

# **Machine Learning assignment report**



Made by  
Anglo Sherif Nabil Soliman  
40-14526

- After importing all the libraries we need .. now it is the time to import the data . and this is done using `pd.read_csv` function

```
data = pd.read_csv('C:\\semester 10\\machine learning\\house_prices_data_training_data.csv')
```

- Then we drop all NAN values using `dropna` function.
- Lets delete the columns of id and date as they dont matter in the dataset
- We preview the data using `data.head()` function.
- Then it is the time to apply the correlation function to see which features are correlated the price of the house , the features that have weak correlation are dropped
- i chose the value 0.3 as a threshold , and the correlation between a feature and the price is less than this value , this feature will be dropped . this is done to simplify the dataset for future analysis.
- `.corr()` function is used to see the correlation between features.
- Then we make normalization . the function of the normalization is the (feature-mean) / standard deviation of these features. normalization is done to reduce the fluctuation in the output.
- I made normalization for both features and the output which is the price

```

17999

In [285]: 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.

            """
            # 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)
            sigma = np.std(X, axis = 0)
            X_norm = (X - mu) / sigma

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

In [286]: data_norm, mu, sigma = featureNormalize(data)

print('Computed mean:', mu)
print('Computed standard deviation:', sigma)
#print(data_norm)
data_norm = np.concatenate([np.ones((m, 1)), data_norm], axis=1)
data_norm=pd.DataFrame(data_norm)
data_norm.head()
```

- In this assignment , we will do model selection , and this is done by splitting the dataset as follows . 60% for the training part , 20% for cross validation and finally 20% for the testing part.

```
host:8888/notebooks/Assignment-1-Part-2.ipynb#
jupyter Assignment-1-Part-2 Last Checkpoint 3 hours ago (autosaved)
File Edit View Insert Cell Kernel Widgets Help
In [287]: def train_validate_test_split(df, train_percent=.6, validate_percent=.2, seed=None):
    np.random.seed(seed)
    perm = np.random.permutation(df.index)
    m = len(df.index)
    train_end = int(train_percent * m)
    validate_end = int(validate_percent * m) + train_end
    train = df.iloc[perm[:train_end]]
    validate = df.iloc[perm[train_end:validate_end]]
    test = df.iloc[perm[validate_end:]]
    return train, validate, test
train, validate, test = train_validate_test_split(data_norm, train_percent=.6, validate_percent=.2, seed=None)
print(validate.shape)
print(train.shape)
print(test.shape)
(3599, 11)
(10799, 11)
(3001, 11)
```

- Now to the compute cost part . this function is used to calculate the error when we apply linear regression technique using thetas.
- The compute cost function is called in the gradient descent function

```
In [292]: 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 points in X and y.
    Parameters
    -----
    X : array_like
        The dataset of shape (m x n+1).
    y : array_like
        A vector of shape (m, ) for the values at a given data point.
    theta : array_like
        The linear regression parameters. A vector of shape (n+1, )
    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
    m = y.shape[0] # number of training examples

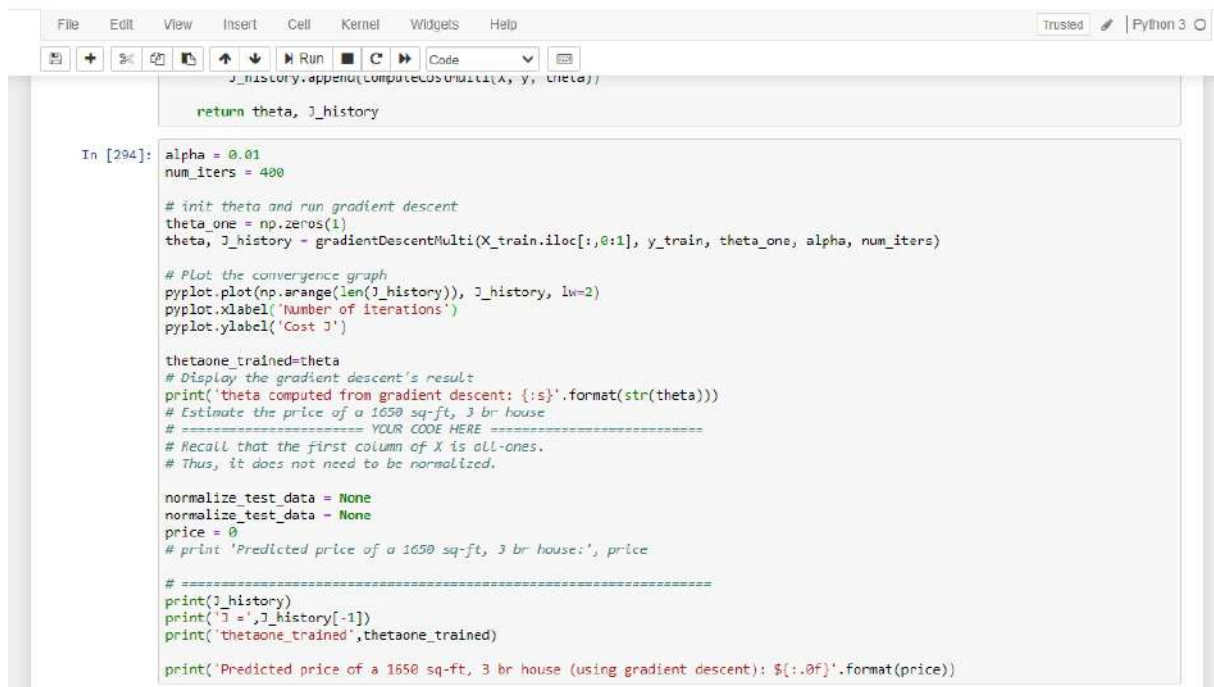
    # You need to return the following variable correctly
    J = 0

    # ===== YOUR CODE HERE =====
    h = np.dot(X, theta)

    J = (1/(2 * m)) * np.sum(np.square(np.dot(X, theta) - y))
    return J
```

- Now to the gradient descent function
- This function is used to learn the theta values.

- The function takes 5 input parameters and they are . X which is the features of the dataset , Y which is the output of the dataset which is the price in our case . alpha , which is the learning rate , and the last parameter is the number of iterations to run gradient descent.
- The gradient descent function returns two values which are the thetas in type of array , and the second value returned is the J value which is calculated using the called function “compute cost”.



```

J_history.append(computeCostMulti(X, y, theta))

return theta, J_history

In [294]: alpha = 0.01
num_iters = 400

# init theta and run gradient descent
theta_one = np.zeros(1)
theta, J_history = gradientDescentMulti(X_train.iloc[:,0:1], y_train, theta_one, 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')

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

normalize_test_data = None
normalize_test_data = None
price = 0
# print 'Predicted price of a 1650 sq-ft, 3 br house: ', price

# =====
print(J_history)
print('J =', J_history[-1])
print('thetaone_trained', thetaone_trained)

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

```

- In the previous figure I chose the learning rate to be equal to 0.01 and the number of iterations = 400
- We keep applying this function to each parameter and the function returns an array of thetas for this parameter.
- So for example .. if we applied this function on the 4<sup>th</sup> feature , we will get a vector that contains 4 values. In addition to the J (error) value,
- The function is applied on the test part only
- The following figure is results for the outputs of the gradient descent function that is used on test part.

```

theta_four = np.zeros(4)
theta, J_history = gradientDescentMulti(X_train.iloc[:,0:4], y_train, theta_four, 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')

thetatfour_trained=theta

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

normalize_test_data = None
normalize_test_data = None
price = 0
# print 'Predicted price of a 1650 sq-ft, 3 br house:'. price
# =====
print(J_history)
print('J = ', J_history[-1])
print('Predicted price of a 1650 sq-ft, 3 br house (using gradient descent): ${:.0f}'.format(price))

```

```

theta computed from gradient descent: 2 -0.058187
3 0.127620
4 0.517108
5 0.219343
Name: 1, dtype: float64
[0.467435474655331, 0.45824935732630107, 0.4494658176293099, 0.4410666163793238, 0.433034347952291, 0.4253524020000311, 0.4180
0492739810066, 0.41087679666616943, 0.40425357356413166, 0.39782148113071006, 0.3916673713546657, 0.3857786965519252, 0.3801434
817659362, 0.3747502865442276, 0.369588239912185, 0.3646460964772814, 0.3599163336425578, 0.35538706906015384, 0.35105005588320
665, 0.3468665496943767, 0.34291862399274564, 0.3391880054211329, 0.33545756012120276, 0.33195985095129305, 0.328608127107050
7, 0.3253956591849157, 0.32231681493056596, 0.319365045640689, 0.3165348731894886, 0.31382087764711114, 0.3112178654695908, 0.3
0872095821796475, 0.30632538185380714, 0.30402665637956977, 0.3018204861478062, 0.29970277046070503, 0.29766959471227383, 0.295

```

J values at each iteration



} theta values

- Now it is the time to apply the compute cost function to the cross validation part to get the J values.

## Now Lets get J values fo the validation part

```

In [303]: j1=computeCostMulti(X_val.iloc[:,0:1], y_val, thetaone_trained)
j2=computeCostMulti(X_val.iloc[:,0:2], y_val, thetatwo_trained)
j3=computeCostMulti(X_val.iloc[:,0:3], y_val, thetathree_trained)
j4=computeCostMulti(X_val.iloc[:,0:4], y_val, thetatfour_trained)
j5=computeCostMulti(X_val.iloc[:,0:5], y_val, thetatfive_trained)
j6=computeCostMulti(X_val.iloc[:,0:6], y_val, thetatsix_trained)
j7=computeCostMulti(X_val.iloc[:,0:7], y_val, thetatseven_trained)
j8=computeCostMulti(X_val.iloc[:,0:8], y_val, thetateight_trained)
j9=computeCostMulti(X_val.iloc[:,0:9], y_val, thetatnine_trained)
#theta_one_val = thetaone_trained
#theta, J_history = gradientDescentMulti(X_val.iloc[:,0:1], y_val, theta_one_val, alpha, num_iters)
print('j1=',j1)
print('j2=',j2)
print('j3=',j3)
print('j4=',j4)
print('j5=',j5)
print('j6=',j6)
print('j7=',j7)
print('j8=',j8)
print('j9=',j9)

```

```

j1= 0.4077958609379826
j2= 0.3208663883956447
j3= 0.22859161568597522
j4= 0.2063628890442122
j5= 0.18810062296397695
j6= 0.18887335860798232
j7= 0.1875871309912664
j8= 0.15788732614573436
j9= 0.158157204845281

```

## Results :

- We notice that from the J values that we got from applying gradient descent function that the least J value that J8 is least error in all J values. With value equal to 0.17350347673138822
- And after applying the compute cost function on each feature we see that also the least J value that J8 is least error in all J value with value equal to 0.157887
- So we can deduce that the feature of J8 which is feature 7 is the best degree in the dataset
- then now is the time to compute the general error in the test part and let's call this error J\_general.

```
print ('j7=',j7)
print ('j8=',j8)
print ('j9=',j9)
```

```
j1= 0.4077958609379826
j2= 0.3208663883956447
j3= 0.22859161568597522
j4= 0.2063628890442122
j5= 0.18810062296397695
j6= 0.18887335860798232
j7= 0.1875871309912664
j8= 0.15788732614573436
j9= 0.158157204845281
```

**So the least J is j8 (feature 7) ... in the training part also feature 7 has the least J ... so Degree 7 is the best degree**

```
In [304]: #thetafinal=[thetaone_trained,thetatwo_trained,thetathree_trained,thetatfour_trained,thetatfive_trained,thetatsix_trained,thetat:

j_general=computeCostMulti(X_test.iloc[:,0:8], y_test, thetateight_trained)
print ('j_general',j_general)
```

```
j_general 0.23966606060854435
```

```
In [ ]:
```