

Homework Assignment 4

COGS 118A: Introduction to Machine Learning I

Due: 11:59pm, Sunday, Nov. 4th, 2018 (Pacific Time).

Instructions: Use **Jupyter Notebook** to answer the questions below, include all your code and figures in the notebook. Export the notebook as PDF file and submit it on Gradescope. You may look up the information on the Internet, but you must write the final homework solutions by yourself.

Please Note: Writing homework without using Jupyter Notebook will lead to point deduction.

Late Policy: 5% of the total points will be deducted on the first day past due. Every 10% of the total points will be deducted for every extra day past due.

Grade: ____ out of 100 points

1 (35 points) Parabola Estimation

We are given the data $S = \{(x_i, y_i), i = 1, \dots, n\}$. Here, $x_i, i = 1, \dots, n$ are one dimensional scalars. We will try to fit the data with a parabola, i.e. $y_i = w_1 x_i + w_2 x_i^2 + w_3$. To solve this problem with matrix manipulation, we represent the data as matrices $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$ and $Y = [y_1, y_2, \dots, y_n]^T$, where \mathbf{x}_i is a feature vector corresponding to the data x_i : $\mathbf{x}_i = [1, x_i, x_i^2]^T$ in this problem. $W = [w_0, w_1, w_2]^T$. We are finding W that minimizes the sum-of-squares error function $g(W)$.

Please download HW4_2018.ipynb from website

1.1 L2 norm (5 points)

Consider L_2 norm as your loss function:

$$g(W) = \|XW - Y\|_2^2 = (XW - Y)^T(XW - Y). \quad (1)$$

Use the **closed form solution** to compute W and plot the scatter graph of data and estimated parabola. Report the parabola function and the figures.

1.2 L1 norm (15 points)

Consider L_1 norm as your loss function:

$$g(W) = \|XW - Y\|_1 = \sum_{i=1}^n |y_i - f(x_i; W)| \quad (2)$$

1. (5 points) Derive the gradient of the cost function $g(W)$ with respect to W and you should have gradient as in the following form: (refer to lecture slides for hints)

$$\frac{\partial g(W)}{\partial W} = \left((\text{sign}(XW - Y))^T X \right)^T$$

where

$$\text{sign}(x) = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0. \end{cases}$$

and if X is a matrix, $\text{sign}(X)$ means performing element-wise $\text{sign}(x_{ij})$ over all element x_{ij} in X .

2. (10 points) Use the **gradient descent method** to find W^* and plot the scatter graph of data and estimated parabola. Report the parabola function and the figures.

Hint 1: In NumPy, you can use `np.sign(x)` to compute the sign of matrix x .

Hint 2: You may need to change the number of iterations to 300000, learning rate to 0.000001 and error threshold for W to 0.00001. Same settings may apply to sub-problem (c).

1.3 L1 and L2 norm (10 points)

Consider L_1 and L_2 norm as your loss function:

$$\begin{aligned} g(W) &= \alpha \|XW - Y\|_2^2 + (1 - \alpha) \|XW - Y\|_1 \\ &= \sum_{i=1}^n \left(\alpha (y_i - f(x_i; W))^2 + (1 - \alpha) |y_i - f(x_i; W)| \right) \end{aligned} \quad (3)$$

Use the **gradient descent method** to find W^* when $\alpha = 0.3$, $\alpha = 0.5$, and $\alpha = 0.7$ respectively and plot the scatter graph of data and estimated parabola. Report the parabola function and the figures.

1.4 Comparison (5 points)

Compare the result from above three subsections (1.1-1.3) and plot all curves in one figure along with scatter graph of data. Try to explain the reason to (1) the position of each curve compared to the position of valid data points and outliers (2) difference between L_2 curve and L_1 curve (3) similarity among L_2 curve and $L_1 + L_2$ curves.

2 (15 points) Logistic Regression (1)

For a logistic regression function with input $x \in \mathbb{R}$ and output $y \in \{0, 1\}$, the probability of $P(y = 1|x) = \frac{e^{\alpha+\beta x}}{1 + e^{\alpha+\beta x}}$.

1. Please write down the formulation of $P(y = 0|x)$.

2. Show that $[P(y = 1|x)]^y \times [P(y = 0|x)]^{1-y} = \frac{1}{1 + e^{-(2y-1)(\alpha+\beta x)}}$

3. What is the decision boundary for classifier:

$$y = \begin{cases} 1, & \alpha + \beta x \geq 0, \\ 0, & \text{else} \end{cases}$$

and how is it related to $P(y = 1|x)$.

3 (10 points) Logistic Regression (2)

For a logistic regression function with input $\mathbf{x} \in \mathbb{R}^m$, i.e. \mathbf{x} is a m-dimensional vector, and output $y \in \{0, 1\}$, the probability of $P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$.

1. Please show that $P(y|\mathbf{x}) = \frac{1}{1 + e^{-(2y-1) \times (\mathbf{w}^T \mathbf{x} + b)}}$

2. Please analyze the decision boundary for this logistic regression classifier. (On what condition of \mathbf{x} , y will be predicated as 1 or 0?)

4 (40 points) Logistic Regression (3)

In logistic regression model, y is the label for each data point, and it can be either 0 or 1. Here, we define our approximate function to

$$p(y = 1|\mathbf{x}) = h(\mathbf{x}; \mathbf{w}, b) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

where $\mathbf{x} \in \mathbb{R}^m$ is the feature, and b is the bias, and $\mathbf{w} \in \mathbb{R}^K$ contains the set of parameters $\mathbf{w} = (w_0, w_1, \dots, w_m)$. The output of $h(\mathbf{x}; \mathbf{w}, b)$ is called the confidence. When $h(\mathbf{x}; \mathbf{w}, b) > 0.5$, the classifier outputs 1 for the given \mathbf{x} ; otherwise, the classifier outputs 0.

Dataset

Download the data file **Q4_data.txt** from the course website. This dataset is modified from **Iris** dataset. The dataset contains 2 classes of 50 instances each, and the two classes are *Iris-versicolour* and *Iris-virginica*. The dataset currently contains 100 instances. The first 15 instances of each class are used for testing, and the rest of data are used for training. The code for the above preprocessing is provided.

Model

In this task, the first 4 attributes and the 5th attribute of i -th instance are considered as its feature $\mathbf{x}_i \in \mathbb{R}^4$ and its label $y_i \in \mathbb{R}$.

The goal is to train a classifier based on logistic regression to predict the correct label y_i from the given feature \mathbf{x}_i . We define a loss function which measures the distance between the correct label and the prediction from the classifier, as is shown below:

$$\mathcal{L}(\mathbf{w}) = - \sum_i \ln p(y_i|\mathbf{x}_i; \mathbf{w}, b)$$

where

$$p(y_i|\mathbf{x}_i; \mathbf{w} + b) = \frac{1}{1 + e^{-(2y_i - 1) \times (\mathbf{w}^T \mathbf{x}_i + b)}},$$

which is called the sigmoid function.

4.1 Gradient Descent (10 pts)

The training procedure is to minimize the loss function on the training set. Consider **gradient descent** method to find the optimal \mathbf{w}^* in $h(\mathbf{x}; \mathbf{w}, b)$.

1. Derive $\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}}$.
2. Suppose that the learning rate is denoted as α . Write down the update rule for \mathbf{w} .

4.2 Training (10 pts)

Implement your own **gradient descent** algorithm to train a binary classifier based on logistic regression. You might need to vectorize your algorithm in order to run efficiently. Note that b is also supposed to be learned in your gradient descent algorithm. **Plot the training curve: cost function vs. # iterations**

4.3 Decision Boundary (10 pts)

Derive the equation of the decision boundary for your trained classifier, and Plot training data and test data separately along with decision boundary

4.4 Test (10 pts)

Classify the data in the test set and report the accuracy.