

CMPSC 448: Machine Learning and Algorithmic AI

Homework 2

Haorui Lyu

September 2022

Instruction

This HW includes both theory and implementation problems:

- You cannot look at anyone else's code
- Your homework must work with Python 3.7 (you may install the Anaconda distribution of Python)
- You need to submit a report including all deliverable and figures (in PDF format), also three files `Problem4.ipynb`, `Problem5.py`, and `Problem5Plot.py`
- The only modules your code can import are: `math`, `numpy`, `matplotlib`, `random`, `pandas`, `sklearn`

Theory and problem solving

Problem 1. [20 points] In the lectures, we showed that the MLE estimator for linear regression when the random noise for each data point is identically and independently distributed (i.i.d.) from a Gaussian distribution $\mathcal{N}(0, \sigma^2)$ reduces to the linear regression with squared loss (OLS). In this problem, we would like to change the distribution of the noise model and derive the MLE optimization problem. In particular, let's assume for a fixed unknown linear model $\mathbf{w}_* \in \mathbb{R}^d$ the response y_i for each data point $\mathbf{x}_i \in \mathbb{R}^d$ in training data $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ is generated by

$$y_i = \mathbf{w}_*^\top \mathbf{x}_i + \epsilon_i$$

where ϵ_i is generated i.i.d. from a Laplace distribution $\epsilon_i \sim \text{Laplace}(0, \sigma)$. A random variable x has a $\text{Laplace}(\mu, \sigma)$ distribution if its probability density function is

$$f(x|\mu, \sigma) = \frac{1}{2\sigma} \exp\left(-\frac{|x - \mu|}{\sigma}\right)$$

Under above assumption on the noise model,

1. Show that each $y_i, i = 1, 2, \dots, n$ is a random variable that follows $\text{Laplace}(\mathbf{w}_*^\top \mathbf{x}_i, \sigma)$ distribution.
2. Write down the MLE estimator for training data in \mathcal{S} and derive the final optimization problem. Note that you just need to state the final minimization problem and not its solution.
3. Compare the obtained optimization problem to one obtained under Gaussian noise model and highlight key differences.

Solution :

1. We find the moment generating function for y_i :

$$\begin{aligned} M_{y_i} &= E[e^{t(\mathbf{w}_*^\top \mathbf{x}_i + \epsilon_i)}] \\ M_{y_i} &= E[e^{t\mathbf{w}_*^\top \mathbf{x}_i + t\epsilon_i}] \\ M_{y_i} &= e^{t\mathbf{w}_*^\top \mathbf{x}_i} E[e^{t\epsilon_i}] \end{aligned}$$

We know that ϵ_i is from a Laplace distribution $\epsilon_i \sim \text{Laplace}(0, \sigma)$. We know that the MGF for Laplace distribution Laplace: $(\mu, b) = \frac{\exp(\mu t)}{1 - b^2 t^2}$ Therefore the Momente generation function for ϵ_i is:

$$M_{\epsilon_i} = \frac{1}{1 - \sigma^2 t^2}$$

Then:

$$M_{y_i} = \frac{e^{t\mathbf{w}_*^\top \mathbf{x}_i}}{1 - \sigma^2 t^2}$$

Which means y_i is also from Laplace distribution. Therefore $y_i, i = 1, 2, \dots, n$ is random variable.

2. What is the distribution of $y_i = \mathbf{w}_*^\top \mathbf{x}_i + \epsilon_i$, y_i from Laplace($\mathbf{w}_*^\top \mathbf{x}_i + \epsilon_i, \sigma$).

The distribution is: $P[y_i | x_i; \mathbf{w}]$.

Let's compute the likelihood for a fixed parameter vector:

$$\begin{aligned} L(S | \mathbf{w}) &= \prod_{i=1}^n P[y_i | x_i; \mathbf{w}] \\ L(S | \mathbf{w}) &= \frac{1}{2\sigma} \exp\left(-\frac{|y_i - \mathbf{w}_*^\top \mathbf{x}_i|}{\sigma}\right) \end{aligned}$$

Let's take the log:

$$L(S | \mathbf{w}) = -\ln 2\sigma - \sum_{i=1}^n \frac{|y_i - \mathbf{w}_*^\top \mathbf{x}_i|}{\sigma}$$

Therefore:

$$\text{MLE} = \arg \max_{\mathbf{w}} - \sum_{i=1}^n \frac{|y_i - \mathbf{w}_*^\top \mathbf{x}_i|}{\sigma}$$

3. We find that the MLE for one obtained under Gaussian noise model is $\arg \max_{\mathbf{w}} - \sum_{i=1}^n \frac{(y_i - \mathbf{w}_*^\top \mathbf{x}_i)^2}{2\sigma^2}$. We can find the key difference is we calculate absolute value for $y_i - \mathbf{w}_*^\top \mathbf{x}_i$ but in Gaussian noise model we calculate the square of it. Also for σ .

Problem 2. [15 points] Suppose we run a ridge regression with regularization parameter λ on a training data with a single variable $\mathcal{S} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, and get coefficient $w_1 \in \mathbb{R}$ (for simplicity, we assumed the data are centered and no bias (intercept) term is needed). We now include an exact copy of first feature to get a new training data as

$$\mathcal{S}' = \{([x_1, x_1]^\top, y_1), ([x_2, x_2]^\top, y_2), \dots, ([x_n, x_n]^\top, y_n)\}$$

where each training example is a 2 dimensional vector with equal coordinates, refit our ridge regression on \mathcal{S}' and get the solution $[w'_1, w'_2]^\top \in \mathbb{R}^2$.

1. Derive the optimal solution for $[w'_1, w'_2]^\top$ and show that $w'_1 = w'_2$.
2. What is relation between w'_1 and w_1 .

Solution :

2.1. We know this problem run a ridge regression.

We derive the optimal solution of RR:

$$w_{RR}^* = \left(\sum_{i=1}^n x_i x_i^T + 2\lambda I \right)^{-1} \sum_{i=1}^n y_i x_i$$

We find $\left(\sum_{i=1}^n x_i x_i^T + 2\lambda I \right)^{-1}$ first :

$$\begin{aligned} \left(\sum_{i=1}^n x_i x_i^T + 2\lambda I \right)^{-1} &= \begin{bmatrix} \sum_{i=1}^n x_i^2 + 2\lambda & \sum_{i=1}^n x_i^2 \\ \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i^2 + 2\lambda \end{bmatrix}^{-1} \\ &= \frac{1}{\left(\sum_{i=1}^n x_i^2 + 2\lambda \right)^2 - \left(\sum_{i=1}^n x_i^2 \right)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 + 2\lambda & -\sum_{i=1}^n x_i^2 \\ -\sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i^2 + 2\lambda \end{bmatrix} \\ &= \frac{1}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 + 2\lambda & -\sum_{i=1}^n x_i^2 \\ -\sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i^2 + 2\lambda \end{bmatrix} \\ &= \begin{bmatrix} \frac{\sum_{i=1}^n x_i^2 + 2\lambda}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} & \frac{-\sum_{i=1}^n x_i^2}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} \\ \frac{-\sum_{i=1}^n x_i^2}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} & \frac{\sum_{i=1}^n x_i^2 + 2\lambda}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \text{then } w_{RR}^* &= \begin{bmatrix} \frac{\sum_{i=1}^n x_i^2 + 2\lambda}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} & \frac{-\sum_{i=1}^n x_i^2}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} \\ \frac{-\sum_{i=1}^n x_i^2}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} & \frac{\sum_{i=1}^n x_i^2 + 2\lambda}{4\lambda^2 + 4\lambda \sum_{i=1}^n x_i^2} \end{bmatrix} \begin{bmatrix} \sum_{i=1}^n y_i x_i \\ \sum_{i=1}^n y_i x_i \end{bmatrix} \\ &= \begin{bmatrix} \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2} \\ \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2} \end{bmatrix} = \begin{bmatrix} w_1' \\ w_2' \end{bmatrix} \end{aligned}$$

Therefore we find $w_1' = w_2'$.

2.2

$$\begin{aligned}
 \mathbf{w}_{RR}^* &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T + 2\lambda \mathbf{I} \right)^{-1} \sum_{i=1}^n y_i \mathbf{x}_i \\
 \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T + 2\lambda \mathbf{I} \right)^{-1} &= \frac{1}{\sum_{i=1}^n x_i^2 + 2\lambda} \\
 w_1 &= \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2 + 2\lambda} \quad w_1' = \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2}
 \end{aligned}$$

Here we can see the relation between w_1 and w_1'

Problem 3. [10 points] As we discussed in the lecture, the Perceptron algorithm will converge only if the data is linearly separable. In particular, for linearly separable data with margin γ , if $\|\mathbf{x}\|_2 \leq R$ for all data points, then it will converge in at most $\left(\frac{R}{\gamma}\right)^2$ iterations as stated in the class.

If the training data $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{-1, +1\}$ is not linearly separable, then there is a simple trick to force the data to be linearly separable and then apply the Perceptron algorithm as follows. If you have n data points in d dimensions, map data point \mathbf{x}_i to the $(d+n)$ -dimensional point $[\mathbf{x}_i, \mathbf{e}_i]^T \in \mathbb{R}^{d+n}$, where $\mathbf{e}_i \in \{0, 1\}^n$ is a n -dimensional vector of zeros, except for the i th position, which is 1 (e.g., $\mathbf{e}_4 = [0, 0, 0, 1, 0, \dots, 0]^T$).

Show that if you apply this mapping, the data becomes linearly separable (you may wish to do so by providing a weight vector \mathbf{w} in $(d+n)$ -dimensional space that successfully separates the data).

Solution :

After the mapping, the data will become

$$\begin{bmatrix} x_1 & 1 & 0 & 0 & \dots & 0 \\ x_2 & 0 & 1 & 0 & \dots & 0 \\ x_3 & 0 & 0 & 1 & \dots & 0 \\ \dots & & & & & \\ x_n & 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

We can process the Perceptron algorithm to separate data. We know that the original \mathcal{S} is not separable. That means the algorithm will loop forever. But after we mapping this data, it should be separable. Here is some example. Suppose $n = 3$, and we have data after mapping $[1, 2, 3, 1, 0, 0]$, $[2, 1, 3, 0, 1, 0]$, $[1, 3, 2, 0, 0, 1]$, we can provide weight vector \mathbf{w} such as $[0, 0, 0, 1, 1, -1]$ to separate this data set.

Programming and experiment

Problem 4. [25 points] In this problem you will use the Pima Indians Diabetes dataset from the UCI repository to experiment with the k -NN algorithm and find the optimal value for the number of neighbors k . You do not need to implement the algorithm and encouraged to use the implementation in `scikit-learn`. Below is a simple code showing the steps to use the NN implementation when the number of neighbors is 3:

To accomplish this task, please do:

- Download the provided `Pima.csv` data file and load it using `pandas`. As a sanity check, make sure there are 768 rows of data (potential diabetes patients) and 9 columns (8 input features including Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, Age, and 1 target output). Note that the data file has no header and you might want to explicitly create the header. The last value in each row contains the target label for that row, and the remaining values are the features. Report the statistics of each feature (min, max, average, standard deviation) and the histogram of the labels (target outputs). For each feature, we need a graph that describes how labels are distributed across the different intervals. For example, in age, we need to draw how many data are labeled 1 and 0, among different intervals (say $[20,30]$, $[30,40]$, etc.).
- Split the data into training and test sets with 80% training and 20% test data sizes. You can easily do this in `scikit-learn`, e.g., Use 5-fold cross-validation on training data to decide the best number of neighbours k . To this end, you can use the built in functionality in `scikit-learn` such as `cross_val_score` (note that this function returns the accuracy for each fold in an array and you have to average them to get the accuracy for all splits). For $k = 1, 2, 3, \dots, 15$ compute the 5-fold cross validation error and plot the results (with values of k on the x -axis and accuracy on the y -axis). Include the plot in your report and justify your decision for picking a particular number of neighbors k .
- Evaluate the k -NN algorithm on test data with the optimal number of neighbours you obtained in previous step and report the test error.
- Process the input data by subtracting the mean (a.k.a. centralization) and dividing by the standard deviation (a.k.a. standardization) over each dimension (feature), repeat the previous part and report the accuracy. Do centralization and standardization affect the accuracy? Why?

Running and Deliverable

You are provided with a `Problem4.ipynb` file to put the implementation code for all parts above. Make sure your notebook is runnable on Anaconda with Python 3.7. The results and discussions should be included in the PDF file.

Problem 5. [30 points] In this problem, we consider a simple linear regression model with a modified loss function and try to solve it with Gradient Descent (GD) and Stochastic Gradient Descent (SGD).

In general setting, the data has the form $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ where \mathbf{x}_i is the d -dimensional feature vector and y_i is a real-valued target. For this regression problem, we will be using linear prediction $\mathbf{w}^\top \mathbf{x}_i$ with the objective function:

$$f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n g_\delta(\mathbf{w}; \mathbf{x}_i, y_i) + \lambda \sum_{j=1}^d w_j^2$$

where $g_\delta(\mathbf{w}; \mathbf{x}, y)$ is the error of linear model with parameter vector $\mathbf{w} \in \mathbb{R}^d$ on a single training pair (\mathbf{x}, y) defined by:

$$g_\delta(\mathbf{w}; \mathbf{x}, y) = \begin{cases} (y - \mathbf{w}^\top \mathbf{x} - \delta)^2 & \text{if } y \geq \mathbf{w}^\top \mathbf{x} + \delta \\ 0 & \text{if } |y - \mathbf{w}^\top \mathbf{x}| < \delta \\ (y - \mathbf{w}^\top \mathbf{x} + \delta)^2 & \text{if } y \leq \mathbf{w}^\top \mathbf{x} - \delta \end{cases}$$

Please note that we simply dropped the intercept term by the simple trick we discussed in the lecture, i.e., adding a constant feature, which is always equal to one to simplify estimation of the “intercept” term.

Gradient Descent

In this part, you are asked to optimize the above objective function using gradient descent and plot the function values over different iterations, which can be done using the python library `matplotlib`.

To this end, in `Problem5.py`, fill in the function `bgd_l2(data, y, w, eta, delta, lam, num_iter)` where `data` is a two dimensional numpy array with each row being a feature vector, `y` is a one-dimensional numpy array of target values, `w` is a one-dimensional numpy array corresponding to the weight vector, `eta` is the learning rate, `delta` and `lam` are parameters of the objective function. This function should return new weight vector, history of the value of objective function after each iteration (python list).

Run this function for the following settings and plot the history of objective function (you should expect a monotonically decreasing function over iteration number):

- $\eta = 0.05, \delta = 0.1, \lambda = 0.001, \text{num_iter} = 50$
- $\eta = 0.1, \delta = 0.01, \lambda = 0.001, \text{num_iter} = 50$
- $\eta = 0.1, \delta = 0, \lambda = 0.001, \text{num_iter} = 100$
- $\eta = 0.1, \delta = 0, \lambda = 0, \text{num_iter} = 100$

Stochastic Gradient Descent

In `Problem5.py` fill in the function `sgd_l2(data, y, w, eta, delta, lam, num_iter, i)` where `data, y, w, lam, delta, num_iter` are same as previous part. In this part, you should use $\frac{1}{\sqrt{t}}$ as a learning rate, where `t` is the iteration number, starting from 1. The variable `i` is for testing the correctness of your function. If `i` set to -1 then you just need to apply the normal SGD (randomly select the data point), which runs for `num_iter`, but if `i` set to something else (other than -1), your code only needs to compute **SGD** for that specific data point (in this case, the `num_iter` will be 1!).

Run this function for the settings below and plot the history of objective function:

- $\eta = 1, \delta = 0.1, \lambda = 0.5, \text{num_iter} = 800$
- $\eta = 1, \delta = 0.01, \lambda = 0.1, \text{num_iter} = 800$
- $\eta = 1, \delta = 0, \lambda = 0, \text{num_iter} = 40$
- $\eta = 1, \delta = 0, \lambda = 0, \text{num_iter} = 800$

Running and Deliverable

You are provided with a sample data sets (`data.npy`). You can load these files into `numpy` array using this syntax: `np.load('data.npy')`. The sample data for this homework has 100 data points (x_i, y_i) in a 100×2 `numpy` array. Please note that you need to add a column of all ones to data to handle the intercept term, making your data a 100×3 `numpy` array. For the plotting part, make sure that your plots have appropriate title, x and y -axis labels. You need to submit two python files `Problem5.py` and `Problem5Plot.py` and a PDF file including all the plots. In `Problem5.py` everything except the imports should be inside the functions definition we mentioned above. The file `Problem5Plot.py` is where you are generating the plots.