



DDA5001 · Homework 1

Due: 23:59, Sep 28.

Instructions:

- Homework problems must be clearly answered to receive full credit. Complete sentences that establish a clear logical progress are highly recommended.
- Submit your answer paper in Blackboard. Please upload a file or a zip file. The file name should be in the format **last name-first name-hw1**.
- The homework must be written in English.
- Late submission will not be graded.
- Each student **must not copy** homework solutions from another student or from any other source. You are encouraged to discuss with others. However, you must write down the answer using your understanding (after discussion) and words.

Problem 1 (12 points). Concepts and Fundamental Knowledge

- (a) (4 points) Concisely state the difference between supervised learning and unsupervised learning.
- (b) (4 points) Which one of the following statements is true?
- 1) Regression is used to fit categorical labels.
 - 2) Least squares must have infinitely many solutions if the number of data points is smaller than feature dimension.
 - 3) The perceptron always converges to a unique linear classifier for a given training dataset.
 - 4) Least squares is a maximum likelihood estimator.
- (c) (4 points) Let the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a full column rank matrix. Explain why $\mathbf{X}^\top \mathbf{X}$ is positive definite (PD). (Hint: Use the definition of linear independence)

Since a PD matrix must be invertible as well, this clarifies why we can take matrix inverse to obtain the unique solution of least squares in case I (where \mathbf{X} has full column rank). From an optimization perspective, this PD property also reveals that the least squares in case I is strongly convex.

Problem 2 (17 points). Least Squares Without Full Column Rank

Consider the least squares problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2,$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\boldsymbol{\theta} \in \mathbb{R}^d$, $\mathbf{y} \in \mathbb{R}^n$.

- (a) (8 points) When $n < d$, the problem yields infinite many solutions. Assume that $\text{rank}(\mathbf{X}) = n$, derive the expression of the infinitely many solutions based on singular value decomposition (SVD). In addition, what is the optimal function value?

Note that the SVD is a notion from linear algebra, and the SVD of the matrix \mathbf{X} can be written as

$$\mathbf{X} = \underbrace{\mathbf{V}}_{\in \mathbb{R}^{n \times n}} \underbrace{\begin{bmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \end{bmatrix}}_{\in \mathbb{R}^{n \times d}} \underbrace{\begin{bmatrix} \mathbf{U}_1^\top \\ \mathbf{U}_2^\top \end{bmatrix}}_{\in \mathbb{R}^{d \times d}} = \mathbf{V}\boldsymbol{\Sigma}_1\mathbf{U}_1^\top.$$

where $\mathbf{V} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix such that $\mathbf{V}^\top \mathbf{V} = \mathbf{V}\mathbf{V}^\top = \mathbf{I}$, $\boldsymbol{\Sigma}_1 \in \mathbb{R}^{n \times n}$ is a diagonal matrix taking the form

$$\boldsymbol{\Sigma}_1(i, j) = \begin{cases} \sigma_i, & i = j \\ 0, & i \neq j \end{cases} \quad \text{with} \quad \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n > 0,$$

$\mathbf{U}_1^\top \in \mathbb{R}^{n \times d}$ is an semi-orthogonal matrix satisfying $\mathbf{U}_1^\top \mathbf{U}_1 = \mathbf{I}$ (but $\mathbf{U}_1 \mathbf{U}_1^\top \neq \mathbf{I}$), and $\mathbf{U}_2^\top \in \mathbb{R}^{(d-n) \times d}$ is an semi-orthogonal matrix satisfying $\mathbf{U}_2^\top \mathbf{U}_2 = \mathbf{I}$ (but $\mathbf{U}_2 \mathbf{U}_2^\top \neq \mathbf{I}$).

(Hint: Let $\mathbf{A} := \mathbf{V}\boldsymbol{\Sigma}_1$ which is square matrix of full rank, $\mathbf{z} = \mathbf{U}_1^\top \boldsymbol{\theta}$, solve $\min_{\mathbf{z}} \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2^2$ first, then solve $\mathbf{U}_1^\top \boldsymbol{\theta} = \mathbf{z}$ by considering that $\mathbf{U}_1^\top \mathbf{U}_2 = \mathbf{0}$.)

- (b) (5 points) Later in the overfitting section, we will study that we can also add a regularizer to solve the issue of infinitely many solutions, resulting in the following problem:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_2^2,$$

where λ is a positive constant. Derive the unique solution of this problem.

Problem 3 (21 points). Robust Linear Regression

Suppose we have the generative linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\epsilon},$$

where $\boldsymbol{\epsilon}$ is the error term and $\boldsymbol{\epsilon} \sim_{i.i.d.} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ and $\mathbf{X} \in \mathbb{R}^{n \times d}$ has full column rank. Under this Gaussian noise setup, the maximum likelihood estimator for $\boldsymbol{\theta}^*$ is given by the least squares:

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{LS} &= \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

- (a) (6 points) Suppose the error term, $\boldsymbol{\epsilon} = [\epsilon_1, \epsilon_2, \dots, \epsilon_n]$ follows the Laplace distribution, i.e., $\epsilon_i \stackrel{i.i.d.}{\sim} L(0, b)$, $i = 1, 2, \dots, n$ and the probability density function is $p(\epsilon_i) = \frac{1}{2b} e^{-\frac{|\epsilon_i - 0|}{b}}$ for some $b > 0$. Under the maximum likelihood estimation principle, derive the machine learning problem formulation for estimating $\boldsymbol{\theta}^*$.

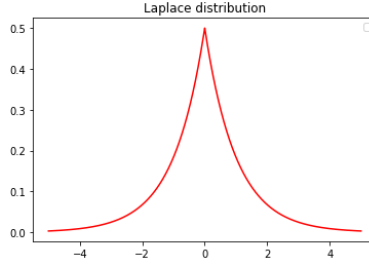


Figure 1: PDF of Laplace distribution

- (b) (5 points) **Huber-smoothing.** ℓ_1 -norm minimization

$$\hat{\boldsymbol{\theta}}_{L1} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_1$$

is a popular approach for robust linear regression. However, it is nondifferentiable. We utilize smoothing technique for approximately solving the ℓ_1 -norm minimization problem. **Huber function is often the choice for smoothing ℓ_1 function.** The definition and sketch map are shown as below. For some $\mu > 0$,

$$h_\mu(z) = \begin{cases} |z|, & |z| \geq \mu \\ \frac{z^2}{2\mu} + \frac{\mu}{2}, & |z| \leq \mu \end{cases}$$

Then,

$$H_\mu(\mathbf{z}) = \sum_{j=1}^n h_\mu(z_j).$$

By using Huber smoothing, the approximation of the above ℓ_1 -norm-based robust linear regression problem can be smoothed as

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} H_\mu(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}).$$

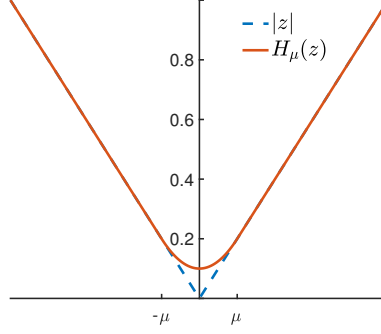


Figure 2: Huber smoothing

Define

$$\mathcal{L}(\boldsymbol{\theta}) = H_{\mu}(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}),$$

find the gradient $\nabla \mathcal{L}(\boldsymbol{\theta})$.

- (c) (10 points) **Gradient descent for minimizing $\mathcal{L}(\boldsymbol{\theta})$.** The procedure of the gradient descent method is shown in the following algorithm pseudocode. (We will study the details of the gradient-based optimization algorithms later)

-
1. **Input:** Training data \mathbf{X}, \mathbf{y} and initialization $\boldsymbol{\theta}_0$
Huber smoothing parameter μ ,
total iteration number T ,
learning rate α .
 2. **for** $k = 0, 1, \dots, T$, **do**
 3. $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \alpha \nabla \mathcal{L}(\boldsymbol{\theta}_k)$
 4. **end for**
 5. **return** $\boldsymbol{\theta}_T$
-

The data set is generated by the linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\epsilon}_1 + \boldsymbol{\epsilon}_2,$$

where each entry of $\boldsymbol{\epsilon}_1 \in \mathbb{R}^n$ follows i.i.d. Gaussian distribution, while $\boldsymbol{\epsilon}_2 \in \mathbb{R}^n$ is a sparse vector that contains *outliers*. Given the observed training data \mathbf{X} and \mathbf{y} ,

- (1) calculate the estimation $\hat{\boldsymbol{\theta}}_{LS}$ returned by least squares and compute $\left\| \hat{\boldsymbol{\theta}}_{LS} - \boldsymbol{\theta}^* \right\|_2$.
- (2) suppose $n = 1000, d = 50$, use Python to implement the gradient descent method to minimize $\mathcal{L}(\boldsymbol{\theta})$ in part (b), the parameters are set as $\mu = 10^{-5}, \alpha = 0.001, T = 1000$, plot the error $\left\| \boldsymbol{\theta}_k - \boldsymbol{\theta}^* \right\|_2$ as a function of iteration number. You can download the data $\{\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^*\}$ on Blackboard.

Problem 4 (24 points). Convergence of The Perceptron for Linearly Separable Data

In lecture, we studied that one iteration of perceptron algorithm will move forward to the currently wrongly classified data. Though such an observation does provide some intuition that perceptron will eventually converge to a linear classifier that correctly separate all the (linearly separable) training data, it is not rigorously justified.

Let us first review the algorithm. The perceptron applies to the linear model

$$f_{\theta}(\mathbf{x}) = \theta^{\top} \mathbf{x}.$$

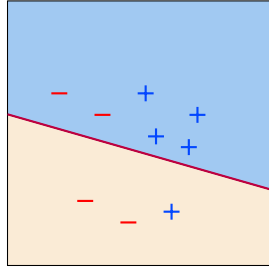
At k -th iteration, the algorithm has parameter θ_{k-1} . Based on θ_{k-1} , the algorithm first pick a wrongly classified data point $(\mathbf{x}_{k-1}, y_{k-1})$ from all the data points $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_n, y_n))$, i.e.,

$$\text{sign}(\theta_{k-1}^{\top} \mathbf{x}_{k-1}) \neq y_{k-1}.$$

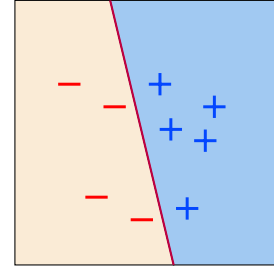
Then, perceptron updates as

$$\theta_k = \theta_{k-1} + y_{k-1} \mathbf{x}_{k-1} \quad \forall k = 1, 2, \dots \quad (1)$$

In this question, you will be guided to show that perceptron will eventually stop, i.e., there exists a \bar{k} such that no wrongly classified data exists after at most \bar{k} number of iterations, and then the perceptron learning algorithm will automatically terminate; see Fig. 3 for an illustration.



(a) Perceptron at the k -th iteration, not all data are correctly classified



(b) Learned perceptron model after at most \bar{k} iterations, all data are correctly classified

Figure 3: Illustration of the convergence progress of the perceptron learning algorithm.

Let θ^* be corresponding to a classifier that correctly separates all the data points like the one in Fig. 3b. For simplicity, we assume that the initial point of the perceptron algorithm is $\theta_0 = \mathbf{0}$.

The full proof procedure is divided into 5 steps.

- (a) (3 points) Let $\rho = \min_{1 \leq i \leq n} y_i(\theta^{*\top} \mathbf{x}_i)$. Show that $\rho > 0$.
- (b) (6 points) Show that $\theta_k^{\top} \theta^* \geq \theta_{k-1}^{\top} \theta^* + \rho$, and hence conclude that $\theta_k^{\top} \theta^* \geq k\rho$.
(Hint: Use the update (1).)
- (c) (6 points) Show that $\|\theta_k\|^2 \leq \|\theta_{k-1}\|^2 + \|\mathbf{x}_{k-1}\|^2$.
(Hint: Use the fact that \mathbf{x}_{k-1} is misclassified by θ_{k-1} .)

(d) (3 points) Show that $\|\boldsymbol{\theta}_k\|^2 \leq kR^2$, where $R = \max_{1 \leq i \leq n} \|\mathbf{x}_i\|$.

(e) (6 points) Using steps 2 and 4, show that

$$\frac{\boldsymbol{\theta}_k^\top \boldsymbol{\theta}^*}{\|\boldsymbol{\theta}_k\|} \geq \sqrt{k} \frac{\rho}{R}$$

and hence show that the perceptron learning algorithm must stop within at most

$$\bar{k} \leq \frac{R^2 \|\boldsymbol{\theta}^*\|^2}{\rho^2}$$

number of iterations.

(Hint: Use the fact $\frac{\boldsymbol{\theta}_k^\top \boldsymbol{\theta}^*}{\|\boldsymbol{\theta}_k\| \|\boldsymbol{\theta}^*\|} \leq 1$. Why?)

In practice, perceptron will converge often faster than the above bound. However, since we do not know ρ and $\boldsymbol{\theta}^*$ in advance, we actually cannot determine the number of iterations to convergence, which does pose a problem if we have non-linearly separable data.

Problem 5 (26 points). Pocket Algorithm for Non-Seperable data

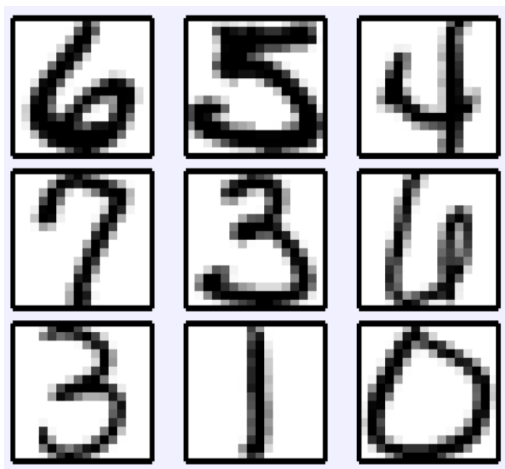
In Problem 4, we have shown that the perceptron algorithm will converge for linearly seperable data. However, for non-linearly separable data, the algorithm will never terminate. In fact, the behavior becomes quite unstable and can jump from a good perceptron to a bad one with only one update (as we will see in this problem). Hence, the quality of Er_{in} cannot be guaranteed.

One simple approach to deal with the unstability is to use the *pocket algorithm*. The algorithm preserves the parameter $\boldsymbol{\theta}$ with the best in-sample error during the update, which ensures monotonically decreasing in-sample error. The algorithm is summarized below.

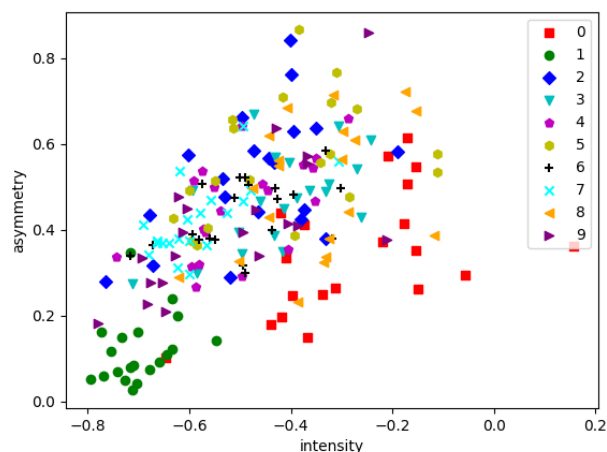
Pocket Algorithm	
1.	set the pocket weight parameter $\hat{\boldsymbol{\theta}}$ to be $\boldsymbol{\theta}_0$
2.	for $k = 0, \dots, T$, do
3.	pick a wrongly classified sample from dataset as (\mathbf{x}_k, y_k)
4.	update as $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + y_k \mathbf{x}_k$
5.	if $\text{Er}_{\text{in}}(\boldsymbol{\theta}_{k+1}) < \text{Er}_{\text{in}}(\hat{\boldsymbol{\theta}})$:
6.	$\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_{k+1}$
7.	end for
8.	return $\hat{\boldsymbol{\theta}}$

In this problem, we will implement the performance of perceptron and pocket algorithms based on the US Post Office Zip Code Data, which is a dataset of 16×16 grayscale images for digits (0-9). The training dataset has 7291 samples and the test dataset has 2007 samples.

We may extract heuristic features to help us do the classification. Let us focus on the binary classification between digits "1" and "6". In this case, the *intensity* may be a good feature: the digit "6" usually occupies larger area than digit "1". The *asymmetry* is another useful feature, as "1" tends to be more symmetric than "6". We can define the asymmetry feature as the average absolute difference between an image and its flipped version.



(a) Sample digit images.



(b) Feature plot of different classes.

Figure 4: Illustration of sample digits and corresponding feature distribution.

1. (10 points) Consider the binary classification problem for digits "1" and "6". Implement the perceptron algorithm and pocket algorithm. The suggested hyperparameters are: threshold = $1e-4$, total iteration = 2000, initial parameter $\theta_0 = \mathbf{0}$. (You may implement the algorithms based on the provided source code, where you should finish the part marked as #TODO)
2. (8 points) Plot the in-sample error and the out-of-sample error of these two algorithms versus the number of iterations (for this problem, we consider the training error as the in-sample error, and the test error as out-of-sample error). Note that for the pocket algorithm, the errors are calculated with respect to $\hat{\theta}$ (instead of θ_k). Briefly discuss the results.
3. (8 points) Plot 500 data points from class "1" and "6", respectively (similar to the Fig. 4b but with only 2 classes). On this figure, plot the final classification boundary of the two algorithms.