

# Homework #1

## Q1) One feature regression for Boston data.

Choose one of the features in Boston data, try to select one that is mostly correlated (or inversely correlated) with the target. Fit a line with linear regression (you can use built-in `LinearRegression()` method). What is the equation of the fitting line? Please plot your data points (chosen feature vs. target) and plot your fitting line on it.

## Q2) Another metric for performance evaluation

We have seen that `LinearRegression()` by `sklearn` provides a method called `score` which outputs a value with a maximum of 1 (the larger the better), which is good to evaluate the performance. Another common performance metric is mean sum of squared error (MSE).

Please implement this metric and compute it when Boston data with all features are used. `lr.score` gave 0.74. What is the MSE?

Important restriction: Do NOT use `mean_squared_error` module from `sklearn.metrics`.

## Q3) Splitting the dataset

In the lecture we did training and evaluation (`score`) on the whole Boston dataset. This is bad. Use the code in Week 1 Python Notebook to:

- shuffle the data,
- select 60% of the data as a training set (do the line fitting), 40% as a test set,
- what is the score on the test data?
- what is the mean squared error?

# Homework #2

## Part 1)

- Load the breast cancer dataset using `datasets.load_breast_cancer()`
- Use scatter plots to look at the data. You can select two features for x and y axes at a time and color the samples according to their classes. You are not expected to try all pairs, which would be too many. But please show a few examples.
- Also use box plots to look at the data. Please show that you are able to examine some features (not all them are necessary).
- Create a training and test set (with shuffling).
- Train a **decision tree classifier** using the training set.
- What are the top 5 most important (discriminative) features?
- Train a **logistic regression model** using the training set.
- Which model (DT or LR) performed better on the test set? Note: You can use `score()` method of built-in classifiers to compare.

## Part 2)

- Get `winequality_white.csv` which uploaded with the Homework file.
- Last column (quality) is the target variable.
- Train a decision tree regressor using the training set. Try different '`min_samples_split`' and '`max_depth`' parameters. Which ones worked best on the test set? What is the MSE on the test set with the best parameters.

## Homework #3

1. Use Iris dataset with only two features, namely petal length (cm) and petal width (cm).
2. Convert your data into a two-class dataset, such that Virginica (class=2) will be one class and Setosa and Versicolor (class=1 and class=0) samples will constitute the other class. (Note: represent both of those classes as a single class)
3. Get your new dataset ready to be fed into the provided plot\_2d\_examples function. (Check week 8 lecture ipynb file and library documentation.)
4. Using plot\_2d\_examples function, evaluate (plot side-by-side) different SVM classifiers with **polynomial kernels of varying degrees**. Also, try several  $C$  values.
5. Do not separate the train/test sets by yourselves, just feed into plot\_2d\_examples. It has its own split inside. Also it prints test set scores on the bottom right corner of each figure. According to these scores, which (degree,  $C$ ) pair produced the best results? Please indicate your answer explicitly.
6. Also comment on  $C$  values. Which  $C$  value range is reasonable do you think? After which value, it looks like an overfit?

## Homework #4

- Download the Mall Customer Segmentation Data from the github repository below (or you can directly use the one attached in the assignment.)
- Use only 3 features (Age, Income and Spending Score) and run k-means algorithm. Determine/choose the optimum "k" value using the elbow method. Show your plot (error vs. k) and comment on plots.
- Again with the same 3 features, scatter your data in a 3D plot. Coloring of samples should denote the determined categories (your optimum k). I.e. we should be able to distinguish the samples of k categories. Comment on optimum k value.

## Homework #5

### Instructions:

Suppose your company is struggling with a series of computer virus attacks for the past several months. The viruses were grouped into a few types with some effort. However, it takes a long time to sort out what kind of virus it is when been hit with. Thus, as a senior IT department member, you undertook a project to classify the virus as quickly as possible. You've been given a dataset of the features that may be handy (or not), and also the associated virus type (target variable).

You are supposed to try different classification methods and apply best practices we have seen in the lectures such as grid search, cross validation, regularization etc. To increase your grade you can add more elaboration such as using ensembling or exploiting feature selection/extraction techniques.

You can download the data (csv file) [here](#).