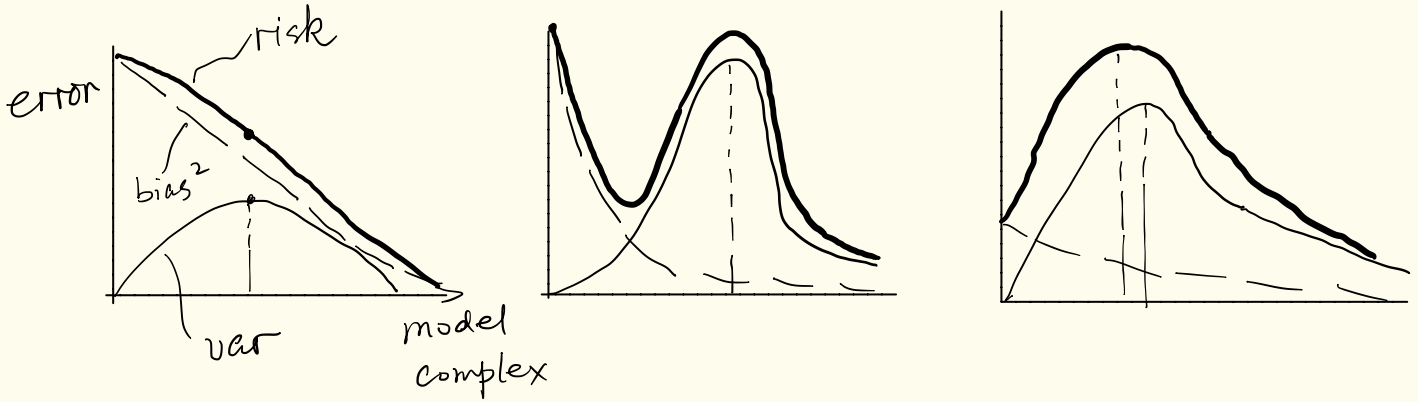


# Bias-Variance Tradeoff for Deep Neural Networks

(Yang et al., ICM 2020)



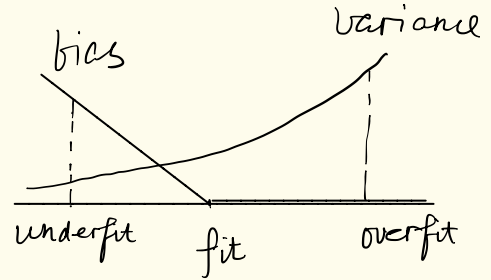
# Bias-Variance Tradeoff (Traditional Statistical Learning Result)

Ex1 (Overfit & Underfit linear model)

1. Underfitting

True model:

$$\underline{y} = \underline{X}\underline{\beta} + \underline{e} = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \underline{e}$$



Fitting / Assumed model:  $\underline{y} = \underline{X}_1 \underline{\beta}_1 + \underline{\varepsilon}$ , w/ LS soln:

$$\hat{\underline{\beta}}_1 = (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T \underline{y}$$

$$\begin{aligned} \text{Bias: } \mathbb{E}[\hat{\underline{\beta}}_1] &= (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \\ &= (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T (x_1 \beta_1 + x_2 \beta_2) \\ &= \underline{\beta}_1 + (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T x_2 \beta_2 \quad (\text{biased}) \end{aligned}$$

Variance:  $\text{Var}(\hat{\beta}_1) = (X_1^T X_1)^{-1} X_1^T \sigma^2 I X_1 (X_1^T X_1)^{-1} = \sigma^2 (X_1^T X_1)^{-1} = B$   
 (has smaller var)

$$\text{Var}(\hat{\beta}_{\text{fit}}) = \sigma^2 (X^T X)^{-1} = \sigma^2 \begin{bmatrix} X_1^T X_1 & X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 \end{bmatrix}^{-1} = \sigma^2 \left[ \begin{array}{c|c} \text{A} & \\ \hline & \end{array} \right]$$

$\begin{bmatrix} X_1^T \\ X_2^T \end{bmatrix} \uparrow \begin{bmatrix} X_1 & X_2 \end{bmatrix}$

The upper corner is  $[X_1^T X_1 - X_1^T X_2 (X_2^T X_2)^{-1} X_2^T X_1]^{-1} \stackrel{\text{def}}{=} A$

Block matrix inversion:  $\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} (A - B D^{-1} C)^{-1} & * \\ * & * \end{bmatrix}$

We have  $A \geq B$  or  $A - B \geq 0$  or  $v^T (A - B) v \geq 0, \forall v \in \mathbb{R}^n$

Def:  $M$  is positive semi-definite (PSD) if  $\tilde{x}^T M \tilde{x} \geq 0, \forall \tilde{x} \in \mathbb{R}^n$

$M$  ... positive definite (PD) if  $\tilde{x}^T M \tilde{x} > 0, \forall \tilde{x} \neq 0$ .

2. Overfitting:

$$\text{True model: } y = X_1 \beta_1 + \underset{\sim}{e}$$

$$\text{Fitting/Assumed model: } y = \underbrace{\begin{bmatrix} X_1 & X_2 \end{bmatrix}}_X \begin{bmatrix} \underset{\sim}{\beta_1} \\ \underset{\sim}{\beta_2} \end{bmatrix} + \underset{\sim}{\varepsilon}, \text{ w/}$$

$$\text{LS soln: } \underset{\sim}{\hat{\beta}} = (X^T X)^{-1} X^T y$$

$$\begin{aligned} \text{Bias: } E[\underset{\sim}{\hat{\beta}}] &= (X^T X)^{-1} X^T (X_1 \beta_1 + X_2 \underset{\sim}{\varepsilon}) = (X^T X)^{-1} X^T X \begin{bmatrix} \beta_1 \\ \underset{\sim}{\varepsilon} \end{bmatrix} \\ &= \begin{bmatrix} \beta_1 \\ \underset{\sim}{\varepsilon} \end{bmatrix} \quad (\text{unbiased (for } \beta_1) \text{ even overfit!}) \end{aligned}$$

$$\text{Variance: } \text{Var}(\underset{\sim}{\hat{\beta}}) = \sigma^2 (X^T X)^{-1} = \sigma^2 \begin{bmatrix} X_1^T X_1 & X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 \end{bmatrix}^{-1} \quad (\text{has larger variance.})$$

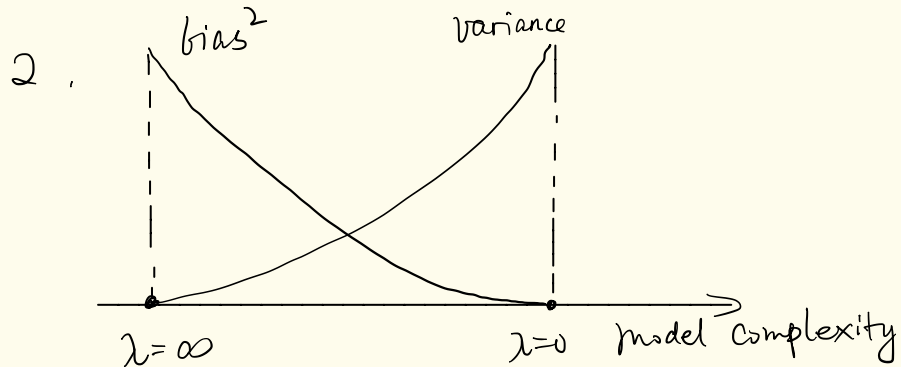
$$\text{Var}(\underset{\sim}{\hat{\beta}}_{\text{fit}}) = \sigma^2 (X_1^T X_1)^{-1}$$

Ex 2 (Ridge Regression)  $\hat{\beta}_{\lambda} = \underset{\beta}{\operatorname{argmin}} \left( \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right),$

$X \in \mathbb{R}^{n \times p}$ ,  $y \in \mathbb{R}^{n \times 1}$ ,  $\beta \in \mathbb{R}^{p \times 1}$ ,  $X$  full column rank,  $\lambda \geq 0$

Assume true model:  $y = X\beta + e$

claims: 1.  $\hat{\beta}_{\lambda} = (X^T X + \lambda I)^{-1} X^T y$  (HW)



$$\text{Bias: } \mathbb{E}[\hat{\beta}] = \mathbb{E}[(X^T X + \lambda I)^{-1} X^T \underset{\sim}{y}] = \underbrace{(X^T X + \lambda I)^{-1} X^T X}_{R} \beta$$

$$R = V \Lambda V^T$$

## Eigenanalysis Review

characteristic

$$1. R = V \Lambda V^T$$

$$= \left[ \underset{\sim}{v_1} \mid \dots \mid \underset{\sim}{v_p} \right] \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_p \end{bmatrix} \begin{bmatrix} \frac{v_1^T}{\lambda_1} \\ \vdots \\ \frac{v_p^T}{\lambda_p} \end{bmatrix}$$

$$= [\lambda_1 \underset{\sim}{v_1} \mid \dots \mid \lambda_p \underset{\sim}{v_p}] \begin{bmatrix} \frac{v_1^T}{\lambda_1} \\ \vdots \\ \frac{v_p^T}{\lambda_p} \end{bmatrix}$$

$$= \sum_{i=1}^p \lambda_i \underset{\sim}{v_i} \underset{\sim}{v_i}^T$$

$$2. R V = V \Lambda = [v_1 \dots v_p] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix}$$

$$3. R v_i = \lambda_i v_i, \quad i=1, \dots, p$$

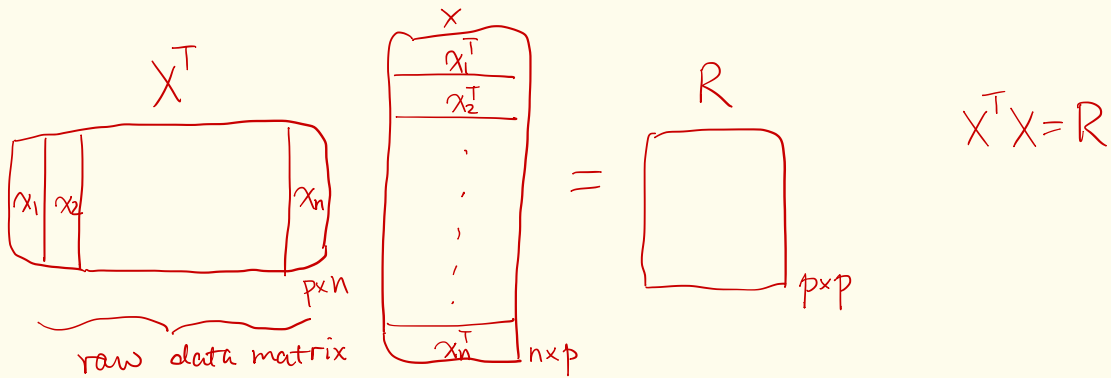
$$X^T X = R \quad (\text{symmetric})$$

$$(X^T X)^T = X^T X = R^T$$

$V$ : orthogonal/orthonormal matrix:

① all vectors are orthogonal, has unit norm, i.e.,  $\underset{\sim}{v_i}^T \underset{\sim}{v_j} = \begin{cases} 0, & \forall i \neq j \\ 1, & \forall i = j \end{cases}$

$$② V^T = V^{-1}$$



$$R = X^T X = \sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T : \text{sample correlation matrix}$$

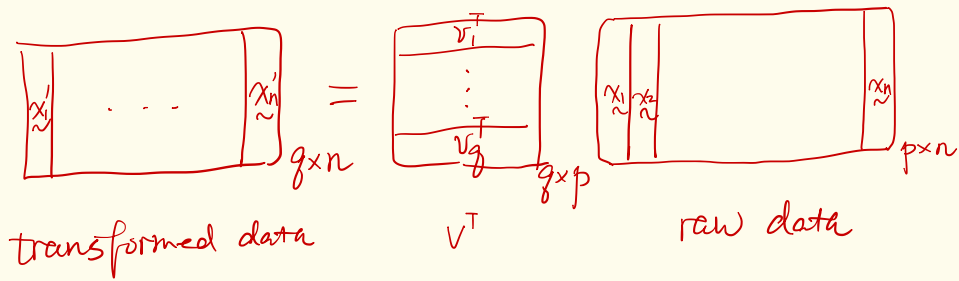
$\mathbb{E}[\tilde{x} \tilde{x}^T] : \text{correlation matrix}$

$$C = \sum_{i=1}^n (\tilde{x}_i - \bar{x})(\tilde{x}_i - \bar{x})^T : \text{sample covariance matrix}$$

$\mathbb{E}[(\tilde{x} - \mathbb{E}\tilde{x})(\tilde{x} - \mathbb{E}\tilde{x})^T] : \text{covariance matrix}$

Principal Component Analysis (PCA), usually  $g \ll p$ ;

Karhunen-Loeve Transform (KLT),  $g = p$ .





$$\text{Bias: } E[\hat{\beta}] = (X^T X + \lambda I)^{-1} X^T X \beta$$

$$= (V \Lambda V^T + \lambda V V^T)^{-1} V \Lambda V^T \beta$$

$$= [V(\Lambda + \lambda I)V^T]^{-1} V \Lambda V^T \beta$$

$$= V(\Lambda + \lambda I)^{-1} V^T V \Lambda V^T \beta$$

$$= V(\Lambda + \lambda I)^{-1} \Lambda V^T \beta$$

$$= V \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{\lambda_p}{\lambda_p + \lambda} \end{bmatrix} V^T \beta$$

$$= V \left( I - \begin{bmatrix} \frac{\lambda}{\lambda_1 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{\lambda}{\lambda_p + \lambda} \end{bmatrix} \right) V^T \beta$$

$$= \beta - \underbrace{\sum_{i=1}^p \frac{1}{\lambda_i/\lambda + 1} v_i v_i^T}_{\neq 0} \beta = \begin{cases} \beta, & \lambda \rightarrow 0 \text{ (unbiased)} \\ 0, & \lambda \rightarrow \infty \text{ (biased)} \end{cases}$$

$$\begin{bmatrix} \lambda_1 + \lambda & & 0 \\ & \ddots & \\ 0 & & \lambda_p + \lambda \end{bmatrix}_{p \times p}$$

$$\text{Variance: } \text{Var}(\hat{\beta}) = \underbrace{(X^T X + \lambda I)^{-1}} \underbrace{X^T} \underbrace{\text{Cov}(y)}_{\sigma^2 I} \underbrace{X (X^T X + \lambda I)^{-1}}$$

$$= \underbrace{V(\Lambda + \lambda I)^{-1} V^T}_{V \Delta V^T} \underbrace{X^T X}_{V \Delta V^T} V(\Lambda + \lambda I)^{-1} V^T \cdot \sigma^2$$

$$= V \begin{bmatrix} \frac{\lambda_1}{(\lambda_1 + \lambda)^2} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{\lambda_p}{(\lambda_p + \lambda)^2} \end{bmatrix} V^T \sigma^2$$

$$= \sigma^2 \sum_{\bar{i}=1}^p \frac{\lambda_{\bar{i}}}{(\lambda_{\bar{i}} + \lambda)^2} \underbrace{v_{\bar{i}}}_{\sim} \underbrace{v_{\bar{i}}^T}_{\sim} = \begin{cases} 0, & \lambda = \infty, \text{ (zero variation)} \\ \sigma^2 (X^T X)^{-1}, & \lambda = 0. \text{ (LS variance)} \end{cases}$$