# BIAS AND VARIANCE ANALYSIS

SUPERVISED $\mathbb{R}$

$$S = \left\{ (x^{(i)}, y^{(i)}) \right\}^{n} \quad , \quad x^{(i)} \in \mathbb{R}^{d} \quad , \quad y \in \mathbb{R}$$

$\Theta^{*}$ "TRUE PARAMETER" UNKNOWN

$X \in \mathbb{R}^{n \times d}$ "DESIGN MATRIX"

← feet →

↑
eq
↓

$\hat{\Theta}_{n}$ — ESTIMATOR OF $\Theta$ USING $n$ TRAINING EXAMPLES

$$\begin{array}{c} x^{(1)}, y^{(1)} \\ x^{(2)}, y^{(2)} \\ \vdots \\ x^{(n)}, y^{(n)} \end{array} \overset{IID}{\sim} P \Rightarrow \boxed{\text{ESTIMATOR}} \Rightarrow \hat{\Theta}_{n}$$

EX

← VARIANCE →

BIAS

$\theta^*$

estim. 1
$\theta_{n_1}$ $\theta_{n_2}$ $\theta_{n_3}$

$\theta^*$

estim. 2
$\theta_{n_1}$ $\theta_{n_3}$ $\theta_{n_2}$

$\theta^*$

estim. 3
$\theta_{n_1}$ $\theta_{n_2}$ $\theta_{n_3}$

$\theta^*$

estim. 4
$\theta_{n_1}$ $\theta_{n_2}$ $\theta_{n_3}$

ESTIMATOR IS A SAMPLING DISTRIBUTION AND THE $\theta$ ARE SAMPLES

$$\text{BIAS}(\hat{\theta}_n) \equiv \mathbb{E}[\hat{\theta}_n - \theta^*] \quad , \quad \text{IF UNBIASED} \quad \mathbb{E}[\hat{\theta}_n] = \theta^*$$

IS A THEORETICAL PROPRIETY THAT CAN ONLY BE OBTAINED ANALITICALLY

$$\text{VARIANCE}(\hat{\theta}_n) \equiv \text{Cov}[\hat{\theta}_n]$$

IF WE USE BOOTSTRAP (ESTIMATION w/ DIFFERENT DATA SETS)
WE CAN GET GOOD ESTIMATES

AND THE MEAN SQUARED ERROR (MSE) WILL BE

$$\text{MSE}(\hat{\theta}_n) = \mathbb{E}[\|\hat{\theta}_n - \theta^*\|^2]$$

$$\text{MSE} = \text{tr}(\text{VAR}(\hat{\theta}_n)) + \|\text{BIAS}(\hat{\theta}_n)\|^2$$

BY LOWERING THE VAR WE CAN INCREASE THE BIAS AND VICEVERSA. WE NEED TO FIND A SWEET SPOT. LET'S TAKE A LOOK AT LINEAR REGRESSION w/ $L_2$ REGULARIZATION. WE DEFINED THE LOSS $J(\Theta)$

$$J(\Theta) = \frac{\lambda}{2} \|\Theta\|_2^2 + \frac{1}{2} \sum_{i=1}^{h} (y^{(i)} - \Theta^T x^{(i)})^2$$

$$\hat{\Theta}_h = \arg\min_{\Theta \in \mathbb{R}^d} J(\Theta)$$

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

$$X^T X = u \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix} u^T \qquad (X^T X \text{ IS PSD BY } \underline{DEF})$$

$$X^T X + \lambda I = u \begin{bmatrix} \sigma_1^2 + \lambda & & \\ & \ddots & \\ & & \sigma_n^2 + \lambda \end{bmatrix} u^T$$

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

$$\hat{\Theta}_n = (X^TX + \lambda I)^{-1} X^T X \theta^* + \left[(X^TX + \lambda I)^{-1} X^T\right] \vec{\varepsilon}$$

$$\mathbb{E}[\hat{\Theta}_n] = \mathbb{E}\left[\underbrace{(X^TX + \lambda I)^{-1} X^T X \theta^*}_{\text{CONST } (*)}\right] + \mathbb{E}\left[\underbrace{\left[(X^TX + \lambda I)^{-1} X^T\right] \vec{\varepsilon}}_{0}\right]$$

$$\mathbb{E}[\hat{\Theta}_n] = (X^TX + \lambda I)^{-1} X^T X \theta^*$$

$$\mathbb{E}(\hat{\Theta}_n) = U \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_n^2}{\sigma_n^2 + \lambda} \end{bmatrix} U^T \Theta^*$$

THIS WAS BIAS, THE MORE $\lambda$ IS USED, THE MORE THE BIAS.

WHAT ABOUT VAR?

$$\text{VAR}(\hat{\Theta}_n) = \text{Cov}\left[\hat{\Theta}_n\right]$$

$$= u \begin{bmatrix} \dfrac{\lambda^2 \sigma_1^2}{(\sigma_1^2 + \lambda)^2} & & \\ & \ddots & \\ & & \dfrac{\lambda^2 \sigma_n^2}{(\sigma_n^2 + \lambda)^2} \end{bmatrix} u^T$$

WE ARE INTRODUCING REGULARIZATION, REDUCING VAR BUT INDUCING BIAS

HOW DO WE FIND THE SWEET SPOT?

FINDING THE SWEET SPOT BY LOOKING AT THE MEAN OF VAR AND BIAS IN TERMS OF PREDICTION

$$S = \left\{ x^{(i)}, y^{(i)} \right\}_{i=1}^{n}$$

$$y = f(x) + \varepsilon \qquad \mathbb{E}[\varepsilon] = 0, \quad V[\varepsilon] = \sigma^2$$

"TRUE" $f$

$$f'(x) = \mathbb{E}\left( y \mid x = x' \right)$$

$$\hat{f}_n(x) = \vec{\hat{y}}$$

THE MSE $_{f_n(x)}$

$$MSE_{f_n(x)} = \mathbb{E}\left[ \left( y_* - \vec{f}_n(x_*) \right)^2 \right] \qquad x_*, y_* \text{ IS}$$

"TEST EXAMPLE"

$$= \mathbb{E}\left[\left(\varepsilon + f(x_*) - \hat{f}_n(x_*)\right)^2\right]$$

$$= \ldots = \gamma^2 + \mathbb{E}\left[\hat{f}_n(x_*) - f(x_*)\right]^2 + V\left[\hat{f}_n(x_*)\right]$$

WILL ALWAYS EXIST,
EVEN W/ TRUE $f$ &larr; IRREDUCIABLE
GIVEN FROM         ERROR
RANDOMNESS         COMPONENT

BIAS

VAR
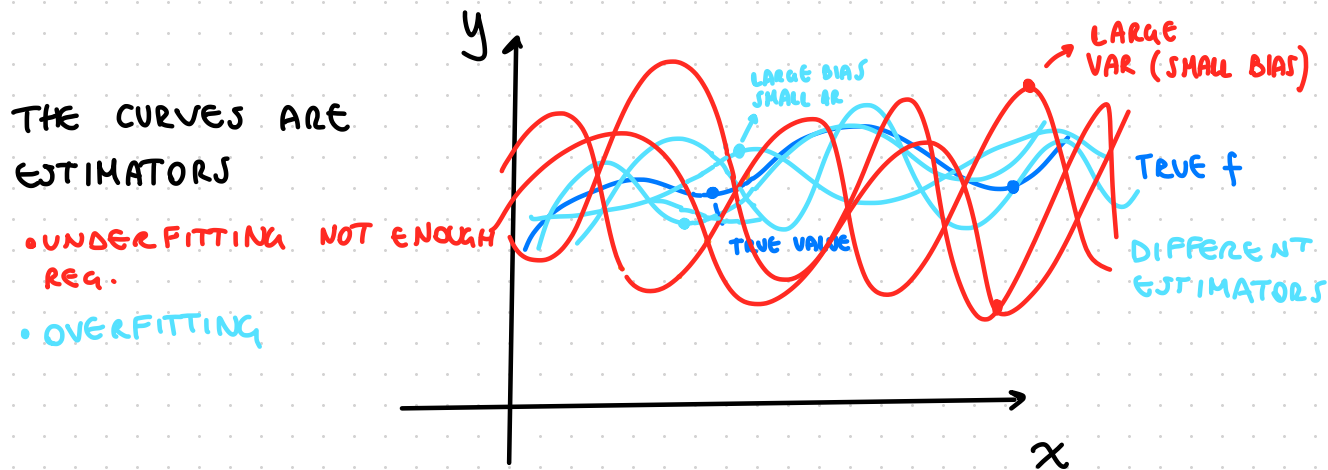HERE LIES RANDOMNESS
(NOISE)

NOTE

$$n \longrightarrow \infty$$

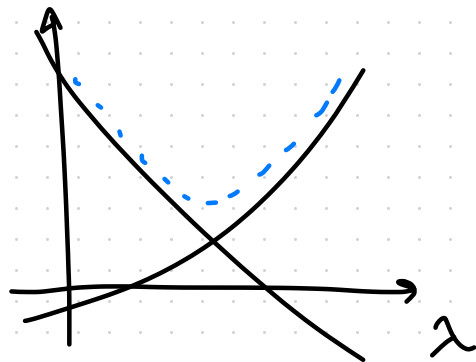$$\text{VAR}\left[\hat{\Theta}_n\right] \longrightarrow 0$$

GOING BACK TO L.R. W/ $L_2$ REG.

$$BIAS(\hat{f}_n) = BIAS(\hat{\Theta}_n)^T x_*$$

$$VAR(\hat{f}_n) = x_*^T VAR[\hat{\Theta}_n] x_*$$



THE CURVES ARE
ESTIMATORS

• UNDERFITTING NOT ENOUGH
  REG.

• OVERFITTING

AS WE INCREASE $\lambda$ INCREAS BIAS → REDUCE VAR
AND THERE IS SOME $\lambda^2$

THE MSE WILL BE



TO KNOW IF WE ARE OVER/UNDER FITTING THERE ARE HEURISTICS

|  | EX 1 | EX 2 | EX 3 |  |
|---|---|---|---|---|
|  | TEST | TEST<br>TRAIN | TRAIN | THROWING MORE DATA<br>LOWERS TRAIN ERROR |
| ERROR |  |  | TEST |  |
|  | TRAIN |  |  | THROW COMPUTING<br>LOWER TEST ERR |
|  | OVERFITTING | UNDERFITTING | NOISE/<br>BUG |  |