

Programming Exercise 1: Linear Regression

Introduction

In this exercise, you will implement linear regression and get to see it work on data. All the information you need for solving this assignment is in this notebook, and all the code you will be implementing will take place within this notebook.

Before we begin with the exercises, we need to import all libraries required for this programming exercise. Throughout the course, we will be using `numpy` [\(http://www.numpy.org/\)](http://www.numpy.org/) for all arrays and matrix operations, and `matplotlib` [\(https://matplotlib.org/\)](https://matplotlib.org/) for plotting.

In [2]:

```
# used for manipulating directory paths
import os

# Scientific and vector computation for python
import numpy as np

# Plotting library
from matplotlib import pyplot
from mpl_toolkits.mplot3d import Axes3D # needed to plot 3-D surfaces

# tells matplotlib to embed plots within the notebook
%matplotlib inline
```

Debugging

Here are some things to keep in mind throughout this exercise:

- Python array indices start from zero, not one (contrary to OCTAVE/MATLAB).
- There is an important distinction between python arrays (called `list` or `tuple`) and `numpy` arrays. You should use `numpy` arrays in all your computations. Vector/matrix operations work only with `numpy` arrays. Python lists do not support vector operations (you need to use for loops).
- If you are seeing many errors at runtime, inspect your matrix operations to make sure that you are adding and multiplying matrices of compatible dimensions. Printing the dimensions of `numpy` arrays using the `shape` property will help you debug.
- By default, `numpy` interprets math operators to be element-wise operators. If you want to do matrix multiplication, you need to use the `dot` function in `numpy`. For example if `A` and `B` are two `numpy` matrices, then the matrix operation `AB` is `np.dot(A, B)`. Note that for 2-dimensional matrices or vectors (1-dimensional), this is also equivalent to `A@B` (requires python ≥ 3.5).

1 Simple python and numpy function

The first part of this assignment gives you practice with python and `numpy` syntax. In the next cell, you will find the outline of a python function. Modify it to return a 5 x 5 identity matrix by filling in the following code:

```
A = np.eye(5)
```

In [3]:

```
def warmUpExercise():
    """
    Example function in Python which computes the identity matrix.

    Returns
    -----
    A : array_like
        The 5x5 identity matrix.

    Instructions
    -----
    Return the 5x5 identity matrix.
    """
    # ===== YOUR CODE HERE =====
    A = np.identity(5)

    # =====
    return A
```

The previous cell only defines the function `warmUpExercise`. We can now run it by executing the following cell to see its output. You should see output similar to the following:

```
array([[ 1.,  0.,  0.,  0.,  0.],
       [ 0.,  1.,  0.,  0.,  0.],
       [ 0.,  0.,  1.,  0.,  0.],
       [ 0.,  0.,  0.,  1.,  0.],
       [ 0.,  0.,  0.,  0.,  1.]])
```

In [4]:

```
warmUpExercise()
```

Out[4]:

```
array([[1., 0., 0., 0., 0.],
       [0., 1., 0., 0., 0.],
       [0., 0., 1., 0., 0.],
       [0., 0., 0., 1., 0.],
       [0., 0., 0., 0., 1.]])
```

2 Linear regression with one variable

Now 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 `Data/ex1data1.txt` contains the dataset for our linear regression problem. The **first column is the population** of a city (in 10,000s) and the **second column is the profit** of a food truck in that city (in \$10,000s). A negative value for profit indicates a loss.

We provide you with the code needed to load this data. The dataset is loaded from the data file into the variables `x` and `y` :

In [5]:

```
# Read comma separated data
data = np.loadtxt(os.path.join('Data', 'ex1data1.txt'), delimiter=',')
X_raw, y = data[:, 0], data[:, 1]

N = y.size # number of training examples
```

Print the shapes of the data

In [6]:

```
print('Dimension of data matrix X_raw:', X_raw.shape)
print('Dimension of the output y:', N)
```

```
Dimension of data matrix X_raw: (97,)
Dimension of the output y: 97
```

2.1 Plotting the Data

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 cannot be plotted on a 2-d plot. There are many plotting libraries in python (see this [blog post \(https://blog.modeanalytics.com/python-data-visualization-libraries/\)](https://blog.modeanalytics.com/python-data-visualization-libraries/) for a good summary of the most popular ones).

In this course, we will be exclusively using `matplotlib` to do all our plotting. `matplotlib` is one of the most popular scientific plotting libraries in python and has extensive tools and functions to make beautiful plots. `pyplot` is a module within `matplotlib` which provides a simplified interface to `matplotlib`'s most common plotting tasks, mimicking MATLAB's plotting interface.

You might have noticed that we have imported the `'pyplot'` module at the beginning of this exercise using the command `'from matplotlib import pyplot'`. This is rather uncommon, and if you look at python code elsewhere or in the `'matplotlib'` tutorials, you will see that the module is named `'plt'`. This is used by module renaming by using the import command `'import matplotlib.pyplot as plt'`. We will not using the short name of `'pyplot'` module in this class exercises, but you should be aware of this deviation from norm.

In the following part, your first job is to complete the `plotData` function below. Modify the function and fill in the following code:

```
pyplot.plot(x, y, 'ro', ms=10, mec='k')
pyplot.ylabel('Profit in $10,000')
pyplot.xlabel('Population of City in 10,000s')
```

In [7]:

```
def plotData(x, y):
    """
    Plots the data points x and y into a new figure. Plots the data
    points and gives the figure axes labels of population and profit.

    Parameters
    -----
    x : array_like
        Data point values for x-axis.

    y : array_like
        Data point values for y-axis. Note x and y should have the same size.

    Instructions
    -----
    Plot the training data into a figure using the "figure" and "plot"
    functions. Set the axes labels using the "xlabel" and "ylabel" functions.
    Assume the population and revenue data have been passed in as the x
    and y arguments of this function.

    Hint
    ----
    You can use the 'ro' option with plot to have the markers
    appear as red circles. Furthermore, you can make the markers larger by
    using plot(..., 'ro', ms=10), where `ms` refers to marker size. You
    can also set the marker edge color using the `mec` property.
    """
    fig = pyplot.figure() # open a new figure

    # ===== YOUR CODE HERE =====

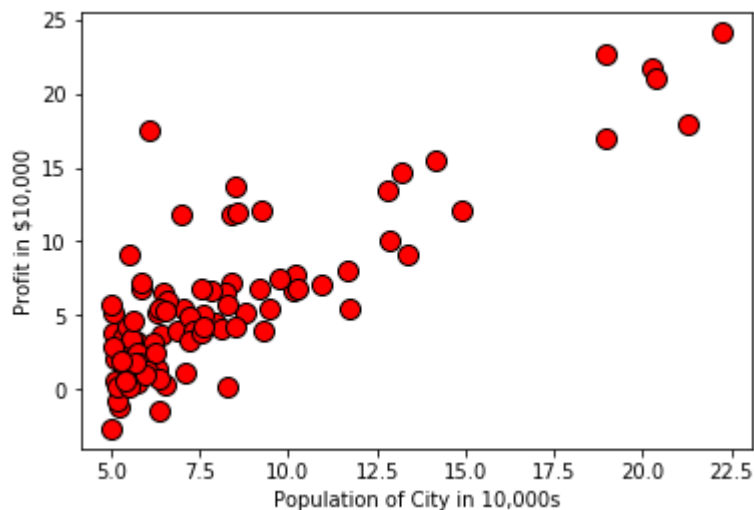
    pyplot.plot(x, y, 'ro', ms=10, mec='k')
    pyplot.ylabel('Profit in $10,000')
    pyplot.xlabel('Population of City in 10,000s')

    # =====
```

Execute the next cell to plot the data

In [8]:

```
plotData(X_raw, y)
```



To quickly learn more about the `matplotlib` `plot` function and what arguments you can provide to it, you can type `?pyplot.plot` in a cell within the jupyter notebook. This opens a separate page showing the documentation for the requested function. You can also search online for plotting documentation.

To set the markers to red circles, we used the option `'or'` within the `plot` function.

In [9]:

```
?pyplot.plot
```

2.2 Gradient Descent

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

2.2.1 Update Equations

The objective of linear regression is to minimize the cost function

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

where the hypothesis $h_{\theta}(x)$ is given by the linear model

$$h_{\theta}(x) = x^T \theta = \theta_0 + \theta_1 x_1$$

So that

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_1(i) - y(i))^2$$

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 batch gradient descent algorithm. In batch gradient descent, each iteration performs the update

$$\theta_j = \theta_j - \alpha \frac{2}{N} \sum_{i=1}^N (h_{\theta}(x(i)) - y(i)) x_j(i) \quad \text{simultaneously update } \theta_j \text{ for all } j$$

In fact, in this univariate case, we have that

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_1(i) - y(i))$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_1(i) - y(i)) x_1(i)$$

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

****Implementation Note:**** We store each example as a row in the the X matrix in Python `numpy`. To take into account the intercept term (θ_0), we add an additional first column to X and set it to all ones. This allows us to treat θ_0 as simply another 'feature'.

2.2.2 Implementation

We have already set up the data for linear regression. In the following cell, we add another dimension to our data to accommodate the θ_0 intercept term.

In [12]:

```
# Add a column of ones to X_raw. The numpy function stack joins arrays along a given axis.
# The first axis (axis=0) refers to rows (training examples)
# and second axis (axis=1) refers to columns (features).
X = np.stack([np.ones(N), X_raw], axis=1)
print('Dimension of data matrix X:', X.shape, '\n')
print('First 4 samples of data matrix X:\n', X[0:4, :])
```

Dimension of data matrix X: (97, 2)

First 4 samples of data matrix X:

```
[[1.    6.1101]
 [1.    5.5277]
 [1.    8.5186]
 [1.    7.0032]]
```

2.2.3 Computing the cost $J(\theta)$

As you perform gradient descent to learn minimize the cost function $J(\theta)$, 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. Remember that the variables X and y are not scalar values. X is a matrix whose rows represent the examples from the training set and y is a vector whose each element represent the value at a given row of X .

In [13]:

```
def computeCost(X, y, theta):
    """
    Compute cost for linear regression. Computes the cost of using theta as the
    parameter for linear regression to fit the data points in X and y.

    Parameters
    -----
    X : array_like
        The input dataset of shape (N x d), where m is the number of examples,
        and d-1 is the number of features. We assume a vector of one's already
        appended to the features so we have d columns.

    y : array_like
        The values of the function at each data point. This is a vector of
        shape (N, ).

    theta : array_like
        The parameters for the regression function. This is a vector of
        shape (d, ).

    Returns
    -----
    J : float
        The value of the regression cost function.

    Instructions
    -----
    Compute the cost of a particular choice of theta.
    You should set J to the cost.
    """

    # initialize some useful values
    N = y.size # number of training examples

    # You need to return the following variables correctly
    J = 0

    # ===== YOUR CODE HERE =====

    J = 1/(N)*np.sum(np.square(np.dot(X, theta) - y)) #single instance implementation
    #J = 1/N * np.linalg.norm( np.dot(X, theta) - y )**2 #matrix implementation

    # =====
    return J
```

Once you have completed the function, the next step will run `computeCost` two times using two different initializations of θ . You will see the cost printed to the screen.

In [14]:

```
J = computeCost(X, y, theta=np.array([0.0, 0.0]))
print('With theta = [0, 0] \nCost computed = %.2f' % J)
print('Expected cost value (approximately) 64.15\n')

# further testing of the cost function
J = computeCost(X, y, theta=np.array([-1, 2]))
print('With theta = [-1, 2]\nCost computed = %.2f' % J)
print('Expected cost value (approximately) 108.48')
```

With theta = [0, 0]
 Cost computed = 64.15
 Expected cost value (approximately) 64.15

With theta = [-1, 2]
 Cost computed = 108.48
 Expected cost value (approximately) 108.48

2.2.4 Gradient descent

Next, you will complete a function which implements gradient descent.

As you program, make sure you understand what you are trying to optimize and what is being updated. Keep in mind that the cost $J(\theta)$ 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.

The starter code for the function `gradientDescent` calls `computeCost` on every iteration and saves the cost to a python list. Assuming you have implemented gradient descent and `computeCost` correctly, your value of $J(\theta)$ should never increase, and should converge to a steady value by the end of the algorithm.

****Vectors and matrices in `numpy`** - Important implementation notes**

A vector in `numpy` is a one dimensional array, for example `np.array([1, 2, 3])` is a vector. A matrix in `numpy` is a two dimensional array, for example `np.array([[1, 2, 3], [4, 5, 6]])`. However, the following is still considered a matrix `np.array([1, 2, 3])` since it has two dimensions, even if it has a shape of 1x3 (which looks like a vector).

Given the above, the function `np.dot` which we will use for all matrix/vector multiplication has the following properties:

- It always performs inner products on vectors. If `x=np.array([1, 2, 3])`, then `np.dot(x, x)` is a scalar.
- For matrix-vector multiplication, so if X is a $N \times d$ matrix and y is a vector of length N , then the operation `np.dot(y, X)` considers y as a $1 \times N$ vector. On the other hand, if y is a vector of length d , then the operation `np.dot(X, y)` considers y as a $d \times 1$ vector.
- A vector can be promoted to a matrix using `y[None]` or `[y[np.newaxis]]`. That is, if $y = \text{np.array}([1, 2, 3])$ is a vector of size 3, then `y[None, :]` is a matrix of shape 1×3 . We can use `y[:, None]` to obtain a shape of 3×1 .

In [15]:

```
def gradientDescent(X, y, theta, alpha, num_iters):
    """
    Performs gradient descent to learn `theta`. Updates theta by taking `num_iters`
    gradient steps with learning rate `alpha`.

    Parameters
    -----
    X : array_like
        The input dataset of shape (N x d).

    y : array_like
        Value at given features. A vector of shape (N, ).

    theta : array_like
        Initial values for the linear regression parameters.
        A vector of shape (d, ).

    alpha : float
        The learning rate.

    num_iters : int
        The number of iterations for gradient descent.

    Returns
    -----
    theta : array_like
        The learned linear regression parameters. A vector of shape (d, ).

    J_history : list
        A python list for the values of the cost function after each iteration.

    Instructions
    -----
    Perform a single gradient step on the parameter vector theta.

    While debugging, it can be useful to print out the values of
    the cost function (computeCost) and gradient here.
    """
    # Initialize some useful values
    N = y.shape[0] # number of training examples

    # make a copy of theta, to avoid changing the original array, since numpy arrays
    # are passed by reference to functions
    theta = theta.copy()

    J_history = [] # Use a python list to save cost in every iteration

    for i in range(num_iters):
        # ===== YOUR CODE HERE =====

        theta = theta - alpha*2/(N) * np.dot( X.T, ( np.dot(X, theta) - y ));

        # =====

        # save the cost J in every iteration
        J_history.append( computeCost(X, y, theta) )

    return theta, J_history
```

After you are finished call the implemented `gradientDescent` function and print the computed θ . We initialize the θ parameters to 0 and the learning rate α to 0.01. Execute the following cell to check your code.

In [17]:

```
# initialize fitting parameters
theta = np.zeros(2)

# some gradient descent settings
iterations = 1500
alpha = 0.005

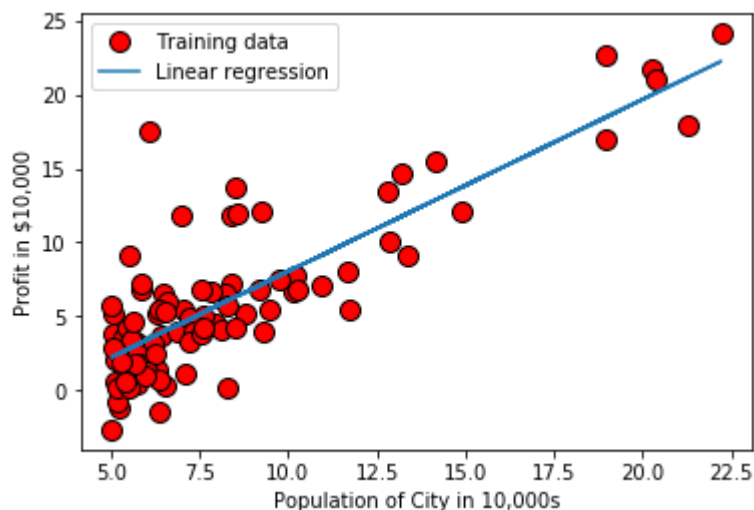
theta, J_history = gradientDescent(X, y, theta, alpha, iterations)
print('Theta found by gradient descent: {:.4f}, {:.4f}'.format(*theta))
print('Expected theta values (approximately): [-3.6303, 1.1664]')
```

Theta found by gradient descent: -3.6303, 1.1664
Expected theta values (approximately): [-3.6303, 1.1664]

We will use your final parameters to plot the linear fit.

In [18]:

```
# plot the linear fit
plotData(X[:, 1], y)
pyplot.plot(X[:, 1], np.dot(X, theta), '-')
pyplot.legend(['Training data', 'Linear regression']);
```



Your final values for θ will also be used to make predictions on profits in areas of 35,000 and 70,000 people.

Note the way that the following lines use matrix multiplication, rather than explicit summation or looping, to calculate the predictions. This is an example of code vectorization in 'numpy'.

Note that the first argument to the 'numpy' function 'dot' is a python list. 'numpy' can internally convert **valid** python lists to numpy arrays when explicitly provided as arguments to 'numpy' functions.

In [19]:

```
# Predict values for population sizes of 35,000 and 70,000
predict1 = np.dot([1, 3.5], theta)
print('For population = 35,000, we predict a profit of {:.2f}\n'.format(predict1*10000))

predict2 = np.dot([1, 7], theta)
print('For population = 70,000, we predict a profit of {:.2f}\n'.format(predict2*10000))
```

For population = 35,000, we predict a profit of 4519.77

For population = 70,000, we predict a profit of 45342.45

2.4 Visualizing $J(\theta)$

To understand the cost function $J(\theta)$ better, you will now plot the cost over a 2-dimensional grid of θ_0 and θ_1 values.

In the next cell, the code is set up to calculate $J(\theta)$ over a grid of values using the `computeCost` function that you wrote. After executing the following cell, you will have a 2-D array of $J(\theta)$ values. Then, those values are used to produce surface and contour plots of $J(\theta)$ using the matplotlib `plot_surface` and `contourf` functions.

The purpose of these graphs is to show you how $J(\theta)$ varies with changes in θ_0 and θ_1 . The cost function $J(\theta)$ is bowl-shaped and has a global minimum. (This is easier to see in the contour plot than in the 3D surface plot). This minimum is the optimal point for θ_0 and θ_1 , and each step of gradient descent moves closer to this point.

In [20]:

```

# grid over which we will calculate J
theta0_vals = np.linspace(-10, 10, 100)
theta1_vals = np.linspace(-1, 4, 100)

# initialize J_vals to a matrix of 0's
J_vals = np.zeros((theta0_vals.shape[0], theta1_vals.shape[0]))

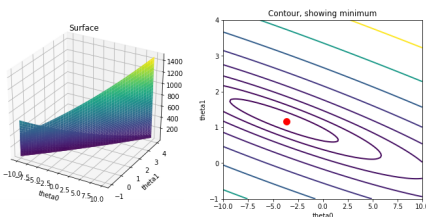
# Fill out J_vals
for i, theta0 in enumerate(theta0_vals):
    for j, theta1 in enumerate(theta1_vals):
        J_vals[i, j] = computeCost(X, y, [theta0, theta1])

# Because of the way meshgrids work in the surf command, we need to
# transpose J_vals before calling surf, or else the axes will be flipped
J_vals = J_vals.T

# surface plot
fig = pyplot.figure(figsize=(12, 5))
ax = fig.add_subplot(121, projection='3d')
ax.plot_surface(theta0_vals, theta1_vals, J_vals, cmap='viridis')
pyplot.xlabel('theta0')
pyplot.ylabel('theta1')
pyplot.title('Surface')

# contour plot
# Plot J_vals as 15 contours spaced logarithmically between 0.01 and 100
ax = pyplot.subplot(122)
pyplot.contour(theta0_vals, theta1_vals, J_vals, linewidths=2, cmap='viridis', levels=np.logspace(0.01, 2, 15))
pyplot.xlabel('theta0')
pyplot.ylabel('theta1')
pyplot.plot(theta[0], theta[1], 'ro', ms=10, lw=2)
pyplot.title('Contour, showing minimum')
pass

```



3 Linear regression with multiple variables

In this part, you will implement linear regression with multiple variables to **predict the prices of houses**. Suppose you are selling your house and you want to know what a good market price would be. One way to do this is to first collect information on recent houses sold and make a model of housing prices.

The file `Data/ex1data2.txt` contains a training set of housing prices in Portland, Oregon. The **first column is the size of the house (in square feet)**, the **second column is the number of bedrooms**, and the **third column is the price of the house**.

3.1 Feature Normalization

We start by loading and displaying some values from this dataset. By looking at the values, note that house sizes are about 1000 times the number of bedrooms. When features differ by orders of magnitude, first performing **feature scaling** can make gradient descent converge much more quickly.

In [21]:

```
# Load data
data = np.loadtxt(os.path.join('Data', 'ex1data2.txt'), delimiter=',')
X = data[:, :2]
y = data[:, 2]
N = y.size

# print out some data points
print('{:>8s}{:>8s}{:>10s}'.format('X[:,0]', 'X[:, 1]', 'y'))
print('-'*26)
for i in range(10):
    print('{:>8.0f}{:>8.0f}{:>10.0f}'.format(X[i, 0], X[i, 1], y[i]))
```

X[:,0]	X[:, 1]	y
2104	3	399900
1600	3	329900
2400	3	369000
1416	2	232000
3000	4	539900
1985	4	299900
1534	3	314900
1427	3	198999
1380	3	212000
1494	3	242500

The code in `featureNormalize` function:

- Subtracts the mean value of each feature from the dataset.
- After subtracting the mean, additionally scales (divide) the feature values by their respective “standard deviations.”

The standard deviation is a way of measuring how much variation there is in the range of values of a particular feature (most data points will lie within ± 2 standard deviations of the mean); this is an alternative to taking the range of values (max-min). In `numpy`, you can use the `std` function to compute the standard deviation.

For example, the quantity `X[:, 0]` contains all the values of x_1 (house sizes) in the training set, so `np.std(X[:, 0])` computes the standard deviation of the house sizes. At the time that the function `featureNormalize` is called, the extra column of 1's corresponding to $x_0 = 1$ has not yet been added to `X`.

You will do this for all the features and your code should work with datasets of all sizes (any number of features / examples). Note that each column of the matrix `X` corresponds to one feature.

****Implementation Note:**** When normalizing the features, it is important to store the values used for normalization - the mean value and the standard deviation used for the computations. After learning the parameters from the model, we often want to predict the prices of houses we have not seen before. Given a new `x` value (living room area and number of bedrooms), we must first normalize `x` using the mean and standard deviation that we had previously computed from the training set.

In [22]:

```

def featureNormalize(X):
    """
    Normalizes the features in X. returns a normalized version of X where
    the mean value of each feature is 0 and the standard deviation
    is 1. This is often a good preprocessing step to do when working with
    learning algorithms.

    Parameters
    -----
    X : array_like
        The dataset of shape (N x d-1).

    Returns
    -----
    X_norm : array_like
        The normalized dataset of shape (N x d-1).

    Instructions
    -----
    First, for each feature dimension, compute the mean of the feature
    and subtract it from the dataset, storing the mean value in mu.
    Next, compute the standard deviation of each feature and divide
    each feature by it's standard deviation, storing the standard deviation
    in sigma.

    Note that X is a matrix where each column is a feature and each row is
    an example. You need to perform the normalization separately for each feature.

    Hint
    ----
    You might find the 'np.mean' and 'np.std' functions useful.
    """
    # You need to set these values correctly
    X_norm = X.copy()
    mu = np.zeros(X.shape[1])
    sigma = np.zeros(X.shape[1])

    # ===== YOUR CODE HERE =====

    mu = np.mean(X, axis=0) # mean of the features
    sigma = np.std(X, axis=0, ddof=1) # mean of the features, divide by N-1
    X_norm = (X - mu)/sigma

    # =====
    return X_norm, mu, sigma

```

Execute the next cell to run the implemented featureNormalize function.

In [23]:

```
# call featureNormalize on the loaded data
X_norm, mu, sigma = featureNormalize(X)

print('Computed mean:', mu)
print('Dimension of the mean mu:', mu.shape)
print('Computed standard deviation:', sigma)
print('Dimension of the standard deviation sigma:', sigma.shape)
```

```
Computed mean: [2000.68085106    3.17021277]
Dimension of the mean mu: (2,)
Computed standard deviation: [7.94702354e+02  7.60981887e-01]
Dimension of the standard deviation sigma: (2,)
```

After the `featureNormalize` function is tested, we now add the intercept term to `X_norm` :

In [24]:

```
# Add intercept term to X
X = np.concatenate([np.ones((N, 1)), X_norm], axis=1)
print('Dimension of data matrix X:', X.shape, '\n')
print('First 4 samples of data matrix X:\n', X[0:4, :])
```

```
Dimension of data matrix X: (47, 3)
```

```
First 4 samples of data matrix X:
[[ 1.          0.13000987 -0.22367519]
 [ 1.         -0.50418984 -0.22367519]
 [ 1.          0.50247636 -0.22367519]
 [ 1.         -0.73572306 -1.53776691]]
```

3.2 Gradient Descent

Previously, you implemented gradient descent on a univariate regression problem. The only difference now is that there is one more feature in the matrix X . The hypothesis function and the batch gradient descent update rule remain unchanged.

The functions `computeCostMulti` and `gradientDescentMulti` implement the cost function and gradient descent for linear regression with multiple variables. Make sure your code supports any number of features and is well-vectorized. You can use the `shape` property of `numpy` arrays to find out how many features are present in the dataset.

****Implementation Note:**** In the multivariate case, the cost function can also be written in the following vectorized form:

$$J(\theta) = \frac{1}{N}(X\theta - Y)^T(X\theta - Y)$$

where

$$X = \begin{pmatrix} -x(1)^T - \\ -x(2)^T - \\ \vdots \\ -x(N)^T - \end{pmatrix} \quad y = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{bmatrix}$$

the vectorized version is efficient when you are working with numerical computing tools like `numpy` . If you are an expert with matrix operations, you can prove to yourself that the two forms are equivalent.

In [25]:

```
def computeCostMulti(X, y, theta):
    """
    Compute cost for linear regression with multiple variables.
    Computes the cost of using theta as the parameter for linear regression to fit the data

    Parameters
    -----
    X : array_like
        The dataset of shape (N x d).

    y : array_like
        A vector of shape (N, ) for the values at a given data point.

    theta : array_like
        The linear regression parameters. A vector of shape (d, )

    Returns
    -----
    J : float
        The value of the cost function.

    Instructions
    -----
    Compute the cost of a particular choice of theta. You should set J to the cost.
    """
    # Initialize some useful values
    N = y.shape[0] # number of training examples

    # You need to return the following variable correctly
    J = 0

    # ===== YOUR CODE HERE =====

    J = 1/N * np.linalg.norm( np.dot(X, theta) - y )**2 #matrix implementation

    # =====
    return J
```

Now define the gradient descent.

In [26]:

```

def gradientDescentMulti(X, y, theta, alpha, num_iters):
    """
    Performs gradient descent to learn theta.
    Updates theta by taking num_iters gradient steps with learning rate alpha.

    Parameters
    -----
    X : array_like
        The dataset of shape (N x d).

    y : array_like
        A vector of shape (N, ) for the values at a given data point.

    theta : array_like
        The linear regression parameters. A vector of shape (d, )

    alpha : float
        The learning rate for gradient descent.

    num_iters : int
        The number of iterations to run gradient descent.

    Returns
    -----
    theta : array_like
        The learned linear regression parameters. A vector of shape (d, ).

    J_history : list
        A python list for the values of the cost function after each iteration.

    Instructions
    -----
    Perform a single gradient step on the parameter vector theta.

    While debugging, it can be useful to print out the values of
    the cost function (computeCost) and gradient here.
    """
    # Initialize some useful values
    m = y.shape[0] # number of training examples

    # make a copy of theta, which will be updated by gradient descent
    theta = theta.copy()

    J_history = []

    for i in range(num_iters):
        # ===== YOUR CODE HERE =====

        theta = theta - alpha*2/(N) * np.dot( X.T, ( np.dot(X, theta) - y ))

        # =====

        # save the cost J in every iteration
        J_history.append( computeCostMulti(X, y, theta) )

    return theta, J_history

```

3.2.1 Selecting learning rates

In this part of the exercise, you will get to try out different learning rates for the dataset and find a learning rate that converges quickly. You can change the learning rate by modifying the following code and changing the part of the code that sets the learning rate.

Use your implementation of `gradientDescentMulti` function and run gradient descent for about 50 iterations at the chosen learning rate. The function should also return the history of $J(\theta)$ values in a vector J .

After the last iteration, plot the J values against the number of the iterations.

If you picked a learning rate within a good range, your plot look similar as the following Figure.

If your graph looks very different, especially if your value of $J(\theta)$ increases or even blows up, adjust your learning rate and try again. We recommend trying values of the learning rate α on a log-scale, at multiplicative steps of about 3 times the previous value (i.e., 0.3, 0.1, 0.03, 0.01 and so on). You may also want to adjust the number of iterations you are running if that will help you see the overall trend in the curve.

****Implementation Note:**** If your learning rate is too large, $J(\theta)$ can diverge and 'blow up', resulting in values which are too large for computer calculations. In these situations, 'numpy' will tend to return NaNs. NaN stands for 'not a number' and is often caused by undefined operations that involve $-\infty$ and $+\infty$.

****MATPLOTLIB tip:**** To compare how different learning rates affect convergence, it is helpful to plot J for several learning rates on the same figure. This can be done by making 'alpha' a python list, and looping across the values within this list, and calling the plot function in every iteration of the loop. It is also useful to have a legend to distinguish the different lines within the plot. Search online for 'pyplot.legend' for help on showing legends in 'matplotlib'.

Notice the changes in the convergence curves as the learning rate changes. With a small learning rate, you should find that gradient descent takes a very long time to converge to the optimal value. Conversely, with a large learning rate, gradient descent might not converge or might even diverge! Using the best learning rate that you found, run the script to run gradient descent until convergence to find the final values of θ . Next, use this value of θ to predict the price of a house with 1650 square feet and 3 bedrooms. You will use this value later to check your implementation of the normal equations. Don't forget to normalize your features when you make this prediction!

In [27]:

```

"""
Instructions
-----
Predict the price of a 1650 sq-ft, 3 br house.

Hint
----
At prediction, make sure you do the same feature normalization.
"""

# Choose some alpha value - change this
alpha = 0.05
num_iters = 400

# init theta and run gradient descent
theta = np.zeros(3)
theta, J_history = gradientDescentMulti(X, y, theta, alpha, num_iters)

# Plot the convergence graph
pyplot.plot(np.arange(len(J_history)), J_history, lw=2)
pyplot.xlabel('Number of iterations')
pyplot.ylabel('Cost J')

# Display the gradient descent's result
print('theta computed from gradient descent: {:s}'.format(str(theta)), '\n')

# Estimate the price of a 1650 sq-ft, 3 br house
# ===== YOUR CODE HERE =====
# Recall that the first column of X is all-ones.
# Thus, it does not need to be normalized.

x_new = np.array([1650, 3])
z_new_norm = (x_new - mu)/sigma
print('New point normalized:', z_new_norm)
print('shape of new point normalized:', z_new_norm.shape, '\n')

z_new = np.concatenate([np.ones(1), z_new_norm], axis=0)
print('New point:', z_new)
print('shape of new point:', z_new.shape, '\n')

price = np.dot(z_new.T, theta) # You should change this

# =====

print('Predicted price of a 1650 sq-ft, 3 br house (using gradient descent): {:.0f}'.format

```

```

theta computed from gradient descent: [340412.65957447 110631.04895815 -664
9.47295013]

```

```

New point normalized: [-0.4412732 -0.22367519]
shape of new point normalized: (2,)

```

```

New point: [ 1.          -0.4412732 -0.22367519]
shape of new point: (3,)

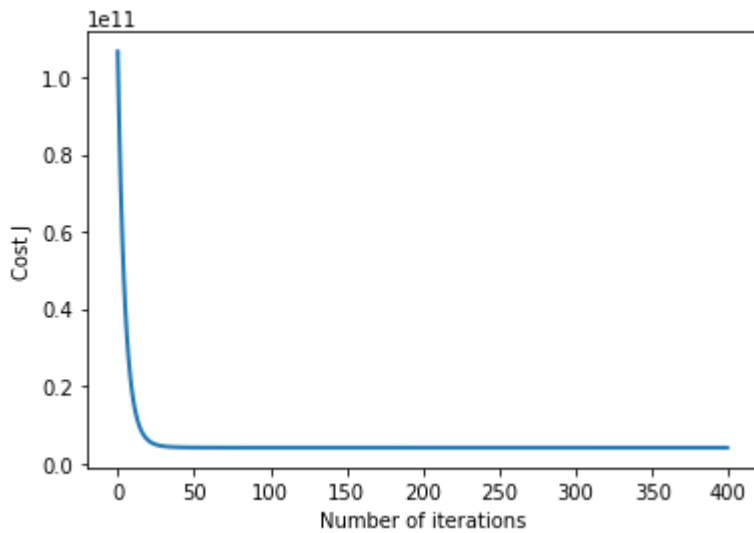
```

```

Predicted price of a 1650 sq-ft, 3 br house (using gradient descent): $29308
1

```





3.3 Normal Equations

We know that the closed-form solution to linear regression is

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

Using this formula does not require any feature scaling, and you will get an exact solution in one calculation: there is no “loop until convergence” like in gradient descent.

First, we will reload the data to ensure that the variables have not been modified. Remember that while you do not need to scale your features, we still need to add a column of 1's to the X matrix to have an intercept term (θ_0). The code in the next cell will add the column of 1's to X for you.

In [28]:

```
# Load data
data = np.loadtxt(os.path.join('Data', 'ex1data2.txt'), delimiter=',')
X = data[:, :2]
y = data[:, 2]
N = y.size
X = np.concatenate([np.ones((N, 1)), X], axis=1)
```

Here is the code for the function `normalEqn` below that uses the formula above to calculate θ .

In [29]:

```

def normalEqn(X, y):
    """
    Computes the closed-form solution to linear regression using the normal equations.

    Parameters
    -----
    X : array_like
        The dataset of shape (N x d).

    y : array_like
        The value at each data point. A vector of shape (N, ).

    Returns
    -----
    theta : array_like
        Estimated linear regression parameters. A vector of shape (d, ).

    Instructions
    -----
    Complete the code to compute the closed form solution to linear
    regression and put the result in theta.

    Hint
    ----
    Look up the function `np.linalg.pinv` for computing matrix inverse.
    """
    theta = np.zeros(X.shape[1])

    # ===== YOUR CODE HERE =====
    X_transpose = X.T
    tmp_1 = np.linalg.pinv( X_transpose.dot(X) )
    tmp_2 = tmp_1.dot(X_transpose)
    theta = tmp_2.dot(y)

    # =====
    return theta

```

Now, once you have found θ using this method, use it to make a price prediction for a 1650-square-foot house with 3 bedrooms. You should find that gives the same predicted price as the value you obtained using the model fit with gradient descent (in Section 3.2.1).

In [30]:

```
# Calculate the parameters from the normal equation
theta = normalEqn(X, y);

# Display normal equation's result
print('Theta computed from the normal equations: {:s}'.format(str(theta)));

# Estimate the price of a 1650 sq-ft, 3 br house
# ===== YOUR CODE HERE =====

price = np.dot( np.array([1, 1650, 3]), theta) # You should change this

# =====

print('Predicted price of a 1650 sq-ft, 3 br house (using normal equations): ${:.0f}'.format
```

```
Theta computed from the normal equations: [89597.90954355  139.21067402 -87
38.01911255]
Predicted price of a 1650 sq-ft, 3 br house (using normal equations): $29308
1
```

In []: