

HW Assignment 2, Implementation

Problems marked with a (*) are mandatory only for ITCS 8156 students. Bonus problems are optional, solving them will result in extra points.

In this assignment, you are asked to run an experimental evaluation of linear regression on the Athens houses dataset, and on an artificial dataset with and without L2 regularization. The input data is available at <https://webpages.uncc.edu/rbunescu/courses/itcs6156/hw02.zip>. Make sure that you organize your code in folders as shown in the table below. Write code only in the Python files indicated in bold.

```
itcs6156/  
  hw02/  
    report.pdf  
    code/  
      simple.py  
      multiple.py  
      polyfit.py  
      train_test_line.png  
    data/  
      simple/  
        train.txt, test.txt  
      multiple/  
        train.txt, test.txt  
      polyfit/  
        train.txt, test.txt, devel.txt
```

All the required results and plots should be included in an appropriately edited homework report, using proper indentation, section titles, and formatting. If you include formulas, make sure that you use appropriate formatting. The PDF of the report `report.pdf` should be submitted on Canvas, together with the code in the folder above. Feel free to create a Jupiter-Notebook for your code for each exercise, which you can save as a PDF using LaTeX, and submit. The assignment will be graded based on **both** the code and the report.

*Note: For exercises 1 and 2, in order to avoid numerical issues, **divide all the house prices by 100.***

1 Simple Regression (20 points)

Train a simple linear regression model to predict house prices as a function of their floor size, based on the solution to the system with 2 linear equations discussed in class (slide 14). Use the dataset from the folder `hw02/data/simple`. Python3 skeleton code is provided in `simple.py`. After training print AND REPORT the parameters AS WELL AS the RMSE and the objective function values on the training and test data. Plot the training using the default blue circles and test examples using lime green triangles. On the same graph also plot the linear approximation.

Hint: You can derive the solution of the system of linear manually, or you can use the `numpy.linalg.solve()` function as shown in the class examples.

2 Multiple Regression (20 points)

Train a multiple linear regression model to predict house prices as a function of their floor size, number of bedrooms, and year. Use the normal equations discussed in class (slide 44), and evaluate on the dataset from the folder `hw02/data/multiple`. Python3 skeleton code is provided in `multiple.py`. After training report the parameters and as well as the RMSE and the objective function values on the training and test data. Compare the test RMSE with the one from the simple case above. Do not forget to use a bias parameter w_0 and the corresponding constant feature set to 1 in each example.

3 Polynomial Curve Fitting (40 points) (*)

In this exercise, you are asked to run an experimental evaluation of a linear regression model, with and without regularization. Use the normal equations discussed in class, and evaluate on the dataset from the folder `hw02/data/polyfit`.

- Select 30 values for $x \in [0, 1]$ uniformly spaced, and generate corresponding t values according to $t(x) = \sin(2\pi x) + x(x + 1)/4 + \epsilon$, where $\epsilon = N(0, 0.005)$ is a zero mean Gaussian with variance 0.005. Save and plot all the values.
Generation and saving was done for you in `dataset.txt`.
- Split the 30 samples (x_n, t_n) in three sets: 10 samples for training, 10 samples for validation, and 10 samples for testing. Save and plot the 3 datasets separately.
Generation and saving was done for you in `train.txt`, `test.txt`, `devel.txt`.
- Consider a linear regression model with polynomial basis functions, trained with the objective shown below:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (h(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Show the closed form solution (vectorized) for the weights \mathbf{w} that minimize $J(\mathbf{w})$.

- Train and evaluate the linear regression model in the following scenarios:
 - Without regularization: Use the training data to infer the parameters \mathbf{w} for all values of $M \in [0, 9]$. For each order M , compute the RMSE separately for the training and test data, and plot all the values on the same graph, as shown in class (slide 29). *Hint: To find the parameters, you can solve the system of linear equations for polynomial curve fitting (slide 22), or the vectorized normal equations (slide 44).*
 - With regularization: Fixing $M = 9$, use the training data to infer the parameters \mathbf{w} , one parameter vector for each value of $\ln \lambda \in [-50, 0]$ in steps of 5. For each parameter vector (lambda value), compute the RMSE separately for the training

and validation data, and plot all the values on the same graph, as shown in class (slide 36, use validation instead of test). Select the regularization parameter that leads to the parameter vector that obtains the lowest RMSE on the validation data, and use it to evaluate the model on the test data. Report and compare the test RMSE with the one obtained without regularization. *Hint: For each value of $\lambda \in \{e^{-50}, e^{-45}, \dots, e^{-5}, e^0\}$, finding the parameters \mathbf{w} is done by using the vectorized normal equations for ridge regression (slide 45).*

3.1 Linear Regression using sklearn (40 bonus points)

Train and evaluate the linear regression models above using `sklearn`.

1. (10 points) Simple Regression on the Houses dataset, without regularization. Compare with the results from your implementation.
2. (10 points) Multiple Regression on the Houses dataset, without regularization. Compare with the results from your implementation.
3. (20 points) Polynomial curve fitting. For ridge regression, add the validation data to the training data and use the `RidgeCV` function to tune the hyper-parameter λ (use 10 folds). Report and compare the test RMSE with vs. without regularization.

4 Submission

Electronically submit on Canvas a `hw02.zip` file that contains the `hw02` folder in which your code is in the 3 required files, as well as the `report.pdf`.

On a Linux system, creating the archive can be done using the command:

```
> zip -r hw02.zip hw02.
```

Please observe the following when handing in homework:

1. Structure, indent, and format your code well.
2. Use adequate comments, both block and in-line to document your code.
3. Make sure your code runs correctly when used in the directory structure shown above. We will not debug your code.
4. For your report, it is recommended to use an editor such as Latex or Word or Jupyter-Notebook that allows editing and proper formatting of equations, plots, and tables with results.
5. Make sure your your Canvas submission contains the correct files by downloading and unzipping it after posting on Canvas.