# Linear regression workbook

This workbook will walk you through a linear regression example. It will provide familiarity with Jupyter Notebook and Python. Please print (to pdf) a completed version of this workbook for submission with HW #1.

ECE 239AS, Winter Quarter 2019, Prof. J.C. Kao, TAs M. Kleinman and A. Wickstrom and K. Liang and W. Chuang

```
In [26]: import numpy as np
         import matplotlib.pyplot as plt

         #allows matlab plots to be generated in line
         %matplotlib inline
```
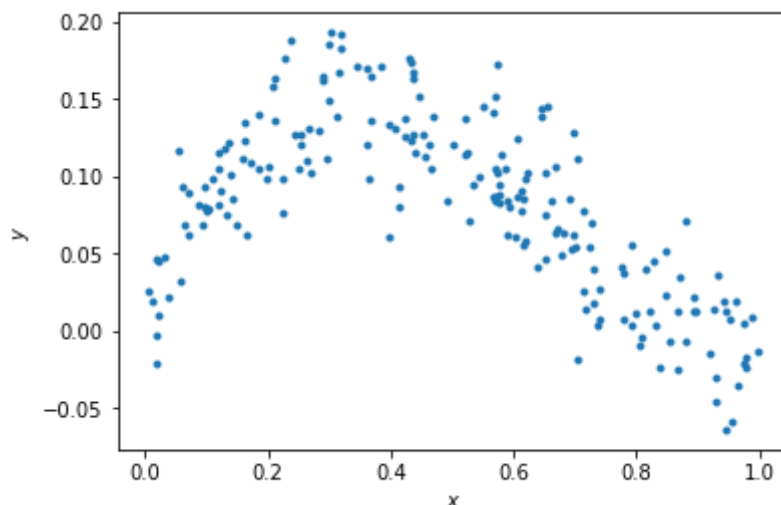
## Data generation

For any example, we first have to generate some appropriate data to use. The following cell generates data according to the model: $y = x - 2x^2 + x^3 + \epsilon$

```
In [27]: np.random.seed(0)   # Sets the random seed.
         num_train = 200      # Number of training data points

         # Generate the training data
         x = np.random.uniform(low=0, high=1, size=(num_train,))
         y = x - 2*x**2 + x**3 + np.random.normal(loc=0, scale=0.03, size=(num_train
         f = plt.figure()
         ax = f.gca()
         ax.plot(x, y, '.')
         ax.set_xlabel('$x$')
         ax.set_ylabel('$y$')
```

Out[27]: Text(0, 0.5, '$y$')

## QUESTIONS:

Write your answers in the markdown cell below this one:

(1) What is the generating distribution of $x$?

(2) What is the distribution of the additive noise $\epsilon$?

### ANSWERS:

(1) x is from a uniform distribution between 0 and 1

(2) $\epsilon$ is from a normal distribution with mean 0 and stddev 0.03

## Fitting data to the model (5 points)

Here, we'll do linear regression to fit the parameters of a model $y = ax + b$.

```
In [28]:  # xhat = (x, 1)
          xhat = np.vstack((x, np.ones_like(x)))

          # ==================== #
          # START YOUR CODE HERE #
          # ==================== #
          # GOAL: create a variable theta; theta is a numpy array whose elements are

          theta = np.linalg.inv((xhat).dot(xhat.T)).dot(xhat.dot(y))

          # ================== #
          # END YOUR CODE HERE #
          # ================== #
```

```
In [29]:  # Plot the data and your model fit.
          f = plt.figure()
          ax = f.gca()
          ax.plot(x, y, '.')
          ax.set_xlabel('$x$')
          ax.set_ylabel('$y$')

          # Plot the regression line
          xs = np.linspace(min(x), max(x),50)
          xs = np.vstack((xs, np.ones_like(xs)))
          plt.plot(xs[0,:], theta.dot(xs))
```

Out[29]:  [<matplotlib.lines.Line2D at 0x113054710>]



## QUESTIONS

(1) Does the linear model under- or overfit the data?

(2) How to change the model to improve the fitting?

```
### ANSWERS

(1) Underfits the data

(2) Increase complexity of the model (add x^2 and x^3 terms)
```

## Fitting data to the model (10 points)

Here, we'll now do regression to polynomial models of orders 1 to 5. Note, the order 1 model is the linear model you prior fit.

```
In [30]:    N = 5
            xhats = []
            thetas = []

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

            # GOAL: create a variable thetas.
            # thetas is a list, where theta[i] are the model parameters for the polynom
            #    i.e., thetas[0] is equivalent to theta above.
            #    i.e., thetas[1] should be a length 3 np.array with the coefficients of
            #    ... etc.
            for i in np.arange(N):
                xhats.append(xhat if i == 0 else np.vstack((x**(i+1), xhats[i-1])))
                thetas.append(np.linalg.inv((xhats[i]).dot(xhats[i].T)).dot(xhats[i].do

            # ================= #
            # END YOUR CODE HERE #
            # ================= #
```

```
In [31]:  # Plot the data
          f = plt.figure()
          ax = f.gca()
          ax.plot(x, y, '.')
          ax.set_xlabel('$x$')
          ax.set_ylabel('$y$')

          # Plot the regression lines
          plot_xs = []
          for i in np.arange(N):
              if i == 0:
                  plot_x = np.vstack((np.linspace(min(x), max(x),50), np.ones(50)))
              else:
                  plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))
              plot_xs.append(plot_x)

          for i in np.arange(N):
              ax.plot(plot_xs[i][-2,:], thetas[i].dot(plot_xs[i]))

          labels = ['data']
          [labels.append('n={}'.format(i+1)) for i in np.arange(N)]
          bbox_to_anchor=(1.3, 1)
          lgd = ax.legend(labels, bbox_to_anchor=bbox_to_anchor)
```
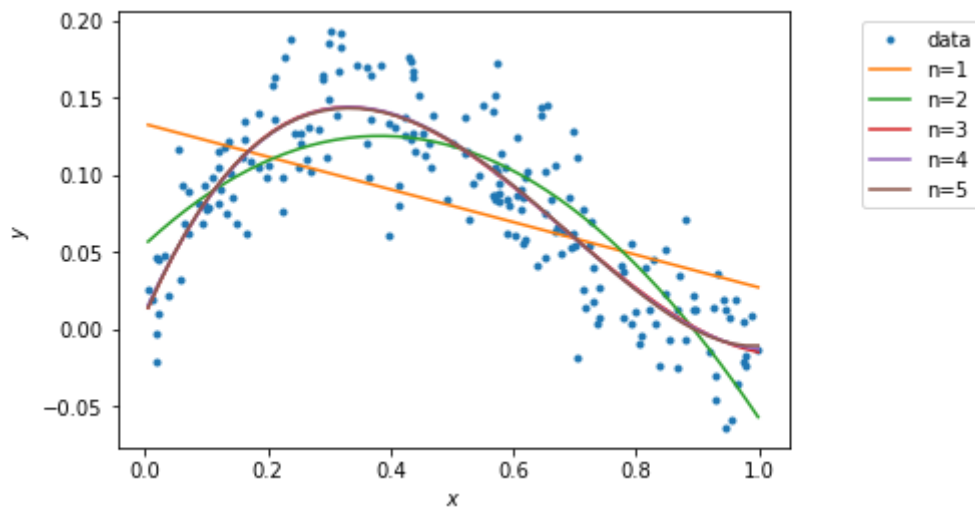


## Calculating the training error (10 points)

Here, we'll now calculate the training error of polynomial models of orders 1 to 5.

```
In [32]: training_errors = []

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

         # GOAL: create a variable training_errors, a list of 5 elements,
         # where training_errors[i] are the training loss for the polynomial fit of
         for i in np.arange(N):
             training_errors.append(0.5*np.sum(np.square(y - (xhats[i].T).dot(thetas

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

         print ('Training errors are: \n', training_errors)
```

```
Training errors are:
 [0.2379961088362701, 0.10924922209268531, 0.08169603801105374, 0.0816535
3735296982, 0.08161479195525295]
```

## QUESTIONS

(1) What polynomial has the best training error?
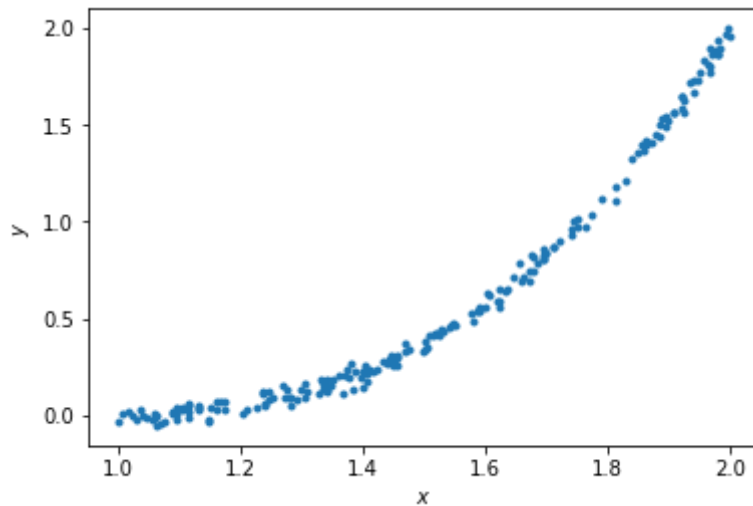
(2) Why is this expected?

## ANSWERS

(1) Polynomial with degree 5

(2) There are more than 5 datapoints, and so training error is strictly decreasing as polynomial degree increases (because polynomial degree 5 encapsulates polynomial degree 4, and so on, thus, it can use the $x^5$ to decrease training error even further)

### Generating new samples and testing error (5 points)

Here, we'll now generate new samples and calculate testing error of polynomial models of orders 1 to 5.

```
In [33]: x = np.random.uniform(low=1, high=2, size=(num_train,))
         y = x - 2*x**2 + x**3 + np.random.normal(loc=0, scale=0.03, size=(num_train
         f = plt.figure()
         ax = f.gca()
         ax.plot(x, y, '.')
         ax.set_xlabel('$x$')
         ax.set_ylabel('$y$')
```

Out[33]: Text(0, 0.5, '$y$')



```
In [34]: xhats = []
         for i in np.arange(N):
             if i == 0:
                 xhat = np.vstack((x, np.ones_like(x)))
                 plot_x = np.vstack((np.linspace(min(x), max(x),50), np.ones(50)))
             else:
                 xhat = np.vstack((x**(i+1), xhat))
                 plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))

             xhats.append(xhat)
```

```python
In [35]:  # Plot the data
          f = plt.figure()
          ax = f.gca()
          ax.plot(x, y, '.')
          ax.set_xlabel('$x$')
          ax.set_ylabel('$y$')

          # Plot the regression lines
          plot_xs = []
          for i in np.arange(N):
              if i == 0:
                  plot_x = np.vstack((np.linspace(min(x), max(x),50), np.ones(50)))
              else:
                  plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))
              plot_xs.append(plot_x)

          for i in np.arange(N):
              ax.plot(plot_xs[i][-2,:], thetas[i].dot(plot_xs[i]))

          labels = ['data']
          [labels.append('n={}'.format(i+1)) for i in np.arange(N)]
          bbox_to_anchor=(1.3, 1)
          lgd = ax.legend(labels, bbox_to_anchor=bbox_to_anchor)
```
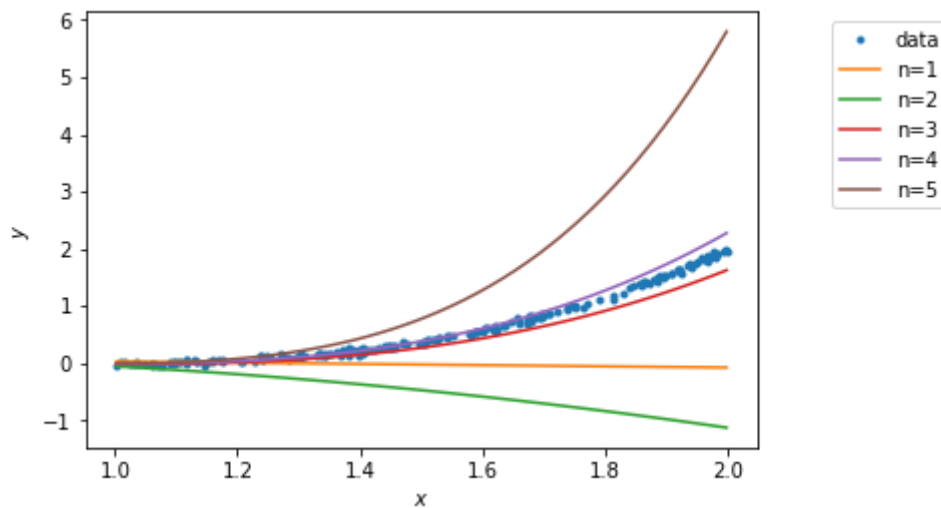
```
In [36]: testing_errors = []

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

         # GOAL: create a variable testing_errors, a list of 5 elements,
         # where testing_errors[i] are the testing loss for the polynomial fit of or
         for i in np.arange(N):
             testing_errors.append(0.5*np.sum(np.square(y - (xhats[i].T).dot(thetas[

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

         print ('Testing errors are: \n', testing_errors)
```

```
Testing errors are:
 [80.86165184550586, 213.19192445057894, 3.1256971082763925, 1.1870765189
474703, 214.91021817652626]
```

## QUESTIONS

(1) What polynomial has the best testing error?

(2) Why polynomial models of orders 5 does not generalize well?

## ANSWERS

(1) Polynomial degree 4

(2) Polynomial degree 5 had overfit the training data which came from an X uniformly distributed over 0 and 1, whereas the testing data has an X uniformly distributed over 1 and 2. The x^5 term then begins to cause very high error, because the testing data seems unlike the training data, despite the underlying function being identical.