

1.

(a) True $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y \quad E(\hat{\beta}_{OLS}) = (X^T X)^{-1} X^T X \beta^* = \beta^*$

(b) True bias-variance trade-off is shown in Q.3.4

Variance = $d_2 \sum_{i=1}^d \frac{d_{ii}}{(d_{ii} + \lambda)^2}$ is monotonically decreasing, when $\lambda \uparrow$

and bias is monotonically increasing when $\lambda \uparrow$

(c) True ① $MSE(\hat{\beta}_{OLS}) = \text{Var}(\hat{\beta}_{OLS})$ Suppose ① = ② at some λ

② $MSE(\hat{\beta}_\lambda^{\text{Lasso}}) = \text{Var}(\hat{\beta}_\lambda^{\text{Lasso}}) + \text{bias}(\hat{\beta}_\lambda^{\text{Lasso}})^2$

$$> \text{Var} + \Delta_r + B - \Delta_B, \text{ where } \Delta_r < \Delta_B$$

$$\textcircled{1} > \textcircled{2}$$

(d) False $\hat{\beta}_\lambda^{\text{Lasso}} = 0$ when $\lambda = \infty$

(e) False $\lambda > 0$, every solution $\hat{\beta}$ has the same L_1 norm.

(f) True $H = VDU^T \quad H^2 = VDU^T H = H$

$$(VDU^T)(VDU^T) = VD^2U^T = VDU^T$$

$$d_{ii}^2 = d_{ii}$$

$$d_{ii}(d_{ii} - 1) = 0 \quad d_{ii} = 1 \text{ or } d_{ii} = 0$$

(g) True

A symmetric matrix with non-negative eigenvalues is a PSD.

$$(X(X^T X)^{-1} X^T)^T = X(X^T X)^{-1} X^T \text{ symmetric}$$

$$X(X^T X)^{-1} X^T \cdot X(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T \text{ idempotent}$$

H has only non-negative eigenvalues, so H is a PSD

(h) False

$$\text{trace}(Q) = \sum_{i=1}^n (I_{ii} - H_{ii}) = n - d = n - d$$

$$\text{trace}(H) = \text{trace}(X^T X (X^T X)^{-1}) = \text{trace}(I_d) = d$$

2.

$$\begin{aligned}
 1. \quad \hat{\beta} &= (X^T X)^{-1} X^T Y \quad E(\hat{\beta}) = (X^T X)^{-1} X^T E(Y) \\
 &= (X^T X)^{-1} X^T X \beta^* \\
 &= \beta^*
 \end{aligned}$$

$$\begin{aligned}
 \text{Cov}(\hat{\beta}) &= \text{Cov}[(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T \text{Cov}(Y) [(X^T X)^{-1} X^T]^T \\
 &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\
 &= \sigma^2 (X^T X)^{-1}
 \end{aligned}$$

Since Y is Normal, $\hat{\beta}$ is a linear function of Y

so $\hat{\beta}$ is also normal $\hat{\beta} \sim N(\beta^*, \sigma^2 (X^T X)^{-1})$

$$\begin{aligned}
 E(X\hat{\beta}) &= X E(\hat{\beta}) = X (X^T X)^{-1} X^T Y = X (X^T X)^{-1} X^T X \beta^* \\
 &= X \beta^*
 \end{aligned}$$

$$2. \quad E(\text{RSS}) = E(\hat{e}^T \hat{e}) = E(Y^T (I-H)^T (I-H) Y) = E(e^T (I-H)e)$$

$$\begin{aligned}
 \text{Since } \hat{e} &= (I-H)Y \\
 &= (I-H)X\beta + e \\
 &= (I-H)e
 \end{aligned}$$

$$= \sum_{i,j} (I_{ij} - H_{ij}) E(e_i e_j)$$

Since $e_i e_j$ uncorrelated when $i \neq j$

$$\begin{aligned}
 \text{tr}(H) &= \text{tr}(X (X^T X)^{-1} X^T) \\
 &= \text{tr}((X^T X)^{-1} X^T X) \\
 &= \text{tr}(I_d) \\
 &= d
 \end{aligned}
 \quad \begin{aligned}
 &= \sum_{i=1}^n \sigma^2 (I_i - h_{ii}) \\
 &= \sigma^2 (n - \text{tr}(H)) \\
 &= \sigma^2 (n - d)
 \end{aligned}$$

3.

We know OLS $\hat{\beta}$ is unbiased

let linear unbiased estimator $\bar{\beta} = DY$, $\bar{\beta}$ is OLS when $D = (X^T X)^{-1} X^T$

now let $D = A + (X^T X)^{-1} X^T$

$$\begin{aligned}\mathbb{E}(\bar{\beta}) &= \mathbb{E}(A Y + (A + (X^T X)^{-1} X^T) Y) \\ &= \mathbb{E}(AY) + \beta^* \\ &= AX\beta^* + \beta^* \quad \text{since } \bar{\beta} \text{ is unbiased, we need } AX=0\end{aligned}$$

$$\text{Cov}(\bar{\beta}) = \text{Cov}(A + (X^T X)^{-1} X^T) Y$$

$$\begin{aligned}&= (A + (X^T X)^{-1} X^T) \text{Cov}(Y) (A + (X^T X)^{-1} X^T)^T \\ &= 6^2 (A + (X^T X)^{-1} X^T) (A^T + X(X^T X)^{-1}) \\ &= 6^2 (A A^T + A X (X^T X)^{-1} + (X^T X)^{-1} X^T A^T + (X^T X)^{-1}) \\ &= 6^2 A A^T + 6^2 (X^T X)^{-1} \\ &= 6^2 A A^T + \text{cov}(\hat{\beta}^{\text{OLS}}) \quad \text{we have } \text{cov}(\bar{\beta}) \geq \text{cov}(\hat{\beta}^{\text{OLS}})\end{aligned}$$

$A A^T$ is a PSD

2.4

We train use $(x_1, y_1) \dots (x_N, y_N)$ to get least square estimate $\hat{\beta}$

this means we minimize $R_{\text{tr}}(\beta) = \frac{1}{N} \sum_{i=1}^N (y_i - \beta^T x_i)^2$ to get $\hat{\beta}$

When we predict, we use the data we've never seen, if we are lucky to have exactly same test data as training data. We have $E(R_{\text{tr}}(\hat{\beta})) = E(R_{\text{te}}(\hat{\beta}))$

But usually we have new data we've never seen, so on average we'll predict worse in test data than in training data. That is $E(R_{\text{tr}}(\hat{\beta})) \leq E(R_{\text{te}}(\hat{\beta}))$

3.

$$\begin{aligned} & \min_{\theta} \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2 \\ &= \min_{\theta} \|X\theta - y\|_2^2 + \|\sqrt{\lambda}\theta + 0\|_2^2 \\ &= \min_{\theta} \|A\theta - Y\|_2^2 \end{aligned}$$

$$\text{let } A = \begin{pmatrix} X_{n \times d} \\ \sqrt{\lambda} I_{d \times d} \end{pmatrix} \quad Y = \begin{pmatrix} y \\ 0 \end{pmatrix}$$

This is similar to LS, we can use normal equation to get $\hat{\theta}_\lambda$

$$\begin{aligned} \hat{\theta}_\lambda &= (A^T A)^{-1} A^T Y \\ &= [(X^T \sqrt{\lambda} I) \cdot \begin{pmatrix} X \\ \sqrt{\lambda} I \end{pmatrix}]^{-1} \cdot (X^T \sqrt{\lambda} I) \begin{pmatrix} y \\ 0 \end{pmatrix} \\ &= (X^T X + \lambda I)^{-1} X^T y \end{aligned}$$

let A be PSD, then $x^T A x \geq 0$ for any vector x

let B be PD, then $x^T B x > 0$

$$x^T (A + B) x = x^T A x + x^T B x > 0$$

$$\geq 0 \quad > 0$$

so $A + B$ is a PD

$X^T X$ is a PSD

λI is a PD for $\lambda > 0$

so $X^T X + \lambda I$ is a PD

and PD is always invertible

so we have a unique solution
for any $\lambda > 0$

$$3.2. \quad E(\hat{\theta}_\lambda) = E((X^T X + \lambda I)^{-1} X^T y)$$

$$= (X^T X + \lambda I)^{-1} X^T E(y)$$

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

$$= W_\lambda \theta^*$$

$$\text{Cov}(\hat{\theta}_\lambda) = (X^T X + \lambda I)^{-1} X^T \text{Cov}(y) X (X^T X + \lambda I)^{-1}$$

$$= \sigma^2 W_\lambda (X^T X + \lambda I)^{-1}$$

Since $\hat{\theta}_\lambda$ is a linear function of y , and y is normal

so $\hat{\theta}_\lambda$ is normal as well.

$$\hat{\theta}_\lambda \sim N(W_\lambda \theta^*, \sigma^2 W_\lambda (X^T X + \lambda I_d)^{-1})$$

3.3

$$\begin{aligned}
\|\mathbb{E}(\hat{\theta}_\lambda) - \theta^*\|_2^2 &= \|\mathbb{W}_\lambda \theta^* - \theta^*\|_2^2 \\
&= \|[(X^T X + \lambda I_d)^{-1} X^T X - I] \theta^*\|_2^2 \\
&= \|[(X^T X + \lambda I_d)^{-1} X^T X - (X^T X + \lambda I_d)^{-1} (X^T X + \lambda I_d)] \theta^*\|_2^2 \\
&= \|-\lambda (X^T X + \lambda I_d)^{-1} \theta^*\|_2^2 \\
&= \lambda^2 \|(V D V^T + \lambda U U^T)^{-1} \theta^*\|_2^2 \\
&= \lambda^2 \|U(D + \lambda I)^{-1} U^T \theta^*\|_2^2 \\
&= \lambda^2 \text{tr}((\theta^*)^T U(D + \lambda I)^{-1} (D + \lambda I)^{-1} U^T \theta^*) \\
&= \lambda^2 \text{tr}(\|(D + \lambda I)^{-1} U^T \theta^*\|_2^2) \quad \text{where } U = \begin{pmatrix} u_1 & \dots & u_d \\ 1 & \dots & 1 \end{pmatrix} \\
&= \sum_{i=1}^d \frac{\lambda^2}{(d+i\lambda)^2} (u_i^T \cdot \theta^*)^2
\end{aligned}$$

$$U^T = \begin{pmatrix} u_1^T \\ \vdots \\ u_d^T \end{pmatrix}$$

u_i^T is the i th column of U

$$\begin{aligned}
&\mathbb{E}\|\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda)\|_2^2 \\
&= \mathbb{E}(\text{tr}((\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))^T (\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda)))) \\
&= \mathbb{E}(\text{tr}((\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda)^T))) \\
&= \text{tr}(\mathbb{E}((\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda)^T))) \\
&= \text{tr}(\text{cov}(\hat{\theta}_\lambda)) \quad \text{where } \text{cov}(\hat{\theta}_\lambda) = 6^2 \mathbb{W}_\lambda (X^T X + \lambda I_d)^{-1} = \\
&= \text{tr}(6^2 (D + \lambda I)^{-1} D (D + \lambda I)^{-1}) \\
&= \text{tr}(6^2 (D + \lambda I)^{-2} D) \\
&= \sum_{i=1}^d 6^2 \cdot \frac{d_i}{(d+i\lambda)^2}
\end{aligned}$$

$$\begin{aligned}
&= 6^2 (V D V^T + \lambda U U^T)^{-1} V D V^T (V D V^T + \lambda I_d)^{-1} \\
&= 6^2 U(D + \lambda I)^{-1} U^T V D V^T U (V D V^T + \lambda I_d)^{-1} U^T \\
&= 6^2 U(D + \lambda I)^{-1} D (D + \lambda I)^{-1} U^T
\end{aligned}$$

3.4

$$\lambda = 0 \text{ squared bias} = 0 \quad \text{scalar variance} = \sigma^2 \sum_{i=1}^n \frac{1}{d_i^2}$$

$$\lambda = \infty \text{ squared bias} = \sum_{i=1}^d \theta_i^{*2} \quad \text{scalar variance} = 0$$

when bias is at its minimum, the variance reaches its maximum

when bias is at its maximum, the variance reaches its minimum.

and squared bias is a monotonic increasing function of λ for $\lambda \in [0, +\infty)$

scalar variance is a monotonic decreasing function of λ for $\lambda \in [0, +\infty)$

Variance = $\sigma^2 \sum_{i=1}^d \frac{1}{(d_i + \lambda)^2}$ as λ increases, variance decreases, bias increases

3.5

$$\text{trace}(M(\lambda)) = \text{trace}(\mathbb{E}[(\hat{\theta}_\lambda - \theta^*)(\hat{\theta}_\lambda - \theta^*)^T])$$

$$= \mathbb{E}(\text{trace}((\hat{\theta}_\lambda - \theta^*)(\hat{\theta}_\lambda - \theta^*)^T))$$

$$= \mathbb{E}(\text{trace}((\hat{\theta}_\lambda - \theta^*)^T(\hat{\theta}_\lambda - \theta^*)))$$

$$= \mathbb{E}(\|\hat{\theta}_\lambda - \theta^*\|_2^2)$$

$$\mathbb{E}(\|\hat{\theta}_\lambda - \theta^*\|_2^2) = \mathbb{E}(\|\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda) + \mathbb{E}(\hat{\theta}_\lambda) - \theta^*\|_2^2)$$

$$= \mathbb{E}(\|\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda)\|_2^2 + \mathbb{E}(\|\mathbb{E}(\hat{\theta}_\lambda) - \theta^*\|_2^2)$$

$$+ 2\mathbb{E}((\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))^T \cdot (\mathbb{E}(\hat{\theta}_\lambda) - \theta^*))]$$

$$= \text{Scalar-variance} + \text{squared bias} \quad \text{where } \mathbb{E}(\hat{\theta}_\lambda) - \theta^* \text{ is constant}$$

$$+ 2(\mathbb{E}(\hat{\theta}_\lambda) - \theta^*) \cdot \mathbb{E}((\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))^T)$$

$$= \text{Scalar-variance} + \text{squared bias} \quad \text{where } \mathbb{E}(\hat{\theta}_\lambda^T) - \mathbb{E}(\mathbb{E}(\hat{\theta}_\lambda)^T) = 0$$

3.6. $\hat{\theta}_{OLS}$ is when $\lambda = 0$

$$\begin{aligned}\text{trace}(M(0)) &= \text{trace}(\mathbb{E}(\hat{\theta}_{OLS} - \theta^*)(\hat{\theta}_{OLS} - \theta^*)^T) \\ &= \mathbb{E}(\text{trace}((\hat{\theta}_{OLS} - \theta^*)^T(\hat{\theta}_{OLS} - \theta^*))) \\ &= \mathbb{E}(\|\hat{\theta}_{OLS} - \theta^*\|^2) \\ &= \text{MSE}(\hat{\theta}_{OLS})\end{aligned}$$

$$\text{MSE}(\hat{\theta}_{OLS}) > \text{MSE}(\hat{\theta}_\lambda)$$

$$\Rightarrow \text{MSE}(\hat{\theta}_{OLS}) - \text{MSE}(\hat{\theta}_\lambda) > 0$$

$$\Rightarrow \text{trace}(M(0)) - \text{trace}(M(\lambda)) > 0$$

$$\Rightarrow \text{trace}(M(0) - M(\lambda)) > 0$$

We need to prove $M(0) - M(\lambda)$ is a \checkmark PD matrix
 Since the trace of a matrix is equal to the sum of its eigenvalues, and a \checkmark P.D matrix has positive eigenvalues so if $M(0) - M(\lambda)$ is a PD, then we can prove

first derive $M(\lambda)$

$$\begin{aligned}M(\lambda) &= \mathbb{E}[(\hat{\theta}_\lambda - \theta^*)(\hat{\theta}_\lambda - \theta^*)^T] \\ &= \mathbb{E}[(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda) + \mathbb{E}(\hat{\theta}_\lambda) - \theta^*)(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda) + \mathbb{E}(\hat{\theta}_\lambda) - \theta^*)^T] \\ &= \mathbb{E}[(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))^T] + \mathbb{E}[(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda))(\mathbb{E}(\hat{\theta}_\lambda) - \theta^*)^T]_{①} \\ &\quad + \mathbb{E}[(\mathbb{E}(\hat{\theta}_\lambda) - \theta^*)(\hat{\theta}_\lambda - \mathbb{E}(\hat{\theta}_\lambda)^T]_{②} + \mathbb{E}[(\mathbb{E}(\hat{\theta}_\lambda) - \theta^*)(\mathbb{E}(\hat{\theta}_\lambda) - \theta^*)^T]\end{aligned}$$

while for ① and ② $\mathbb{E}(\hat{\theta}_\lambda) - \theta^*$ is a constant, and $\mathbb{E}(X - \mathbb{E}(X)) = 0$

$$= \text{Var}(\hat{\theta}_\lambda) + 0 + 0 + (\mathbb{E}(\hat{\theta}_\lambda) - \theta^*)(\mathbb{E}(\hat{\theta}_\lambda) - \theta^*)^T$$

$$= \sigma^2 W_\lambda (X^T X + \lambda I_d)^{-1} + (W_\lambda \theta^* - I_d) \theta^* \theta^{*T} (W_\lambda - I_d)^T$$

next page

$$M(0) - M(\lambda)$$

when $\lambda = 0 \quad W_\lambda = I$

$$= \underbrace{6^2(X^T X)^{-1} - 6^2 W_\lambda (X^T X + \lambda I_d)^{-1}}_{\textcircled{1}} + \underbrace{(W_\lambda \Theta^* - I_d) \Theta^* \Theta^{*T} (W_\lambda - I_d)^T}_{\textcircled{2}}$$

$$\textcircled{1} = 6^2 (W_\lambda W_\lambda^{-1} (X^T X)^{-1} (W_\lambda^T)^{-1} W_\lambda^T - W_\lambda (X^T X)^{-1} W_\lambda^T)$$

$$= 6^2 W_\lambda (W_\lambda^T (X^T X)^{-1} (W_\lambda^T)^{-1} - (X^T X)^{-1}) W_\lambda^T \quad (\text{let } S = (X^T X))$$

$$= 6^2 W_\lambda (S^{-1}(S + \lambda I) S^{-1}(S + \lambda I) S^{-1} - S^{-1}) W_\lambda^T$$

$$= 6^2 W_\lambda (S^{-1}(SS^{-1}S + SS^{-1}\lambda + \lambda S^{-1}S + \lambda S^{-1}\lambda) S^{-1} - S^{-1}) W_\lambda^T$$

$$= 6^2 W_\lambda (S^{-1}(\lambda(2I + \lambda S^{-1})S^{-1}) W_\lambda^T$$

$$= 6^2 W_\lambda (2\lambda(X^T X)^{-2} + \lambda^2(X^T X)^{-3}) W_\lambda^T$$

$$= 6^2 (X^T X + \lambda I_d)^{-1} (2\lambda I_d + \lambda^2 (X^T X)^{-1}) (X^T X + \lambda I_d)^{-1}$$

$$\textcircled{2} \text{ use the fact } W_\lambda - I_d = -\lambda (X^T X + \lambda I_d)^{-1}$$

$$\begin{aligned} \text{Because } (X^T X + \lambda I)^{-1} X^T X - I_d &= (X^T X + \lambda I)^{-1} X^T X - (X^T X + \lambda I)^{-1} (X^T X + \lambda I) \\ &= (X^T X + \lambda I)^{-1} (-\lambda I) \\ &= -\lambda (X^T X + \lambda I_d)^{-1} \end{aligned}$$

$$\textcircled{2} = -\lambda (X^T X + \lambda I_d)^{-1} \Theta^* \Theta^{*T} (-\lambda (X^T X + \lambda I_d)^{-1})$$

$$= \lambda^2 (X^T X + \lambda I_d)^{-1} \Theta^* \Theta^{*T} (X^T X + \lambda I_d)^{-1}$$

next page

$$M(0) - M(\lambda)$$

$$= \textcircled{1} - \textcircled{2}$$

$$= \lambda (X^T X + \lambda I_d)^{-1} \underbrace{[26^2 I_d + 6^2 \lambda (X^T X)^{-1} - \lambda \theta^* \theta^{*T}]}_B (X^T X + \lambda I_d)^{-1} \underbrace{A}_{\lambda}$$

A is invertible and $A^T = ((X^T X + \lambda I_d)^{-1})^T = (X^T X + \lambda I_d)^{-1} = A$

We have $M(0) - M(\lambda)$

$$= \lambda A^T B A$$

This is a note:

$$B = 26^2 I_d + 6^2 \lambda (X^T X)^{-1} - \lambda \theta^* \theta^{*T}$$

$X^T X$ is symmetric for any ^{non-zero} vector $z \in \mathbb{R}^n$

$$z^T X^T X z = (Xz)^T (Xz) = \|Xz\|_2^2 \geq 0$$

Thus $X^T X$ is a PSD, this means $(X^T X)^{-1}$ is a PSD

We only need B to be PD

because let z be any non-zero vector in \mathbb{R}^n

$$z^T A^T B A z = (Az)^T B A z$$

$$= w^T B w > 0 \text{ iff } B \text{ is PD}$$

Since A is invertible $Az \neq 0$

PSD + PD = PD, so we need $26^2 I_d - \lambda \theta^* \theta^{*T}$ to be a PD

Similarly let z be any nonzero vector $\in \mathbb{R}^n$

$$26^2 z^T I_d z - \lambda z^T \theta^* \theta^{*T} z$$

$$= 26^2 \|z\|_2^2 - \lambda (\theta^{*T} z)^T \theta^{*T} z \geq 0 \quad \text{note } \theta^{*T} z \text{ is a scalar}$$

$$\lambda < \frac{26^2 \|z\|_2^2}{\|\theta^{*T} z\|_2^2} = \frac{26^2 \|z\|_2^2}{\|\theta^*\|_2^2 \|z\|_2^2} = \frac{26^2}{\|\theta^*\|_2^2}$$

Conclude $MSE(\hat{\theta}_\lambda) < MSE(\theta^{OLS})$
for $\lambda \in (0, \frac{26^2}{\|\theta^*\|_2^2})$

$$4.1 \quad ① f(\lambda x + (1-\lambda)y) = \frac{L}{2} (\lambda x + (1-\lambda)y)^2 = \frac{L}{2} (\lambda^2 x^2 + 2\lambda x(1-\lambda)y + (1-\lambda)^2 y^2)$$

$$② \lambda f(x) + (1-\lambda)f(y) = \lambda \frac{L}{2} x^2 + (1-\lambda) \frac{L}{2} y^2$$

To prove convexity of $f(x) = \frac{L}{2} x^2$

let $② > ①$

$$\lambda f(x) + (1-\lambda)f(y) - f(\lambda x + (1-\lambda)y) = \lambda(1-\lambda) \frac{L}{2} (x^2 + y^2 - 2xy)$$

$$= \lambda(1-\lambda) \frac{L}{2} (x-y)^2 \geq 0$$

Since $\lambda \in [0,1]$, $\frac{L}{2} \geq 0$

Thus $f(x)$ is convex

$$\nabla_x f(x) = Lx \quad \text{gradient update}$$

$$x_{t+1} = x_t - \gamma L x_t$$

$$= x_t (1 - \gamma L)$$

$$= x_{t-1} (1 - \gamma L) (1 - \gamma L)$$

$$= x_0 (1 - \gamma L)^{t+1}$$

4.2

when $\gamma = \frac{1}{L}$, $X_t = X_{(0)}(1 - \frac{1}{L}L)^t$ $X_{(1)}$ becomes 0 right away
 $= 0$

when $\gamma = \frac{2}{L}$, $X_t = X_{(0)}(1 - \frac{2}{L} \cdot L)^t$
 $= X_{(0)}(-1)^t$ X_t jumps between $\{-X_0, X_0\}$

$\gamma \in (0, \frac{2}{L})$ where $\gamma \neq \frac{1}{L}$

$$\gamma L \in (0, 1) \cup (1, 2)$$

$$1 - \gamma L \in (-1, 0) \cup (0, 1)$$

We have $X_t = X_{(0)}(1 - \gamma L)^t$

Since $1 - \gamma L$ shrink $X_{(0)}$ as t increases

We get a convergence

let $|X_t - X_*| \leq \epsilon$ X_* is a point where $f(x)$ gets optimum (gradient is zero)

$$|X_t| - |X_*| \leq |X_t - X_*| \leq \epsilon$$

$$|X_{(0)}(1 - \gamma L)^t| \leq \epsilon + |X_*|$$

$$|(1 - \gamma L)^t| \leq \frac{\epsilon + |X_*|}{|X_{(0)}|}$$

$$t \log(|1 - \gamma L|) \leq \log\left(\frac{\epsilon + |X_*|}{|X_{(0)}|}\right)$$

$$t \leq \frac{\log\left(\frac{\epsilon}{|X_{(0)}|}\right)}{\log(|1 - \gamma L|)}$$

$X_* = 0$ in this case

4.3

graphically $f(x) = \frac{L}{2}(x-c)^2$ is a horizontally shifted $f(x) = \frac{L}{2}x^2$

so we'll get the same steps

$$4.4. \nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{pmatrix} = \begin{pmatrix} L & 0 \\ 0 & m \end{pmatrix}$$

Let $z \in \mathbb{R}^n z \neq 0$, $z^T \begin{pmatrix} L & 0 \\ 0 & m \end{pmatrix} z = (z_1, z_2) \begin{pmatrix} L & 0 \\ 0 & m \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$

$$= Lz_1^2 + mz_2^2 > 0 \quad \nabla^2 f(x) \text{ is PD}$$

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix} = \begin{pmatrix} LX_1 \\ mX_2 \end{pmatrix}$$

Thus $f(x)$ is convex

$$\begin{aligned} x^t &= x^{t-1} - r \nabla f(x) \\ &= \begin{pmatrix} x_1^{t-1} \\ x_2^{t-1} \end{pmatrix} - r \begin{pmatrix} LX_1 \\ mX_2 \end{pmatrix} \\ &= \begin{pmatrix} x_1^{t-1}(1-rL) \\ x_2^{t-1}(1-rm) \end{pmatrix} \\ &= \begin{pmatrix} x_1^{t-2}(1-rL)^2 \\ x_2^{t-2}(1-rm)^2 \end{pmatrix} \\ &\vdots \\ &= \begin{pmatrix} x_1^0(1-rL)^t \\ x_2^0(1-rm)^t \end{pmatrix} \end{aligned}$$

5.1 when $d > n$, $\text{rank}(X) = n$

while the number of columns are $d > n$

X is not full rank, so X is not invertible

$X^T X$ is also not invertible. Thus OLS is not unique.

$P_{X^T X}$

$$\text{rank}(X) = \text{rank}(X^T X) \leq n < p$$

residuals are still orthogonal to $\text{span}(X)$

$$f(\theta) = \frac{1}{2} \|X\theta - Y\|_2^2$$

and their sum is 0

$$\nabla_{\theta} f(\theta) = X^T(X\theta - Y) = 0$$

$$X^T(X\hat{\theta}_{OLS} - Y) = 0$$

$$X^T \cdot \hat{e} = 0$$

5.2 ridge solution is unique, because $X^T X + \lambda I_d$ is a PD.

$X^T X$ is PSD, λI_d is PD for $\lambda > 0$, so $X^T X + \lambda I_d$ is PD

PD is invertible, so $\hat{\beta}(\lambda) = (X^T X + \lambda I_d)^{-1} X^T Y$ has a unique solution

5.3

$$\lim_{\lambda \rightarrow 0} (X^T X + \lambda I_d)^{-1} X^T Y = \hat{\beta}_{OLS} \text{ and this OLS is not unique}$$

```
## Warning: package 'glmnet' was built under R version 3.4.4

## Warning: package 'Matrix' was built under R version 3.4.4

## Warning: package 'ggplot2' was built under R version 3.4.4

## Warning: package 'caret' was built under R version 3.4.4

#5.4

n=1000
d=5000
x = matrix(rnorm(5000*1000, mean = 0, sd = 1), nrow = 1000)
err = rnorm(1000, 0, sqrt(0.25))
beta = c(rep(1, 15), rep(0, 4985))
y = x %*% beta + err
```

```
#5.5
```

```
training = x[1:800,]
training.label = y[1:800]
test = x[801:1000,]
test.label = y[801:1000]
```

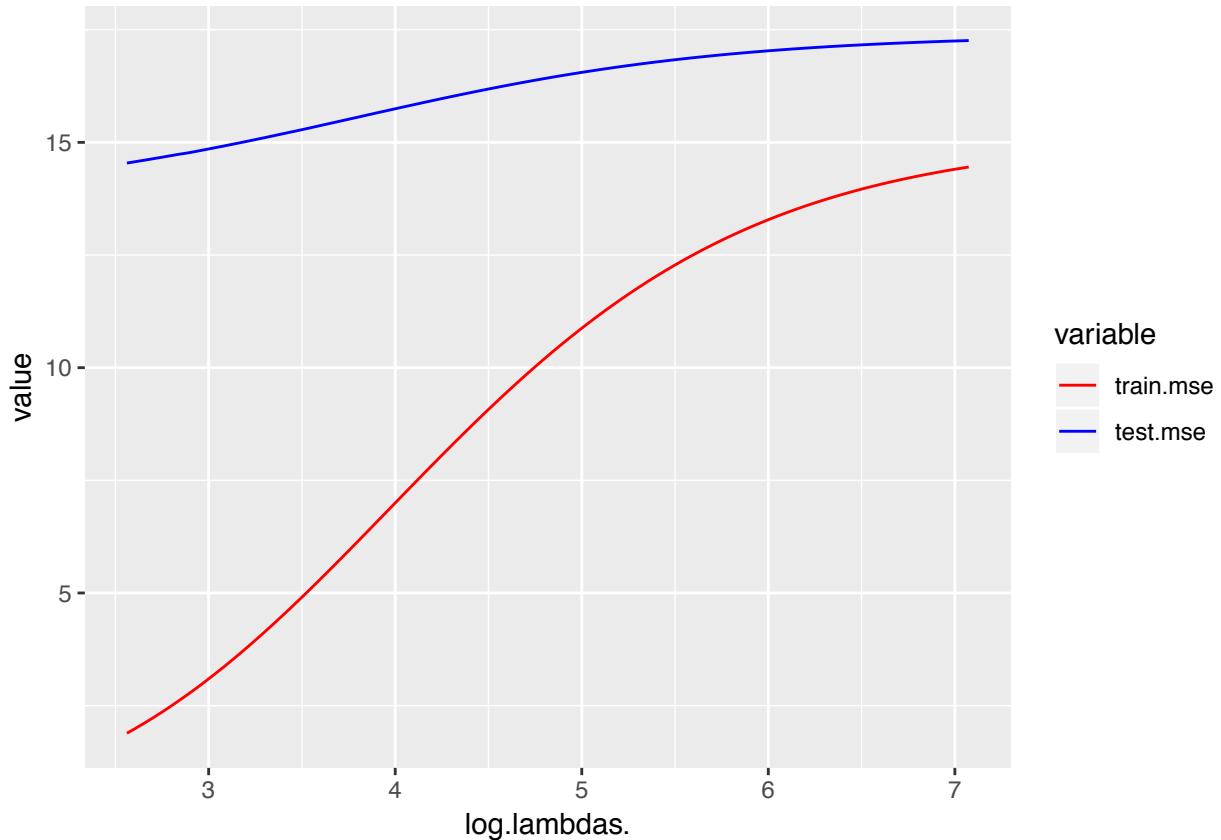
```
#5.6
```

```
mod = cv.glmnet(training, training.label, alpha=0, type.measure = "mse")

lambdas = mod$lambda
training.error = rep(0,length(lambdas))
test.error = rep(0,length(lambdas))

for (i in 1:length(lambdas)) {
  fit.ridge= glmnet(training, training.label, alpha=0, lambda = lambdas[i])
  training.error[i] = mean((training.label - predict(fit.ridge, training))^2)
  test.error[i] = mean((test.label - predict(fit.ridge, test))^2)
}

dat <- data.frame("log(lambda)" = log(lambdas), train.mse = training.error, test.mse = test.error)
dat.m <- melt(dat, id.vars = "log.lambda")
ggplot(dat.m, aes(log.lambda., value, colour = variable)) +
  geom_line() +
  scale_colour_manual(values = c("red", "blue"))
```



```
ridge.mse = test.error[which(test.error == min(test.error))]
ridge.lambda.min = lambdas[which(test.error == min(test.error))]
```

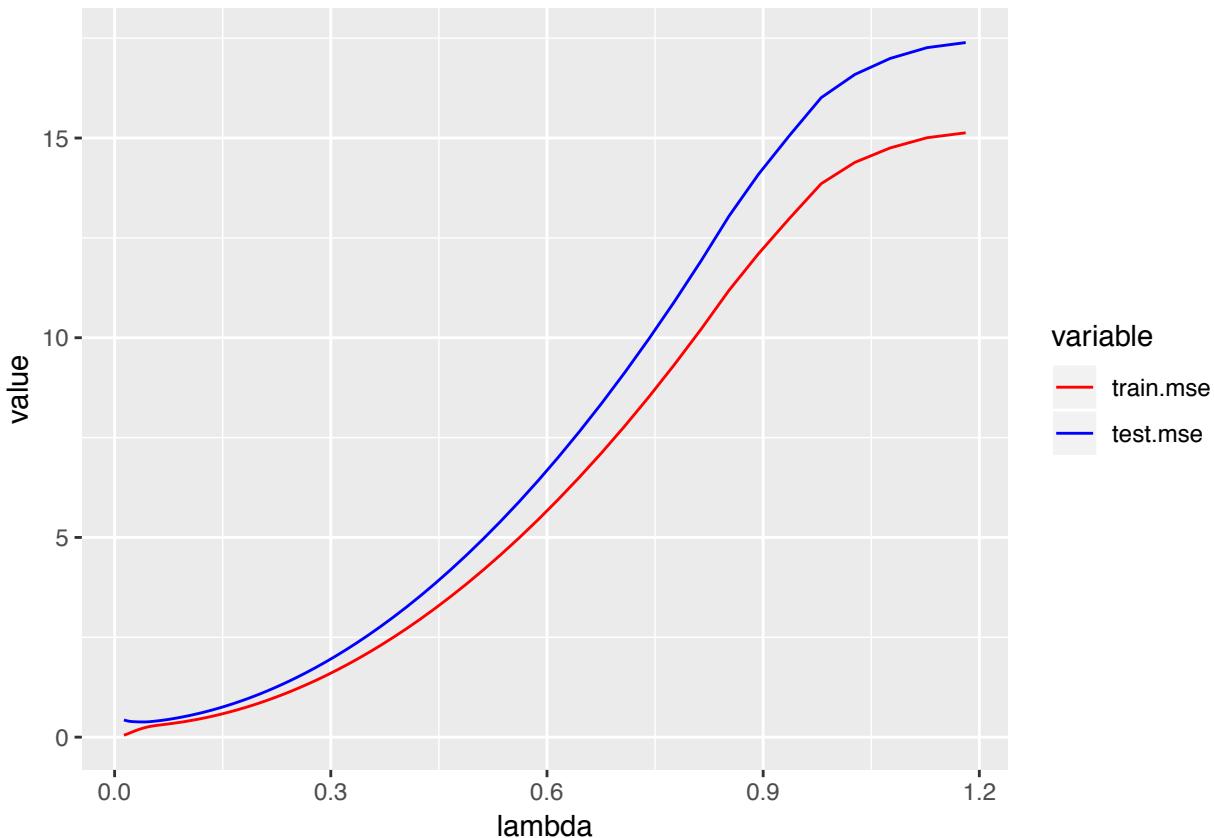
```
#5.7
```

```
mod = cv.glmnet(training, training.label, alpha=1, type.measure = "mse")

lambdas = mod$lambda
training.error = rep(0,length(lambdas))
test.error = rep(0,length(lambdas))

for (i in 1:length(lambdas)) {
  fit.lasso= glmnet(training, training.label, alpha=1, lambda = lambdas[i])
  training.error[i] = mean((training.label - predict(fit.lasso, training))^2)
  test.error[i] = mean((test.label - predict(fit.lasso, test))^2)
}

dat <- data.frame(lambda = lambdas, train.mse = training.error, test.mse = test.error)
dat.m <- melt(dat, id.vars = "lambda")
ggplot(dat.m, aes(lambda, value, colour = variable)) +
  geom_line() +
  scale_colour_manual(values = c("red", "blue"))
```



```
lasso.mse = test.error[which(test.error == min(test.error))]
lasso.lambda.min = lambdas[which(test.error == min(test.error))]
```

```
#5.8  
ridge.lambda.min  
  
## [1] 12.96238  
  
ridge.mse  
  
## [1] 14.54114  
  
lasso.lambda.min  
  
## [1] 0.03958523  
  
lasso.mse  
  
## [1] 0.3805818  
  
#6
```

#6.1-6.2

```
MSE = function(y, x, beta){  
  sum((y-(x %*% beta))^2)/length(y)  
}  
R2 = function(y, x, beta){  
  cor(y, x %*% beta)^2  
}
```

```

#6.3

test.mse= rep(0, 11)
modelQuality = matrix(rep(0,22),nrow=11)
modelQualityRImp =matrix(rep(0,22),nrow=11)
for (i in 1:length(Xnames)) {
  formula = as.formula(paste("log(SalePrice + 1) ~ ", paste(Xnames[1:i], collapse= "+")))
  lmod = lm(formula, data = AmesTinyTrain)
  training.label = log(AmesTinyTrain$SalePrice+1)
  mse = MSE(training.label, model.matrix(lmod),lmod$coefficients)
  r2 = R2(training.label,model.matrix(lmod),lmod$coefficients)
  modelQuality[i,1] = mse
  modelQuality[i,2] = r2
  r.out = summary(lmod)
  modelQualityRImp[i,1] = sum(r.out$residuals^2)/nrow(model.matrix(lmod))
  modelQualityRImp[i,2] = r.out$r.squared

  test.label = log(AmesTinyTest$SalePrice+1)
  test.mse[i] = sum((test.label-predict(lmod, AmesTinyTest))^2)/nrow(AmesTinyTest)
}

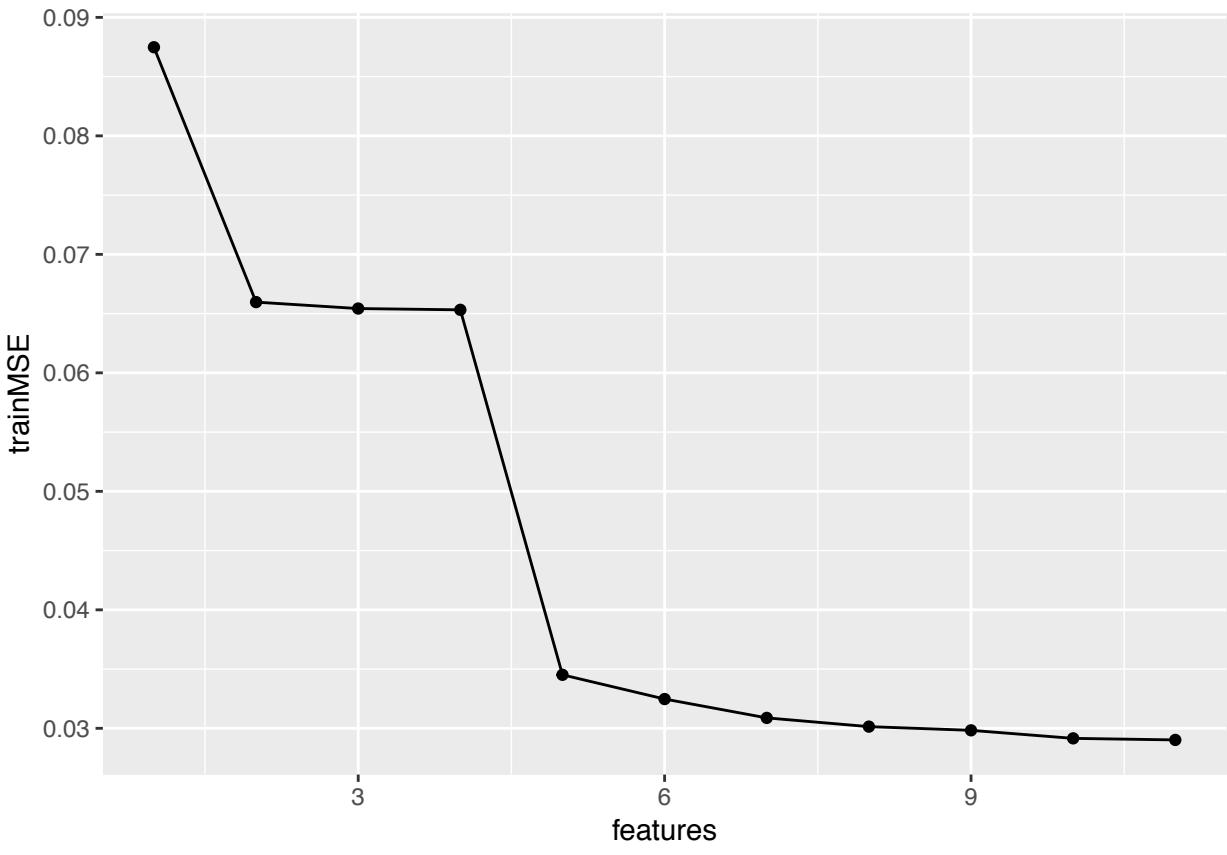
modelQuality = as.data.frame(modelQuality)
modelQualityRImp = as.data.frame(modelQualityRImp)
colnames(modelQuality) <- c("MSE", "R2")
colnames(modelQualityRImp) <- c("MSE", "R2")
modelQuality

##          MSE         R2
## 1  0.08748009 0.4745942
## 2  0.06596899 0.6037900
## 3  0.06542675 0.6070466
## 4  0.06531829 0.6076980
## 5  0.03451592 0.7926972
## 6  0.03247054 0.8049818
## 7  0.03087920 0.8145394
## 8  0.03014235 0.8189649
## 9  0.02982983 0.8208419
## 10 0.02915284 0.8249079
## 11 0.02902088 0.8257005

modelQualityRImp

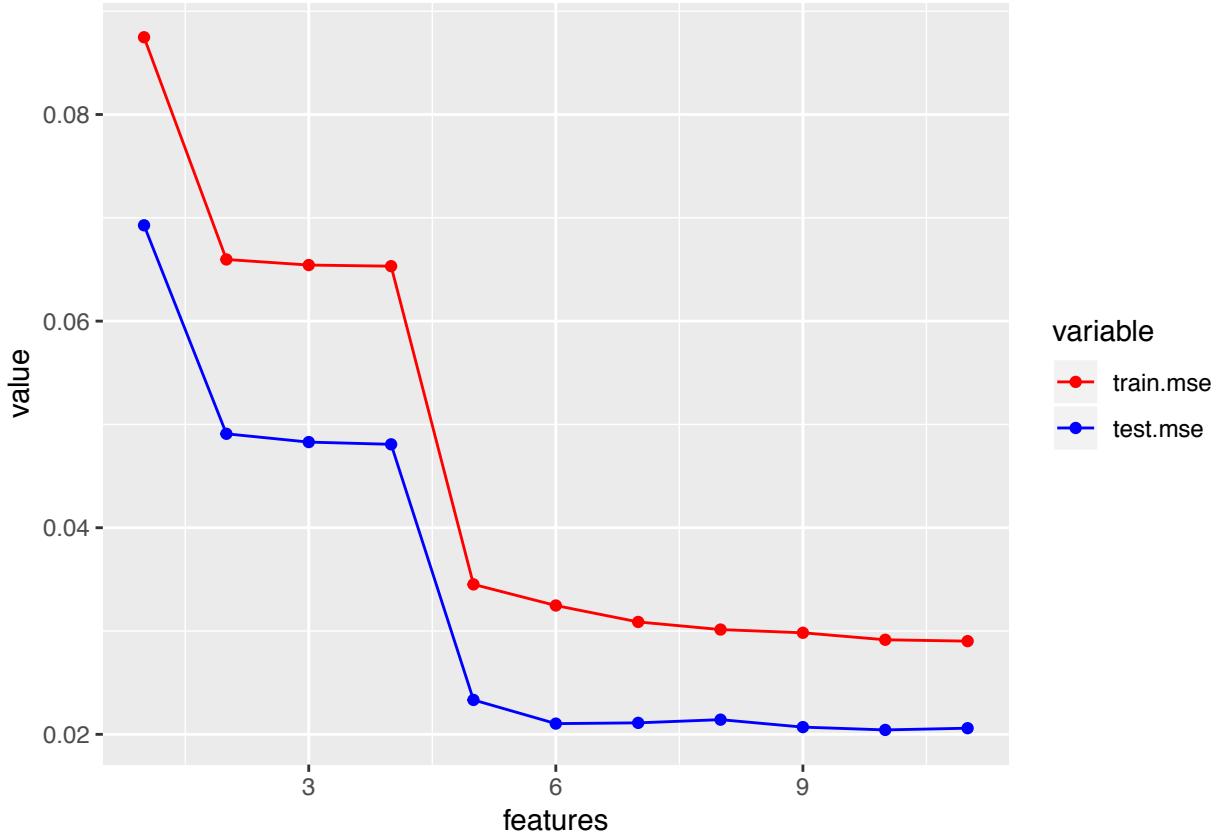
##          MSE         R2
## 1  0.08748009 0.4745942
## 2  0.06596899 0.6037900
## 3  0.06542675 0.6070466
## 4  0.06531829 0.6076980
## 5  0.03451592 0.7926972
## 6  0.03247054 0.8049818
## 7  0.03087920 0.8145394
## 8  0.03014235 0.8189649
## 9  0.02982983 0.8208419
## 10 0.02915284 0.8249079
## 11 0.02902088 0.8257005

```



#6.4

the training MSE is decreasing as the number of predictors increase



#6.5

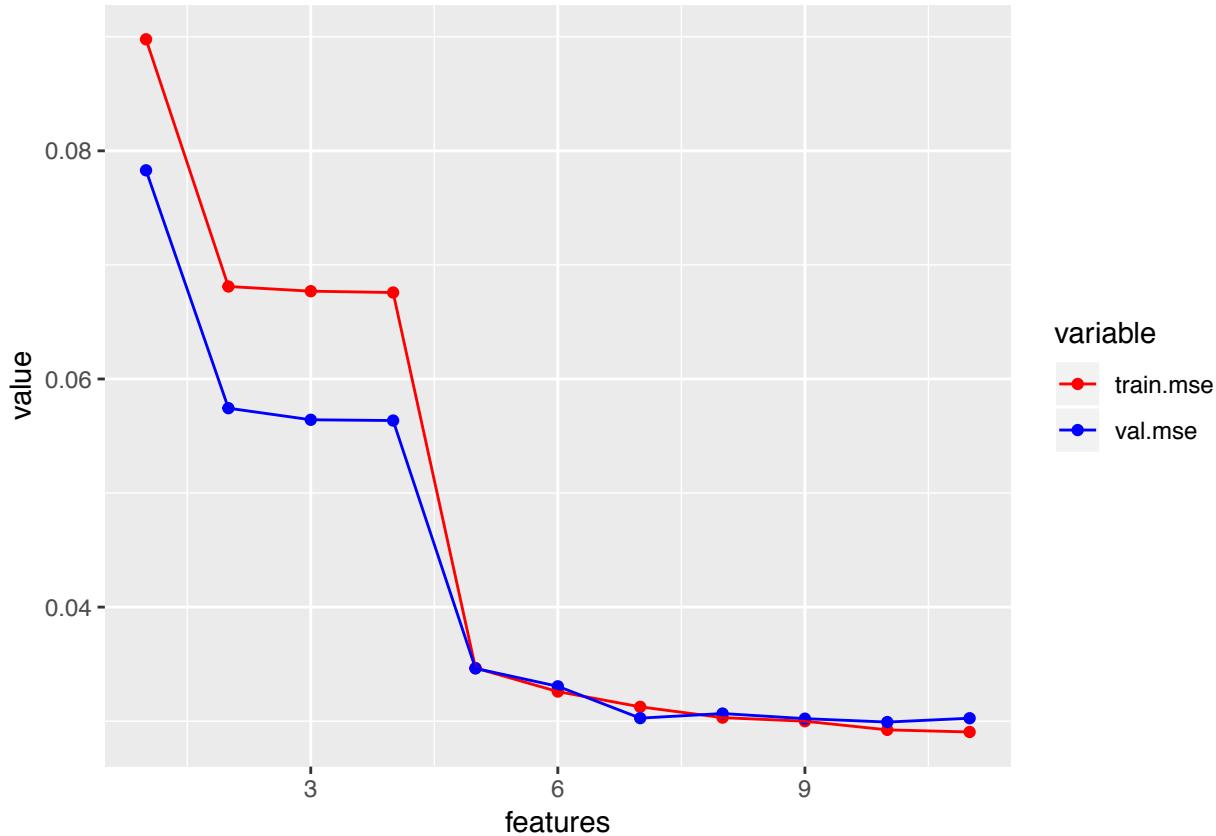
test MSE is decreasing in [1:10] and starts increasing at 11 and the 10th model gives the lowest test MSE

```
#6.2.2
```

```
modelQualitySingleVal = matrix(rep(0,22),nrow=11)
for (i in 1:length(Xnames)) {
  formula = as.formula(paste("log(SalePrice + 1) ~ ", paste(Xnames[1:i], collapse= "+")))
  lmod = lm(formula, data = AmesTinyActTrain)
  training.label = log(AmesTinyActTrain$SalePrice+1)
  mse = MSE(training.label, model.matrix(lmod), lmod$coefficients)
  modelQualitySingleVal[i,1] = mse

  val.label = log(AmesTinyActVal$SalePrice+1)
  modelQualitySingleVal[i, 2] = sum((val.label-predict(lmod, AmesTinyActVal))^2)/nrow(AmesTinyActVal)
}

dat <- data.frame(features = 1:length(Xnames), train.mse = modelQualitySingleVal[,1], val.mse = modelQualitySingleVal[,2])
dat.m <- melt(dat, id.vars = "features")
ggplot(dat.m, aes(features, value, colour = variable)) +
  geom_point() + geom_line(aes(features, value, colour = variable))+
  scale_colour_manual(values = c("red", "blue"))
```



```
which(modelQualitySingleVal[,2] == min(modelQualitySingleVal[,2]))
```

```
## [1] 10
```

the 10th model gives the lowest validation MSE

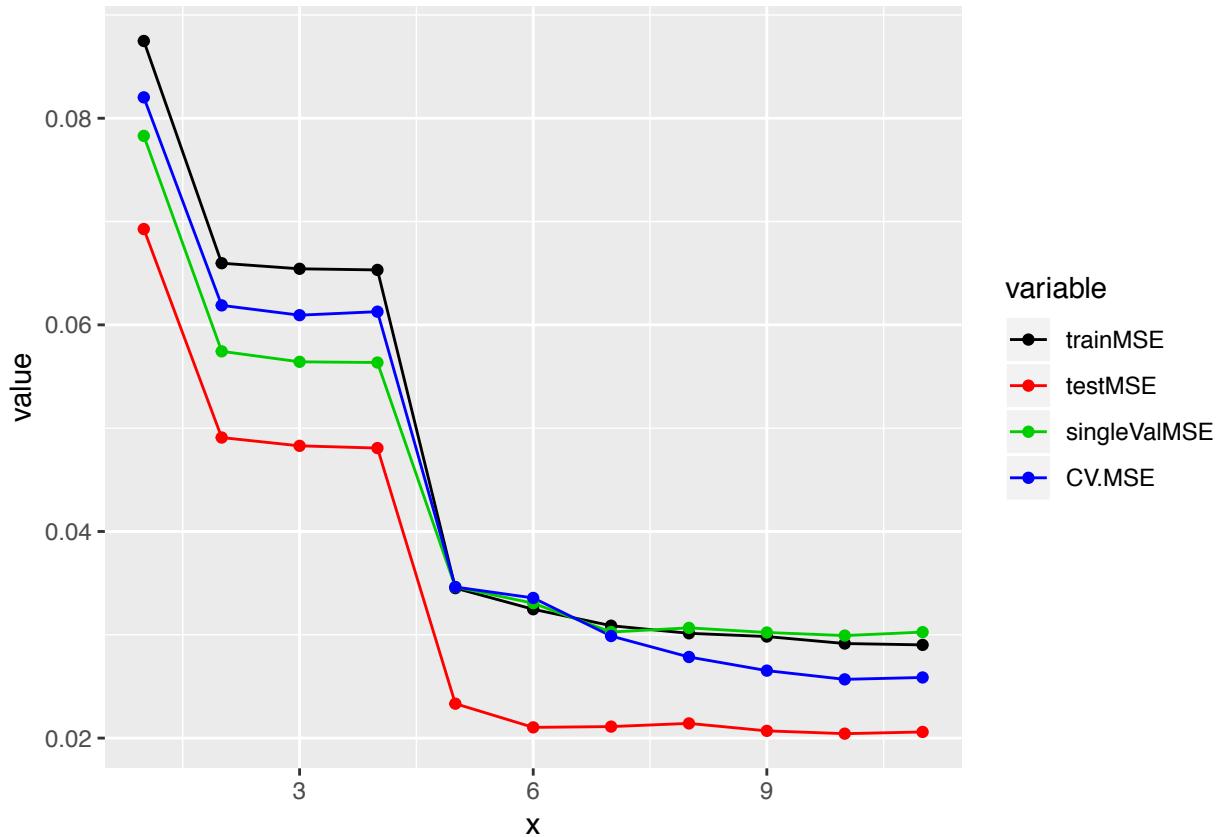
#6.3.1

```
set.seed(123)
folds <- createFolds(AmesTinyTrain$SalePrice, k = 5)
cvMSE = matrix(rep(0,55),nrow=11)

for (i in 1:length(Xnames)) {
  formula = as.formula(paste("log(SalePrice + 1) ~ ", paste(Xnames[1:i], collapse= "+")))
  for (f in 1:5) {
    train = AmesTinyTrain[-folds[[f]], ]
    lmod = lm(formula, data = train)
    test = AmesTinyTrain[folds[[1]], ]
    pred<- predict(lmod, test)
    true_y<- log(test$SalePrice + 1)
    mse1 = 1/length(folds[[1]]) * sum((pred-true_y)^2)
    cvMSE[i,f] = mse1
  }
}
```

```
#6.3.2
```

```
dat <- data.frame(x = 1:length(Xnames), trainMSE = modelQuality$MSE, testMSE = test.mse, singleValMSE =  
dat.m <- melt(dat, id.vars = "x")  
ggplot(dat.m, aes(x, value, colour = variable)) +  
  geom_line() + geom_point(aes(x, value, colour = variable)) +  
  scale_colour_manual(values = 1:5)
```



```
#10th model  
# which(dat$CV.MSE == min(dat$CV.MSE))
```

the 10th model gives the lowest CV-MSE

#6.4.1

```
train = AmesTiny[setdiff(names(AmesTiny), c("SalePrice"))]

x_train <- model.matrix(~ .-1, train)
train.label = log(AmesTiny$SalePrice + 1)
mod.ridge = glmnet(x_train, log(AmesTiny$SalePrice + 1), alpha = 0, lambda = 1)
```

```
#6.4.2
```

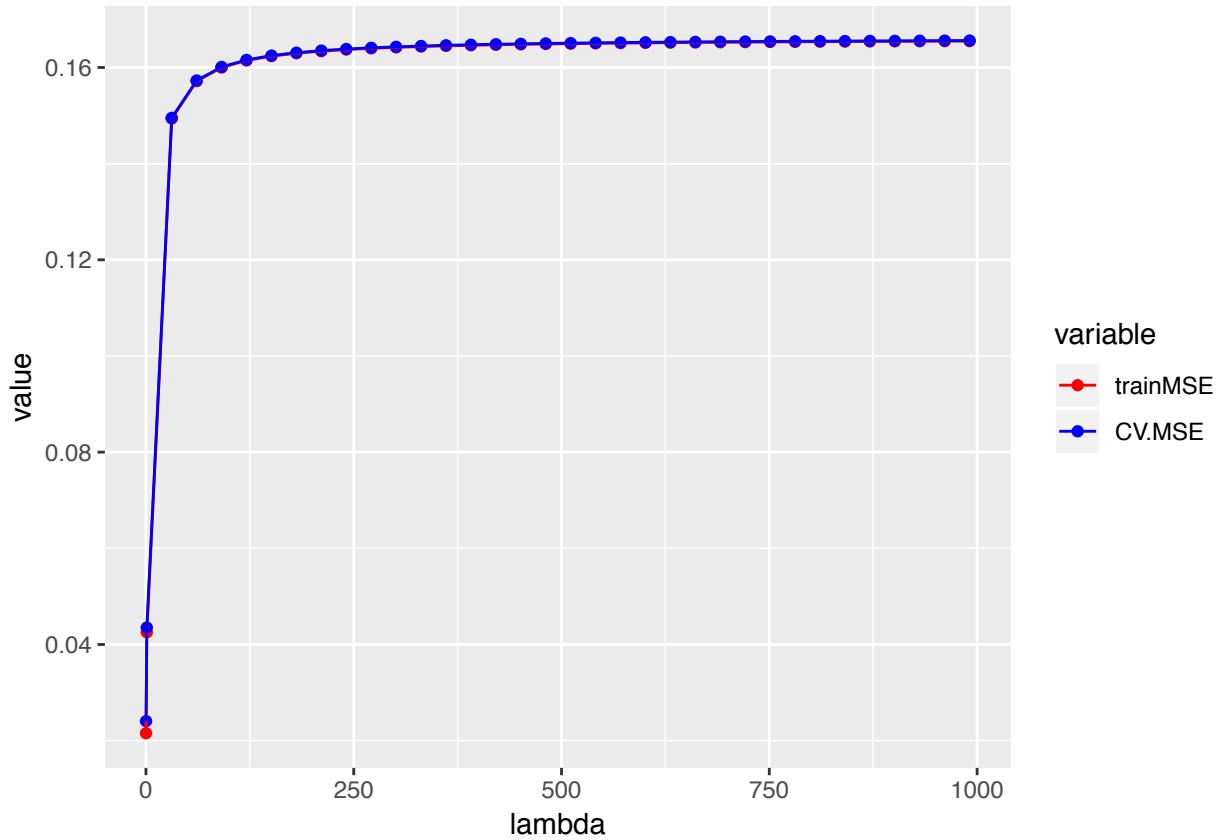
```
lambdas = c(0.1, seq(1,1000,30))

training.error = rep(0,length(lambdas))
test.error = rep(0,length(lambdas))

for (i in 1:length(lambdas)) {
  fit.ridge= glmnet(x_train, train.label, alpha=0, lambda = lambdas[i])
  training.error[i] = mean((train.label - predict(fit.ridge, x_train))^2)
}

cvmod = cv.glmnet(x_train, train.label, alpha=0, lambda = lambdas, type.measure = 'mse', nfolds = 5)

dat <- data.frame(lambda = lambdas, trainMSE = training.error, CV.MSE = rev(cvmod$cvm))
dat.m <- melt(dat, id.vars = "lambda")
ggplot(dat.m, aes(lambda, value, colour = variable)) +
  geom_point() + geom_line(aes(lambda, value, colour = variable))+
  scale_colour_manual(values = c("red", "blue"))
```



when $\lambda = 0.1$, we have both minimum training MSE and CV-MSE