

Homework2: Linear Regression

In this assignment, we will start with utilizing scikit-learn to implement a linear regression model. Afterwards, we will be dropping scikit-learn and implementing these algorithms from scratch without the use of machine learning libraries. While you would likely never have to implement your own linear regression algorithm from scratch in practice, such a skill is valuable to have as you progress further into the field and find many scenarios where you actually may need to perform such implementations manually. Additionally, implementing algorithms from scratch will help you better understand the underlying mathematics behind each model.

Import Libraries

We will be using the following libraries for this homework assignment. For the questions requiring manual implementation, the pre-existing implementations from scikit-learn should *not* be used.

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import operator
%matplotlib inline

from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import mean_squared_error
```

Preparing Data

The file named **dataset1.csv** includes data that was generated from an n-degree polynomial with some gaussian noise. The data has 2 columns - first column is the feature (input) and the second column is its label (output). The first step is to load the data and split them into training, validation, and test sets. A reminder that the purpose of each of the splitted sets are as follows:

- **Training Set:** The sample of data used to fit the model
- **Validation Set:** The sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyperparameters.
- **Test Set:** The sample of data used to provide an unbiased evaluation of a final model fit on the training dataset.

In the section below, we load the csv file and split the data randomly into 3 equal sets.

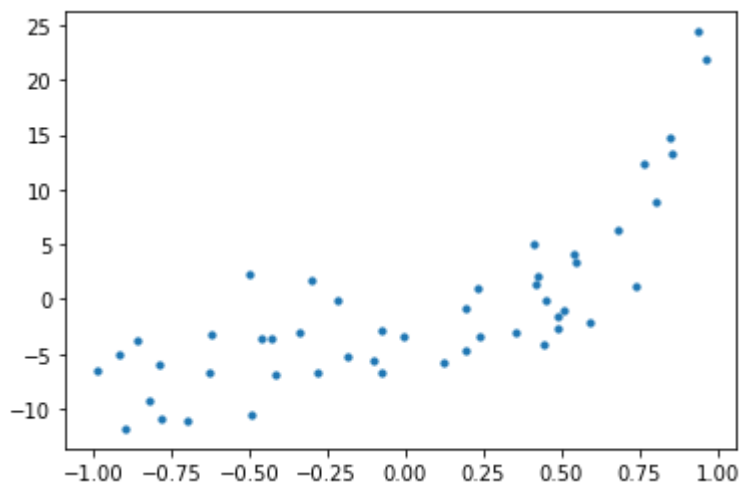
Note that in practice, we usually aim for around a 70-20-10 split for train, valid, and test respectively, but due to limited data in our case, we will do an even split in order to have sufficient data for evaluation

```
In [2]: # Load the data and split into 3 equal sets
data = pd.read_csv('./datasets/dataset1.csv', header=None)
data = data.iloc[:, :-1]
train, valid, test = np.split(data, [int(.33*len(data)), int(.66*len(data))])

# We sort the data in order for plotting purposes later
train.sort_values(by=[0], inplace=True)
valid.sort_values(by=[0], inplace=True)
test.sort_values(by=[0], inplace=True)
```

Let's take a look at what our data looks like

```
In [3]: plt.scatter(train[0], train[1], s=10)
plt.show()
```



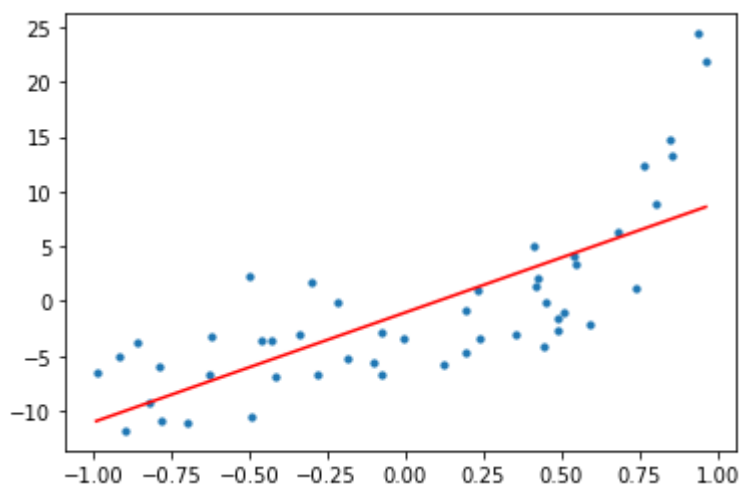
Let's apply a linear regression model using scikit-learn and see what the results look like.

```
In [4]: # Reshape arrays since scikit-learn only takes in 2D arrays
train_x = np.array(train[0])
train_y = np.array(train[1])
valid_x = np.array(valid[0])
valid_y = np.array(valid[1])

train_x = train_x.reshape(-1,1)
train_y = train_y.reshape(-1,1)
valid_x = valid_x.reshape(-1,1)
valid_y = valid_y.reshape(-1,1)

# Apply linear regression model
model = LinearRegression()
model.fit(train_x, train_y)
y_pred = model.predict(train_x)

# Plot the results
plt.scatter(train_x, train_y, s=10)
plt.plot(train_x, y_pred, color='r')
plt.show()
```



By analyzing the line of best fit above, we can see that a straight line is unable to capture the patterns of the data. This is an example of underfitting. As seen in the latest lecture, we can generate a higher order equation by adding powers of the original features as new features.

The linear model:

$$y(x) = w_1x + w_0$$

can be transformed to a polynomial model such as:

$$y(x) = w_2x^2 + w_1x + w_0$$

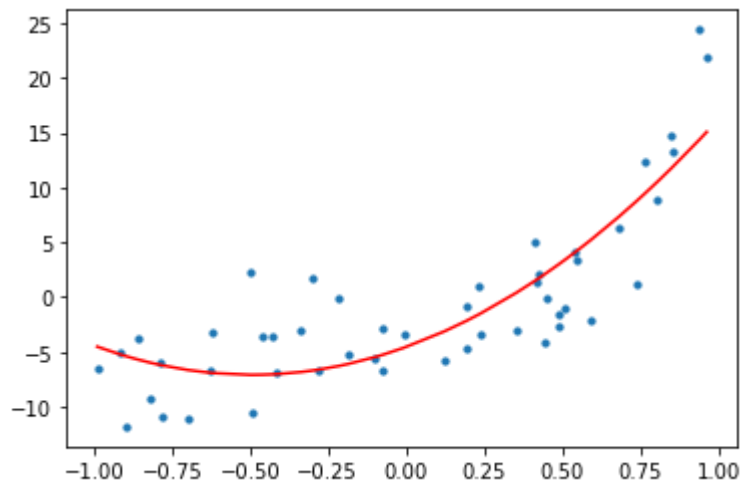
Note that this is still considered to be linear model as the coefficients/weights associated with the features are still linear. x^2 is only a feature. However the curve that we would be fitting in this case is quadratic in nature.

Below we show an example of a quadratic curve being fit to the data

```
In [5]: # Create polynomial features with degree 2
polynomial_features = PolynomialFeatures(degree=2)
x_poly = polynomial_features.fit_transform(train_x)

# Apply linear regression
model = LinearRegression()
model.fit(x_poly, train_y)
y_poly_pred = model.predict(x_poly)

# Plot the results
plt.scatter(train_x, train_y, s=10)
plt.plot(train_x, y_poly_pred, color='r')
plt.show()
```



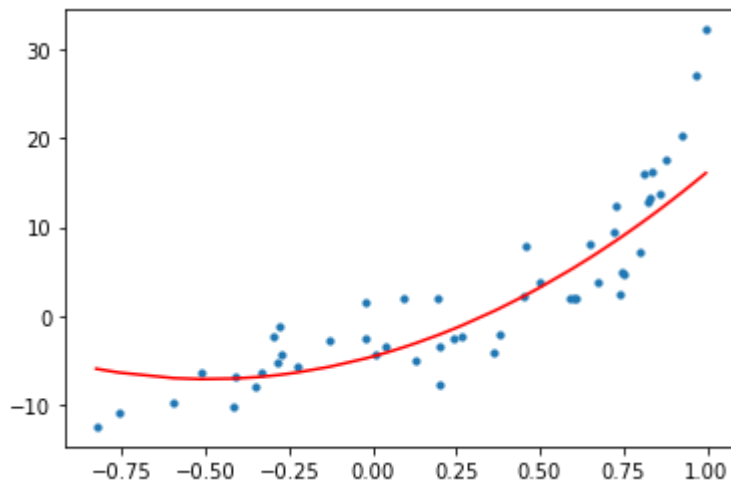
As you can see, we get a slightly better fit with a quadratic curve. Let's use the model to make predictions on our validation set and compute the mean squared error, which is the error which we wish to minimize.

```
In [6]: # Make predictions using pretrained model
valid_y_poly_pred = model.predict(polynomial_features.fit_transform(valid_x))

# Calculate mean squared error
mse = mean_squared_error(valid_y, valid_y_poly_pred)
print("Mean Squared Error: {}".format(mse))

# Plot the prediction results
plt.scatter(valid_x, valid_y, s=10)
plt.plot(valid_x, valid_y_poly_pred, color='r')
plt.show()
```

Mean Squared Error: 20.485214511024225



Question 1: Polynomial Regression Using Scikit-learn [10pts]

Now it is your turn! Following the same format as above, implement a 5-degree polynomial regression model on the training data and plot your results. Use your model to predict the output of the validation set and calculate the root mean square error. Report and plot the results.

Grading policy:

- Q1.1 [4pts]
 - Fit a 5-degree polynomial using scikit-learn [1pts]
 - Use model to predict output of validation set [1pts]
 - Calculate and report the MSE [1pt]
 - Plot curves on the training set and the validation set [1pt]
- Q1.2 [1pts]
- Q1.3 [4pts]
 - Fit a 10-degree polynomial using scikit-learn [1pts]
 - Use model to predict output of validation set [1pts]
 - Calculate and report the MSE [1pts]
 - Plot curves on the training set the validation set [1pt]
- Q1.4 [1pts]

Q1.1

```
In [7]: ### YOUR CODE HERE - Fit a 5-degree polynomial using scikit-learn
# Create polynomial features with degree 5
polynomial_features_5 = PolynomialFeatures(degree=5)
x_poly_5 = polynomial_features_5.fit_transform(train_x)

# Apply linear regression
model = LinearRegression()
model.fit(x_poly_5, train_y)
y_poly_pred_5 = model.predict(x_poly_5)

### YOUR CODE HERE - Plot your the curve on the training data set
# Plot the training data set
##Although in the document above, we donnot require to plot this image
plt.figure(figsize=[15,10])
plt.subplot(2,1,1)

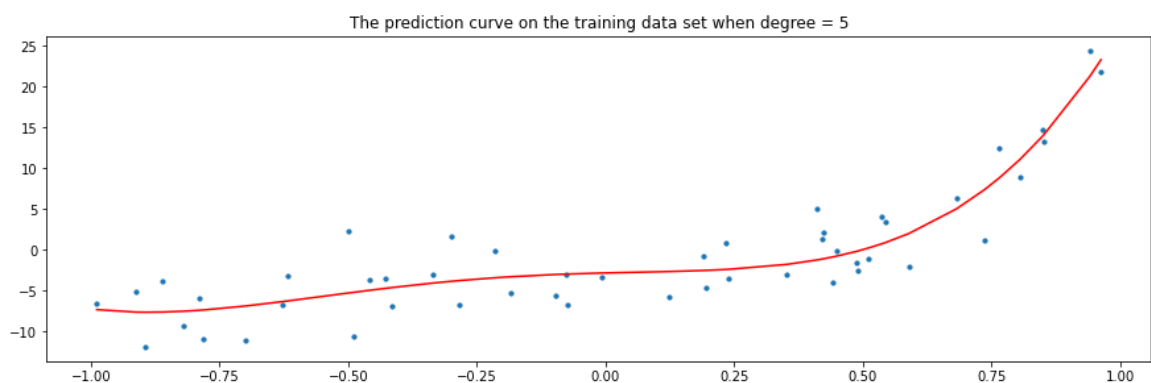
plt.scatter(train_x, train_y, s=10)
plt.plot(train_x, y_poly_pred_5, color='r')
plt.title("The prediction curve on the training data set when degree = 5")
plt.show()

### YOUR CODE HERE - Use model to predict output of validation set
valid_y_poly_pred_5 = model.predict(polynomial_features_5.fit_transform(valid_x))

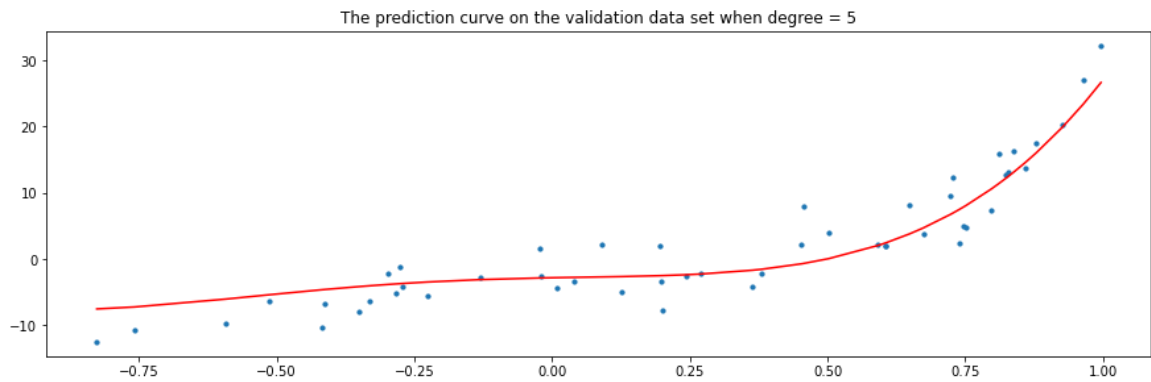
### YOUR CODE HERE - Calculate the MSE. Report and plot the curve on the validation s
# Calculate mean squared error
mse_5 = mean_squared_error(valid_y, valid_y_poly_pred_5)
print("Mean Squared Error: {}".format(mse_5))

# Plot the prediction results on the validation set
plt.figure(figsize=[15,10])
plt.subplot(2,1,2)

plt.scatter(valid_x, valid_y, s=10)
plt.plot(valid_x, valid_y_poly_pred_5, color='r')
plt.title("The prediction curve on the validation data set when degree = 5")
plt.show()
```



Mean Squared Error: 10.363005107697864



Q1.2 Did the mean squared error go up or down as compared to the 2-degree polynomial curve? Why do you think this is the case?

----- ANSWER HERE -----

The mean squared error goes down when we set the degree of polynomial curve equal to 5. The reason is that when the degree is equal to 5, the curve is more suitable to the validation data set, which means we can relieve the underfitting problem by appropriately increasing model complexity.

Q1.3

Now repeat the above for a 10-degree polynomial regression model.

```

In [8]: ### YOUR CODE HERE - Fit a 10-degree polynomial using scikit-learn
# Create polynomial features with degree 10
polynomial_features_10 = PolynomialFeatures(degree=10)
x_poly_10 = polynomial_features_10.fit_transform(train_x)

# Apply linear regression
model = LinearRegression()
model.fit(x_poly_10, train_y)
y_poly_pred_10 = model.predict(x_poly_10)

### YOUR CODE HERE - Plot your the curve on the training data set
# Plot the training data set
##Although in the document above, we donnot require to plot this image
plt.figure(figsize=[15,10])
plt.subplot(2,1,1)

plt.scatter(train_x, train_y, s=10)
plt.plot(train_x, y_poly_pred_10, color='r')
plt.title("The prediction curve on the training data set when degree = 10")
plt.show()

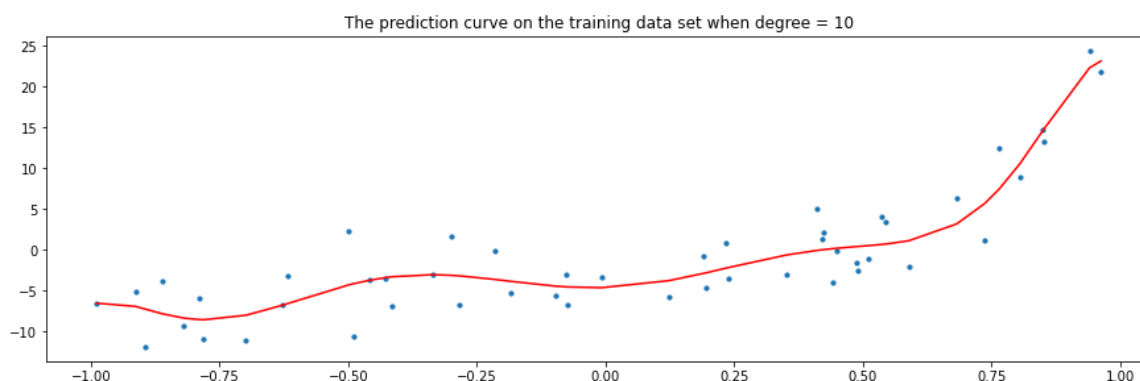
### YOUR CODE HERE - Use model to predict output of validation set
valid_y_poly_pred_10 = model.predict(polynomial_features_10.fit_transform(valid_x))

### YOUR CODE HERE - Calculate the MSE. Report and plot the curve on the validation s
# Calculate mean squared error
mse_10 = mean_squared_error(valid_y, valid_y_poly_pred_10)
print("Mean Squared Error: {}".format(mse_10))

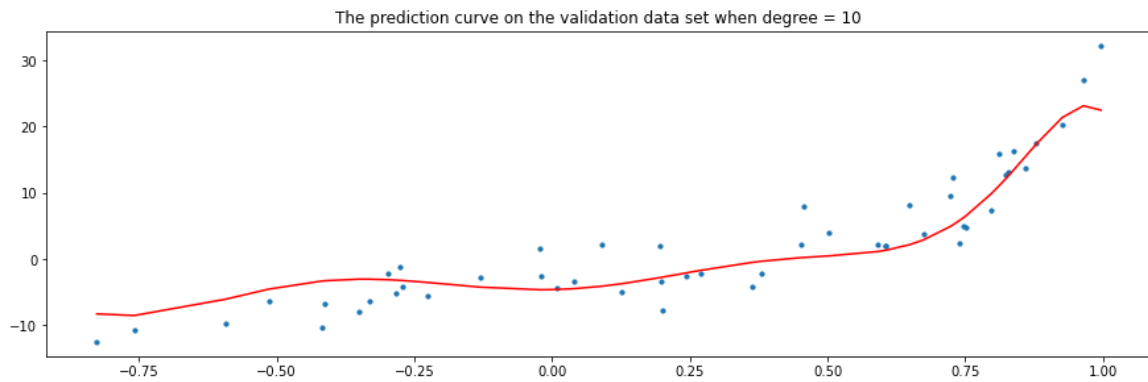
# Plot the prediction results on the validation set
plt.figure(figsize=[15,10])
plt.subplot(2,1,2)

plt.scatter(valid_x, valid_y, s=10)
plt.plot(valid_x, valid_y_poly_pred_10, color='r')
plt.title("The prediction curve on the validation data set when degree = 10")
plt.show()

```



Mean Squared Error: 13.115190778263633



Q1.4 How does the mean square error compare to the polynomial regression with degree 5? Why do you think this is the case?

----- ANSWER HERE -----

The mean squared error goes up compared to the case with degree 5 when we set the degree of polynomial curve equal to 10. The reason is that when the degree is equal to 10, the curve is more suitable to the training data set but less suitable to the validation data set. That is because when we increase the complexity of the model too much, it will produce the overfitting problem.

Question 2: Manual Implementation [10pts]

Now it's time to appreciate the hard work that open source developers have put, in order to allow you to implement machine learning models without doing any math! No more scikit-learn (or any other libraries like Tensorflow, Pytorch, etc) for the rest of this assignment!

Your first step is to fit a **10-degree polynomial** to the dataset we have been used above. Then using your results, calculate the mean squared error on both the training and validation set.

A reminder that in polynomial regression, we are looking for a solution for the equation:

$$Y(X) = W^T * \phi(X),$$

where

$$\phi(X) = [1, X, X^2, X^3, \dots, X^n]^T.$$

Let $\phi(\mathbf{X}) = [\phi(X_1)^T; \dots; \phi(X_n)^T]$ denote the data matrix after polynomial transformation and \mathbf{Y} denote the target vector. Recall the the closed-form solution for linear regression is given by normal equation, which is:

$$W = (\phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^T \mathbf{Y}.$$

Make sure to review the slides, do some research, and/or ask for clarification if this doesn't make sense. You must understand the underlying math before being able to implement this properly.

Suggestion - Use the original pandas dataframes variables named train, valid, and test instead of the reshaped arrays that were used specifically for scikit-learn. It will make your computations cleaner and more intuitive.

Grading policy:

- Q2.1 [4pts]
 - Create the polynomial matrix $\phi(X)$ [1pts]
 - Find the weighted matrix W by normal equation [1pts]
 - Make predictions on the training set, calculate and report the mean squared error [1pts]
 - Make predictions on the validation set, calculate and report the mean squared error [1pts]
- Q2.2 [6pts]
 - Implement gradient decent manually [6pts]
 - Correctness of the gradient [2pts]
 - Correctness of the update in each iteration [2pts]
 - Convergence of the algorithm [1pts]
 - Calculate and report the mean squared error on both training and validation set [1pts]

Q2.1

```
In [9]: ### YOUR CODE HERE - Create the polynomial matrix \phi(X), which is a numpy array
# phi_x = [] # modify this line
def create_matrix(x, degree):
    phi_x = []
    for i in range(degree + 1): ##Because it includes constant term
        phi_x.append(np.power(x, i))
    return np.concatenate(phi_x,-1)

train_poly_matrix_x = create_matrix(train_x, 10)
##print(train_poly_matrix_x)
### YOUR CODE HERE - Find the weighted matrix W by normal equation
W = np.matmul(np.matmul(np.linalg.inv(np.matmul(train_poly_matrix_x.T, train_poly_mat
##np.matmul matrix multiply
# np.linalg is always using when we need to deal with linear algebra problem.

### YOUR CODE HERE - Make predictions on the training set and calculate the mean square
##print(W)
##print(train_poly_matrix_x)
train_predict_y = np.matmul(train_poly_matrix_x, W)
mse_train = np.mean(np.square(train_predict_y - train_y))

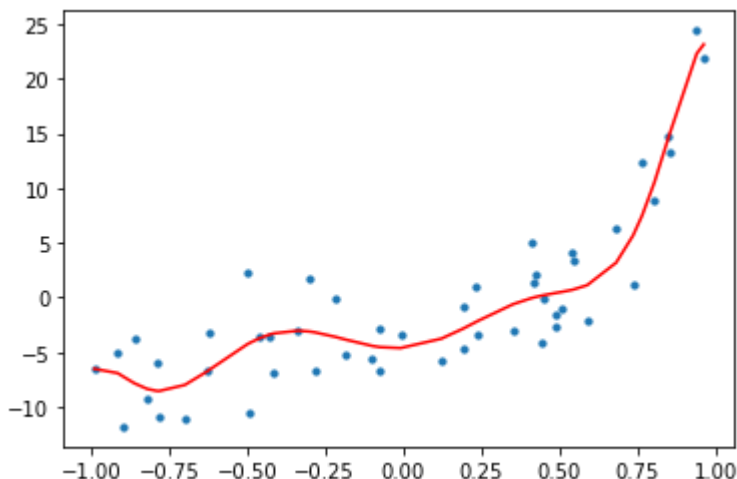
# =====

print("Mean Squared Error (Training): {}".format(mse_train))
plt.scatter(train_x, train_y, s=10)
plt.plot(train_x, train_predict_y, color='r')
plt.show()

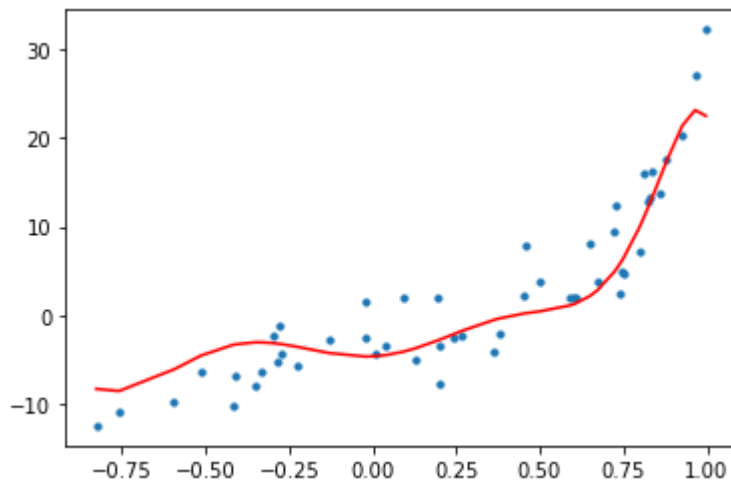
### YOUR CODE HERE - Make predictions on the validation set and calculate the mean square
valid_poly_matrix_x = create_matrix(valid_x, 10) ##fit transform
valid_predict_y = np.matmul(valid_poly_matrix_x, W) ##fit model
mse_valid = np.mean(np.square(valid_predict_y - valid_y))
# =====

print("Mean Squared Error (Validation): {}".format(mse_valid))
plt.scatter(valid_x, valid_y, s=10)
plt.plot(valid_x, valid_predict_y, color='r')
plt.show()
```

Mean Squared Error (Training): 8.60066434175743



Mean Squared Error (Validation): 13.11519077813295



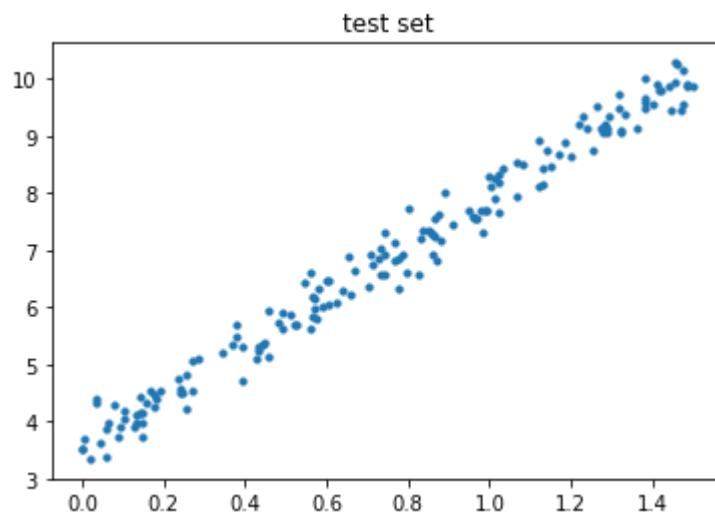
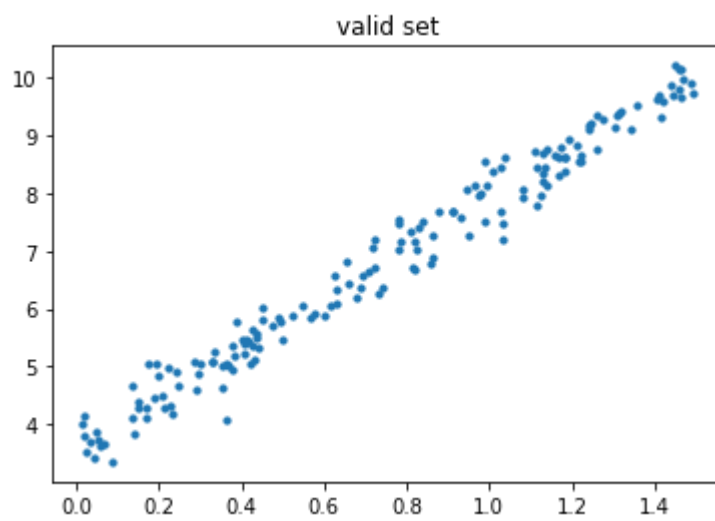
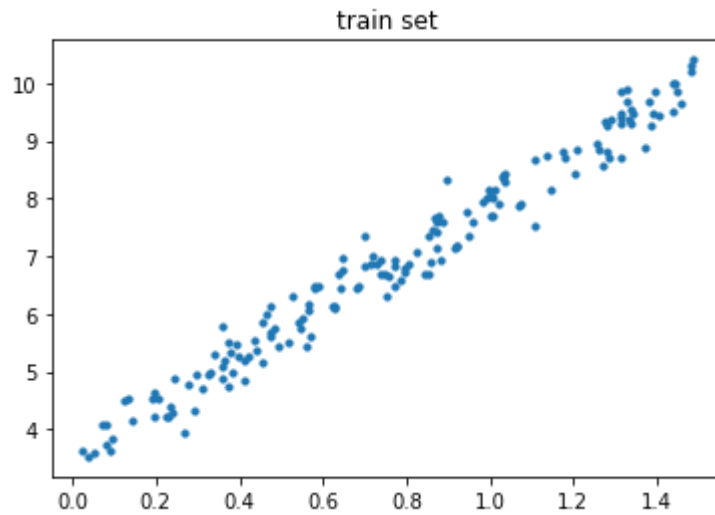
For the rest of the assignment, we will use the other dataset named **dataset2.csv**. First load the csv and split the model into train, valid, and test sets as shown earlier in the assignment.

```
In [10]: # Load dataset2.csv and split into 3 equal sets
data = pd.read_csv('./datasets/dataset2.csv', header=None)
data = data.iloc[:, :-1]
train, valid, test = np.split(data, [int(.33*len(data)), int(.66*len(data))])
# Sort the data in order for plotting purposes later
train.sort_values(by=[0], inplace=True)
valid.sort_values(by=[0], inplace=True)
test.sort_values(by=[0], inplace=True)

train_x = np.array(train[0])
train_y = np.array(train[1])
valid_x = np.array(valid[0])
valid_y = np.array(valid[1])
test_x = np.array(test[0])
test_y = np.array(test[1])
train_x = train_x.reshape(-1, 1)
train_y = train_y.reshape(-1, 1) ## -1 means fuzzy processing, which means we just fo
valid_x = valid_x.reshape(-1, 1)
valid_y = valid_y.reshape(-1, 1)
test_x = test_x.reshape(-1, 1)
test_y = test_y.reshape(-1, 1)
```

Plot the data below to see what it looks like

```
In [11]: plt.title('train set')
plt.scatter(train_x, train_y, s=10)
plt.show()
plt.title('valid set')
plt.scatter(valid_x, valid_y, s=10)
plt.show()
plt.title('test set')
plt.scatter(test_x, test_y, s=10)
plt.show()
```



Q2.2

If done properly, you should see that the points fall under a relatively straight line with minor deviations. Looks like a perfect example to implement a linear regression model using the **gradient descent** method without the use of any machine learning libraries!

Since the data falls along a straight line, we can assume the solution follows the form:

$$y(x) = wx + b$$

A reminder that in gradient descent, we essentially want to iteratively get closer to the minimum of our objective function (the mean squared error), such that:

$$MSE(w_0) > MSE(w_1) > MSE(w_2) > \dots$$

The algorithm is as follows:

**** 1) Pick initial w_0 randomly. ****

**** 2) For $k = 1, 2, \dots \Rightarrow w_{k+1} = w_k - \alpha g(w_k)$ where $\alpha > 0$ is the learning rate and $g(w_k)$ is the gradient. ****

**** End when $|w_{k+1} - w_k| < \epsilon$ ****

There are many resources online for gradient descent. You must understand the underlying math before being able to implement this properly.

Now once you understand, it is time to implement the gradient descent below. You may set the learning rate to 1e-6 or whatever value you think is best. As usual, calculate the mean squared error and plot your results. This time, training should be done using the training and validation sets, while the final mean squared error should be computed using the testing set.

Feel Free to modify the existing code, or just ignore it and provide your own implementation!

```

In [23]: ### Implement gradient decent
m = 0 # initial w (slope)
b = 1 # initial b (intercept)
lr = 0.001
epsilon = 1e-9
epsilon_valid = 1e-12
lr_decay = 0.9
epsilon_lr = 0.2
symbol_m = 1
symbol_b = 1
decay_steps = 1000

def get_mse(m, b, x, y):
    return np.sum(np.multiply((y - m * x - b), (y - m * x - b))) / len(x)

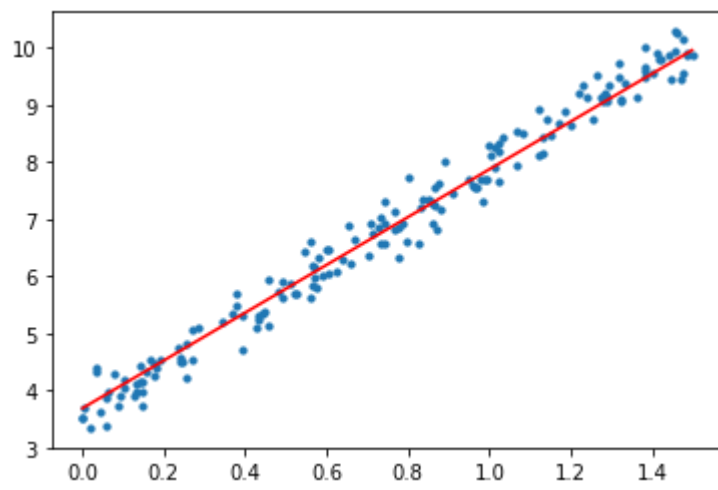
# ===== YOUR CODE HERE =====
finish_m = 0
finish_b = 0
gra_m = 0
gra_b = 0
step = 0
##print(train_x.shape)
##print(train_y.shape)
##print(train_x * m+1)
##print(train_y)
while(True):
    #y = f(x) = wx + b
    y_f_x = m * train_x + b
    #descent gradient
    gra_m = 2 * lr * np.mean(train_x * y_f_x - train_x * train_y)
    gra_b = 2 * lr * np.mean(y_f_x - train_y)
    ##print("123", gra_m, gra_b)
    if(finish_m == 0):
        m -= gra_m
    if(finish_b == 0):
        b -= gra_b
    step += 1
    # shrink lr every 1000 steps
    if step % decay_steps == 0:
        lr *= lr_decay
    if(abs(gra_m) <= epsilon):
        finish_m = 1
    if(abs(gra_b) <= epsilon):
        finish_b = 1
    if(finish_m == 1 and finish_b == 1):
        break
train_y_predict = m * train_x + b
# =====

### Calculate the mean squared error on both training and validation set and plot the
print("Mean Squared Error (Training): {}".format(get_mse(m, b, train_x, train_y)))
print("Mean Squared Error (Testing): {}".format(get_mse(m, b, test_x, test_y)))
pred_y_test = m * test_x + b
plt.scatter(test_x, test_y, s=10)
plt.plot(test_x, pred_y_test, color='r')
plt.show()

```

Mean Squared Error (Training): 0.09862300510208004

Mean Squared Error (Testing): 0.07978290571616456



In []: