

CS 434: Assignment 1

Due April 14th 11:59PM, 2017

General instructions.

1. The following languages are acceptable: Java, C/C++, Matlab, Python and R.
2. You can work in team of up to 3 people. Each team will only need to submit one copy of the source code and report.
3. You need to submit your source code (self contained, well documented and with clear instruction for how to run) and a report via TEACH. In your submission, please clearly indicate your team members' information.
4. Be sure to answer all the questions in your report. Your report should be typed, submitted in the pdf format. You will be graded based on both your code as well as the report. In particular, the clarity and quality of the report will be worth 10 % of the pts. So please write your report in clear and concise manner. Clearly label your figures, legends, and tables.

1 Linear regression

In this assignment you will use the Boston Housing dataset from the CMU StatLib Library that concerns the housing prices in Boston suburbs. The data set contains 13 attributes describing each area (e.g., crime rate, accessibility to major highways) and the target variable is the median value of housing (in thousands) for that area. The description of the data is in the file housing desc.txt associated with this assignment on canvas. The goal is to predict the median value of housing of an area based on 13 attributes. For your convenience the data has been divided into two datasets: (1) a training dataset housing train.txt you should use in the learning phase, and (2) a testing dataset housing test.txt to be used for testing. Your task is to implement the linear regression learning algorithm presented in class and explore some variations with it. In particular:

1. Given the training data, load the data into the corresponding X and Y matrices, where X stores the features and Y stores the desired outputs. The rows of X and Y correspond to the examples and the columns of X correspond to the features. Introduce the dummy variable to X by adding an extra column of ones to X (You can make this extra column to be the

first column. Adding extra column to any other position will only change the order of the learned weights and does not matter in practice).

2. Compute the optimal weight vector w using $w = (X^T X)^{-1} X^T Y$. Most programming languages have numerical packages that you can directly use to perform the computation. You don't need to implement your own matrix inversion function. Report the learned weight vector.
3. Apply the learned weight vector to the training and testing data and compute the sum of squared error (SSE) on training and testing data sets. Report the SSE values.
4. Consider the situation where we do not introduce the dummy variable to X , repeat parts [2] and [3]. Apply the learned weight vector to the training and testing data and compute the sum of squared error (SSE) on training and testing data sets. Report the SSE values. How does the dummy variable impact training and testing SSEs?
5. Modify the data by artificially generating additional random features. In particular, you need to generate several (consider for example values 2, 4, 6, 8, and 10, feel free to explore more choices) random features, each from a specified uniform distribution (i.e., generate a feature vector of uniformly distributed numbers in the interval $[0, a]$, where parameter a is different for each feature). For each case, you need to find the optimal weight vector using formula in part [2] and compute both training SSE and testing SSE. Plot the results across different numbers of random features. What trends do you observe? Do more features lead to better performance? Explain your observations.
6. Consider a variant of linear regression, where the optimal weight is computed as

$$w = (X^T X + \lambda I)^{-1} X^T Y$$

where I is the identity matrix of the same size as $X^T X$ and λ is a user specified parameter for learning. Compute the optimal w using this formula with different values of λ (e.g., 0.01, 0.05, 0.1, 0.5, 1, 5. Feel free to explore more choices.) Evaluate each of the learned w by computing the SSE on both training and testing data. Plot the SSE values as a function of λ . What behavior do you observe? What do you think is the best λ value for this problem? NOTE: this part should be done on the original data (without artificially generated features).

7. Compare the different w 's that you got in [6]. As the λ value gets bigger, what impact do you observe it has on the weight values?
8. This variant that we introduce is the solution for minimizing the following modified objective:

$$\sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$$

where the first term is the regular SSE and the second term is called a regularization term and computes the norm of the weight vector w . Can you use this objective to explain the behavior that you observe in [7]?