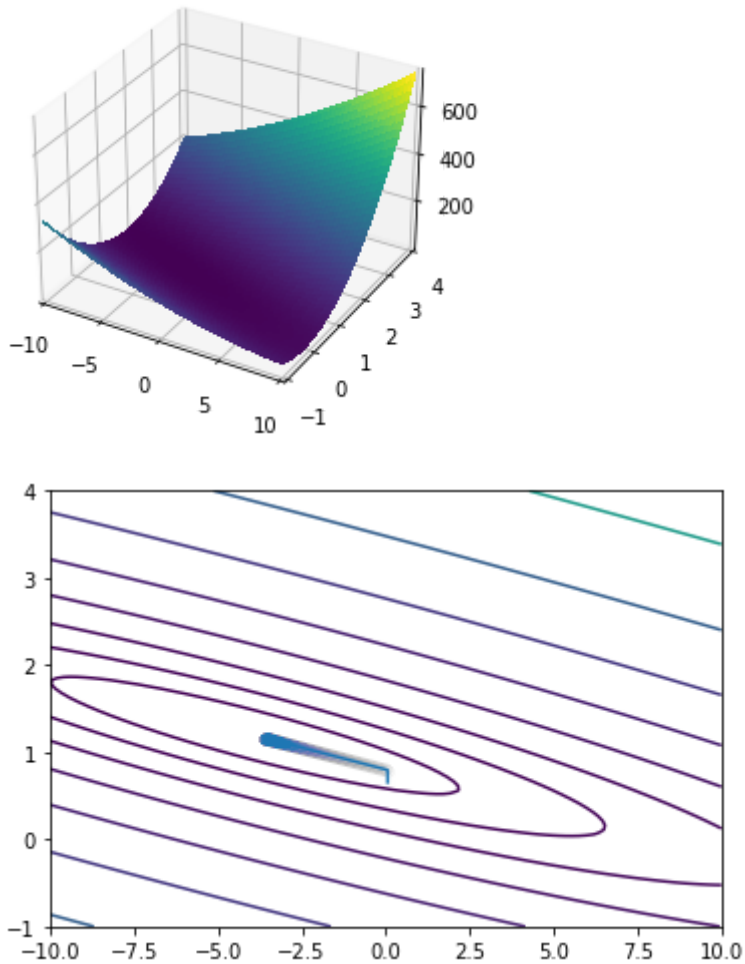
In [392]: from mpl_toolkits.mplot3d import Axes3D
import matplotlib.cm as cm
limits = [(-10, 10), (-1, 4)]
space = [np.linspace(*limit, 100) for limit in limits]
theta_1_grid, theta_2_grid = np.meshgrid(*space)
theta_meshgrid = np.vstack([theta_1_grid.ravel(), theta_2_grid.ravel()])
loss_test_vals_flat = (((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.virid
is,
                                     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.1
ogspace(-2, 3, 20))
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 [393]: data_new = np.loadtxt('ex1data2.txt', delimiter=',')
X_new, y_new = data_new[:, :-1], data_new[:, -1, np.newaxis]
n_new = data.shape[0]
print(X_new.shape, y_new.shape, n_new)
print(X_new[:10], '\n', y_new[:10])
```

```
(47, 2) (47, 1) 97
[[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 [394]: def add_column_new(X):
            a = np.insert(X, 0, 1, 1)
            return a

def predict_new(X, theta):
    """ Computes  $h(x; \theta)$  """
    X_prime = add_column_new(X)
    pred = np.matmul(X_prime, theta)
    return pred

def loss_new(X, y, theta):
    X_prime = add_column_new(X)
    predictedvals = predict_new(X, theta)
    diffsquared = (y - predictedvals) ** 2
    summation = np.sum(diffsquared)
    loss = summation / (2 * len(predictedvals))

    return loss

def loss_gradient_new(X, y, theta):
    X_prime = add_column_new(X)
    mat = X_prime.T

    predictedvals = predict_new(X, theta)

    diff = predictedvals - y
    diffmat = diff.transpose()
    arr = diffmat[0]

    numrows = theta.shape[0]

    diffmat = np.tile(arr, (numrows, 1))

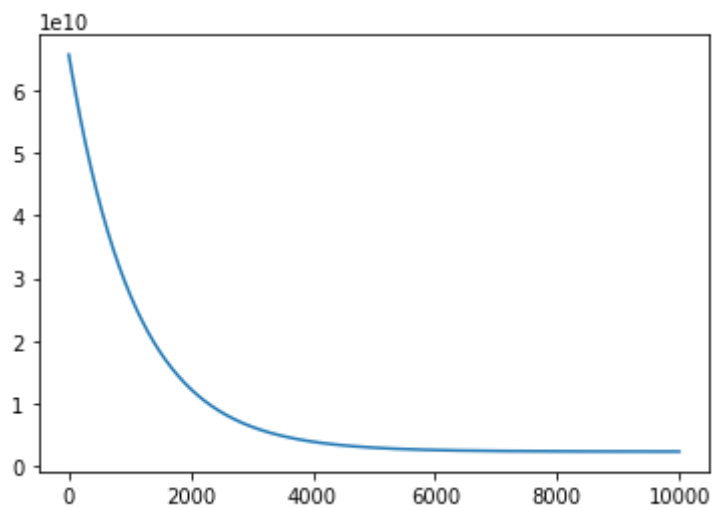
    difftimesx = np.multiply(mat, diffmat)

    summation = difftimesx.sum(axis=1)

    summation = np.array([summation]).T
    # print(summation)
    loss_grad = summation / len(y)
    # print(loss_grad)
    # print('AAAA')
    return loss_grad

theta_init_new = np.zeros((3, 1))
result = run_gd(loss_new, loss_gradient_new, X_new, y_new, theta_init_new,
                n_iter=10000, lr=1e-10)
theta_est, loss_values, theta_values = result
plt.plot(loss_values)
plt.show()

```

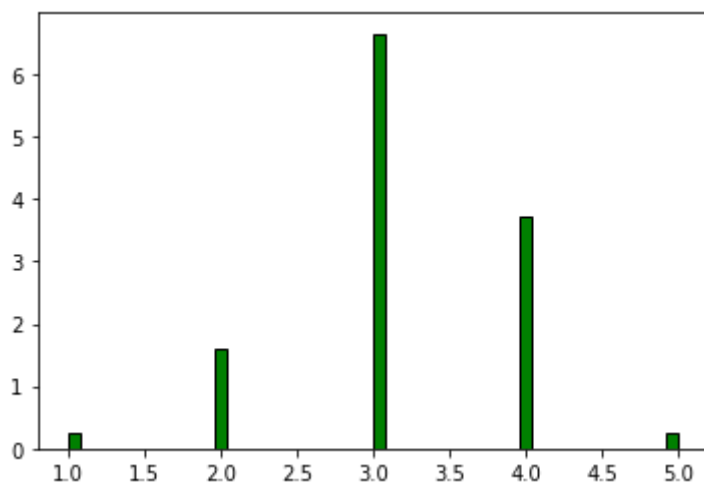
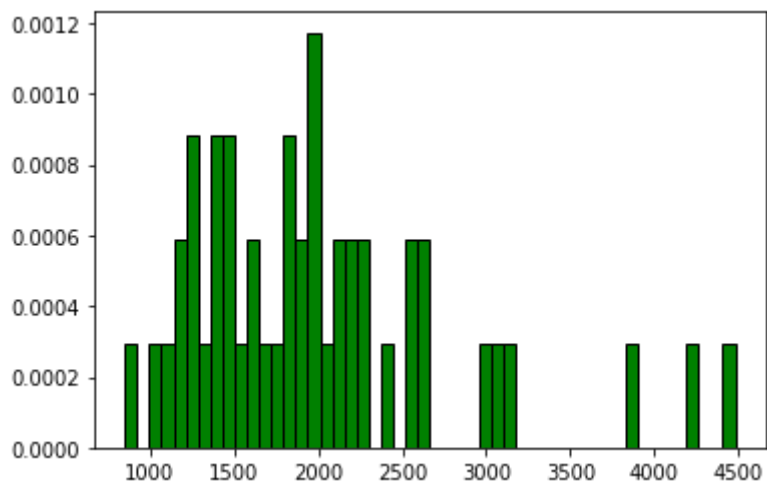


3.2 [20pt] Draw a histogram of values for the first and second feature. Why is feature normalization important? Normalize features and re-run the gradient decent. Compare loss plots that you get with and without feature normalization.


```
In [416]: feature1 = X_new[:, 0]
feature2 = X_new[:, 1]

plt.hist(feature1, bins=50, density = 'true', color = 'green', edgecolor
='black')
plt.show()

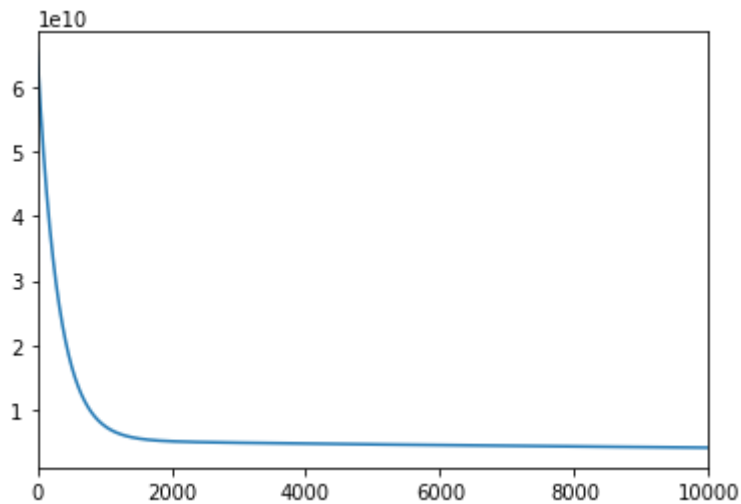
plt.hist(feature2, bins=50, density='true', color = 'green', edgecolor=
'black', align='mid')
plt.show()
```



```
In [446]: theta_init_norm = np.zeros((3, 1))
X_normed = X_new / X_new.max(axis=0)

result = run_gd(loss_new, loss_gradient_new, X_normed, y_new, theta_init_norm, n_iter=10000, lr=1e-3)
theta_est, loss_values, theta_values = result

plt.plot(loss_values)
plt.xlim([0, 10000])
plt.show()
```



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 [496]: print('NOT NORMALIZED CASES:')
          print('test 1: learning rate too small')

          result = run_gd(loss_new, loss_gradient_new, X_new, y_new, theta_init_new, n_iter=10000, lr=1e-24)
          theta_est, loss_values, theta_values = result
          plt.plot(loss_values)
          plt.xlim([0, 10000])
          plt.show()

          print('test 2: learning rate too big')

          result = run_gd(loss_new, loss_gradient_new, X_new, y_new, theta_init_new, n_iter=10000, lr=1e-6*.45)
          theta_est, loss_values, theta_values = result
          plt.plot(loss_values)
          plt.xlim([0, 10000])
          plt.show()

          print('NORMALIZED CASES:')
          print('test 1: learning rate too small')

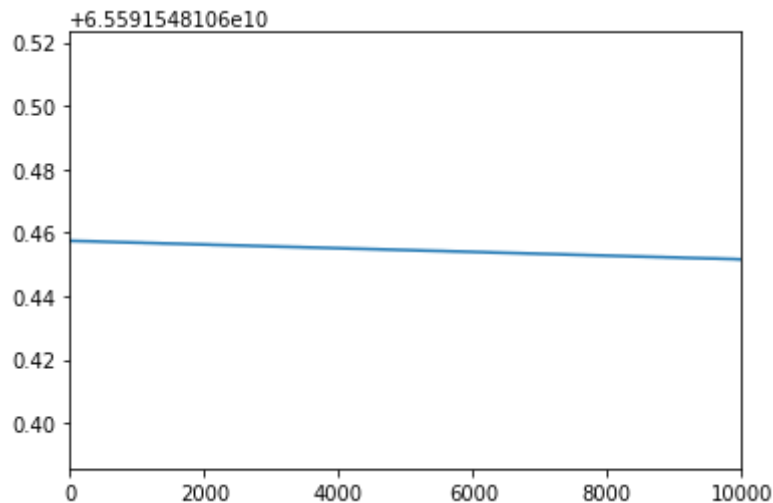
          result = run_gd(loss_new, loss_gradient_new, X_normed, y_new, theta_init_norm, n_iter=10000, lr=1e-18)
          theta_est, loss_values, theta_values = result
          plt.plot(loss_values)
          plt.xlim([0, 10000])
          plt.show()

          print('test 2: learning rate too big')

          result = run_gd(loss_new, loss_gradient_new, X_normed, y_new, theta_init_norm, n_iter=10000, lr=1.5)
          theta_est, loss_values, theta_values = result
          plt.plot(loss_values)
          plt.xlim([0, 10000])
          plt.show()
```

NOT NORMALIZED CASES:

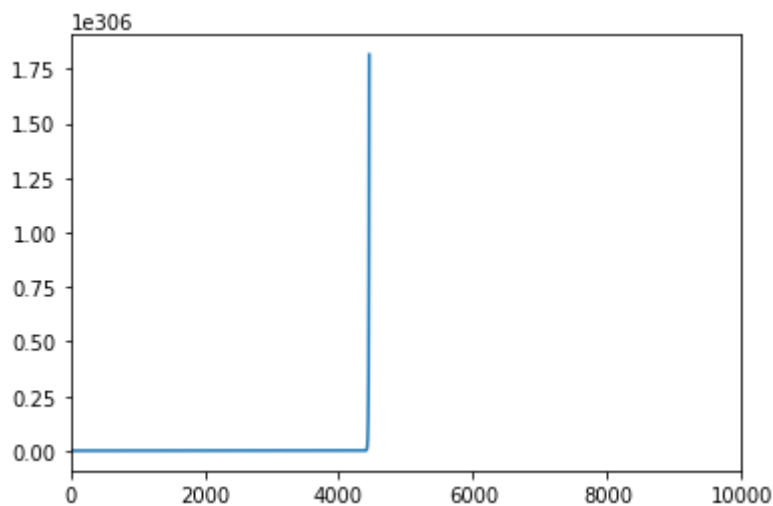
test 1: learning rate too small



test 2: learning rate too big

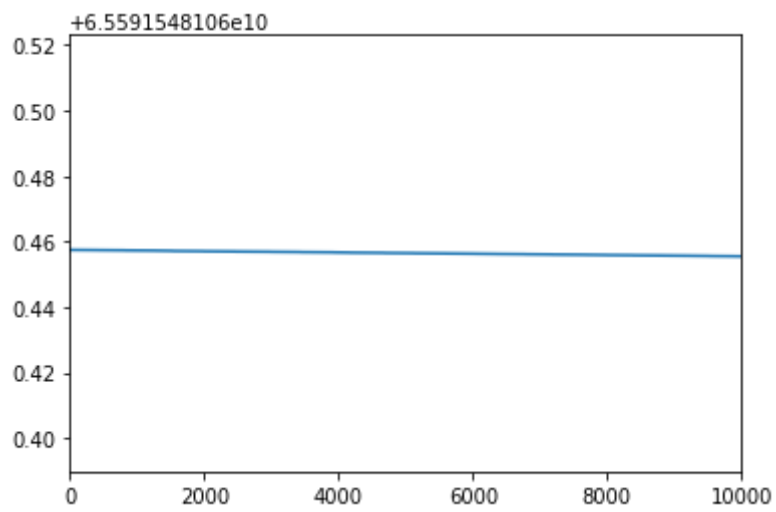
/opt/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:14: RuntimeWarning: overflow encountered in square

/opt/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:63: RuntimeWarning: invalid value encountered in subtract

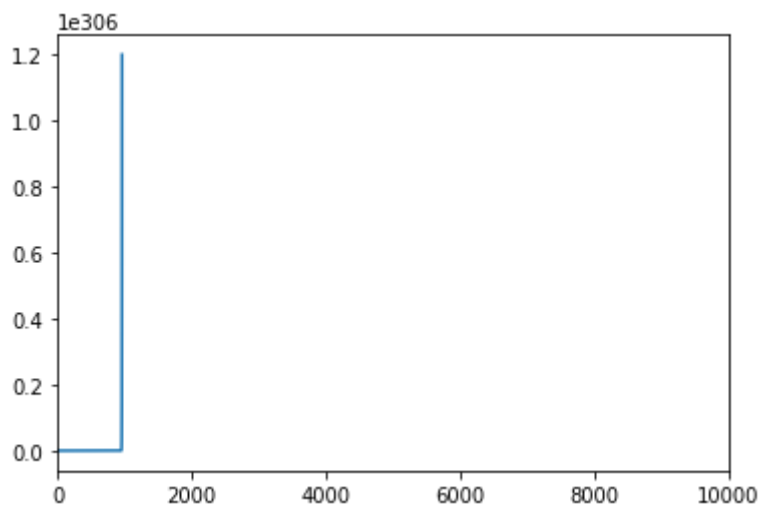


NORMALIZED CASES:

test 1: learning rate too small



test 2: learning rate too big



In []: