# Introduction to Data Science

## 8 - Scikit Learn Basics

## This Lecture

- Basic usage of the Scikit learn package

## Setup

We'll typically start a notebook from now on with a set of standard imports, and any relevant Jupyter notebook "magic" directives:

```
In [1]:   1  import numpy as np
          2  import matplotlib.pyplot as plt
          3  import sklearn as sk
          4
          5  %matplotlib inline
```

## Example Problem Setup

We'll continue to use our synthetic problem to illustrate linear regression in Scikit learn:

$$f(x) = 3 + 0.5n - n^2 + 0.15n^3$$
$$y = f(x) + \epsilon$$

```
In [2]:   1  # "ground truth" function
          2  def f(x):
          3      return 3 + 0.5 * x - x**2 + 0.15 * x**3
          4
```

Some more functions:

```
In [3]:    1  # ensure repeatability of this notebook
           2  # (comment out for new results each run)
           3  np.random.seed(12345)
           4
           5  # convenience function for generating samples
           6  def sample(n, fn, limits, noise=1):
           7      width = limits[1] - limits[0]
           8      x = np.random.random(n) * width + limits[0]
           9      y = fn(x) + np.random.randn(n) * noise
          10      return x, y
          11
          12  # there's a scikit learn tool for generating
          13  # polynomial features - this will be more useful
          14  # when working with multivariate inputs, so we'll
          15  # stick with this simpler solution for now
          16  def phi(x, k):
          17      return np.array([x ** p for p in range(k)]).T
```

## Regression Workflow

- Obtain training samples (data and target)
- [optional] Do some initial visualization, statistics
- [optional] Preprocess data (generate features, dimensionality reduction)
- Initialize a model object
- Split data into training and test sets (or use cross validation, more on this another time)
- Train the model
- Use the trained model to make predictions
- Evaluate approximation quality (e.g., examine MSE)
- Visualize results
- Repeat steps as needed to refine model

## Obtain Training Samples

Usually this involves getting data from an external source: the internet, your own research, etc.
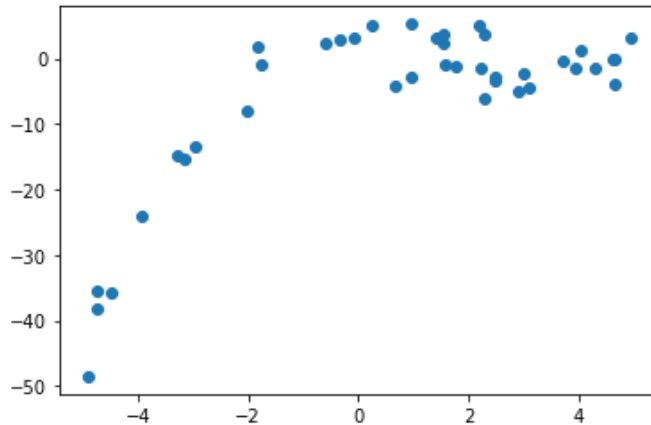
For today, we'll simply sample from our synthetic problem.

```
In [4]:    1  n = 40
           2
           3  # we'll start using scikit learn's names for things
           4  data, target = sample(n, f, [-5,5], 3)
```

## Initial Visualization

This varies a lot. One common visualization is scatter plots showing correlations between pairs of inputs and/or the training data:

```
In [5]:  1  plt.scatter(data, target)
         2  plt.show()
```



## Preprocess Data

Our initial visualization suggests non-linearity. Let's use some polynomial features.

```
In [6]:  1  Phi = phi(data, 5)
         2  Phi.shape
```

Out[6]:  (40, 5)

## Initialize a Model Object

Now we get to some actual Scikit learn code.

There are a bunch of regression models; we're going to use the linear_model.LinearRegression (http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html#sklearn.linear_model.LinearRegress one.

The basic process of obtaining a model, training, and then using it for prediction is uniform-ish across learning methods. Yay!

```
In [7]:  1  from sklearn.linear_model import LinearRegression
         2
         3  lr = LinearRegression(fit_intercept=False)
```

The `fit_intercept` parameter above defaults to True, but we already generated an intercept term in our design matrix.

## Split Data into Training/Test Sets

There's a Scikit learn [function (http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html#sklearn.model_selection.train_test](http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html#sklearn.model_selection.train_test) for this:

```
In [8]:   1  from sklearn.model_selection import train_test_split
          2
          3  X_train, X_test, y_train, y_test = train_test_split(
          4      Phi, target, test_size = 0.5)
```

The `test_size` parameter is used to control what percentage of the data to hold out for testing.
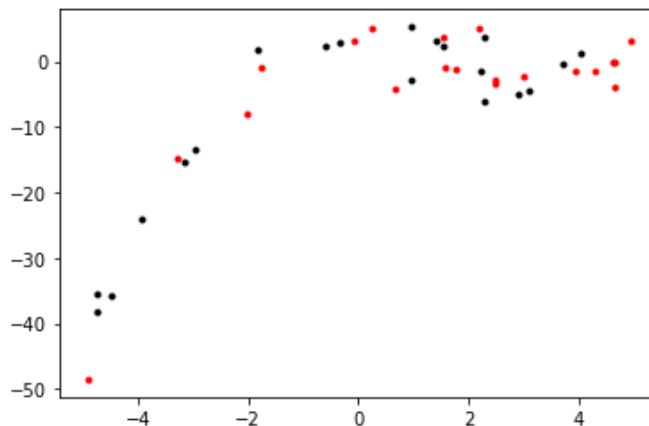
We can check to see that our data matrices/vectors are the size expected:

```
In [9]:   1  print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
```

```
(20, 5) (20, 5) (20,) (20,)
```

```
In [10]:  1  plt.plot(X_train[:,1], y_train, 'k.', X_test[:,1], y_test, 'r.')
```

```
Out[10]:  [<matplotlib.lines.Line2D at 0x1bc496843d0>,
           <matplotlib.lines.Line2D at 0x1bc49684370>]
```



## Train the Model
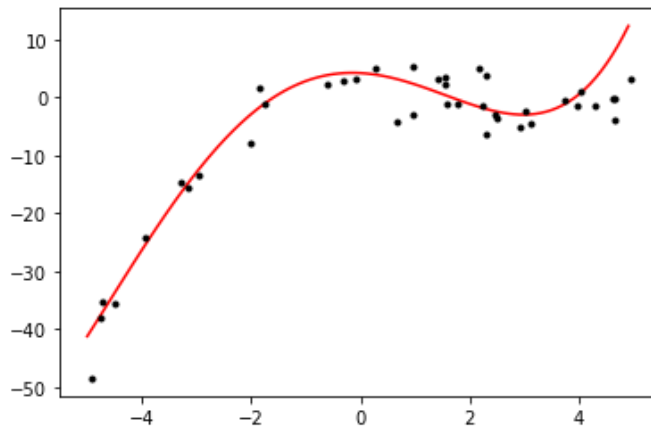
This is done using the `fit` method of the model object:

```
In [11]:  1  lr.fit(X_train, y_train)
```

```
Out[11]:  LinearRegression(fit_intercept=False)
```

## Predict (and Visualize)
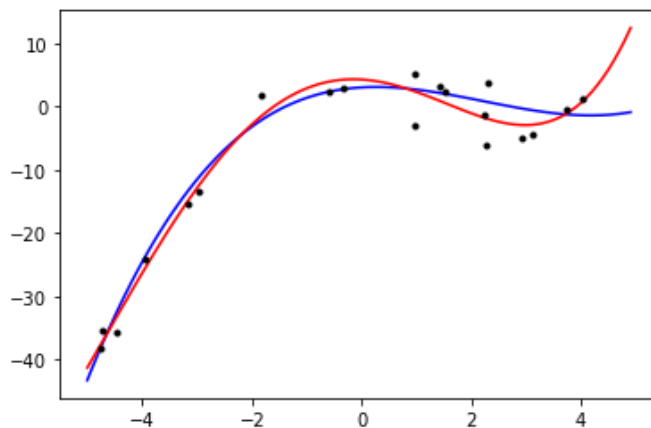
Let's see what we got:

```
In [12]:    1  x = np.arange(-5, 5, 0.1)
            2  yhat = lr.predict(phi(x, 5))
            3  plt.plot(x, yhat, 'r-', data, target, 'k.')
            4  plt.show()
```



In our case, we also know ground truth, so let's add that in:

```
In [13]:    1  plt.plot(x, f(x), 'b-', x, yhat, 'r-', X_train[:,1], y_train, 'k.')
            2  plt.show()
```



## Evaluate

Let's compute MSE and RMSE on our test set:

```
In [14]:    1  MSE = ((y_test - lr.predict(X_test)) ** 2).mean()
            2  RMSE = np.sqrt(MSE)
            3  print("MSE: ", MSE)
            4  print("RMSE:", RMSE)
```

```
MSE:   30.201240859904964
RMSE: 5.495565563243238
```

## Refine Model

A lot we could do here! For now, let's repeat our work on evaluating RMSE for different orders.

I have to redo some work from above to get all the powers up to 12...

In [15]:
```
1  Phi = phi(data, 12)
2  X_train, X_test, y_train, y_test = train_test_split(
3      Phi, target, test_size = 0.5)
```

Now I can just pare down my feature matrices using NumPy's array slicing capabilities.
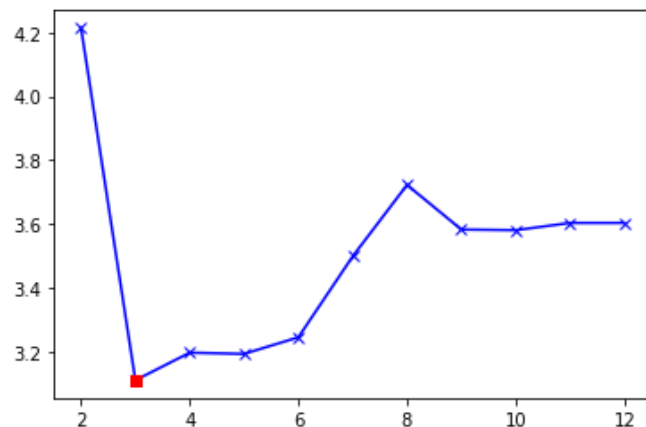
This works somewhat like list slicing, but with multi-dimensional support:

In [16]:
```
1  print(Phi.shape)
2
3  print(Phi[:,:3].shape)
```

```
(40, 12)
(40, 3)
```

Compute RMSEs across a range of orders and find the minimum:

In [18]:
```
1  RMSEs = []
2  orders = range(2,13)
3  for p in orders:
4      lr.fit(X_train[:,:(p+1)], y_train)
5      MSE = ((y_test - lr.predict(X_test[:,:(p+1)])) ** 2).mean()
6      RMSEs.append(np.sqrt(MSE))
7  RMSEs = np.array(RMSEs)
8  plt.plot(orders, RMSEs, 'b-x')
9  plt.plot(orders[RMSEs.argmin()], RMSEs.min(), 'rs') # marks the lowest RMSE value
10 plt.show()
11 print(RMSEs.min(), "at order", orders[RMSEs.argmin()])
```



```
3.1110201095243624 at order 3
```

In [19]:
```
1  lr.fit(X_train[:,:5], y_train)
```

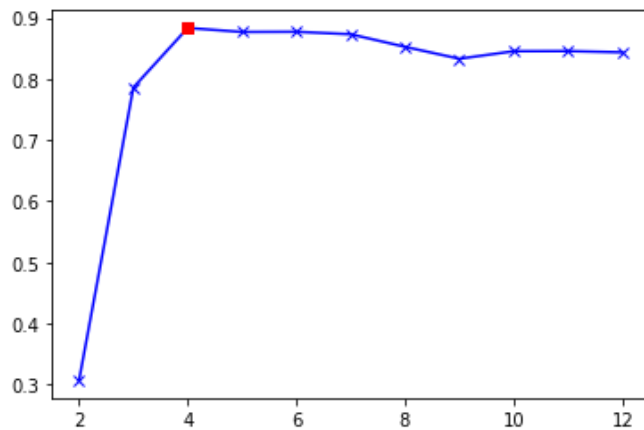Out[19]:  LinearRegression(fit_intercept=False)

```
In [20]:    1  lr.score(X_test[:,:5], y_test)
```

Out[20]:  0.8768295733190168

```
In [21]:    1  X_test.shape
```

Out[21]:  (20, 12)

```
In [23]:    1  COEFFs = []
            2  for p in orders:
            3      lr.fit(X_train[:,:p], y_train)
            4      COEFFs.append(lr.score(X_test[:,:p], y_test))
            5  COEFFs = np.array(COEFFs)
            6  plt.plot(orders, COEFFs, 'b-x')
            7  plt.plot(orders[COEFFs.argmax()], COEFFs.max(), 'rs') # plots the max value with a
            8  plt.show()
            9  print(COEFFs.max(), "at order", orders[COEFFs.argmax()])
```



0.8833925877665701 at order 4

```
In [ ]:     1
```