## $\alpha$ · Regression

$$y = x^T \beta^* + \xi$$

$$\hat{\beta}(\lambda) = \underset{\beta}{\text{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^{n} (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_2^2 \right\}$$

$$Y = X^T \beta^* + \xi$$

$$\hat{\beta}(\lambda) = \underset{\beta}{\text{argmin}} \left( |Y - X^T \beta|^2 + \lambda \|\beta\|^2 \right)$$

$$f(\beta) = \underset{\beta}{\text{argmin}} \left( Y - X\beta \right)^T (Y - X\beta) + \lambda \|\beta\|^2$$

$$\hat{\beta}(\lambda) \sim N(\beta^* \mathbb{E}(\beta|x) \, \text{Var})$$

$$\frac{\partial f(\beta)}{\partial \beta} = -2 X^T (Y - X\beta) + 2\lambda \beta$$

$$\boxed{\hat{\beta}_\lambda^{Ridge} = (X^T X + \lambda I_p)^{-1} X^T Y}$$

$$\hat{E}(\hat{\beta}_\lambda^{Ridge}|X) = E\left\{ (I_p + \lambda X^T X) \hat{\beta}^{ls} \right)$$

$$= (I_p + \lambda (X^T X)^{-1}) \beta^*$$

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

$$\text{Var}[\hat{\beta}|X] = \sigma^2 (X^T X + \lambda I)^{-1} (X^T X)(X^T X + \lambda I)^{-1}$$

$$= x^T \beta^* - E_D\left[ x^T \cdots \right.$$

$$\text{bias} = x^T \beta^* - E_D\left[ x^T (X^T X + \lambda I_p)^{-1} X^T Y \right]$$

$$\text{bias} = x^T \beta^* - E_D\left[ x^T (X^T X + \lambda I_p)^{-1} X^T (X\beta^* + \xi) \right]$$

$$\text{bias} = x^T \beta^* - E_D\left[ x^T (X^T X + \lambda I_p)^{-1} X^T X \beta^* + x^T (X^T X + \lambda I_p)^{-1} X^T \xi \right]$$

$$\text{bias} = x^T \beta^* - E_D\left[ x^T \beta^* + x^T (\lambda I_p)^{-1} X^T X \beta^* + x^T (X^T X)^{-1} X^T \xi + x^T (\lambda I_p) X^T \xi \right]$$

$$\text{bias} = x^T \beta^* - x^T \beta^* - x^T (X^T X)^{-1} X^T E_D[\xi]$$
$$- E_D\left[ x^T (\lambda I_p)^{-1} X^T (X\beta^* + \xi) \right]$$

$$\text{bias} = -x^T (X^T X)^{-1} X^T E_D[\xi] - E\left[ x^T (\lambda I_p)^{-1} X^T Y \right]$$

$$\text{bias} = -\underbrace{x^T (X^T X)^{-1} X^T E_D[\xi]}_{} - \underbrace{x^T (\lambda I_p)^{-1} X (X\beta^* + E[\xi])}_{}$$

Constant with Sample $x$ and given set $(X, Y)$

$f(\lambda)$ with $|f(\lambda)|$ ~~is~~ proportional to $\lambda$ value

Variance

$$E\left[\left(x^T\hat{\beta}(\lambda) - E\left(x^T\hat{\beta}(\lambda)\right)\right)^2\right]$$

$$= E\left[\left(x^T\hat{\beta}(\lambda) - x^T\beta^*\right)^2\right]$$

$$= E\left[\left(x^T(X^TX + \lambda I_p)^{-1}X^Ty - x^T\beta^*\right)^2\right]$$

$$= E\left[\left(x^T(X^TX)^{-1}X^Ty + x^T(\lambda I_p)^{-1}X^Ty - x^T\beta^*\right)^2\right]$$

$$= E\left[\left(x^T\beta^* + x^T(\lambda I_p)^{-1}X^Ty - x^T\beta^* \right.\right.$$
$$\left.\left. + x^T(X^TX)^{-1}X^T\varepsilon\right)^2\right]$$

$$= E\left[\left(x^T(\lambda I_p)^{-1}X^T(X^T\beta^* + \varepsilon) + x^T(X^TX)^{-1}X^T\varepsilon\right)^2\right]$$

$$= E\left[\sigma^2 x^T(X^TX + \lambda I_p)^{-1}x\right].$$

$\lambda$ is high $\rightarrow$ bias is high, variance is low $\rightarrow$ underfit

$\lambda$ is low $\rightarrow$ bias is low, variance is high $\rightarrow$ overfit