**Midterm exam2, High-dimensional Data Analysis, 2018 Spring [+32 points]**

**Not only answer but also calculation**      Name: Dong, Yi-Shian

+6

**Q1 [+6]** Consider a model

$$Y_i = \beta_0 + \sum_{j=1}^{p} \beta_j x_{ij} + \varepsilon_i, \quad \varepsilon_i \overset{iid}{\sim} N(0, \sigma^2), \quad i = 1, \ldots, n.$$

+3

(1) [+3] Derive the Lasso estimator ( $\hat{\boldsymbol{\beta}}_\lambda, \lambda = ?$ ) as a <u>posterior mode</u>.

(define a prior density and derive a posterior density)

Let $\beta_j \sim$ double exponential $(b)$ $\forall j = 1, 2, \ldots, p \Rightarrow f(\beta_j) = \frac{1}{2b} \exp\left(\frac{|\beta_j|}{b}\right)$ $\forall j = 1, 2, \ldots, p$

And because of $Y_i | \beta \sim N(\beta_0 + \sum_{j=1}^{p} \beta_j x_{ij}, \sigma^2)$ $\forall i = 1, 2, \ldots, n$.

$$\Rightarrow f(\beta | Y) \propto f(Y | \beta) f(\beta)$$

$$\propto \left[\prod_{i=1}^{n} \exp\left(\frac{-1}{2\sigma^2}(Y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij})^2\right)\right] \times \left[\prod_{j=1}^{p} \exp\left(\frac{-|\beta_j|}{b}\right)\right]$$

$$= \exp\left[\frac{-1}{2\sigma^2}(y - X\beta)^T(y - X\beta) - \frac{1}{b}\sum_{j=1}^{p}|\beta_j|\right] \quad \left(\text{where } y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}, X = \begin{bmatrix} 1 & x_{11} \cdots x_{1p} \\ \vdots \\ 1 & x_{n1} \cdots x_{np} \end{bmatrix}\right)$$

$$= \exp\left[\frac{-1}{2\sigma^2}\left((y - X\beta)^T(y - X\beta) + \frac{2\sigma^2}{b}\sum_{j=1}^{p}|\beta_j|\right)\right]$$

$$\Rightarrow \hat{\beta}_\lambda = \arg\max_\beta f(\beta | Y) = \arg\min_\beta\left[(y - X\beta)^T(y - X\beta) + \frac{2\sigma^2}{b}\sum_{j=1}^{p}|\beta_j|\right]$$

+3

(2) [+3] Derive the ridge estimator ( $\hat{\boldsymbol{\beta}}_\lambda, \lambda = ?$ ) as a <u>posterior mode</u>.   $\Rightarrow$ When $\lambda = \frac{2\sigma^2}{b}$, $\hat{\beta}_\lambda$ are equivalent to the Lasso estimator

(define a prior density and derive a posterior density)

Let $\beta_j \sim N(0, c)$ $\forall j = 1, 2, \ldots, p$

$\because Y_i | \beta \sim N(\beta_0 + \sum_{j=1}^{p} \beta_j x_{ij}, \sigma^2)$

$\therefore f(\beta | Y) \propto f(Y | \beta) f(\beta) \propto \left[\prod_{i=1}^{n} \exp\left(-\frac{(Y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij})^2}{2\sigma^2}\right)\right] \times \left[\prod_{j=1}^{p} \exp\left(-\frac{\beta_j^2}{2c}\right)\right]$

$$= \exp\left(-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(Y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij})^2 - \frac{1}{2c}\sum_{j=1}^{p}\beta_j^2\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2}(y - X\beta)^T(y - X\beta) - \frac{1}{2c}\sum_{j=1}^{p}\beta_j^2\right) \quad \left(\text{where } y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, X = \begin{bmatrix} 1 & x_{11} \cdots x_{1p} \\ \vdots \\ 1 & x_{n1} \cdots x_{np} \end{bmatrix}\right)$$

$$= \exp\left(\frac{-1}{2\sigma^2}\left[(y - X\beta)^T(y - X\beta) + \frac{\sigma^2}{c}\sum_{j=1}^{p}\beta_j^2\right]\right)$$

$$\Rightarrow \hat{\beta}_\lambda = \arg\max_\beta f(\beta | Y) = \arg\max_\beta \exp\left(\frac{-1}{2\sigma^2}\left[(y - X\beta)^T(y - X\beta) + \frac{\sigma^2}{c}\sum_{j=1}^{p}\beta_j^2\right]\right)$$

$$= \arg\min_\beta\left[(y - X\beta)^T(y - X\beta) + \frac{\sigma^2}{c}\sum_{j=1}^{p}\beta_j^2\right]$$

$\Rightarrow$ When $\lambda = \frac{\sigma^2}{c}$ ( or $c = \frac{\sigma^2}{\lambda}$ ), this $\hat{\beta}_\lambda$ will

equivalent to the ridge estimator.

Moreover, $\hat{\beta}_\lambda$ can be represented as $(X^T X + \lambda I)^{-1} X^T y$.

**+|2** **Q2 [+12]** Consider a model <u>without an intercept</u>:

$$Y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i, \quad \varepsilon_i \overset{iid}{\sim} N(0, \sigma^2), \quad i = 1, \ldots, n. \text{ Assume } \sum_{i=1}^{n} x_{i1} = \sum_{i=1}^{n} x_{i2} = 0 \text{ and } \sum_{i=1}^{n} x_{i1}^2 = \sum_{i=1}^{n} x_{i2}^2 = 1.$$

Answer the questions by using $r_{12} = \sum_{i=1}^{n} x_{i1} x_{i2}$, $r_{1y} = \sum_{i=1}^{n} x_{i1} Y_i$, and $r_{2y} = \sum_{i=1}^{n} x_{i2} Y_i$.

$$X = \begin{bmatrix} x_{11} & x_{12} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} \rightarrow X^T X = \begin{bmatrix} x_{11} \cdots x_{n1} \\ x_{12} \cdots x_{n2} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} = \begin{bmatrix} \sum x_{i1}^2 & \sum x_{i1} x_{i2} \\ \sum x_{i1} x_{i2} & \sum x_{i2}^2 \end{bmatrix}$$

**+|** **(1) [+1]** $X^T X = \begin{bmatrix} 1 & r_{12} \\ & \\ r_{12} & 1 \end{bmatrix}_{2\times 2}$

$$\Rightarrow (X^T X)^{-1} = \frac{1}{1 - r_{12}^2} \begin{bmatrix} 1 & -r_{12} \\ -r_{12} & 1 \end{bmatrix}$$

**+2** **(2) [+2]** Derive the LSE

$$\hat{\beta}_1 = \frac{r_{1y} - r_{12} r_{2y}}{1 - r_{12}^2}$$

$$\Rightarrow \text{LSE:} \quad (X^T X)^{-1} X^T y = \frac{1}{1 - r_{12}^2} \begin{bmatrix} 1 & -r_{12} \\ -r_{12} & 1 \end{bmatrix} \begin{bmatrix} r_{1y} \\ r_{2y} \end{bmatrix}$$

$$\hat{\beta}_2 = \frac{r_{2y} - r_{12} r_{1y}}{1 - r_{12}^2}$$

$$= \frac{1}{1 - r_{12}^2} \begin{bmatrix} r_{1y} - r_{12} r_{2y} \\ r_{2y} - r_{12} r_{1y} \end{bmatrix}$$

---

Below, we assume $r_{12} = r_{1y} = r_{2y} = 1/2$. $\Rightarrow X^T X = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$ and LSE: $\hat{\beta} = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$.

**12** **(3) [+2]** Derive the ridge estimators $\hat{\beta}_{1\lambda}$ and $\hat{\beta}_{2\lambda}$
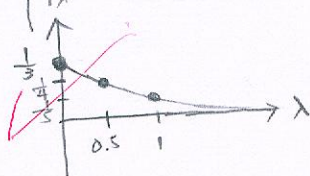
$$\hat{\beta}_\lambda = (X^T X + \lambda I)^{-1} X^T y = \begin{bmatrix} 1+\lambda & r_{12} \\ r_{12} & 1+\lambda \end{bmatrix}^{-1} \begin{bmatrix} r_{1y} \\ r_{2y} \end{bmatrix} \overset{r_{12} = r_{1y} = r_{2y} = \frac{1}{2}}{=} \begin{bmatrix} 1+\lambda & 0.5 \\ 0.5 & 1+\lambda \end{bmatrix}^{-1} \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

$$= \frac{4}{(2\lambda+1)(2\lambda+3)} \begin{bmatrix} 1+\lambda & -0.5 \\ -0.5 & 1+\lambda \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \frac{1}{2} = \frac{2}{(2\lambda+1)(2\lambda+3)} \begin{bmatrix} \lambda + 0.5 \\ \lambda + 0.5 \end{bmatrix} = \frac{1}{2\lambda+3} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\Rightarrow \hat{\beta}_{1\lambda} = \hat{\beta}_{2\lambda} = \frac{1}{2\lambda+3} \text{ when } r_{12} = r_{1y} = r_{2y} = \frac{1}{2}$$

**+2** **(4) [+2]** Draw the ridge trace for $\hat{\beta}_{1\lambda}$

| $\lambda$ | 0 | 0.5 | 1 | $\cdots$ | $\rightarrow \infty$ |
|---|---|---|---|---|---|
| $\hat{\beta}_{1\lambda}$ | $\frac{1}{3}$ | $\frac{1}{4}$ | $\frac{1}{5}$ | | 0 |

$\hat{\beta}_{LSE}$



this curve is decreasing and bounded at 0.

**+2** **(5) [+2]** Derive the degrees of freedom $df_\lambda$

$$df_\lambda = tr[X(X^T X + \lambda I)^{-1} X^T] = tr[X^T X (X^T X + \lambda I)^{-1}] = tr\left[ \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \cdot \frac{4}{(2\lambda+1)(2\lambda+3)} \begin{bmatrix} 1+\lambda & -0.5 \\ -0.5 & 1+\lambda \end{bmatrix} \right]$$

$$= \frac{4}{(2\lambda+1)(2\lambda+3)} tr\begin{bmatrix} \lambda+0.75 & 0.5\lambda \\ 0.5\lambda & \lambda+0.75 \end{bmatrix} = \frac{4}{(2\lambda+1)(2\lambda+3)} \times (2\lambda+1.5) = \frac{8\lambda+6}{(2\lambda+1)(2\lambda+3)}$$

**+1** **(6) [+1]** Derive the degrees of freedom $df_\lambda$ when $\lambda = 0$

$$df_{\lambda=0} = \left. \frac{8\lambda+6}{(2\lambda+1)(2\lambda+3)} \right|_{\lambda=0} = \frac{6}{1 \times 3} \checkmark = 2$$

**+1** **(7) [+1]** Derive the degrees of freedom $df_\lambda$ when $\lambda \rightarrow \infty$

$$df_{\lambda\rightarrow\infty} = \lim_{\lambda \rightarrow \infty} \frac{8\lambda+6}{(2\lambda+1)(2\lambda+3)} = \lim_{\lambda \rightarrow \infty} \frac{8\lambda+6}{4\lambda^2+8\lambda+3} = \lim_{\lambda \rightarrow \infty} \frac{\frac{8}{\lambda} + \frac{6}{\lambda^2}}{4 + \frac{8}{\lambda} + \frac{3}{\lambda^2}} \checkmark = \frac{0}{4} = 0$$

**+1** **(8) [+1]** Choose $\lambda$ by setting $df_\lambda = 1$ and $r_{12} = \frac{1}{2}$.

That is, solving that $\frac{8\lambda+6}{4\lambda^2+8\lambda+3} = 1 \Rightarrow 4\lambda^2 + 8\lambda + 3 = 8\lambda + 6$

$$\Rightarrow 4\lambda^2 = 3 \Rightarrow \lambda = \frac{\sqrt{3}}{2} \quad (\text{negative value is improper})$$

2

**+14**

**Q3 [+14]** Let $f(x) = \sum_{j=1}^{K} \beta_j M_j(x)$, where $M_j(x)$'s are the spline basis functions. For a knot

sequence $\xi_1 < \xi_2 < \xi_3$ with $\Delta = \xi_2 - \xi_1 = \xi_3 - \xi_2$, let $z_i(t) = (t - \xi_i)/\Delta$ for $i = 1, 2,$ and $3$.

**+2** **(1) [+2]** Show that $\int f''(x)^2 dx$ can be written as a quadratic form for $\boldsymbol{\beta} = (\beta_1, \ldots, \beta_K)^T$.

$$\int f''(x)^2 dx = \int \left( \sum_{i=1}^{K} \beta_i M_i''(x) \right) \left( \sum_{j=1}^{K} \beta_j M_j''(x) \right) dx = \sum_{i=1}^{K}\sum_{j=1}^{K} \beta_i \beta_j \int M_i''(x) M_j''(x) dx$$

$$= [\beta_1 \cdots \beta_K] \begin{bmatrix} \int M_1''(x) M_1''(x) dx & \cdots & \int M_1''(x) M_K''(x) dx \\ \vdots & & \\ \int M_K''(x) M_1''(x) dx & \cdots & \int M_K''(x) M_K''(x) dx \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_K \end{bmatrix} = \boldsymbol{\beta}^T \underbrace{\begin{bmatrix} \int M_1''(x) M_1''(x) dx \cdots & \int M_1''(x) M_K''(x) dx \\ \vdots & \\ \int M_K''(x) M_1''(x) dx \cdots & \int M_K''(x) M_K''(x) dx \end{bmatrix}}_{A} \boldsymbol{\beta}$$

**+4** **(2) [+4]** Calculate $\int M_2''(t)^2 dt$,

where $M_2(t) = \dfrac{I(\xi_1 \le t < \xi_2)}{2\Delta} \{ 7z_1(t)^3 - 18 z_1(t)^2 + 12 z_1(t) \} - \dfrac{I(\xi_2 \le t < \xi_3)}{2\Delta} z_3(t)^3 \qquad \left( \dfrac{d}{dt} z_i(t) = \dfrac{1}{\Delta} \right)$

$\to M_2'(t) = \dfrac{I(\xi_1 \le t < \xi_2)}{2\Delta^2} \left( 21 z_1(t)^2 - 36 z_1(t) + 12 \right) - \dfrac{3 I(\xi_2 \le t < \xi_3)}{2 \Delta^2} z_3(t)^2$ ✓

$\to M_2''(t) = \dfrac{I(\xi_1 \le t < \xi_2)}{2\Delta^3} \left( 42 z_1(t) - 36 \right) - 3 \dfrac{I(\xi_2 \le t < \xi_3)}{\Delta^3} z_3(t)$ ✓

$\Rightarrow \int M_2''(t)^2 dt = \int \dfrac{I(\xi_1 \le t < \xi_2)}{4\Delta^6} (42 z_1(t) - 36)^2 dt + 9 \int \dfrac{I(\xi_2 \le t < \xi_3)}{\Delta^6} z_3^2(t) dt$

$\qquad - 3 \underbrace{\int \dfrac{I(\xi_1 \le t < \xi_2) I(\xi_2 \le t < \xi_3)}{\Delta^6} z_3(t)(42 z_1(t) - 36) \, dt}_{\text{this term equal to zero} \quad \because I(\xi_1 \le t < \xi_2) I(\xi_2 \le t < \xi_3) = 0 \ \forall t}$

$= \dfrac{1}{4\Delta^6} \int_{\xi_1}^{\xi_2} (42 z_1(t) - 36)^2 dt + \dfrac{9}{\Delta^6} \int_{\xi_2}^{\xi_3} z_3^2(t) dt$

$\left( \text{Let } x = z_1(t) = \dfrac{t - \xi_1}{\Delta} \quad \text{and} \quad \text{let } y = z_3(t) = \dfrac{t - \xi_3}{\Delta} \right.$
$\left. \qquad \to dt = \Delta dx \qquad\qquad \qquad \to dt = \Delta dy \right)$

$= \dfrac{1}{4\Delta^6} \int_0^1 \Delta (42x - 36)^2 dx + \dfrac{9}{\Delta^6} \int_{-1}^0 \Delta y^2 dy$

$= \dfrac{1}{4\Delta^5} \times \left. \dfrac{(42x - 36)^3}{3 \times 42} \right|_{x=0}^{x=1} + \dfrac{9}{\Delta^5} \times \left. \dfrac{y^3}{3} \right|_{y=-1}^{y=0}$

$= \dfrac{1}{4\Delta^5} \times \dfrac{6^3 + 36^3}{126} + \dfrac{9}{\Delta^5} \times \dfrac{1}{3} = \dfrac{1}{4\Delta^5} \times 372 + \dfrac{3}{\Delta^5}$

✓ $= \dfrac{93}{\Delta^5} + \dfrac{3}{\Delta^5} = \dfrac{96}{\Delta^5}$

3

**+2** (3) [+2] Let $\Omega$ be a matrix with elements $\Omega_{jk} = \int M_j''(t)M_k''(t)dt$. Let $\{(Y_i, x_i), \ i=1,\ldots,n\}$ be data.

Let $\hat{f}_\lambda(x) = \sum_{j=1}^{K} \hat{\beta}_{j\lambda} M_j(x)$ be a smoothing spline. Derive the formula of $\hat{\boldsymbol{\beta}}_\lambda = (\hat{\beta}_{1\lambda}, \cdots, \hat{\beta}_{1\lambda})^T$.

RSS with penalty of roughness : $RSS(\beta) = (y-X\beta)^T(y-X\beta) + \lambda \int f''(x)^2 dx$

$\checkmark$ $(y-X\beta)^T(y-X\beta) + \lambda \beta^T \Omega \beta$ where $y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$, $\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$, $X = \begin{bmatrix} M_1(x_1) & \cdots & M_K(x_1) \\ \vdots & & \vdots \\ M_1(x_n) & \cdots & M_K(x_n) \end{bmatrix}$

$\Rightarrow \hat{\beta}_\lambda = \underset{\beta}{\arg\min} \ RSS(\beta) = \underset{\beta}{\arg\min} \ (y-X\beta)^T(y-X\beta) + \lambda \beta^T \Omega \beta$

That's equivalent to $\quad \frac{\partial}{\partial \beta} RSS(\beta)\Big|_{\hat{\beta}_\lambda} = -2X^T(y-X\hat{\beta}_\lambda) + 2\lambda \Omega \hat{\beta}_\lambda \overset{set}{=\!=\!=} 0$

$\Rightarrow (X^TX + \lambda \Omega)\hat{\beta}_\lambda = X^Ty \quad \checkmark \Rightarrow \hat{\beta}_\lambda = \checkmark (X^TX + \lambda \Omega)^{-1} X^Ty$

**+2** (4) [+2] Define $CV(\lambda)$.

$CV(\lambda) = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{Y_i - \hat{f}_\lambda(x_i)}{1 - h_i} \right)^2$,
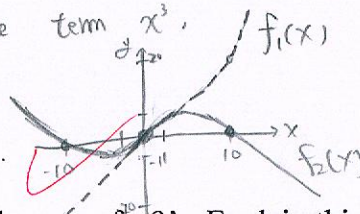
where $h_i$ is the $i$-th diagonal element of $\checkmark X(X^TX + \lambda \Omega)^{-1} X^T$

**+2** (5) [+2] The polynomial regression is sensitive to changes of $\beta$'s. Explain this by a figure.

Here compare these to function $\begin{cases} f_1(x) = x + \frac{1}{100} x^3 \\ f_2(x) = x - \frac{1}{100} x^3 \end{cases}$
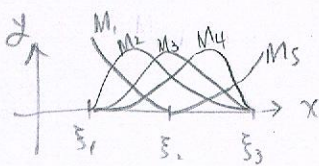
The coefficient $\frac{1}{100}$ can be treated as a small value, but it has large influence of the term $x^3$.

| $x$ | -10 | -1 | 0 | 1 | 10 |
|-----|-----|-----|-----|-----|-----|
| $f_1$ | -20 | -1.01 | 0 | 1.01 | 20 |
| $f_2$ | 0 | -0.99 | 0 | 0.99 | 0 |



$\rightarrow$ The performance of $f_1$ and $f_2$ are very different, this show that polynomial regression is sensitive to the changes of $\beta$'s.
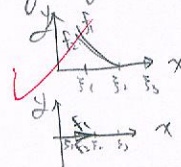
**+2** (6) [+2] The cubic spline is not sensitive to changes of $\beta$'s. Explain this by a figure.
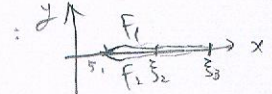


We see the figure marginally,

1. $\begin{cases} f_1(x) = 1.01 \ M_1(x) \\ f_2(x) = 0.99 \ M_1(x) \end{cases}$ ;

2. $\begin{cases} f_1(x) = 0.01 \ M_1(x) \\ f_2(x) = -0.01 \ M_2(x) \end{cases}$ ;

3. $\begin{cases} f_1(x) = 0.01 \ M_2(x) \\ f_2(x) = -0.01 \ M_2(x) \end{cases}$

$\cdots$ etc.

we can see the small change of $\beta_j$'s does not have obviously influence on the different model marginally, and by the continuity of cubic spline, a very small difference of each $\beta$'s also no influence the result obviously (not sensitive to changes of $\beta$'s)

4