

HOMWORK 1:

LINEAR REGRESSION AND LOGISTIC REGRESSION

ELEC 400M @ UBC

TAs: Sadegh Mahdavi (Primary), Yinjia Huo, Chun-Yin Huang

Instructions

- **Homework Submission:** Submit your code and report to Canvas. You will use Co-lab to implement the coding tasks. Please check Piazza for updates about the homework.
 - Upload a zip file containing two files: Your report in .PDF format and your notebook in .Ipynb format.
 - To ensure the reproducibility of your results: (1) set a seed for numpy and python random modules on top of your Colab notebook (2) restart and run all the cells of your notebook once before submission.
- **Collaboration policy:** The purpose of student collaboration is to facilitate learning, not to circumvent it. Studying the material in groups is strongly encouraged. It is also allowed to seek help from other students in understanding the material needed to solve a particular homework problem, provided no written notes (including code) are shared, or are taken at that time, and provided learning is facilitated, not circumvented. The actual solution must be done by each student alone. The presence or absence of any form of help or collaboration, whether given or received, must be explicitly stated and disclosed in full by all involved.
- **Description:** In Assignment 1, you will calculate the analytic solution of a linear model on a training dataset and report the error on both training and testing set. You will use Python and Sklearn to confirm your results. Next, you will practice using gradient descent to train a logistic regression model by coding the gradient descent algorithm by yourself. Again, you will use Python and Sklearn to confirm your results.

1 Linear Regression [7 pts]

Assume we have a training set and a testing set as follows:

Sample ID	x	y
1	5.86	0.74
2	1.34	1.18
3	3.65	0.51
4	4.69	-0.48
5	4.13	-0.07
6	4.87	0.37
7	7.91	1.35
8	5.57	0.30
9	7.30	1.64
10	7.89	1.75

Table 1: Training set

Sample ID	x	y
1	5.80	0.93
2	0.57	1.87
3	4.30	-0.06
4	6.55	1.60
5	0.82	1.22
6	3.72	0.90
7	5.80	0.93
8	3.26	1.53
9	6.75	1.73
10	4.77	-0.51

Table 2: Testing Set

Please fit the linear model $\hat{Y} = f_{\Theta}(\mathbf{X})$ using RSS objective $J(\Theta) = \|\hat{Y} - Y\|_2^2$.

(a) Calculate the analytic solution of the linear model $f_{\Theta}(x) = \theta_0 + \theta_1 x$. Then, use scatter plot to plot the training data points and draw the fitted line on the same figure.

(b) Suppose we want to increase the model complexity, by considering y as a linear function of both x and x^2 . Namely $f_{\Theta}(x) = \theta_0 x + \theta_1 x + \theta_2 x^2$. In this case, calculate the analytic solution of model and plot the curve of the model, together with the data points in training set.

(hint: in this case, $\mathbf{X} = \begin{pmatrix} 1 & x^{(1)} & x^{(1)^2} \\ \vdots & \vdots & \vdots \\ 1 & x^{(10)} & x^{(10)^2} \end{pmatrix}$)

(c) Let us further increase the model complexity by assuming y is related to higher-order forms of x , i.e., $f_{\Theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$. Again, calculate the analytic solution of the model and plot the curve of the function, together with the data points in training set.

(d) Observe the above three functions, please point out which could be faced with underfitting, which could be faced with overfitting, and which one is relatively a good one? Then, you can calculate the values of prediction error on the test data to verify your thoughts.

(e) For question (a-c), please verify your optimal Θ^* using the linear regression function in *sklearn*. The example code is provided in ELEC_400M_HW1_Example_Codes.ipynb.

2 Logistic Regression [8 pts]

Assume that we have a training set and a test set as follows: Use these data to implement a logistic classifier.

Sample ID	x_1	x_2	y
1	0.346	0.780	0
2	0.303	0.439	0
3	0.358	0.729	0
4	0.602	0.863	1
5	0.790	0.753	1
6	0.611	0.965	1

Table 3: Training set

Sample ID	x_1	x_2	y
1	0.959	0.382	0
2	0.750	0.306	0
3	0.395	0.760	0
4	0.823	0.764	1
5	0.761	0.874	1
6	0.844	0.435	1

Table 4: Testing Set

We use the linear model $f_{\Theta}(x_1, x_2) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ and the logistic regression function is written as $\sigma_{\Theta}(x_1, x_2) = \frac{1}{1 + e^{-f_{\Theta}(x_1, x_2)}}$. We use cross-entropy error as the loss function.

As introduced in the lecture, we will use gradient descent method to update the model based on the data points in the training set. The model parameters are initialized as $\theta_0 = -1, \theta_1 = -1.5, \theta_2 = 0.5$. We set step size $\alpha = 0.1$.

- (a) Please write down the logistic model $P(\hat{y} = 1|x_1, x_2)$ and its cross-entropy error function.
- (b) Use gradient descent to update θ_0, θ_1 and θ_2 for ONE iteration. Write down the gradients and updated parameters. Then please implement the iterative updating using your own Python algorithm till convergence. The example code is provided in ELEC_400M_HW1_Example_Codes.ipynb.
- (c) Use *Sklearn* to find the best θ_0, θ_1 and θ_2 for one iteration. The example code is provided in ELEC_400M_HW1_Example_Codes.ipynb.
- (d) Use the above new model to make predictions for all the samples in the test dataset. Please compare your results in (b) and the results generated by *Sklearn* in (c). Then, calculate the accuracy, precision and recall to evaluate the models.

Note

1. Remember to submit your assignment by 23:59pm of the due date. Late submission will affect your scores.
2. If you submit multiple times, ONLY the content and time-stamp of the latest one would be considered.
3. We strictly follow the rules of UBC Academic Misconduct.