

Assignment #2 (Linear Model)

Instructor: Beilun Wang

Name: Qipeng Zhu, ID: 58119304

Problem Description:

Problem 1: Linear Regression

Give data set $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)})^\top$ and $\mathbf{y} = (y^{(1)}, y^{(2)}, \dots, y^{(n)})^\top$ where $(\mathbf{x}^{(i)\top}, y^{(i)}) = (x_1^{(i)}, x_2^{(i)}, \dots, x_p^{(i)}, y^{(i)})$ is the i -th observation. We focus on the model $y = \boldsymbol{\theta}^\top \mathbf{x} + \varepsilon$.

- (1) Assuming $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, write down the log-likelihood function of \mathbf{y} . You can ignore any unnecessary constants.
- (2) Based on your answer to (1), show that finding Maximum Likelihood Estimate of $\boldsymbol{\theta}$ is equivalent to solving $\operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$.
- (3) Prove that $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$ with $\lambda > 0$ is Positive Definite (Hint: definition).
- (4) Show that $\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$ is the solution to $\operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 + \lambda \|\boldsymbol{\theta}\|^2$.
- (5) Assuming $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ and $\theta_i \sim \mathcal{N}(0, \tau^2)$ for $i = 1, 2, \dots, p$ in $\boldsymbol{\theta}$ ($\boldsymbol{\theta}$ does not vary in each sample), write down the estimate of $\boldsymbol{\theta}$ by maximizing the conditional distribution $f(\boldsymbol{\theta} | \mathbf{y})$ (Hint: Bayes' formula). You can ignore any unnecessary constants. Also find out the relationship between your estimate and the solution in (4).

Problem 2: Gradient Descent

Continuously differentiable function $f : \mathbb{R} \mapsto \mathbb{R}$ is called β -**smooth** when its derivative f' is β -**Lipschitz**, which for $\beta > 0$ implies that

$$|f'(x) - f'(y)| \leq \beta |x - y|.$$

Now suppose f is β -**smooth** and **convex** as a loss function in a gradient descent problem.

- (1) Prove that

$$f(y) - f(x) \leq f'(x)(y - x) + \frac{\beta}{2}(y - x)^2.$$

(Hint: Newton-Leibniz formula.)

- (2) Give $x_{k+1} = x_k - \eta f'(x_k)$ as one step of GD. Prove that

$$f(x_{k+1}) \leq f(x_k) - \eta(1 - \frac{\eta\beta}{2})(f'(x_k))^2.$$

- (3) Based on (2), let $\eta = 1/\beta$ and assume the unique global minimum point x^* of f exists. Prove that

$$\lim_{k \rightarrow \infty} f'(x_k) = 0, \quad \lim_{k \rightarrow \infty} x_k = x^*.$$

(Hint: show that for $K \in \mathbb{N}_+$, $\sum_{k=1}^K (f'(x_k))^2 \leq 2\beta(f(x_1) - f(x_{K+1}))$.)

(4) Recall one of the properties of convex function: $f(y) \geq f(x) + f'(x)(y - x)$. Prove that

$$f(y) - f(x) \geq f'(x)(y - x) + \frac{1}{2\beta}(f'(y) - f'(x))^2.$$

(Hint: let $z = y - \frac{1}{\beta}(f'(y) - f'(x))$.)

Problem 3: Kernel functions

Kernel function $k : \mathbb{R}^p \times \mathbb{R}^p \mapsto \mathbb{R}$ is called **Positive Semi-Definite(PSD)** when its Gramian matrix K is PSD, where $K_{ij} = k(\mathbf{u}_i, \mathbf{u}_j)$ for any group of vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n \in \mathbb{R}^p$. Let k_1 and k_2 be two PSD kernels.

- (1) Give a function $f : \mathbb{R}^p \mapsto \mathbb{R}$. Show that the kernel k defined by $k(\mathbf{u}, \mathbf{v}) = f(\mathbf{u})f(\mathbf{v})$ is PSD.
- (2) Show that the kernel k defined by $k(\mathbf{u}, \mathbf{v}) = k_1(\mathbf{u}, \mathbf{v})k_2(\mathbf{u}, \mathbf{v})$ is PSD. (Hint: consider about the Hadamard product and eigendecomposition.)
- (3) Give P as a polynomial with non-negative coefficients(e.g., $P(x) = \sum_i a_i x^i$ with $a_i \geq 0$). Show that the kernel k defined by $k(\mathbf{u}, \mathbf{v}) = P(k_1(\mathbf{u}, \mathbf{v}))$ is PSD.
- (4) Show that the kernel k defined by $k(\mathbf{u}, \mathbf{v}) = \exp(k_1(\mathbf{u}, \mathbf{v}))$ is PSD. (Hint: use the series expansion.)

Answer:**Problem 1: Linear Regression**

(1) Because $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, and $y = \theta^\top x + \varepsilon$, we can find $(y - \theta^\top x) \sim \mathcal{N}(0, \sigma^2)$. So

$$\begin{aligned} f(y|x, \theta) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y - \theta^\top x)^2}{2\sigma^2}} \\ f(y|X, \theta) &= \prod_{i=1}^n f(y_i|x_i, \theta) \\ &= \frac{1}{(\sqrt{2\pi}\sigma)^n} e^{-\frac{\sum_{i=1}^n (y_i - \theta^\top x_i)^2}{2\sigma^2}} \\ &= C e^{-\frac{(X\theta - y)^\top (X\theta - y)}{2\sigma^2}} \quad (C \text{ is a constant}) \end{aligned}$$

So the log-likelihood function of y is:

$$\begin{aligned} \ln f(y|X, \theta) &= -\frac{(X\theta - y)^\top (X\theta - y)}{2\sigma^2} \\ &= (X\theta - y)^\top (X\theta - y) \end{aligned}$$

(2) Suppose

$$L(\theta) = \ln f(y|X, \theta)$$

So

$$\nabla_{\theta} L(\theta) = 2X^\top X\theta - 2X^\top y$$

We let

$$\nabla_{\theta} L(\theta) = 0$$

then we can get

$$\theta_{ML} = (X^\top X)^{-1} X^\top y$$

So θ_{ML} is the solution of $\operatorname{argmin}_{\theta} \|y - X\theta\|^2$

(3) For any $v \neq 0$:

$$v^\top (X^\top X + \lambda I) v = v^\top X^\top X v + \lambda v^\top v = \|Xv\|_0^2 > 0$$

So $X^\top X + \lambda I$ with $\lambda > 0$ is Positive Definite.

(4) Suppose

$$J(\theta) = \|y - X\theta\|^2 + \lambda \|\theta\|^2$$

First we can simplify the optimization:

$$J(\theta) = (y - X\theta)^\top (y - X\theta) + \lambda \theta^\top \theta = \theta^\top X^\top X \theta - 2\theta^\top X^\top y + y^\top y + \lambda \theta^\top \theta$$

so its gradient is:

$$\nabla_{\theta} J(\theta) = 2X^\top X\theta - 2X^\top y + 2\lambda\theta$$

let

$$\nabla_{\theta} J(\theta) = 0$$

we can get

$$\theta^* = (X^\top X + \lambda I)^{-1} X^\top y$$

Because

$$\operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 + \lambda\|\boldsymbol{\theta}\|^2$$

is a convex optimization, $\boldsymbol{\theta}^*$ is its solution.

(5) From problem(1), we can know

$$f(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \frac{1}{(2\sqrt{2\pi}\sigma)^n} e^{-\frac{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})}{2\sigma^2}}$$

Because

$$\theta_i \sim \mathcal{N}(0, \tau^2) \text{ for } i=1, 2, \dots, p$$

we can know

$$f(\boldsymbol{\theta}) = \prod_{i=1}^p f(\theta_i) = \frac{1}{(2\sqrt{2\pi}\tau)^p} e^{-\frac{\boldsymbol{\theta}^T \boldsymbol{\theta}}{2\tau^2}}$$

Then

$$f(\boldsymbol{\theta}|\mathbf{y}) = \frac{f(\mathbf{y}|\boldsymbol{\theta})f(\boldsymbol{\theta})}{\int_{-\infty}^{+\infty} f(\mathbf{y}|\boldsymbol{\theta})f(\boldsymbol{\theta})d\mathbf{x}} = \alpha e^{-\left(\frac{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})}{2\sigma^2} + \frac{\boldsymbol{\theta}^T \boldsymbol{\theta}}{2\tau^2}\right)}$$

$$\ln f(\boldsymbol{\theta}|\mathbf{y}) = \ln \alpha - \left(\frac{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})}{2\sigma^2} + \frac{\boldsymbol{\theta}^T \boldsymbol{\theta}}{2\tau^2} \right)$$

Let

$$\nabla_{\boldsymbol{\theta}} \ln f(\boldsymbol{\theta}|\mathbf{y}) = -2(\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^T \mathbf{y} + \frac{\sigma^2}{\tau^2} \boldsymbol{\theta}) = 0$$

we can get

$$(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I}) \boldsymbol{\theta} = \mathbf{X}^T \mathbf{y}$$

Finally, we can get

$$\boldsymbol{\theta}^* = (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

When we suppose $\lambda = \frac{\sigma^2}{\tau^2}$, the solution is the same as the solution of problem(4)

Problem 2: Gradient Descent

(1) From Lagrange theorem: $\forall x, y \in \mathbb{R}, \exists \xi$ between x and y , s.t

$$f''(\xi) = \frac{f'(y) - f'(x)}{y - x}$$

So

$$|f''(\xi)| = \left| \frac{f'(y) - f'(x)}{y - x} \right| \leq \beta$$

For $f(y)$, we can use Taylor expansion at x , we can get

$$f(y) = f(x) + f'(x)(y - x) + \frac{f''(\xi)}{2}(y - x)^2 (\xi \in (x, y))$$

So

$$f(y) - f(x) \leq f'(x)(y - x) + \frac{|f''(\xi)|}{2}(y - x)^2 \leq f'(x)(y - x) + \frac{\beta}{2}(y - x)^2$$

(2) From problem(1), we can easily get:

$$f(x_{k+1}) - f(x_k) \leq f'(x_k)(x_{k+1} - x_k) + \frac{\beta}{2}(x_{k+1} - x_k)^2$$

$$= f'(x_k)(-\eta f'(x_k)) + \frac{\beta}{2}\eta^2(f'(x_k))^2 = -\eta(1 - \frac{\eta\beta}{2})(f'(x_k))^2$$

(3) From problem to we can get,

$$\begin{aligned} \exists n \in N_+, s.t. \quad & f(x_{n+1}) \leq f(x_n) - \frac{1}{2\beta}(f'(x_n))^2 \\ \left\{ \begin{array}{ll} f(x_{n+1}) \leq f(x_n) - \frac{1}{2\beta}(f'(x_n))^2 & (n) \\ f(x_n) \leq f(x_{n-1}) - \frac{1}{2\beta}(f'(x_{n-1}))^2 & (n-1) \\ \vdots & \\ f(x_2) \leq f(x_1) - \frac{1}{2\beta}(f'(x_1))^2 & (1) \end{array} \right. \end{aligned}$$

Add up these n formulas, we can get:

$$\sum_{k=1}^n (f'(x_k))^2 \leq 2\beta(f(x_1) - f(x_{n+1}))$$

because $f(x^*)$ is the minimum, so $f(x_{n+1}) \geq f(x^*)$, then

$$\sum_{k=1}^n (f'(x_k))^2 \leq 2\beta(f(x_1) - f(x^*)) = c \text{ (c is a constant)}$$

It's easy to know $\sum_{k=1}^n (f'(x_k))^2$ is monotonic increasing, and bounded, which means it can converge. So

$$\lim_{k \rightarrow \infty} f'(x_k) = 0$$

And $\lim_{k \rightarrow \infty} x_k$ is solution of $f'(x)=0$,

we can get

$$\lim_{k \rightarrow \infty} x_k$$

is extreme point. Because $f(x)$ is convex, so its local minimum is global minimum. So $f(x)$

$$\lim_{k \rightarrow \infty} x_k = x^*$$

.

(4) Suppose:

$$z = y - \frac{1}{\beta}(f'(y) - f'(x))$$

Because $f(x)$ is convex:

$$f(x) - f(z) \leq f'(x)(x - z) = f'(x)(x - y) + f'(x)(y - z) \quad (1)$$

From proble(1), we can know:

$$f(z) - f(y) \leq f'(y)(z - y) + \frac{\beta}{2}(z - y)^2 = -f'(y)(y - z) + \frac{\beta}{2}(y - z)^2 \quad (2)$$

Let (1) + (2), we can get:

$$f(x) - f(y) \leq f'(x)(x - y) + (f'(x) - f'(y))(y - z) + \frac{\beta}{2}(y - z)^2$$

Substituting $y - z = \frac{1}{\beta}(f'(y) - f'(x))$

$$f(x) - f(y) \leq f'(x)(x - y) + \frac{1}{\beta}(f'(x) - f'(y))(f'(y) - f'(x)) - \frac{1}{2\beta}(f'(y) - f'(x))^2$$

So

$$f(y) - f(x) \leq f'(x)(y - x) + \frac{1}{2\beta}(f'(y) - f'(x))^2$$

Problem 3: Kernel functions

(1) Suppose \mathbf{K} is the Gramian matrix of $k(\mathbf{u}, \mathbf{v})$, We can easily know:

$$\begin{aligned} \mathbf{K} &= \begin{pmatrix} f(\mathbf{u}_1)f(\mathbf{u}_1) & f(\mathbf{u}_1)f(\mathbf{u}_2) & \dots \\ f(\mathbf{u}_2)f(\mathbf{u}_1) & f(\mathbf{u}_2)f(\mathbf{u}_2) & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \\ &= (f(\mathbf{u}_1), f(\mathbf{u}_2), \dots, f(\mathbf{u}_n))^T \cdot (f(\mathbf{u}_1), f(\mathbf{u}_2), \dots, f(\mathbf{u}_n)) \end{aligned}$$

For any given $\mathbf{v} \in \mathbb{R}^p$, we have:

$$\begin{aligned} \mathbf{v}\mathbf{K}\mathbf{v}^T &= \mathbf{v}(f(\mathbf{u}_1), f(\mathbf{u}_2), \dots, f(\mathbf{u}_n))^T \cdot (f(\mathbf{u}_1), f(\mathbf{u}_2), \dots, f(\mathbf{u}_n))\mathbf{v}^T \\ &= \|\mathbf{v} \cdot (f(\mathbf{u}_1), f(\mathbf{u}_2), \dots, f(\mathbf{u}_n))^T\|_2^2 \geq 0 \end{aligned}$$

So $k(\mathbf{u}, \mathbf{v})$ is PSD kernels.

(2) Suppose \mathbf{A}, \mathbf{B} and \mathbf{C} is respectively the Gramian matrix of k_1, k_2, k , because:

$$k(\mathbf{u}, \mathbf{v}) = k_1(\mathbf{u})k_2(\mathbf{v})$$

We can easily know $c_{ij} = a_{ij}b_{ij}$, so we can get:

$$\mathbf{C} = \mathbf{A} \circ \mathbf{B}$$

According the the nature of Hadamard product:

$$\lambda_{\min}(\mathbf{C}) = \lambda_{\min}(\mathbf{A} \circ \mathbf{B}) \geq \lambda_{\min}(\mathbf{A})\lambda_{\min}(\mathbf{B})$$

Because \mathbf{A} and \mathbf{B} are PSD, so $\lambda_{\min}(\mathbf{A}) \geq 0, \lambda_{\min}(\mathbf{B}) \geq 0$, so we can get:

$$\lambda_{\min}(\mathbf{C}) \geq 0$$

So \mathbf{C} is PSD, and $k(\mathbf{u}, \mathbf{v})$ is PSD.

(3) Suppose $m_i(\mathbf{u}, \mathbf{v}) = k_1^i(\mathbf{u}, \mathbf{v})$:

$$\begin{cases} m_1(\mathbf{u}, \mathbf{v}) = k_1(\mathbf{u}, \mathbf{v}) \\ m_2(\mathbf{u}, \mathbf{v}) = k_1^2(\mathbf{u}, \mathbf{v}) \\ \vdots \\ m_n(\mathbf{u}, \mathbf{v}) = k_1^n(\mathbf{u}, \mathbf{v}) \end{cases}$$

When $i = 1$, we can easily get $m_1(\mathbf{u}, \mathbf{v})$ is PSD.

Suppose $m_n(\mathbf{u}, \mathbf{v})$ is PSD;

When $i = n+1$; For $m_{n+1}(\mathbf{u}, \mathbf{v})$, we have:

$$m_{n+1}(\mathbf{u}, \mathbf{v}) = m_n(\mathbf{u}, \mathbf{v}) * m_1(\mathbf{u}, \mathbf{v})$$

From problem(2) and hypothesis, we can know $m_{n+1}(\mathbf{u}, \mathbf{v})$ is PSD. So for any $i \in \mathbb{N}$, we have $m_i(\mathbf{u}, \mathbf{v})$ is PSD. Suppose \mathbf{M}_i is the Gramian matrix of m_i , so \mathbf{M}_i is PSD, for any $i \in \mathbb{N}$. Then we have:

$$\lambda_{\min}(\mathbf{M}_i) \geq 0, \text{ for any } i \in \mathbb{N}$$

Because

$$k(\mathbf{u}, \mathbf{v}) = P(k_1(\mathbf{u}, \mathbf{v})) = \sum_i a_i m_i(\mathbf{u}, \mathbf{v})$$

Suppose \mathbf{K} is the Gramian matrix of $k(\mathbf{u}, \mathbf{v})$, so

$$\mathbf{K} = \sum_i a_i \mathbf{M}_i \implies \lambda(\mathbf{K}) = \sum_i a_i \lambda(\mathbf{M}_i) \geq 0$$

So \mathbf{K} is PSD, and $k(\mathbf{u}, \mathbf{v})$ is PSD.

(4) From Maclaurin series expansion,

$$\exp(k_1(\mathbf{u}, \mathbf{v})) = \sum_{i=0}^{\infty} \frac{f^{(i)}(x^*)}{i!} x^i = \sum_{i=0}^{\infty} \frac{\exp(x^*)}{i!} x^i \quad (x^* = k_1(\mathbf{u}, \mathbf{v}))$$

It's clear that $\exp(x^*) \geq 0$ for any $i \in \mathbb{N}$. So from problem(3), we can get $k(\mathbf{u}, \mathbf{v})$ is PSD.