# Homework Q's

4.2 $\quad f(\vec{u}, \vec{v}) = \vec{u}^T A \vec{u} + \vec{v}^T \underbrace{(B\vec{u})} + c$

use defns.

$$\vec{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \qquad A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} & \ddots \end{bmatrix}$$

12.3

$$\frac{1}{n-1} \sum_{i=1}^{n} \underbrace{(\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T} \qquad \vec{x}_i \quad 15\text{-dim}$$

$\qquad\qquad\qquad$ outer product $= 15 \times 15$

$\begin{bmatrix} \frac{1}{\vec{\mu}} & \cdots & \frac{1}{\vec{\mu}} \end{bmatrix}$ $\qquad$ broadcasting $\qquad$ A = np.array($\cdots$)

X $\quad$ - mu[:, np.newaxis]

$\qquad\qquad$ np.repeat

A.shape $\quad$ (100,) $\quad$ (100,1)

# Problem 11 → new PDF



error

under fitting

overfitting → high var

high bias

test error

var

unbiased

bias

model complexity

big coefficients on polynomials

$\beta$'s $\gg 100$

0   1   polynomial degree

model complexity

Ridge          OLS          all linear models

"size of $\beta$'s"

# Ridge Regression

# Ordinary least squares (OLS) is <u>unbiased</u>

$$\vec{y} = X\vec{\beta}^* + \vec{\varepsilon}$$

Assume:
$X^T X$ not singular

↑ true model

np.random.randn
noise Gaussian random vars

$$\mathbb{E}[\vec{\varepsilon}] = 0$$

$$\vec{\beta} = (X^T X)^{-1} X^T \vec{y} = (X^T X)^{-1} X^T \left( \underline{X\vec{\beta}^*} + \underline{\vec{\varepsilon}} \right)$$

$$= \underline{(X^T X)^{-1} X^T X \vec{\beta}^*} + \underline{(X^T X)^{-1} X^T \vec{\varepsilon}}$$

$$= \vec{\beta}^* + (X^T X)^{-1} X^T \vec{\varepsilon}$$

$$\boxed{\mathbb{E}[\vec{\beta}] = \vec{\beta}^*} + (X^T X)^{-1} X^T \cancel{\mathbb{E}[\vec{\varepsilon}]}$$

averaging over noise,
OLS estimate = truth

0

$$\text{Bias} = \mathbb{E}[\vec{\beta}] - \vec{\beta}^*$$

# Observations

$\left( \begin{array}{l} \text{Cholesky decomposition} \\ \text{generating correlated} \\ \text{Gaussians} \end{array} \right)$

$\sigma_1 \approx \sigma_2$

$\sigma_1 \approx 5\sigma_2$
same
rough size

$\sigma_1 \gg \sigma_2$
$\sigma_1 \approx 100\,\sigma_2$
small
$\sigma_2$

$$\vec{\beta}_{OLS} = \sum_{i=1}^{rank(X)} \vec{V}_i \left( \frac{\overbrace{\vec{u}_i^T \vec{y}}^{\text{has noise}}}{\underbrace{\sigma_i}} \right)$$
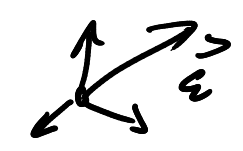
$W_i$

$\dfrac{1}{small} = big$

noise = equally shared
among the components

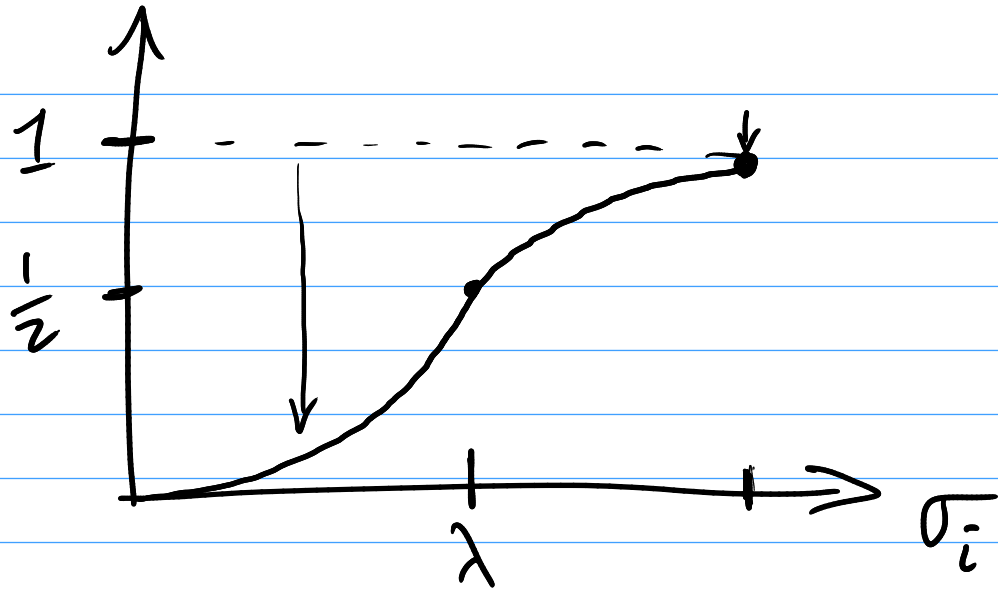correlations in X
$\iff$ small s.v.'s

If linear model is true:

$$w_i = \vec{v}_i^T \vec{\beta}^* + \frac{\vec{u}_i^T \vec{\varepsilon}}{\sigma_i} \xleftarrow{\text{i.i.d.}}$$

$\underbrace{\qquad}$ $i^{th}$ component $\vec{\beta}^*$ in $V$ basis

$\underbrace{\frac{\vec{u}_i^T \vec{\varepsilon}}{\sigma_i}}_{\text{effect of noise}}$

$$\vec{\beta}_{ridge} = \sum_{i=1}^{rank(X)} \vec{v}_i \underbrace{\left( \frac{\vec{u}_i^T \vec{y}}{\sigma_i + \lambda} \right)}_{w_i}, \quad w_i = \underbrace{\left( \frac{\sigma_i}{\sigma_i + \lambda} \right)}_{\text{plot}} \vec{v}_i^T \vec{\beta}^* + \frac{\vec{u}_i^T \vec{\varepsilon}}{\sigma_i + \lambda}$$

$\lambda$ ridge parameter, must be tuned for dataset / scenario

$\sigma_i > \lambda$ : not much effect

$\sigma_i < \lambda$ : have strong effect on those components

$$\frac{\sigma_i}{\sigma_i + \lambda} = \frac{1}{1 + \frac{\lambda}{\sigma_i}}$$

Ridge: add bias ($\lambda = 0$ no bias, OLS)
reduce variance

Extra: How to get $w_i$ formula

$$w_i = \frac{\vec{u_i}^T \vec{y}}{\sigma_i} = \frac{\vec{u_i}^T (X\vec{\beta}^* + \vec{\varepsilon})}{\sigma_i}$$ <span style="color:blue">using linear assumption</span>

$$= \frac{\vec{u_i}^T (USV^T \vec{\beta}^* + \vec{\varepsilon})}{\sigma_i}$$ <span style="color:blue">using $X = USV^T$</span>

$$= \underbrace{\vec{u_i}^T USV^T \vec{\beta}^* \left(\frac{1}{\sigma_i}\right)}_{\vec{u_i}^T U = \begin{pmatrix} 0 & 0 & & \\ & & 1 & \\ & & \uparrow & 0 & 0 \end{pmatrix}} + \underbrace{\frac{\vec{u_i}^T \vec{\varepsilon}}{\sigma_i}} = \underbrace{\vec{v_i}^T \vec{\beta}^*}_{i^{th} \text{ component in } V \text{ basis}} + \overbrace{\frac{\vec{u_i}^T \vec{\varepsilon}}{\sigma_i}}^{\text{noise term}}$$

$i^{th}$ position

So it selects $i^{th}$ singular vector, i.e.
$$\vec{u_i}^T USV^T \vec{\beta}^* = \sigma_i \vec{v_i}^T \vec{\beta}^*$$