

CS189 HW3

(a)

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \log P(x_i | \theta) + n \log P(\theta)$$

$$= \arg \min_{\theta} \|X\theta - y\|_2^2 + \frac{\sigma^2}{\sigma_h^2} \|\theta - \mu_0\|_2^2$$

$$\begin{aligned} \nabla_{\theta} \left(\|X\theta - y\|_2^2 + \frac{\sigma^2}{\sigma_h^2} \|\theta - \mu_0\|_2^2 \right) \\ = \nabla_{\theta} \left((X\theta - y)^T (X\theta - y) + \frac{\sigma^2}{\sigma_h^2} (\theta - \mu_0)^T (\theta - \mu_0) \right) \\ = \nabla_{\theta} \left((X\theta)^T X\theta - 2y^T X\theta + y^T y + \frac{\sigma^2}{\sigma_h^2} (\theta^T \theta - \theta^T \mu_0 - \mu_0^T \theta + \mu_0^T \mu_0) \right) \end{aligned}$$

$$= 2X^T X\theta - 2X^T y + \frac{2\sigma^2}{\sigma_h^2} (\theta - \mu_0) = 0$$

$$X^T X\theta^* - X^T y + \frac{\sigma^2}{\sigma_h^2} (\theta^* - \mu_0) = 0$$

$$\theta^* = \left(X^T X + \frac{\sigma^2}{\sigma_h^2} I \right)^{-1} (X^T y + \frac{\sigma^2}{\sigma_h^2} \mu_0)$$

(b) The two figures for $\sigma_h = 1$ and $\sigma_h = 10$ are included in Appendix

The code is in `problem1-starter.py` from # Generating Data and Labels with Random Gaussian Noise to the end. After changing σ_h to 10, the range and variation of prior obviously increases. This is because σ_h is the standard deviation of prior distribution. σ_h increases means the variance of prior increases.

Code also includes the implementation of MLE, MAP and prior.

(c) The two figures for $\mu_0 = \begin{bmatrix} -5 \\ -5 \end{bmatrix}$ and $\begin{bmatrix} 3 \\ 3 \end{bmatrix}$ are included in Appendix. The code is in `problem1-starter.py`.

When change to $\begin{bmatrix} -5 \\ -5 \end{bmatrix}$, the center of prior plot moves from (0, 0) to (-5, -5). When change to $\begin{bmatrix} 3 \\ 3 \end{bmatrix}$, the center of prior plot moves from (-5, -5) to (3, 3). θ_{MLE} doesn't change according to the plots with respect to prior.

$$(d) \theta_{MAP} = \arg \max_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \log P(x_i | \theta) + \log P(\theta)$$

$$= \arg \max_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - x_i \theta)^2}{2\sigma^2}} \right) + \log \left(\frac{1}{2\sigma_h} e^{-\frac{|\theta - \mu_0|}{\sigma_h}} \right)$$

$$= \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \frac{(y_i - x_i \theta)^2}{2\sigma^2} + \frac{|\theta - \mu_0|}{\sigma_h}$$

$$= \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - x_i \theta)^2 + \frac{2\sigma^2}{\sigma_h} |\theta - \mu_0|$$

$$= \arg \min_{\theta \in \mathbb{R}^d} \|y - X\theta\|_2^2 + \frac{2\sigma^2}{\sigma_h} |\theta - \mu_0|$$

(e) The two plots for $p(\theta_i) \sim L(0, 1)$ and $p(\theta_i) \sim L(0, 0.0001)$
 The code is in problem1-starter.py. Which implement of Laplacian-prior and MAP-contour.

After change σ_h from 1 to 0.0001 the weight of prior in θ_{MAP} increases a lot.

This is because: According to the result in (d)

$$\theta_{MAP} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|y - X\theta\|^2 + 2 \underbrace{\frac{\sigma^2}{\sigma_h^2}}_{\text{weight of prior}} \|\theta - \mu\|$$

This is the prior term in θ_{MAP} . Its weight is $\frac{\sigma^2}{\sigma_h^2}$. If I change σ_h from 1 to 0.0001, the weight changes to 10^8 times before.

(f) If prior of θ is $\text{Unit}(0, 1)$ then $p(\theta) = 1$

$$\because y = X\theta + \varepsilon \quad \varepsilon \sim N(0, 1) \quad \therefore y \sim N(X\theta, 1)$$

$$\therefore \theta_{MAP} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|y - X\theta\|^2 + 2\sigma^2 \times \log p(\theta)$$

$$= \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|y - X\theta\|^2 + 2\sigma^2(\log 1)$$

$$= \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|y - X\theta\|^2 = \theta_{MLE}$$

So, now θ_{MAP} and θ_{MLE} are equivalent.

2. Probabilistic Model of Linear regression

(a) $\therefore Y = XW_1 + W_0 + Z, \quad Z \sim \mathcal{N}(0, 1)$

$\therefore Y|X=x \sim \mathcal{N}(XW_1 + W_0, 1)$

$$P(Y|X) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(Y - (XW_1 + W_0))^2}{2}}$$

(b) log likelihood function

$$\operatorname{argmax}_W \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi}} e^{-\frac{(Y_i - (X_i W_1 + W_0))^2}{2}}\right)$$

$$= \operatorname{argmin}_W \sum_{i=1}^n (Y_i - (X_i W_1 + W_0))^2$$

$$= \operatorname{argmin}_W \|Y - XW_1 - W_0 \times \mathbf{1}\|^2 \quad \mathbf{1} \in \mathbb{R}^n \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$= \operatorname{argmin}_W (Y - XW_1 - W_0 \mathbf{1})^T (Y - XW_1 - W_0 \mathbf{1})$$

$$= \operatorname{argmin}_W (Y^T - W_1 X^T - W_0 \mathbf{1}^T) (Y - XW_1 - W_0 \mathbf{1})$$

$$= \operatorname{argmin}_W (Y^T Y - Y^T X W_1 - Y^T W_0 \mathbf{1} - W_1 X^T Y + W_1 X^T X W_1 + W_0 X^T W_0 \mathbf{1} - W_0 \mathbf{1}^T Y + W_0 \mathbf{1}^T X W_1 + W_0 \mathbf{1}^T W_0 \mathbf{1})$$

$$= \operatorname{argmin}_W \sum_{i=1}^n Y_i^2 - 2 Y^T X W_1 - 2 W_0 \sum_{i=1}^n Y_i + W_1^2 X^T X + 2 W_0 W_1 \sum_{i=1}^n X_i + n W_0^2$$

$$\begin{cases} \nabla W_1 = -2 X^T Y + 2 X^T X W_1 + 2 W_0 \sum_{i=1}^n X_i \end{cases}$$

$$\nabla W_0 = -2 \sum_{i=1}^n Y_i + 2 W_1 \sum_{i=1}^n X_i + 2 n W_0$$

$$\begin{cases} -2 X^T Y + 2 X^T X W_1 + 2 \sum_{i=1}^n X_i W_0 = 0 \\ -2 \sum_{i=1}^n Y_i + 2 W_1 \sum_{i=1}^n X_i + 2 n W_0 = 0 \end{cases}$$

$$\begin{cases} W_1 = \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i) + n \sum_{i=1}^n X_i Y_i}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \\ W_0 = \frac{(\sum_{i=1}^n Y_i)(\sum_{i=1}^n X_i^2) + (\sum_{i=1}^n X_i Y_i)(\sum_{i=1}^n X_i)}{(\sum_{i=1}^n X_i)^2 - n \sum_{i=1}^n X_i^2} \end{cases}$$

(c) $\therefore Z \sim \mathcal{U}[-0.5, 0.5] \quad Y = XW + Z$

$\therefore Y|X \sim \mathcal{U}[-0.5 + XW, 0.5 + XW]$

(d) $Y_i = X_i W + Z_i$ log likelihood function

$$\log P(Y=Y_1 | X=X_1) \times \dots \times P(Y=Y_n | X=X_n)$$

$$Y_i = X_i W + Z_i = \sum_{i=1}^n \log P(Y=Y_i | X=X_i)$$

\therefore when $-0.5 + X_i W \leq Y_i \leq 0.5 + X_i W$, $P(Y=Y_i | X=X_i) = 1$ otherwise it is 0

To maximize log likelihood function, I need to make $-0.5 + X_i W \leq Y_i \leq 0.5 + X_i W$ hold for as many (X_i, Y_i) as possible

$$\therefore -0.5 + X_i W \leq Y_i \leq 0.5 + X_i W \Rightarrow \frac{Y_i - 0.5}{X_i} \leq W \leq \frac{Y_i + 0.5}{X_i}$$

\therefore I think one appropriate estimate for W_{MLE}

$$\text{is } \frac{1}{n} \sum_{i=1}^n \frac{Y_i}{X_i}$$

(e) $Y_i = X_i W + Z_i$ $Z_i \sim N(0, 1)$

$$Y_i \sim N(X_i W, 1) \quad W \sim N(0, \sigma^2)$$

$$\log P(W|Y, X) \propto \log P(Y|W, X) + \log P(W)$$

$$= \sum_{i=1}^n \log P(Y_i | X_i, W) + \log P(W)$$

$$\propto - \sum_{i=1}^n \frac{(Y_i - X_i W)^2}{2} - \frac{W^2}{2\sigma^2}$$

$$\propto - \frac{\sigma^2 \sum_{i=1}^n (Y_i - X_i W)^2 + W^2}{2\sigma^2}$$

$$\propto \frac{(\sum_{i=1}^n X_i^2 \sigma^2 + 1) W^2 - \sum_{i=1}^n 2\sigma^2 X_i Y_i W + \sum_{i=1}^n \sigma^2 Y_i^2}{2\sigma^2}$$

$$\propto \frac{\left(W - \frac{\sum_{i=1}^n \sigma^2 X_i Y_i}{\sum_{i=1}^n \sigma^2 X_i^2 + 1} \right)^2}{2\sigma^2}$$

$\therefore P(W|Y, X) \sim N(\mu', \sigma'^2)$ Posterior mean

where $\mu' = \frac{\sum_{i=1}^n \sigma^2 X_i Y_i}{\sum_{i=1}^n \sigma^2 X_i^2 + 1}$ $\sigma'^2 = \frac{\sigma^2}{\sum_{i=1}^n \sigma^2 X_i^2 + 1}$ is $\mu' = \frac{\sum_{i=1}^n \sigma^2 X_i Y_i}{\sum_{i=1}^n \sigma^2 X_i^2 + 1}$

(f) $Y_i = W^T X_i + Z_i$ $Z_i \sim N(0, 1)$ $\therefore Y_i \sim N(W^T X_i, 1)$

log likelihood function

$$\sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi}} e^{-\frac{(y_i - w^T x_i)^2}{2}} \right]$$

$$W_{MLE} = \arg \max_w \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi}} e^{-\frac{(y_i - w^T x_i)^2}{2}} \right]$$

$$= \arg \min_w \sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2} + \frac{1}{2} \log 2\pi$$

has no w

$$= \arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$= \arg \min_w \|Y - W^T X\|^2 = W_{LS}$$

\therefore In this setting, the maximum likelihood estimator for w is the solution to a least square problem

(9) $y_i = w^T x_i + z_i$ $z_i \sim N(0, 1)$ $w_i \sim N(0, \sigma^2)$

$$p(w|x, Y) \propto P(x|w) P(w)$$

$$p(w_i|x, Y) \propto \left(e^{-\frac{\sum_{i=1}^n (y_i - w^T x_i)^2}{2}} - \frac{w_i^2}{2\sigma^2} \right)$$

I didn't manage to solve this problem.

3. Simple Bias-Variance Tradeoff

(a) ① $E\left(\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right) = \frac{1}{n} \times n\mu - \mu = 0$
 ② $E\left(\frac{x_1 + x_2 + \dots + x_{n+1}}{n+1} - \mu\right) = \frac{1}{n+1} \times n\mu - \mu = -\frac{\mu}{n+1}$
 ③ $E\left(\frac{x_1 + x_2 + \dots + x_{n+n_0}}{n+n_0} - \mu\right) = \frac{1}{n+n_0} \times n\mu - \mu = -\frac{n_0\mu}{n+n_0}$
 ④ $E(0 - \mu) = -\mu$

(b) X_1, X_2, \dots, X_n are iid

① $\text{Var}(\hat{X}) = \frac{1}{n^2} [\text{Var}(X_1) + \dots + \text{Var}(X_n)] = \frac{1}{n^2} \times n\sigma^2 = \frac{\sigma^2}{n}$
 ② $\text{Var}(\hat{X}) = \frac{1}{(n+1)^2} \times n\sigma^2 = \frac{n\sigma^2}{(n+1)^2}$
 ③ $\text{Var}(\hat{X}) = \frac{1}{(n+n_0)^2} \times n\sigma^2 = \frac{n\sigma^2}{(n+n_0)^2}$
 ④ $\text{Var}(\hat{X}) = 0$

(c) $E[(\hat{X} - X')^2] = \text{Var}(\hat{X} - X') + [E(\hat{X} - X')]^2$
 $= \text{Var}(\hat{X}) + \text{Var}(X') + \{E[(\hat{X} - \mu) + (\mu - X')]\}^2$
 $= \text{Var}(\hat{X}) + \sigma^2 + [E(\hat{X} - \mu) + E(\mu - X')]^2$
 $= \text{Var}(\hat{X}) + \sigma^2 + [E(\hat{X} - \mu)]^2$
 $= \text{Var}(\hat{X}) + \sigma^2 + [\text{bias}(\hat{X})]^2$

$E[(\hat{X} - \mu)^2] = \text{Var}(\hat{X} - \mu) + [E(\hat{X} - \mu)]^2$
 $= \text{Var}(\hat{X}) + [\text{bias}(\hat{X})]^2$

Compare them, there is a σ^2 in $E[(\hat{X} - X')^2]$ but not in $E[(\hat{X} - \mu)^2]$

\hat{X} is calculated based on training samples

so \hat{X} will perform worse on fresh sample X' , therefore the error will be larger.

(d) ① $= \text{Var}(\hat{X}) + \text{bias}^2(\hat{X}) = \frac{\sigma^2}{n} + 0^2 = \frac{\sigma^2}{n}$
 ② $= \frac{n}{(n+1)^2} \sigma^2 + \left(\frac{\mu}{n+1}\right)^2 = \frac{n\sigma^2 + \mu^2}{(n+1)^2}$
 ③ $= \frac{n\sigma^2}{(n+n_0)^2} + \frac{n_0^2\mu^2}{(n+n_0)^2} = \frac{n\sigma^2 + n_0^2\mu^2}{(n+n_0)^2}$

④ $= 0 + \mu^2 = \mu^2$

(e) ① is when $n_0 = 0$ ② is when $n_0 = 1$ ④ is when $n_0 = \infty$

(f) As n_0 increases, |bias| increases (since bias ≤ 0 here, actually bias is decreasing), variance decreases.

(g) $n_0 = \alpha n$, $\text{Var}(\hat{X}) = \frac{n\sigma^2}{(n+\alpha n)^2}$ bias $(\hat{X}) = -\frac{\alpha n \mu}{n + \alpha n}$

$$\text{expected total error} = \frac{n\sigma^2}{(n+\alpha n)^2} + \frac{(\alpha n \mu)^2}{(n+\alpha n)^2}$$

$$= \frac{n\sigma^2 + \alpha^2 n^2 \mu^2}{(n+\alpha n)^2}$$

$$\frac{\partial}{\partial \alpha} \text{expected total error} = \frac{n^2 \mu^2 \times 2\alpha (n+\alpha n)^2 - (n\sigma^2 + \alpha^2 n^2 \mu^2) \times 2(n+\alpha n) \times n}{(n+\alpha n)^4}$$

$$= 0$$

$$2n^2 \mu^2 \alpha (n+\alpha n)^2 = 2n(n+\alpha n)(n\sigma^2 + \alpha^2 n^2 \mu^2)$$

$$n\mu^2 (n+\alpha n) \alpha = n\sigma^2 + \alpha^2 n^2 \mu^2$$

$$n^2 \mu^2 \alpha + n^2 \mu^2 \alpha^2 = n^2 \mu^2 \alpha + n\sigma^2$$

$$n^2 \mu^2 \alpha^2 = n\sigma^2$$

$$\alpha = \frac{n\sigma^2}{n^2 \mu^2} = \frac{\sigma^2}{n\mu^2}$$

(h) When σ is large and μ is small (close to 0)

$$\alpha = \frac{\sigma^2}{n\mu^2} \rightarrow \infty \text{ (very large)}$$

$$(i) X' = X - \mu_0$$

$$E(X') = E(X - \mu_0) = \mu - \mu_0$$

$$\text{Var}(X') = \text{Var}(X - \mu_0) = \text{Var}(X) = \sigma^2$$

(j) In ridge regression, as λ increases, model bias \uparrow , model variance \downarrow .

This is very similar to α 's influence of $\hat{\beta}$'s bias and variance.

So in ridge regression, we can use cross validation to select the λ value with the smallest validation error.

4. Robotic Learning of Controls from Demonstrations and Images.

(a) The 0th, 10th and 20th images in the crashing set are plotted in the appendix.

Their corresponding control vectors are shown in appendix.

The code is in robotic_ridge_code_starter.py under section # 4(a).

(b) The code to do this is under section # 4(b).

When I attempt to do this, it will raise singular matrix error.

This is because $\det(X^T X) = 0$, $X^T X$ is not invertible.

(c) For each λ in $[0.1, 1.0, 10.0, 100.0, 1000.0]$ the result is shown in appendix.

The code is under section # 4(c).

(d) The result is shown in appendix.

The code is under section # 4(d).

(e) The result is shown in appendix.

The code is under section # 4(e).

By increasing λ value, the bias will increase because $\lambda \uparrow$ means the model complexity will decrease, therefore bias will increase.

By increasing λ value, the variance will decrease because $\lambda \uparrow$, model complexity \downarrow , therefore variance \downarrow .

(f) The result is shown in appendix.

The code is under section # 4(f).