# 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 [1]:

```python
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: $y = x - 2x^2 + x^3 + \epsilon$
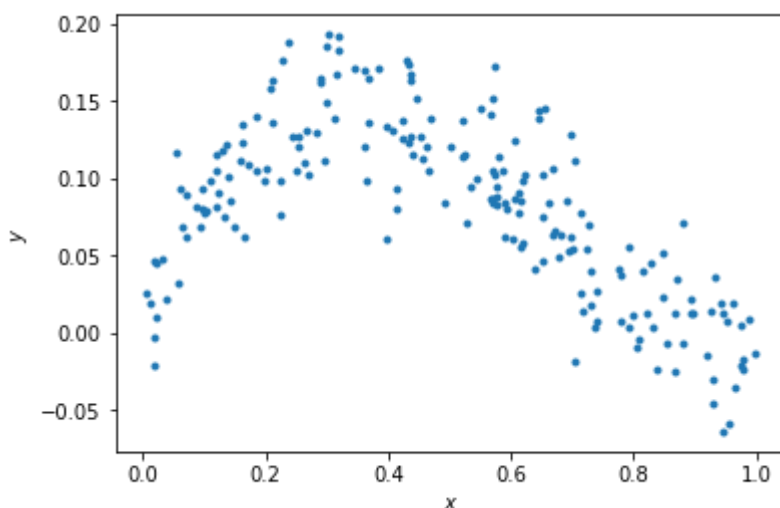
In [2]:

```python
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[2]:

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) The distribution of x is uniform distribution.

(2) The distribution of noise is normal distribution.

## Fitting data to the model (5 points)

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

In [3]:

```python
# xhat = (x, 1)
xhat = np.vstack((x, np.ones_like(x)))
xhat = xhat.transpose()
y = y.transpose()
# =================== #
# START YOUR CODE HERE #
# =================== #
# GOAL: create a variable theta; theta is a numpy array whose elements are [a, b]

theta = np.matmul(np.matmul(np.linalg.inv(np.matmul(xhat.transpose(),xhat)),xhat.transp
ose()),y)
# ================== #
# END YOUR CODE HERE #
# ================== #
```
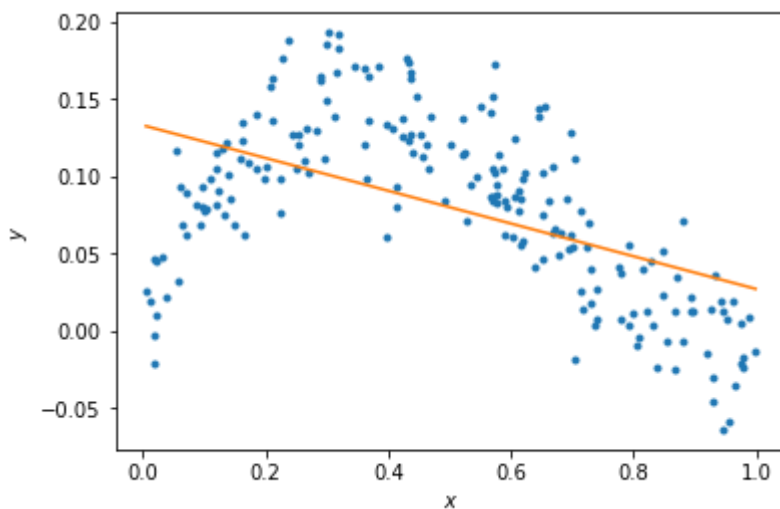
In [4]:

```
# 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[4]:

```
[<matplotlib.lines.Line2D at 0x2234766a390>]
```



## QUESTIONS

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

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

## ANSWERS

(1) The linear model under-fit the data.

(2) Use a polinomial model with higher orders.

## 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 [5]:

```python
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 polynomial fit of o
rder i+1.
#    i.e., thetas[0] is equivalent to theta above.
#    i.e., thetas[1] should be a length 3 np.array with the coefficients of the x^2, x,
 and 1 respectively.
#    ... etc.

for i in range(N):
    if i == 0:
        xhat = np.vstack((x,np.ones_like(x)))
        xhats.append(xhat)
    else:
        xhat = np.vstack((np.power(x,i+1),xhats[-1]))
        xhats.append(xhat)
    xhat = xhat.transpose()
    theta = np.matmul(np.matmul(np.linalg.inv(np.matmul(xhat.transpose(),xhat)),xhat.tr
anspose()),y)
    thetas.append(theta)

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

In [6]:

```python
thetas
```

Out[6]:

```
[array([-0.10599633,  0.13315817]),
 array([-0.48023061,  0.36743967,  0.05521084]),
 array([ 0.8843808 , -1.82077417,  0.91178032,  0.00979068]),
 array([ 0.14080037,  0.60466289, -1.64250929,  0.87250485,  0.01175321]),
 array([ 0.52432592, -1.164568  ,  1.76052438, -2.07430275,  0.93373916,
         0.009716  ])]
```
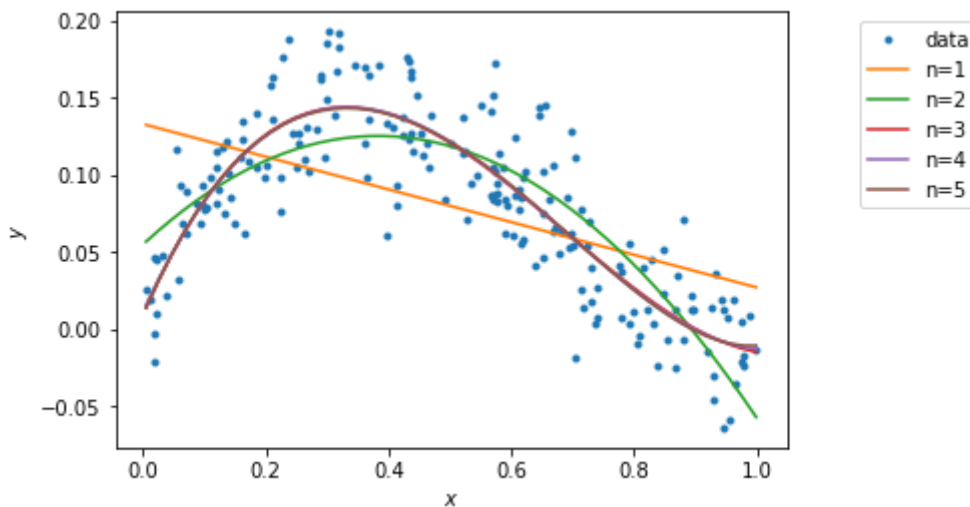
In [7]:

```python
# 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 [8]:

```
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 order i+1.

for xhat, theta in zip(xhats,thetas):
    error = np.linalg.norm(y - np.matmul(theta,xhat), ord=2)**2
    training_errors.append(error)


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

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

```
Training errors are:
 [0.4759922176725402, 0.2184984441853706, 0.16339207602210748, 0.163307074
70593958, 0.1632295839105059]
```

## QUESTIONS

(1) What polynomial has the best training error?

(2) Why is this expected?

## ANSWERS

(1) Polynomial with order 5 has the best training error.

(2) A higher order of polynomial is expectd to fit the training data better than lower order.

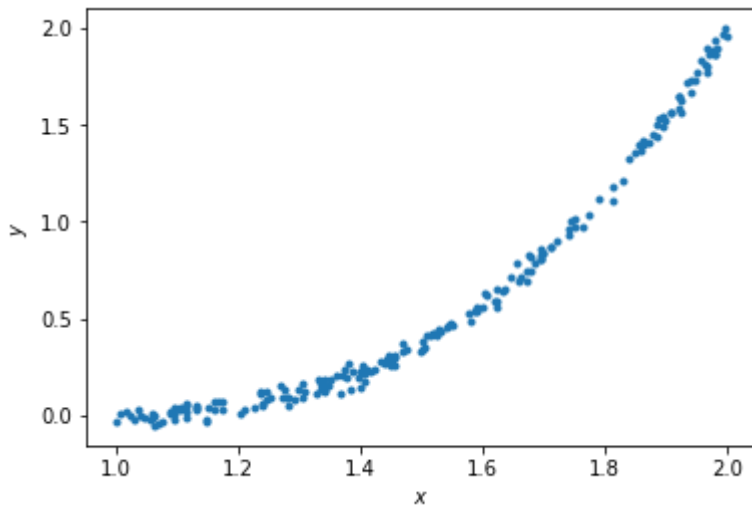### Generating new samples and testing error (5 points)

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

In [9]:

```
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[9]:

Text(0,0.5,'$y$')



In [10]:

```
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)
```
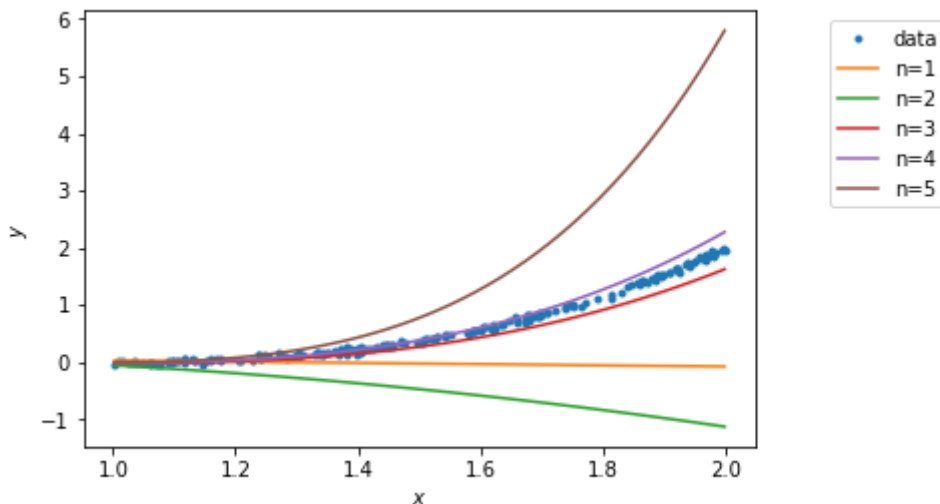
In [11]:

```python
# 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 [12]:

```
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 order i+1.
for xhat, theta in zip(xhats,thetas):
    error = np.linalg.norm(y - np.matmul(theta,xhat), ord=2)**2
    testing_errors.append(error)

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

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

Testing errors are:
 [161.72330369101172, 426.38384890115896, 6.251394216547027, 2.37415303785
4201, 429.8204363359659]

## QUESTIONS

(1) What polynomial has the best testing error?

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

## ANSWERS

(1) Polynomial with order 4 has the best testing error.

(2) Becasue a high order polynomial may overfit the data. It will have low training error but a high testing error.