# Homework #1

MTH 9899 Baruch College

DATA SCIENCE II: Machine Learning

Due: April 6, 2016 - 18:00

**Notes**

- Code for this **MUST** be written in Python.

- Do NOT use $3^{rd}$ Party Packages for the regression functions.

- One thing not mentioned in class is the relationship between $\lambda$ and the size of the dataset. As the dataset grows in size, you will generallly need to include higher $\lambda$ values.

- For this assignment, please submit the relevant graphs and a short paragraph, as asked for below. Also, please submit a copy of your python code.

**Problem 1** In our first lecture, we spoke at length about Ridge Regression and the tradeoff between bias and variance. In this question, we'll run some simulations to see how $\lambda$ can affect the variance of $\beta$. Below, you will find python code to generate 4 different datasets for testing. Based on this code, you need to:

    i For datasets 1-4 in the code, generate each dataset 1000 times. For each of these 1000 times, perform simple OLS regression and record the $\beta$ values. Plot a histogram of the $\beta$ values and report the $\mu_\beta$ and $\sigma_\beta^2$.

    ii Repeat the above trials with Ridge Regression instead, using reasonable $\lambda$ values. Prepare a graph of how $\mu_\beta$ and $\sigma_\beta^2$ change as a function of $\lambda$ for each of the datasets - you do NOT need to include histograms of all of your distributions. Also, please calculate the effective degrees of freedom, to make sure that the $\lambda$ values you are using are reasonable, you should see effective DOFs from 2 down to less than 1.

    iii Calculate the expected $var(\beta^R)$ using the formulas from class. How do these compare to the simulated distributions for $\beta$? Please explain any differences.

```
import numpy as np

def generate_data(num_rows, num_features, true_betas, sigma_2 = 1, seed = None):
    """
    Args:
        num_rows (int): The number of sample rows of data
```

```
        num_features ( int ): The number of features
        true_betas ( array ): The true beta values used to generate y
        sigma_2 ( float ): The multiplier for the random noise
    """
    if seed:
        np.random.seed(seed)
    X = np.random.randn(num_rows, num_features)
    Y = X.dot(true_betas) + np.random.randn(num_rows) * sigma_2
    return X,Y


def get_dataset(set_num):
    if set_num == 1:
        return generate_data(1000, 2, np.array([1,1]))
    elif set_num == 2:
        return generate_data(50, 2, np.array([1,0]));
    elif set_num == 3:
        return generate_data(250, 2, np.array([1,0]));
    elif set_num == 4:
        return generate_data(100000, 2, np.array([1,0]));
    assert "Shouldn't be here"
```