

Problem 1

(1) The main difference is that supervised learning uses labeled input and output data, while unsupervised learning does not.

(2) training data: used for learning to fit the parameters of the chosen model.

Test data: used to assess the performance of the trained model.

Validation data: used to tune parameters. The specific reason that why we need validation data is that in the process of adjusting parameters, the final model may have the problem of overfitting, which may lead to bad generalization on new data. So we use part of training data to approximate out of sample error, that is validation, and we test the model and further, "validate" our hyper-parameter base on validation data.

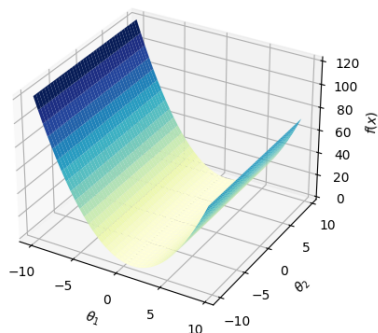
(3) not necessarily. From the concentration inequality, we know $P(|h_i - \mathbb{E}[h_i]| \geq t) \leq 2e^{-\frac{t^2}{2\sigma^2}}$ so if $\sigma^2 > 0.02$, the probability of right arm's length leaves the height by amount t is enlarged, which means it will concentration inequality for sub-Gaussian couldn't hold. If $\sigma^2 = 0.02$, the left arms and right arms are from the same distribution, it's better to use both of them.

(4). $X = V\Sigma U^T$, $V \in \mathbb{R}^{n \times n}$, $U \in \mathbb{R}^{d \times d}$, $\Sigma \in \mathbb{R}^{n \times d}$ and $V^T V = I$, $U^T U = I$
 $X^T X = U \Sigma^T V^T V \Sigma U^T = U \Sigma^T \Sigma U^T$
 \because the X is full-rank, $\Sigma = \begin{pmatrix} \lambda_1 & \lambda_2 & \dots & \lambda_n \end{pmatrix}$ and $\lambda_i \neq 0$
 \therefore so $\Sigma^T \Sigma > 0$, and $U \Sigma^T \Sigma U^T = \sum_{i=1}^d \lambda_i^2 U_i^T U_i > 0 \Rightarrow X^T X$ is positive definite

$$\begin{aligned} (5) \text{ set } f(\theta) &= \|X\theta - y\|^2 + \lambda \|\theta\|^2 \\ &= (\theta^T X^T - y^T)(X\theta - y) + \lambda \theta^T \theta \\ &= \theta^T X^T X \theta - \theta^T X^T y - y^T X \theta + y^T y + \lambda \theta^T \theta \\ \frac{\partial f(\theta)}{\partial \theta} &= 2X^T X \theta - 2X^T y + 2\lambda \theta \\ &= 0 \\ \Rightarrow (X^T X + \lambda I) \theta^* &= X^T y \end{aligned}$$

Problem 2

1. $\min_{\theta} \|\theta - 1\|^2$, the 3D-figure for $f_0 = (0, -1)^T$ is:



$$(2) \min_{\theta} \|x\theta - y\|_2^2 = \min_z \|Az - y\|_2^2$$

$$\nabla f_z = z \|Az - y\| = 0 \Rightarrow A^T A \hat{z} = A^T y$$

$\because n < d \quad \therefore \text{rank}(A^T A) = n < d \Rightarrow \hat{z}$ is not unique

$\because z := U_1^T \theta$, $U \in \mathbb{R}^{d \times d}$ are orthonormal $\therefore U^T = U^{-1}$
and $U_1^T \theta = \hat{z}$, $\hat{\theta} = U^{-T} \hat{z} \Rightarrow \hat{\theta}$ is not unique.

Problem 3

(a) $\because p(\epsilon_i) = \frac{1}{2b} e^{-|\frac{\epsilon_i - 0}{b}|}$ and $\epsilon_i \stackrel{iid}{\sim} L(0, b)$

$$\therefore p(\epsilon|\theta) = \pi p(\epsilon_i|\theta) = \prod_{i=1}^n \frac{1}{2b} e^{-|\frac{\epsilon_i}{b}|} = (\frac{1}{2b})^n e^{-\sum_{i=1}^n |\frac{\epsilon_i}{b}|}$$

and $\epsilon = y - x\theta$, so further $p(\epsilon|\theta) = (\frac{1}{2b})^n e^{-\sum_{i=1}^n \frac{|y_i - x_i^T \theta|}{b}}$

$$L(\theta|\epsilon) = -n \log 2b - \sum_{i=1}^n |y_i - x_i^T \theta|$$

$$\text{learning problem: } \hat{\theta}_b = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \theta)^2$$

$$(b). \quad H_n(x_i^T \theta - y_i) = \begin{cases} \|x_i^T \theta - y_i\|, & \|x_i^T \theta - y_i\| \geq \mu \\ \frac{\|x_i^T \theta - y_i\|^2}{2\mu} + \frac{\mu}{2}, & \|x_i^T \theta - y_i\| \leq \mu \end{cases}$$

$$\nabla H_n(x_i^T \theta - y_i) = \begin{cases} \|x_i\|, & \|x_i^T \theta - y_i\| \geq \mu \\ \frac{x_i(x_i^T \theta - y_i)}{\mu}, & \|x_i^T \theta - y_i\| \leq \mu \end{cases}$$

$$\therefore \nabla f(\theta) = \nabla \sum_i f_i(\theta) = \sum_{i=1}^n \nabla H_n(x_i^T \theta - y_i)$$

$$\begin{aligned}
 (c) \quad (1) \quad \hat{\theta}_k &= (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X \theta^* + \varepsilon_1 + \varepsilon_2) \\
 &= (X^T X)^{-1} (X^T X \theta^* + X^T \varepsilon_1 + X^T \varepsilon_2) \\
 &= \theta^* + (X^T X)^{-1} X^T (\varepsilon_1 + \varepsilon_2) \\
 \|\hat{\theta}_k - \theta^*\|_2 &= \|(X^T X)^{-1} X^T (\varepsilon_1 + \varepsilon_2)\|_2
 \end{aligned}$$

(2) The function of error with respect to iteration number is below:

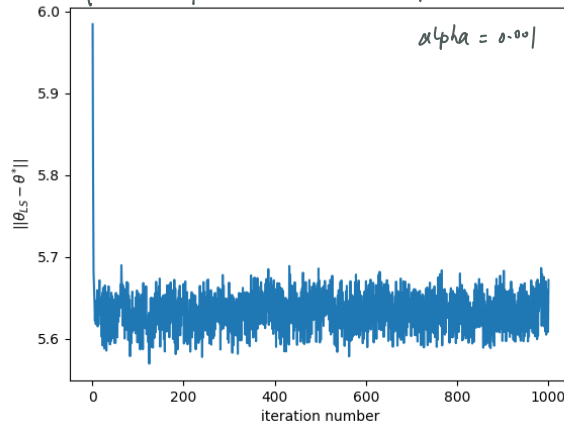
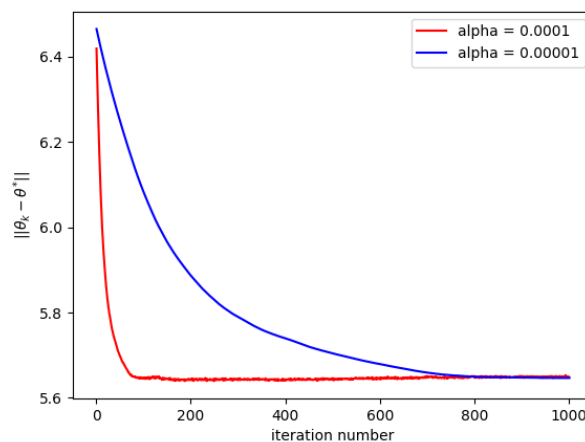


Figure 3-1

but we can see that the result doesn't converge. So I increase the iteration number to 10000, it still fails to converge. Therefore, I think it's the problem of our optimization strategy. The step size α is fixed, which maybe too large to reach the optimal value in the end.

Further I decrease the alpha to $1e-4$ and $1e-5$ respectively, but result still doesn't change, just speed up the descent as follow:



In conclusion, this has to do nothing with our optimization method. The key is that $\varepsilon_1 \in \mathbb{R}^n$ follows Gaussian distribution, and the figure 3-1 is a white Gaussian noise. We actually make it to converge, but maybe because the standard deviation of ε_1 is large, our convergence is not ideal.