

7.a.

$$P(Y|X) = \frac{P(Y, X)}{P(X)}$$

$$P(Y = Xw_1 + w_0 + z | X=x)$$

$$P(z | Y = Xw_1 + w_0 | X=x) = P(z | Y = Xw_1 + w_0)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (Y - Xw_1 - w_0)^2\right)$$

7.b. MLE

$$\text{Like}(w, w_0) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (Y_i - X_i w_1 - w_0)^2\right)$$

$$L = \log(\text{Like}(w, w_0)) = \text{Constant} - \frac{1}{2} \sum (Y_i - X_i w_1 - w_0)^2$$

$$\arg \min_{w, w_0} \frac{1}{2} \sum (Y_i - X_i w_1 - w_0)^2$$

$$\frac{\partial L}{\partial w_1} = \left\{ \begin{array}{l} \delta L / \delta w_1 = \frac{1}{2} \cdot 2 \sum (Y_i - X_i w_1 - w_0) \cdot (-x_i) \\ = \sum_{i=1}^n X_i (Y_i - X_i w_1 - w_0) \end{array} \right. \quad ①$$

$$\frac{\delta L}{\delta w_0} = - \sum_{i=1}^n (Y_i - X_i w_1 - w_0) \quad ②$$

Set each to zero

$$① \quad \sum x_i y_i + \sum x_i w_1 + \sum w_0 = 0$$

$$w_0 = \frac{1}{n} \sum (y_i - x_i w_1)$$

$$① - \sum_{i=1}^n (x_i y_i - x_i^2 w) + \sum x_i w_0 = 0$$

$$\sum x_i w_0 = \sum x_i y_i + x_i^2 w, \quad w_0 = \frac{\sum (x_i y_i - x_i^2 w)}{\sum x_i}$$

set each sum for w_0

$$\frac{1}{n} \sum (y_i - w \sum x_i) = \sum (x_i y_i) - w \sum x_i^2$$

$$\frac{1}{n} \sum x_i^2 (y_i - w \sum x_i) = \sum (x_i y_i) - w \sum x_i^2$$

$$\frac{1}{n} \sum x_i^2 y_i - \frac{w}{n} \sum x_i \sum x_i^2 = \sum (x_i y_i) - w \sum x_i^2$$

$$\frac{1}{n} \sum x_i^2 y_i - \sum (x_i y_i) = \frac{w}{n} \sum x_i \sum x_i^2 - w \sum x_i^2$$

$$= w \left(\frac{\sum x_i \sum x_i^2}{n} - \sum x_i^2 \right)$$

$$\left. \begin{aligned} w_1 &= \frac{1}{n} \sum x_i^2 y_i - \sum x_i y_i \\ &\frac{\sum x_i \sum x_i^2}{n} - \sum x_i^2 \end{aligned} \right\}$$

$$\left. \begin{aligned} w_0 &= \frac{1}{n} \sum y_i - \frac{1}{n} \sum x_i \cdot \frac{\frac{1}{n} \sum x_i^2 y_i - \sum x_i y_i}{\frac{\sum x_i \sum x_i^2}{n} - \sum x_i^2} \end{aligned} \right]$$

11

2.6.

$$z \sim \mathcal{U}[-0.5, 0.5]$$

$$\frac{1}{N}$$

-c

$$0.5 - (-0.5) = 1$$

□

$$P(y|x=x) = P(z=y-xw | x=x)$$

$$= P(z=y-xw) \quad P(y_1 \dots y_n | x_1 \dots x_n)$$

$$= \frac{1}{b-a} = \frac{1}{0.5+0.5} = 1 \quad z_i = y_i - x_i w$$

|z| \leq 0.5

$$f(z) = \begin{cases} 1 & -0.5 \leq z_i \leq 0.5 \\ 0 & \text{else} \end{cases} \quad \begin{cases} 1 & \text{if } z_i \in [-0.5, 0.5] \\ 0 & \text{otherwise} \end{cases}$$

7. d.

$$\arg \max_w \text{lik}(w) = \prod_{i=1}^n 1 \cdot I(z_i \in [-0.5, 0.5])$$

$$\text{lik} = 1 \text{ if } -0.5 \leq \max(|z_1|, \dots, |z_n|) \leq 0.5$$

$$\text{lik} = 0 \text{ if } \max(|z_1|, \dots, |z_n|) > 0.5$$

$$\prod_{i=1}^n 1 \{ -0.5 \leq y_i - x_i w \leq 0.5 \}$$

$$\arg \max_w \prod_{i=1}^n 1 \left\{ \frac{0.5 + y_i}{x_i} \leq w \leq \frac{y_i - 0.5}{x_i} \right\} \quad \begin{cases} 0 & i=1, \dots, n \\ 1 & x_i \neq 0 \end{cases}$$

• y_i must be in range $\in x_i = 0 \{ -0.5 \leq y_i \leq 0.5 \}$

7. D cont.

• Need one particular estimate for w

$$\frac{0.5 + y_i}{x_i} \leq w \leq \frac{0.5 + y_i}{x_i}$$

Want tightest bound for w as our particular estimate

(A) largest $\frac{0.5 + y_i}{x_i}$ from our Data

B Smallest $\frac{0.5 + y_i}{x_i}$ from our Data, such that
 $A < B$.

Our estimate should be within these tightest bounds.

2.F. One \rightarrow Ridge Regression $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

$$Y_i = \gamma_i w + z_i \quad z_i \sim N(0, 1) \text{ iid}$$

$$\left. \begin{array}{l} w \sim N(0, \sigma^2) \\ x_i \\ z_i \sim N(0, 1) \end{array} \right\} \text{Independent}$$

use Bayes

$$y_i | x_i, w \sim N(w^T x_i, \sigma^2)$$

$$P(w | \{(x_i, y_i), \dots, (x_n, y_n)\})$$

$$[w \sim N(0, \sigma^2) \Rightarrow \tilde{w} = \sigma \tilde{v} \quad w | \tilde{v} \sim N(0, 1)]$$

$$\text{Bayes. } \left[P(w | \text{data} = \{(x_i, y_i)\}_n) = \frac{P(\text{data} = D | w) P(w)}{P(\text{data} = D)} \right]$$

$$\left[\text{of } P(\text{data} = D | w) P(w) \right] \text{ (posterior)}$$

"posterior" is $P(w | x_1, y_1, x_2, y_2, \dots, x_n, y_n)$ & $P(w) \cdot P(x_1, y_1, \dots, x_n, y_n | w)$

proportional to "prior likelihood" $\left\{ \left(\prod_{i=1}^n P(y_i | x_i, w) P(w | \sigma) \right) \right\}$ Due to independence

$$= P(w) \cdot \prod_{i=1}^n P(y_i | x_i, w)$$

$$\left[P(y_i | x_i, w) = \frac{1}{(2\pi)^{1/2}} \exp \left\{ -\frac{1}{2} (y_i - x_i^T w)^2 \right\} \right]$$

$$= P(z_i = y_i - x_i^T w | x_i, w) =$$

$$\left[P(w | \sigma) = P(w) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left(-\frac{1}{2\sigma^2} (w)^2 \right) \right]$$

$$\star \text{ as } e^{(y_1^2 - 2y_1x_i w + x_i^2 w)} e^{(y_2^2 - 2y_2x_i w + x_i^2 w)} \dots = e^{2y_1^2} e^{-2x_i^2 w} e^{2x_i^2 w}$$

2.F (cont.) :

$$\prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y_i - x_i w)^2\right) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{w^2}{2\sigma^2}\right) \right]$$

$$= \frac{1}{(2\pi)^{n/2}} \left[\frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \prod_{i=1}^n \left[\exp\left(-\frac{1}{2}\left((y_i - x_i w)^2 - \frac{w^2}{\sigma^2}\right)\right) \right] \right]$$

$$\exp\left[\frac{-w^2}{2\sigma^2}\right] \exp\left[-\frac{1}{2}\left(\sum_{i=1}^n y_i^2 - 2w \sum_{i=1}^n x_i y_i + w^2 \sum_{i=1}^n x_i^2 - \frac{w^2}{\sigma^2}\right)\right]$$

$$\exp\left[-\frac{1}{2}\left(\sum_{i=1}^n y_i^2 - 2w \sum_{i=1}^n x_i y_i + w^2 \sum_{i=1}^n x_i^2 - \frac{w^2}{\sigma^2}\right)\right]$$

$$-\frac{1}{2\sigma^2} \left(\sigma^2 \sum_{i=1}^n y_i^2 - 2w \sigma^2 \sum_{i=1}^n x_i y_i + w^2 \sigma^2 \sum_{i=1}^n x_i^2 - w^2 \right)$$

$$-\frac{1}{2\sigma^2} \left(w^2 (\sigma^2 \sum_{i=1}^n x_i^2 - 1) + w (-2\sigma^2 \sum_{i=1}^n x_i y_i) + \sigma^2 \sum_{i=1}^n y_i^2 \right)$$

$$\underline{a(x+\delta)^2 + e} \quad d = \frac{b}{2a} = \frac{-2\sigma^2 \sum_{i=1}^n x_i y_i}{2\sigma^2 \sum_{i=1}^n x_i^2 - 2}$$

$$e = c - \frac{b^2}{4a} = \sigma^2 \sum_{i=1}^n y_i^2 - \frac{(2\sigma^2 \sum_{i=1}^n x_i y_i)^2}{4\sigma^2 \sum_{i=1}^n x_i^2 - 4}$$

$$\frac{-1}{2\sigma^2} \left[\left(\sigma^2 \sum_{i=1}^n x_i^2 - 1 \right) \left(w + \frac{-2\sigma^2 \sum_{i=1}^n x_i y_i}{2\sigma^2 \sum_{i=1}^n x_i^2 - 2} \right)^2 + \sigma^2 \sum_{i=1}^n y_i^2 - \frac{(2\sigma^2 \sum_{i=1}^n x_i y_i)^2}{(4\sigma^2 \sum_{i=1}^n x_i^2 - 4)} \right]$$

$$\left[\frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{1}{2\sigma^2} \sum_i (y_i - \mu)^2\right) \cdot \exp\left[-\frac{1}{2\sigma^2} \left(\sum_i x_i^2 - 1\right) \left(w - \frac{\sum_i x_i y_i}{\sum_i x_i^2 - 1}\right)^2\right] \right]$$

$$\boxed{\mu = \frac{\sigma^2 \sum x_i y_i}{(\sigma^2 \sum x_i^2 - 1)}}$$

Gaussian

$$N(\mu, \sigma^2)$$

e.g. $y_i = w^T x_i + \epsilon_i$ $y_i \in \mathbb{R}$, $w, x_i \in \mathbb{R}^d$
 $\epsilon_i \sim N(0, 1)$
 w fixed

MLE gives least squares

$$\begin{aligned}
 \text{Lik}(w) &= P(y_i | w, x_i) = P(\epsilon_i = y_i - w^T x_i | w, x_i) \\
 &= P(\epsilon_1 | \epsilon_2, \dots, \epsilon_n, w, x_i) \cdot P(\epsilon_2 | \epsilon_1, \dots, \epsilon_n, w, x_i) \cdots \\
 &= \prod_{i=1}^n P(\epsilon_i | x_i, w) \\
 &= \prod_{i=1}^n P(\epsilon_i = y_i - w^T x_i | x_i, w) \\
 &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (y_i - w^T x_i)^2\right),
 \end{aligned}$$

$$\arg \max_w \log(\text{lik}(w)) = \text{const} + \sum_{i=1}^n -\frac{1}{2} (y_i - w^T x_i)^2$$

$$\begin{aligned}
 \text{arg min}_w &= \hat{\sum}_{i=1}^n (y_i - w^T x_i)^2 \\
 &\quad \left[\begin{array}{l} \text{"Has same objective function as w/ OLS"} \\ \therefore \text{it is clearly also the OLS estimator } \hat{w}. \end{array} \right]
 \end{aligned}$$

currently: No Bias

2.6 Multi-Dimensional Ridge Regression

\vec{w} random vector

$$Y_i = \vec{w}^T x_i + \varepsilon_i$$

$$Y_i \in \mathbb{R}, \vec{w}, x_i \in \mathbb{R}^d$$

$$\varepsilon_i \sim N(0, 1)$$

$$w_j \sim N(0, \sigma^2) P(\vec{w}) = N(0, I\sigma^2)$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_d^T \end{bmatrix} \quad Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_d \end{bmatrix}$$

$$\Rightarrow X^T X = \sum x_i x_i^T$$

$$X^T Y = \sum x_i Y_i$$

observe $P(\vec{w}_j | (x_i, Y_i)_{i=1, \dots, d}) \propto P((x_i, Y_i)_{i=1, \dots, d} | w_j) P(w_j)$
but

$$w^T x_i = x_i^T w \Rightarrow P(Y_i | x_i, w_j) P(w_j) \quad \text{"by independence"}$$

$$= P(\varepsilon_i = Y_i - w^T x_i | x_i, w_j) P(w_j)$$

$$= \prod_{i=1}^d P(w_j) P(\varepsilon_i = Y_i - w^T x_i)$$

$$= \prod_{i=1}^d \left[\frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{1}{2} \frac{w_j^2}{\sigma^2}\right\} \cdot \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} (Y_i - w^T x_i)^2\right\} \right]$$

$$= \frac{1}{\sqrt{2\pi}\sigma^2} \prod_{i=1}^d \exp\left\{-\frac{1}{2} \left(\frac{w_j^2}{\sigma^2} + (Y_i - w^T x_i)^2 \right)\right\}$$

$$\exp\left\{-\frac{1}{2} \left(\frac{w_j^2}{\sigma^2} + Y_1^2 - 2w^T x_1 Y_1 + w^T x_1^2 \right)\right\} \left(\frac{w_2^2}{\sigma^2} + Y_2^2 - 2w^T x_2 Y_2 + w^T x_2^2 \right) \cdots$$

2. H cont.

$$\{w_1, w_2, \dots, w_d\} \quad [d \times 2]$$

$$\exp \left\{ -\frac{1}{2} \left(\sum y_i^2 - 2 \sum w^T x_i y_i + \sum w^T x_i^2 + \frac{1}{\sigma^2} \sum w_i^2 \right) \right\}$$

$$\exp \left\{ -\frac{1}{2} \left(w^T \left(\sum x_i^2 \right) + w^T \left(-2 \sum x_i y_i \right) + \sum y_i^2 + \frac{1}{\sigma^2} \sum w_i^2 \right) \right\}$$

$$d = \frac{b}{2a}, e = c = \frac{b^2}{4a}$$
$$a(x+d)^2 + e$$

$$d = \frac{-2 \sum x_i y_i}{2 \sum x_i^2}$$
$$e = (\sum y_i^2 + \frac{1}{\sigma^2} \sum w_i^2) = \frac{4 \sum x_i^2 y_i^2}{4 \sum x_i^2}$$

$$\exp \left\{ -\frac{1}{2} \left(\sum x_i^2 \left(w^T - \frac{\sum x_i y_i}{\sum x_i^2} \right)^2 + \sum y_i^2 + \frac{1}{\sigma^2} \sum w_i^2 - \frac{\sum x_i^2 y_i^2}{\sum x_i^2} \right) \right\}$$

$$M = \frac{\sum y_i y_i}{\sum y_i^2} = \frac{\mathbf{x}^T \mathbf{y}}{\sum x_i^2}$$

3. a.

$$\text{bis } E(\hat{x} - \mu) \quad \text{vom } \text{Var}[\hat{x}]$$

$$\text{i)} \quad \hat{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \mu$$

$$E(\hat{x} - \mu) = E(\hat{x}) - E(\mu) \Rightarrow E(\hat{x}) - \mu = 0 \\ \frac{\mu - \mu}{n} = 0$$

$$\text{ii)} \quad \hat{x} = \frac{x_1 + \dots + x_n}{n+1}$$

$$E\left(\frac{x_1 + \dots + x_n}{n+1}\right) = \frac{1}{n+1}(E(x_1) + \dots + E(x_n))$$

$$\frac{\Delta \cdot M}{n+1}$$

$$\frac{\Delta \cdot M}{n+1} - \mu = M\left(\frac{\Delta}{n+1} - 1\right)$$

$$\text{iii)} \quad \hat{x} = \frac{x_1 + \dots + x_n}{n+n_0} \quad M\left(\frac{\Delta}{n+n_0} - 1\right)$$

$$\text{iv)} \quad \hat{x} = 0 \quad 0 = E(\mu) = -\mu$$

$$\text{3.b. i)} \quad \text{Var}(\hat{x}) = \text{Var}\left(\frac{1}{n}(x_1 + \dots + x_n)\right) = \frac{1}{n} \text{Var}(x_i) = \frac{\sigma^2}{n}$$

$$\text{ii)} \quad \text{Var}\left(\frac{1}{n+1}(x_1 + \dots + x_n)\right) = \frac{1}{(n+1)^2} n \text{Var}(x_i) = \frac{n \sigma^2}{(n+1)^2}$$

$$\text{iii)} \quad \text{Var}\left(\frac{1}{n+n_0}(x_1 + \dots + x_n)\right) = \frac{1}{(n+n_0)^2} n \sigma^2 = \frac{n \sigma^2}{(n+n_0)^2}$$

$$\text{iv)} \quad \text{Var}(\lambda) = E(\lambda^2) - E(\lambda)^2 = 0$$

x' is an iid copy of x

$$3.C. E[(\hat{x} - x')^2]$$

$$\text{Var}(x) = E(x^2) - E(x)^2$$

$$E(\hat{x}^2 - 2\hat{x}x' + x'^2)$$

$$E(\hat{x}^2) - 2E(\hat{x}x') + E(x'^2)$$

$$\sigma_{\hat{x}}^2 + \hat{\mu}^2 - 2E(\hat{x})E(x') + \sigma_x^2 + \mu^2$$

$$\sigma_{\hat{x}}^2 + \hat{\mu}^2 - 2\hat{\mu}\mu + \sigma_x^2 + \mu^2 = \hat{\sigma}^2 + (\hat{\mu} - \mu)^2$$

$$E[(\hat{x} - \mu)^2] = E[\hat{x}^2 - 2\hat{x}\mu + \mu^2]$$

$$E(\hat{x}^2) = 2E(\hat{x})E(\mu) + E(\mu^2)$$

$$= \sigma_{\hat{x}}^2 + \hat{\mu}^2 - 2E(\hat{x})E(\mu) + \mu^2$$

$$= \sigma_{\hat{x}}^2 + \hat{\mu}^2 - 2\hat{\mu}\mu + \mu^2 = \hat{\sigma}^2 + (\hat{\mu} - \mu)^2$$

- Our second derivation does not have variance. The x' introduces variance to the system.

$$\underline{E(\hat{\mu} - \mu) = 0}$$

$$E[(\hat{\mu} - \mu)^2] = \sigma_{\hat{\mu}}^2 + (\hat{\mu} - \mu)^2 \quad \frac{\Delta - 1}{n+1} \frac{n-1}{n+1}$$

3.e. i) $\left(\frac{\sigma^2}{n} + 0^2 \right) \quad \frac{\Delta - 1}{n+1} = \frac{1 - \frac{1}{n+1}}{n+1} = \frac{1}{n+1}$

ii) $\left(\frac{\Delta \sigma^2}{(n+1)^2} + \left(\frac{-\mu}{n+1} \right)^2 \right) \quad \frac{1+\frac{1}{n+1}}{n+1}$

iii) $\left(\frac{\Delta \sigma^2}{(n+n_0)^2} + \left(\frac{-\mu n_0}{n+n_0} \right)^2 \right)$

iv) $0 + (-\mu)^2 = \mu^2$

3.e $\hat{\mu} = \frac{x_1 + \dots + x_n}{n+n_0}$ to samples of $n+n_0$
"dummies"

1) has $n_0 = 0$

2) has $n_0 = 1$

3) has $n_0 \neq n_0$

4) has $n_0 = \infty$

3.f. Bias as n_0 increase: our bias approaches $-\mu$
 $F(x) \notin \{0, 1\}$ the model gets worse approach zero
 Var as n_0 increase: our variance approaches 0

3.g. $\lambda_0 = \alpha\lambda$ Setting λ to minimize

$$\text{Min } E[(\lambda - \mu)^2] = \frac{\lambda\sigma^2}{(\lambda + \alpha\lambda)^2} + \left(\frac{\mu - \lambda\alpha}{\lambda + \alpha\lambda} \right)^2$$

$$\text{Min} \left[\frac{\lambda\sigma^2}{(\lambda + \alpha\lambda)^2} + \left(\frac{\mu - \lambda\alpha}{\lambda + \alpha\lambda} \right)^2 \right]$$
$$\mu^2\lambda^2\sigma^2(\lambda + \alpha\lambda)^{-2}$$

differentiation

$$\frac{d}{d\lambda} = \lambda\sigma^2 \left[(-2)(\lambda + \alpha\lambda)^{-3} \cdot 1 \right]$$

$$+ \mu^2\lambda^2 \left[2\alpha(\lambda + \alpha\lambda)^{-2} + \alpha^2(-2)(\lambda + \alpha\lambda)^{-3} \cdot 1 \right]$$

set to zero

$$2 \left(\frac{\alpha\lambda\mu^2}{(\lambda + \alpha\lambda)^3} - \frac{\sigma^2}{\lambda} \right) = 0$$

solve for λ

$$\frac{\alpha\lambda\mu^2}{(\lambda + \alpha\lambda)^3} - \frac{\sigma^2}{\lambda} = 0$$

$$\frac{\alpha\lambda\mu^2}{(\lambda + \alpha\lambda)^3} = \frac{\sigma^2}{\lambda}$$

$$\frac{\alpha\lambda\mu^2}{(\lambda + \alpha\lambda)^3} \cdot (\lambda + \alpha\lambda)^3 = \frac{\sigma^2}{\lambda}, \quad \lambda = \frac{\sigma^2}{\alpha^2\mu^2}$$

3.h

$$\lambda = \frac{\sigma^2}{\lambda^2 \mu^2}$$

μ small
 $(\lambda \text{ approaches } \infty)$ fall fast
as μ is small $\geq \sigma^2$
approaches a large #

3.i

$$\boxed{\lambda' = \lambda - \mu_0}$$

s.t.

$$E(\lambda') \text{ small } \approx 0$$

$E(\lambda) - E(\mu_0) \approx 0$ implies
it is close to μ_0
 $\lambda - \mu_0 \approx 0$

Bias ↙

$$\begin{aligned} \text{Variance} \rightarrow \text{Var}(\lambda') &= \text{Var}(\lambda) - \text{Var}(\mu_0) \\ &\quad 0 \\ &= \text{Var}(\lambda) \end{aligned}$$

3.j relation λ & λ'

- The regularization parameter reduces variance & also simultaneously increases bias.
- The greater our α the less impact our data will have on our GLS fit's.
- λ seems to do the same, regularizing our data
- Reg. most effective when we can't fit our data.
- Suggests that large values of λ can reduce the variance as our prev. q showed

4.a.

Show $E[\hat{\omega}] = \omega^*$

$$\hat{\omega} = \arg \min_{\omega} \|y - X\omega\|_2^2$$

$$\hat{\omega} = (X^T X)^{-1} X^T y$$

$$y = y^* + z$$

$$\hat{\omega} = (X^T X)^{-1} X^T (y^* + z)$$

$$E[(X^T X)^{-1} X^T y^*] + E[(X^T X)^{-1} X^T z]$$

$$E[(X^T X)^{-1} X^T y^*] + 0 = \omega^* = (X^T X)^{-1} X^T y^*$$

Show that

(A) $E(\|y^* - X\hat{\omega}\|^2) = y^{*\top} y^* - y^{*\top} X \omega^* - \hat{\omega}^T X^T y^* + \hat{\omega}^T X^T X \omega^*$

$$E[(y^* - X\hat{\omega})^T (y^* - X\hat{\omega})] =$$

$$E((y^{*\top} - \hat{\omega}^T X^T)(y^* - X\hat{\omega})) = E(y^{*\top} y^* - y^{*\top} X \omega^* - \hat{\omega}^T X^T y^* + \hat{\omega}^T X^T X \hat{\omega})$$

(B) $\|y^* - E(X\hat{\omega})\|_2^2 = (y^* - E(X\hat{\omega}))^T (y^* - E(X\hat{\omega}))$

$$(y^{*\top} - \hat{\omega}^{*\top} X^T)(y^* - X\hat{\omega}^*) = \left[y^{*\top} y^* - y^{*\top} X \omega^* - \hat{\omega}^{*\top} X^T y^* + \hat{\omega}^{*\top} X^T X \hat{\omega} \right]$$

4.a. (cont.)

$$\textcircled{C} \rightarrow E\left(\|X\hat{\omega} - E(X\hat{\omega})\|\right)$$

$$= E[(X\hat{\omega} - E(X\hat{\omega}))^T (X\hat{\omega} - E(X\hat{\omega}))]$$

$$= E[(\hat{\omega}^T X^T - E(\hat{\omega}^T))^T (X\hat{\omega} - E(X\hat{\omega}))]$$

$$= E[\hat{\omega}^T X^T X \hat{\omega} - \hat{\omega}^T X^T E(X\hat{\omega}) - E(\hat{\omega}^T)^T X \hat{\omega} + E(\hat{\omega}^T)^T E(X\hat{\omega})]$$

$$= E(\hat{\omega}^T X^T X \hat{\omega}) -$$

$$w^T X^T X w^* = w^{*T} X^T X w^* = w^T X^T X w^* + w^{*T} X^T X w^*$$

$$= 0$$

$$A = B + C$$

$$C = 0$$

∴

$A = B$ which is sum to be true
on the previous page.

$$4.6. \quad v \sim N(0, \Sigma) \quad \sum \in \mathbb{R}^{d \times d}, A \in \mathbb{R}^{k \times d}$$

$$Av \sim N(0, A\Sigma A^T)$$

then $\hat{\omega} = \mathcal{N}(\omega^*, \sigma^2 (x^T x)^{-1})$

$$z \sim N(0, \sigma^2) \quad \hat{\omega} = (x^T x)^{-1} x^T (y^* + z)$$

$$\hat{\omega} = (x^T x)^{-1} x^T y^* + (x^T x)^{-1} x^T z$$

$$\hat{\omega} = \omega^* + (x^T x)^{-1} x^T z$$

$$x^T z \sim N(0, x^T \sigma^2 x)$$

$$\omega^* + x^T z \sim N(\omega^*, x^T \sigma^2 x)$$

$$\omega^* + (x^T x)^{-1} x^T z \sim N(\omega^*, (x^T x)^{-1} x^T \sigma^2 x (x^T x)^{-1})$$

$$\therefore (x^T x)^{-1} x^T \sigma^2 x (x^T x)^{-1} = \sigma^2 (x^T x)^{-1}$$

$$x^{-1} x^T \sigma^2 x (x^{-1} x^T)^{-1}$$

$$x^{-1} \underbrace{x^T \sigma^2 x}_{I} (x^T)^{-1} \underbrace{x^{-1}}_{I} = \sigma^2 (x^T x)^{-1}$$

$$x^{-1} \sigma^2 x^{-1} = \sigma^2 x^{-1} x^T$$

σ^2 is a scalar thus the above expression is true equivalent. \therefore
 $N \sim (\omega^*, \sigma^2 (x^T x)^{-1})$

$$\text{Var}(\hat{\omega}) = (\sigma^2)(\mathbf{x}^\top \mathbf{x})^{-1} \quad \frac{1}{n} \mathbb{E}(\|\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*\|^2_2) = \underline{\frac{\sigma^2}{n}}$$

4.c. If $\beta \sim N(\mu, \Sigma)$

Math $\Rightarrow \mathbb{E}(\|\beta\|_2^2) = \|\mu\|_2^2 + \text{tr}(\Sigma)$

Statistical \therefore But we find $0 = \mathbb{E}(\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*)$,

$$= \mathbf{x}\mathbb{E}(\hat{\omega}) - \mathbf{x}\omega^* = \mathbf{x}\omega^* - \mathbf{x}\omega^* = 0$$

$$\text{so } \mathbb{E}(\|\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*\|_2^2) = \|0\|_2^2 + \text{tr}(\Sigma) = \text{tr}(\Sigma)$$

where ϵ is variance of $\underbrace{\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*}_{\mathbf{x}\omega^* \text{ is fixed}}$

$$\text{Var}(\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*) = \text{Var}(\mathbf{x}\hat{\omega}) \quad \begin{matrix} \text{"xw* is fixed"} \\ \mathbf{x}^\top \mathbf{x} \sigma^2 (\mathbf{x}^\top \mathbf{x})^{-1} \\ \sigma^2 \mathbf{x}^\top \mathbf{x} \end{matrix}$$

$$\text{Var}(\mathbf{x}\hat{\omega}) = \mathbf{x}^\top \mathbf{x} \text{Var}(\hat{\omega}) = \mathbf{x}^\top \mathbf{x} \sigma^2 (\mathbf{x}^\top \mathbf{x})^{-1}$$

$$\text{trace}(\mathbf{x}^\top \mathbf{x} \sigma^2 (\mathbf{x}^\top \mathbf{x})^{-1}) = \sigma^2 \text{trace}(\mathbf{x}^\top \mathbf{x} (\mathbf{x}^\top \mathbf{x})^{-1})$$

$$= \sigma^2 \cdot \text{tr}(\mathbf{I}) = \sigma^2 \cdot d$$

$$\frac{1}{n} \mathbb{E}(\|\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*\|_2^2) = \frac{1}{n} (\|0\|_2^2 + \text{tr}(\text{Var}(\mathbf{x}\hat{\omega} - \mathbf{x}\omega^*)))$$

$$= \boxed{\frac{\sigma^2 d}{n}}$$

$$\{\alpha_i\}_{i=1}^n$$

$$(1) \quad \underline{y = y^* + \varepsilon}$$

$$u.d. \quad \{\alpha_i, y_i\}_{i=1}^n$$

$$y_i = w_i \alpha_i + w_0 + \varepsilon_i$$

χ w/ D+2 polynomial Features $P_D(x_i)$

$$= (1, \alpha_1, \dots, \alpha_D)$$

2.2

$$w^* = (\chi^T \chi)^{-1} \chi^T y^*$$

$$y_i = \underbrace{w_i \alpha_i + w_0 + \varepsilon_i}_{{= y^*}} \quad y_i = \underbrace{y^* + \varepsilon_i}$$

$$y_i^* = w_i \alpha_i + w_0$$

$$\hat{w} = (\chi^T \chi)^{-1} \chi^T y = (\chi^T \chi)^{-1} \chi^T (w_i \alpha_i + w_0)$$

$$= (\chi^T \chi)^{-1} \chi^T$$

bias

$$\|y^* - \chi w^*\|_2^2 = 0$$

$$\text{b/c } \chi w^* = y^*$$

$$w^* = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix} \quad \text{a } (D+2) \times 1 \text{ vector}$$

avg expected prediction squared error

$$\frac{1}{n} E \left(\| \chi \hat{w} - \chi w^* \|_2^2 \right) \stackrel{\text{using u.c.}}{=} \frac{\sigma^2(D+1)}{n} \leq \epsilon$$

$$\lambda \geq \frac{\sigma^2(D+1)}{\epsilon}$$

HW03 - CS 189

1.

Who else did you work with on this homework? In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

I worked on this homework with Ehimare Okoyomon, Prashanth Ganeth, and Daniel Mockaitis. We worked by getting together throughout the week and communicating on facebook.

I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up}

Nicholas Lorio, 26089160

Question 6.

What time is it?

9:57PM

Sources

<https://ocw.mit.edu/courses/mathematics/18-443-statistics-for-applications-fall-2006/lecture-notes/lecture2.pdf>

```
In [167]: import numpy as np
import matplotlib.pyplot as plt
import pickle
```

HW03

Question 2

2.E

```
In [55]: import numpy as np
import matplotlib.pyplot as plt

sample_size = [5,25,125,625]
plt.figure(figsize=[12, 10])
low_bound = -0.5
high_bound = 0.5
N = 10001
W = np.linspace(40, 60, num=N)
w_true = 50
print(W)

for k in range(4):
    n = sample_size[k]

    # generate data
    # np.linspace, np.random.normal and np.random.uniform might be useful functions
    Xs = np.linspace(1, 100, num=n)
    Ys = np.array([x*w_true + np.random.uniform(low_bound, high_bound)
for x in Xs])

    likelihood = np.ones(N) # likelihood as a function of w

    for i in range(N):
        w_i = W[i]
        in_bound = True
        for j in range(n):
            y_j = Ys[j]
            x_j = Xs[j]

            if w_i > (y_j + 0.5)/x_j or w_i < (y_j - 0.5)/x_j:
                in_bound = False
                break
        if in_bound:
            likelihood[i] = 1
            print("Likelihood is 1 for w equal to " + str(w_i) + " in sample" + str(n))
        else:
            likelihood[i] = 0

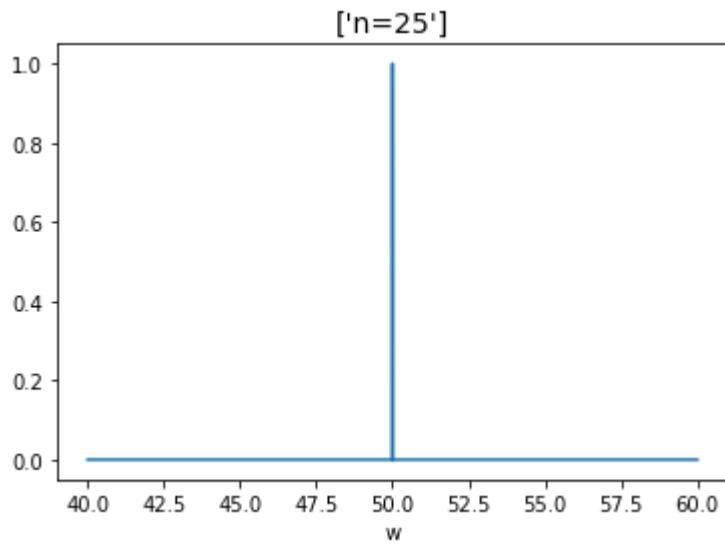
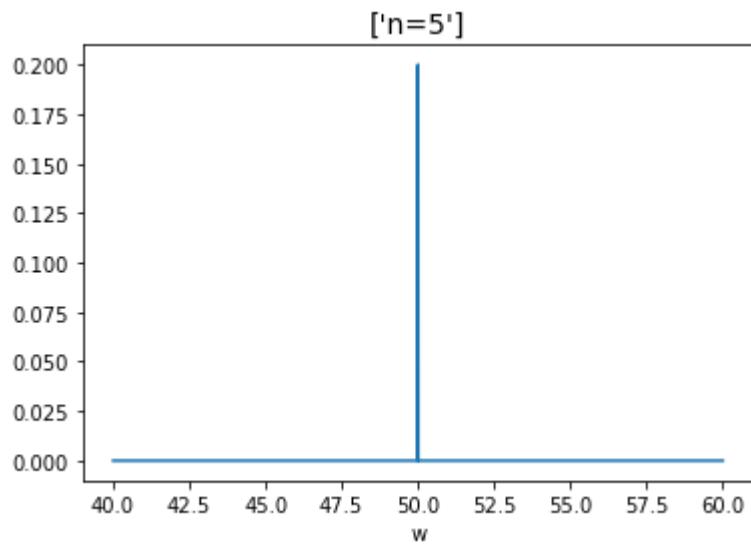
        # compute likelihood
        #print(sum(likelihood))
        likelihood /= sum(likelihood) # normalize the likelihood

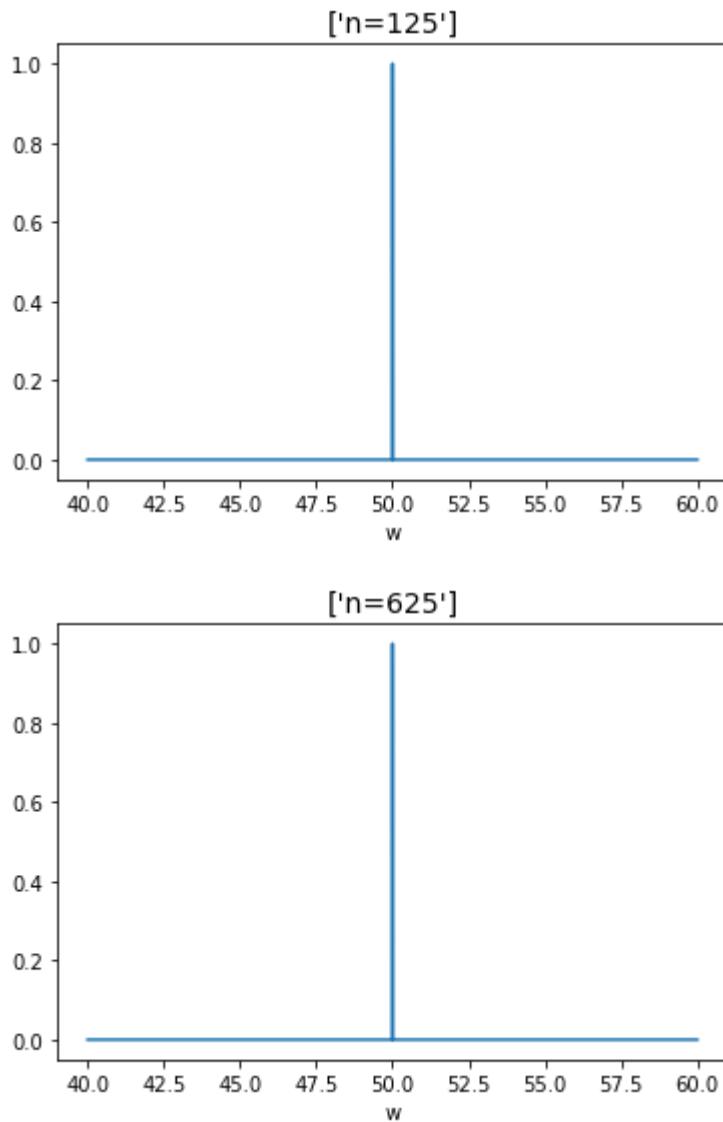
    plt.figure()
    # plotting likelihood for different n
    plt.plot(W, likelihood)
    plt.xlabel('w', fontsize=10)
    plt.title(['n=' + str(n)], fontsize=14)

plt.show()
```

```
[40.      40.002 40.004 ... 59.996 59.998 60.      ]  
Likelihood is 1 for w equal to 49.994 in sample5  
Likelihood is 1 for w equal to 49.996 in sample5  
Likelihood is 1 for w equal to 49.998 in sample5  
Likelihood is 1 for w equal to 50.0 in sample5  
Likelihood is 1 for w equal to 50.002 in sample5  
Likelihood is 1 for w equal to 50.0 in sample25  
Likelihood is 1 for w equal to 50.0 in sample125  
Likelihood is 1 for w equal to 50.0 in sample625
```

```
<matplotlib.figure.Figure at 0x7f7c3fea72e8>
```





As n gets large, the amount of w estimates that accurately fit within our bounds decreases. The MLE of the uniform distribution can be either 1 or zero. As n increases our estimation for the model becomes less general and the variance of our MLE parameter decreases.

2.I

```
In [174]: import numpy as np
import matplotlib.pyplot as plt

sample_size = [5,25,125]
w0_true = 20
w1_true = 50
N = 10001
w0 = np.linspace(10,30, num=N)
w1 = np.linspace(40, 60, num=N)
Xs = np.random.rand(n,2)
Ys = np.array([x[0]*w0_true + x[1]*w1_true + np.random.normal(0, 1) for
   x in Xs])

for k in range(4):
    n = sample_size[k]

    # generate data
    # np.linspace, np.random.normal and np.random.uniform might be useful
    # functions

# compute likelihood

N = 1001
# W0s =
# W1s =
likelihood = np.ones([N,N]) # likelihood as a function of w_1 and w_0

for i1 in range(N):
    # w_1 = W1s[i1]
    for i2 in range(N):
        # w_2 = W2s[i2]
        for i in range(n):
            # compute the likelihood here

# plotting the likelihood
plt.figure()
# for 2D likelihood using imshow
plt.imshow(likelihood, cmap='hot', aspect='auto', extent=[0,4,0,4])
plt.xlabel('w0')
plt.ylabel('w1')
plt.show()
print(n)
```

File "<ipython-input-174-edcd92a9c715>", line 38
 plt.figure()

IndentationError: expected an indented block

Question 4

4.E

```
In [173]: # assign problem parameters
w0=1
w1=1
N = [10, 100, 1000]
# generate data
# np.random might be useful
alpha = []
z = []
D = 10
for i in range(D):
    alpha.append(np.random.uniform(-1, 1))
    z.append(np.random.normal(0, 1))

y = np.array(alpha)*w1 + w0 + z

err = np.zeros(D)
for d in range(D):
    Xs = np.array([[x**i for i in range(d)] for x in alpha])
    w = np.linalg.lstsq(Xs, y)[0]
    err[d] = np.sum(np.square(Xs.dot(w) - y)) / len(y)

plt.figure()
# plotting likelihood for different n
plt.plot(D, err)
plt.xlabel('D', fontsize=10)
plt.ylabel('Err', fontsize=10)
plt.title(['n=' + str(n)], fontsize=14)

# fit data with different models
# np.polyfit and np.polyval might be useful

# plotting figures
# sample code

# plt.figure()
# plt.subplot(121)
# plt.semilogy(np.arange(1, deg+1), error[:, -1])
# plt.xlabel('degree of polynomial')
# plt.ylabel('log of error')
# plt.subplot(122)
# plt.semilogy(np.arange(n_s, n_s+step), error[-1, :])
# plt.xlabel('number of samples')
# plt.ylabel('log of error')
# plt.show()
```

```
-----
-----  
LinAlgError                                Traceback (most recent call  
    last)  
<ipython-input-173-6e9650117cd2> in <module>()  
    17 for d in range(D):  
    18     Xs = np.array([[x**i for i in range(d)] for x in alpha])  
--> 19     w = np.linalg.lstsq(Xs, y)[0]  
    20     err[d] = np.sum(np.square(Xs.dot(w) - y)) / len(y)  
    21  
  
/usr/lib/python3.6/site-packages/numpy/linalg/linalg.py in lstsq(a, b,  
rcond)  
    1976         b = b[:, newaxis]  
    1977     _assertRank2(a, b)  
-> 1978     _assertNoEmpty2d(a, b) # TODO: relax this constraint  
    1979     m = a.shape[0]  
    1980     n = a.shape[1]  
  
/usr/lib/python3.6/site-packages/numpy/linalg/linalg.py in _assertNoEm  
pty2d(*arrays)  
    223     for a in arrays:  
    224         if _isEmpty2d(a):  
--> 225             raise LinAlgError("Arrays cannot be empty")  
    226  
    227 def transpose(a):  
  
LinAlgError: Arrays cannot be empty
```

Question 5

```
In [75]: class HW3_Sol(object):
```

```
    def __init__(self):
        pass

    def load_data(self):
        self.x_train = pickle.load(open('x_train.p','rb')), encoding='latin1')
        self.y_train = pickle.load(open('y_train.p','rb')), encoding='latin1')
        self.x_test = pickle.load(open('x_test.p','rb')), encoding='latin1')
        self.y_test = pickle.load(open('y_test.p','rb')), encoding='latin1')

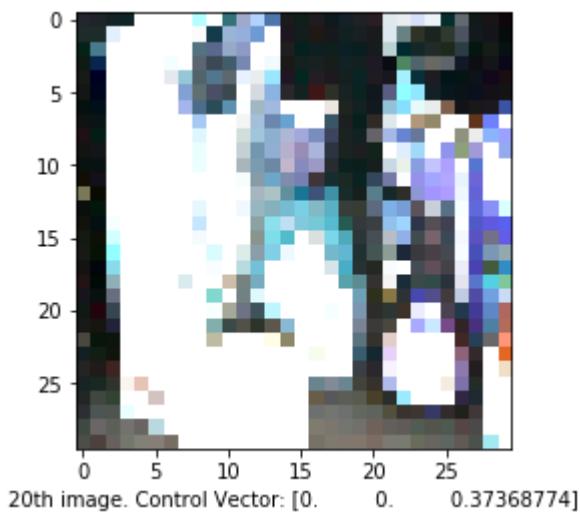
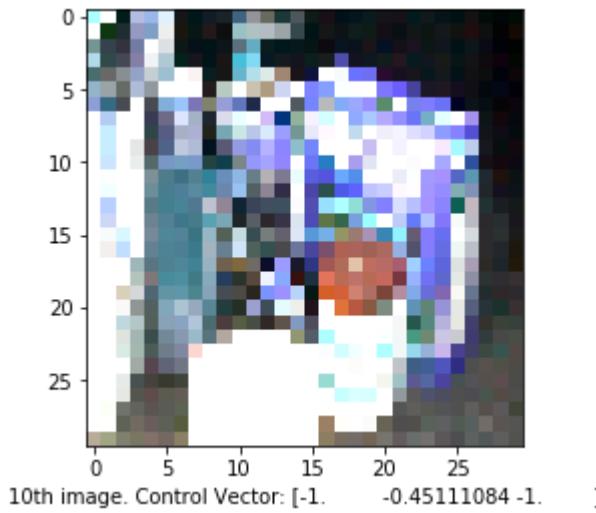
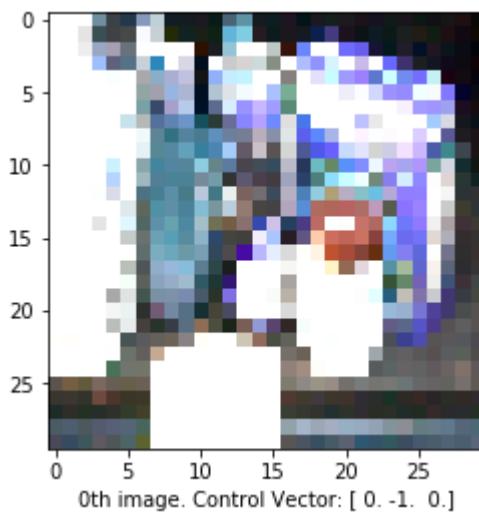
hw3_sol = HW3_Sol()
hw3_sol.load_data()

#please visualize the 0th, 10th and 20th images in the training dataset.
# Also find out what's their corresponding control vectors.

#Data
    # Tuples of n values
    # First index is sample of n
    # Second index is row of sample
    # Third index is column of a row
    # Last index is an entry in RGB pixel
    #Shape of x_train (91, 30, 30, 3)

#5.a
def visualize_control_vectors(index):
    sample = hw3_sol.x_train[index]
    y_sample = hw3_sol.y_train[index]
    plt.imshow(sample)
    plt.xlabel(str(index) + "th image. Control Vector: " + str(y_sample))
    plt.show()

visualize_control_vectors(0)
visualize_control_vectors(10)
visualize_control_vectors(20)
```



```
In [156]: #5.b
X = np.vstack([sample.flatten() for sample in hw3_sol.x_train])
U = hw3_sol.y_train

w = np.linalg.lstsq(X, U, rcond=None)
print("PI = ")
print(w[0])

PI =
[[2.89426558e-05 6.73236427e-05 9.28172481e-05]
 [4.33752719e-05 6.97139943e-05 6.69960177e-05]
 [4.24503660e-05 5.92797801e-05 6.84843436e-05]
 ...
 [1.36854182e-05 8.81609378e-05 5.15078903e-05]
 [1.96098789e-06 8.02702569e-05 3.78611891e-05]
 [7.46474147e-07 7.25823830e-05 3.73556513e-05]]
```

5.b

The weights are on the order of e^{-5} to e^{-7} . However, some of the weights are very large and some of the weights are very small in relation to each other. There is a large difference in weight values due to the relatively large differences in the pixel values of our data. We can amend this through the implementation of standardization and ridge regression.

```
In [98]: #5.c
def lstsqL(A, b, lambda_):
    return np.linalg.solve((A.T @ A) + np.eye(len(A[0]))*lambda_, A.T
@ b)
lambda_array = [0.1, 1.0, 10, 100, 1000]

for l in lambda_array:
    w_hat = lstsqL(X, U, l)
    err = np.average(np.square(X.dot(w_hat) - U))
    print(str(err) + "      Traning Error for Lambda: " + str(l))

print()
print("This is different than Ehimares, Confirm")
```

29320754873.53885	Traning Error for Lambda: 0.1
134705132893.04636	Traning Error for Lambda: 1.0
2439924413.431048	Traning Error for Lambda: 10
3207855552.0886855	Traning Error for Lambda: 100
7188177320.168516	Traning Error for Lambda: 1000

In [154]: #5.d

```
X_S = (X/255)*2 - 1 # BUT THESE DATA POINTS ARE NOT ALL BETWEEN 0 and 1
?!
print("Errors for standardized X")
print()
for l in lambda_array:
    w_hat = lstsqL(X_S, U, l)
    err = np.average(np.square(X_S.dot(w_hat) - U))
    print(str(err) + "      Traning Error for Lambda: " + str(l))
```

Errors for standardized X

1.1413557810626031e-07	Traning Error for Lambda: 0.1
1.0168794432272311e-05	Traning Error for Lambda: 1.0
0.0005420243232873813	Traning Error for Lambda: 10
0.011143580331285194	Traning Error for Lambda: 100
0.0814696917144622	Traning Error for Lambda: 1000

In [155]: #5.e

```
Xtest = np.vstack([sample.flatten() for sample in hw3_sol.x_test])
Utest = hw3_sol.y_test

for l in lambda_array:
    w_hat = lstsqL(Xtest, Utest, l)
    err = np.average(np.square(Xtest.dot(w_hat) - Utest))
    print(str(err) + "      Test Error for Lambda: " + str(l))
print()
print()
X_S_test = (Xtest/255)*2 - 1
for l in lambda_array:
    w_hat = lstsqL(X_S_test, Utest, l)
    err = np.average(np.square(X_S_test.dot(w_hat) - Utest))
    print(str(err) + "      Standardized Test Error for Lambda: " + str(l))
```

1808276541.284668	Test Error for Lambda: 0.1
30271423502.21572	Test Error for Lambda: 1.0
2270580651.7596126	Test Error for Lambda: 10
2041647982.9062707	Test Error for Lambda: 100
954418744.3903918	Test Error for Lambda: 1000

1.82861668664678e-08	Standardized Test Error for Lambda: 0.1
1.7546293832295183e-06	Standardized Test Error for Lambda: 1.0
0.00012689114816606716	Standardized Test Error for Lambda: 10
0.0037786682640806486	Standardized Test Error for Lambda: 100
0.04386417414236367	Standardized Test Error for Lambda: 1000

5.e.

Minimum error in test data OLS Ridge Regression with lambda = 100 for both standardized and unstandardized test data.

Lambda affect on performance in terms of bias:

As lambda increases our bias increases.

Lambda affect on performance in terms of variance:

As lambda increases our variance decreases.

```
In [157]: #5.f

#With training data

without = \
np.linalg.svd((X.T @ X) + 100*np.eye(len(X[0])), compute_uv = False)
k = np.max(without) / np.min(without)
print(k)

with_ = \
np.linalg.svd((X_S.T @ X_S) + 100*np.eye(len(X_S[0])), compute_uv = False)
k = np.max(with_) / np.min(with_)
print(k)

1626465.0894416121
440.7329578145653
```

Question 6.

Q: Regarding question 5, robotic learning of control from demonstrations and images, how would higher resolution visuals in our data affect our prediction error and the performance of Huey?

A: We have already seen and proven how more data can improve our model. Increasing the resolution of our visuals would increase the amount of good features that we have in our model. Sharper images would provide greater distinction between the data points, thus we would have more singular values that are discrete. In particular, we can note that better features would be those that better define the outlines of the objects. As such sharper images could potentially overcomplicate\ our model as our goal is to recognize the boundaries between object and not the specific details of each object. Features which give greater detail to the inner boundary portions of the objects would not serve as good features.

We can interpret this as something to keep in mind when designing machine learning models, our task can determine what data we need to collect for training.