

BIAS & VARIANCE: CLASSICAL & MODERN ELEMENTS

+ CLASSICAL theory

→ REGULARIZATION

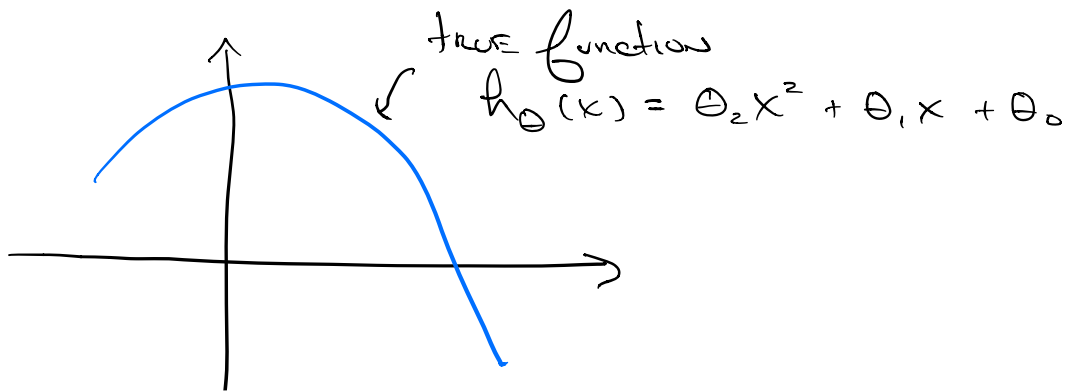
→ PARAMETER SELECTION

Computer efficient: Successive Halving

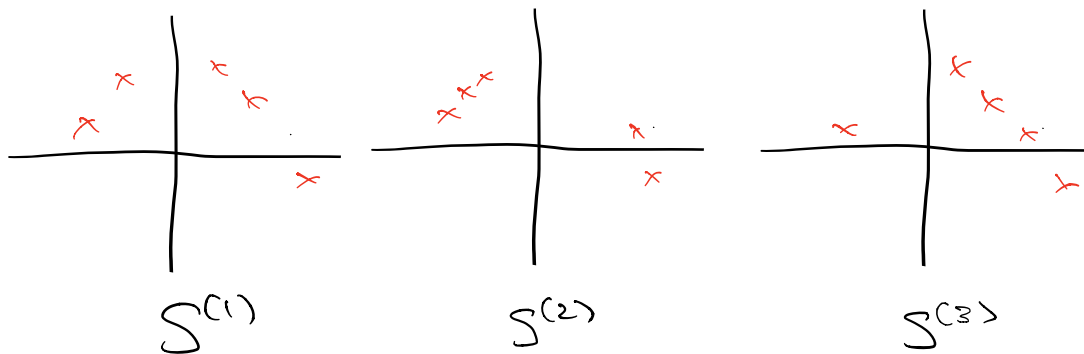
DATA

: k-fold

+ MODERN theory (BONUS)

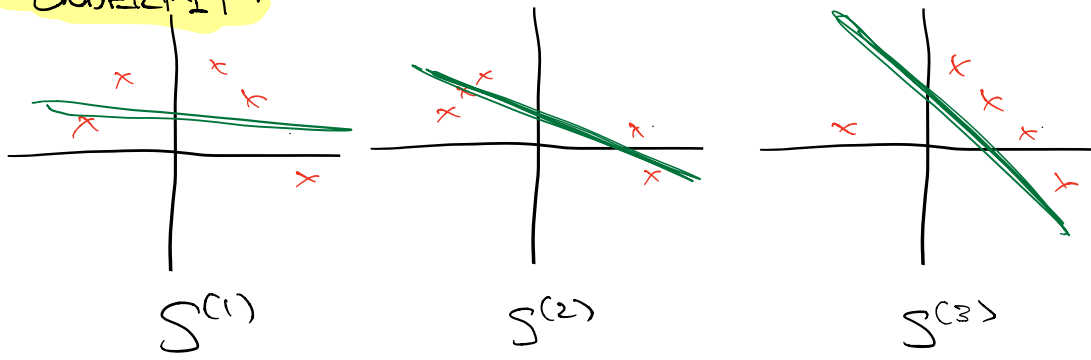


WE DON'T GET TO SEE h_0 DIRECTLY - ONLY SAMPLES



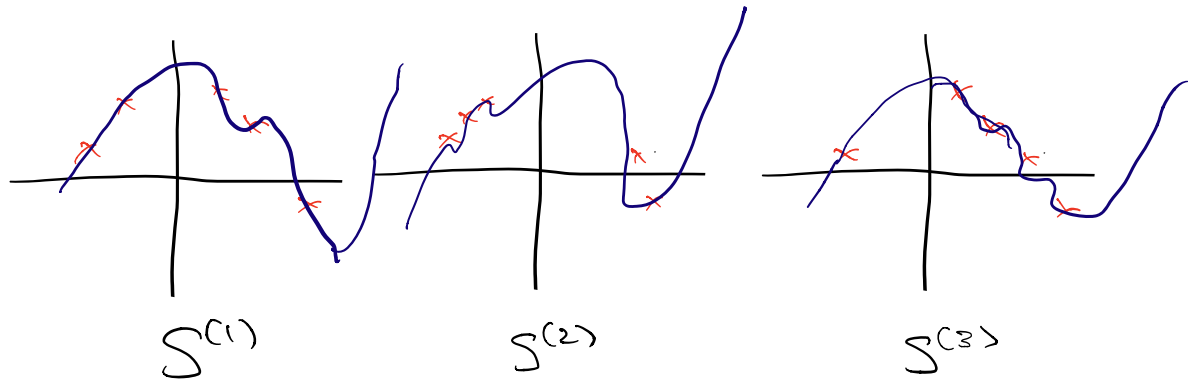
WHAT HAPPENS IF WE FIT A LINE TO THESE SAMPLES?

"UNDERFIT"



WE INFORMALLY CALL THIS underfit the error is pretty high (WE USE FOR DET.)

WHAT HAPPENS IF WE USE DEGREE 5 POLYNOMIAL?



"OVERFIT"

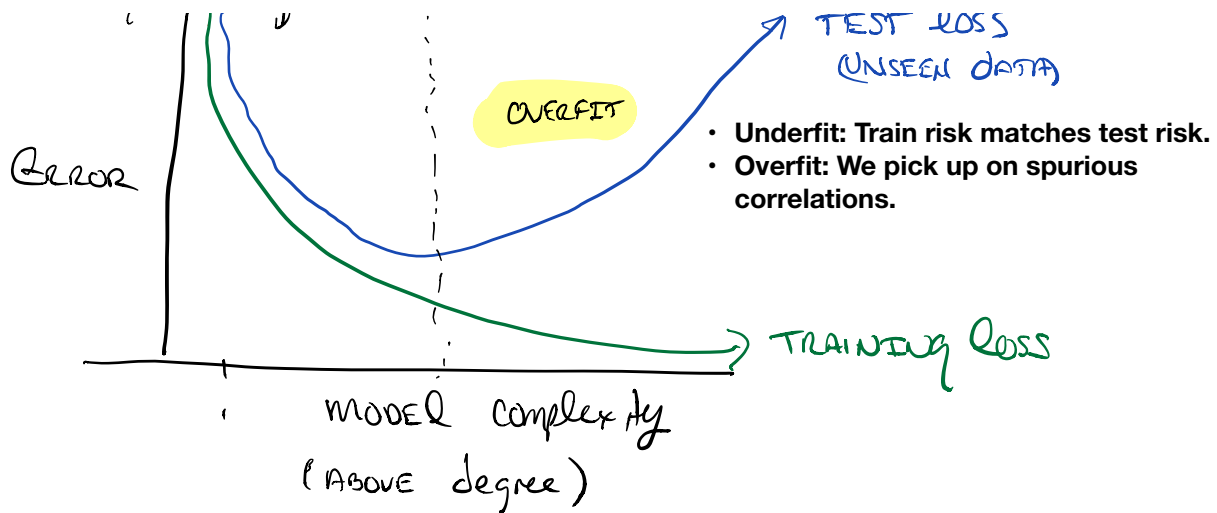
This fits EACH SAMPLE well. BUT the function totally CHANGES PER SAMPLE! (HIGH VARIANCE)

WHAT IF WE USE QUADRATICS?

→ low error & low variance.

→ IT FITS

UNDERFIT → optimal complexity



NB: This is CLASSICAL BIAS VARIANCE.

→ helpful to understand many ML ideas.

→ Incomplete for modern models in important ways (more later)

MORE FORMAL BIAS-VARIANCE

CONSIDER LINEAR REGRESSION

$$\text{Output } y \in \mathbb{R} = \theta \cdot x + \epsilon$$

$\epsilon \sim \mathcal{N}(0, \sigma^2)$
 $x \in \mathbb{R}^d$ (FEATURES (DATA/INPUT))
 $\theta \in \mathbb{R}^d$ (PARAMETERS IN \mathbb{R}^d)

PROCEDURE

Fix $x \in \mathbb{R}^d$, A TEST POINT (REASON ABOUT ERROR HERE)

1. DRAW n points $(x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$

2. TRAIN A LINEAR REGRESSOR $h_S: \mathbb{R}^d \rightarrow \mathbb{R}$

3. DRAW TEST SAMPLE (x, y) such that

$$h_S(x) + \epsilon = y$$

$\epsilon \sim \mathcal{N}(0, \sigma^2)$ NOISE

4. MEASURE $(h_S(x) - y)^2$

WE EXAMINE

$$\mathbb{E}[(h_S(x) - y)^2]$$

ϵ 's

→ TWO SOURCES OF RANDOMNESS

Goal: DECOMPOSE THIS ERROR

$$\mathbb{E}[(h_S(x) - (h_0(x) + \epsilon))^2] =$$

$$\underbrace{\mathbb{E}[\epsilon^2]}_{\sigma^2} + \underbrace{\mathbb{E}[(h_S(x) - h_0(x))^2]}_S + 2\mathbb{E}[\epsilon(h_S(x) - h_0(x))]$$

$\mathbb{E}[\epsilon] = 0$

DEPENDS ON TRAIN SET S & ϵ 'S INDEPENDENT

UNAVOIDABLE

ERROR TERM

$$h_{\text{avg}}(x) \triangleq \mathbb{E}_S[h_S(x)] \quad \text{"long run AVERAGE of many } S \text{ x"}$$

$$\begin{aligned}
 &= \mathbb{E}_S [(h_\theta(x) - h_{avg}(x) + h_{avg}(x) - h_S(x))^2] \\
 &= \mathbb{E}_S [(h_\theta(x) - h_{avg}(x))^2] + \mathbb{E} [(h_{avg}(x) - h_S(x))^2] + \overset{\text{cross term}}{0} \\
 &= \underbrace{(h_\theta(x) - h_{avg}(x))^2}_{\substack{\text{does NOT DEPEND ON } S \\ \text{ON CLASS of hypothesis}}} + \text{VAR}_S(h) \quad \text{"VARIANCE OVER TRAINING SET"} \\
 &\rightarrow \boxed{\text{BIAS}} \qquad \qquad \qquad \boxed{\text{VARIANCE}}
 \end{aligned}$$

RECAP

$$\mathbb{E}_{S, \epsilon} [(y - h_S(x))^2] = \sigma^2 \quad \downarrow \text{NOISE TERM IN TEST DATA} + \boxed{\text{BIAS}} + \boxed{\text{VARIANCE}}$$

Examples

	<u>LINEAR</u>	<u>degree 5</u>
BIAS	<u>LARGE</u>	○ (fits EVERY point!)
VARIANCE	lower	higher
	just right fit combines both!	

NOTE:

- Having different DEV / HOLDOUT SET allows us to assess VARIANCE (and hence stability)
- IF WE USE model class that is expressive MAY NEED TO "trust points less" (REDUCE VARIANCE)
 - Regularization

Regularization IS AT HEART of BOTH classical

AND modern theory. SPEND A little bit of TIME ON this...

REGULARIZATION

REDUCE VARIANCE TO GET ROBUST MODEL

→ CAN BE EXPLICIT (CHANGE MODEL)

IMPLICIT (PROCEDURE)

MOST CLASSICAL LINEAR REGRESSION

$$\underset{\Theta \in \mathbb{R}^d}{\text{ARGMIN}} \frac{1}{2} \sum_{i=1}^n (x^{(i)} \cdot \Theta - y^{(i)})^2 + \frac{\lambda}{2} \|\Theta\|_2^2$$

PARAMETER
↗ $\in \mathbb{R}_+$

→ PENALTY FOR REALLY COMPLEX MODEL (MINIMUM NORM SOLUTION)

$\lambda = 0 \rightarrow$ ORDINARY LEAST SQUARES

$\lambda = 10^{100} \rightarrow \Theta = 0$ probably looks pretty good!

SET λ TO SOME VALUE TO BALANCE LOSS & VAR.

HOW WILL STAY GOOD!

Solution? Fix $\lambda > 0$.

TAKE DERIVATE WRT TO Θ

$$X^T X \Theta - X^T y + \lambda \Theta = 0$$

$$(X^T X + \lambda I) \Theta = X^T y$$

UNDETERMINED CASE

IF $X^T X$ IS NOT FULL RANK $\nexists \lambda = 0$ MAY NOT HAVE UNIQUE SOLUTION ($x \in \mathbb{R}^{n \times d}$ $n < d$)

IF $X^T X$ IS NOT FULL RANK $\exists v$ s.t.

$$v \neq 0 \text{ \underline{w}h } X^T X v = 0$$

$\nexists X^T X \theta = X^T y$ then $X^T X(\theta + v) = X^T y$ AS WELL \Rightarrow NO UNIQUE SOLUTION

\Rightarrow if $\lambda > 0$ then IT DOES HAVE A UNIQUE SOLUTION, SINCE $X^T X + \lambda I$ IS FULL RANK.

that is eigenvalues of $X^T X$ $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2 \geq 0$

then $X^T X + \lambda I$ has eigenvalues

$$\sigma_1^2 + \lambda, \sigma_2^2 + \lambda, \dots, \sigma_n^2 + \lambda > 0$$

IN THIS CASE, $\theta_\lambda = (X^T X + \lambda I)^{-1} X^T y$

BACK TO VARIANCE

$$\mathbb{E}_S [(f_S(x) - f_{\text{avg}}(x))^2]$$

VARIANCE IN θ_λ the solutions for fixed λ

$$\text{VAR}_S = \mathbb{E}_S [(\theta_\lambda \cdot x - \mathbb{E}_S[\theta_\lambda \cdot x])^2]$$

$$\approx \mathbb{E} [\|\theta_\lambda - \mathbb{E}[\theta_\lambda]\|^2] \quad (*)$$

TO Simplify Analysis,

AND ONLY RANDOMNESS IN DATA IS THE TRAIN
POINTS ERROR

$$y = X \cdot \theta + v \quad v \in N(0, c^2 I_n)$$

RANDOM NOISE PER POINT FIXED X

$$\textcircled{\otimes} = \mathbb{E} \left[\left\| \underbrace{(X^T X + \lambda I)}_A^{-1} X^T v \right\|^2 \right]$$

then $Av \sim N(0, c^2 A A^T)$

$$\text{Hence,} \quad \leq c^2 \frac{\sigma_{\max}^2}{(\sigma_{\max}^2 + \lambda)^2}$$

SO AS λ INCREASES, VARS DECREASES

Bonus Observation CAN sometimes implicitly
regularize as well. (Surprisingly important!)
(is modern theory)

Thought Experiment, WE RUN gradient descent with $\lambda=0$
IS UNDERDETERMINED CASE

Claim IF WE INITIALIZE TO $\theta^{(0)}$ THEN OUR SOLUTION IS θ_n

$$\Theta_{\text{GD}} = \underbrace{P_{\text{NUL}(X)}(\Theta_{\text{GD}})} + P_{\text{SPAN}(X)}(\Theta_{\text{GD}}) \\ = \Theta_0!$$

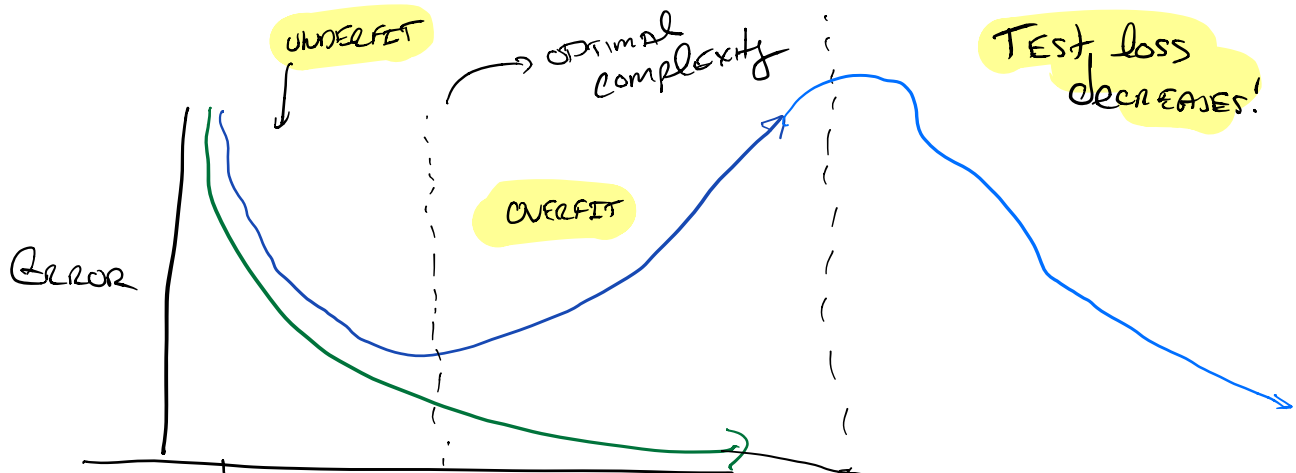
why? $\Theta^{(k+1)} = \Theta^{(k)} - \alpha X^T (X^T \Theta - y)$
 only CHANGES IN SPAN(X)!

OBSERVATION: We can Regularize by Initialization!
 SET $\Theta^{(0)} = 0$ has good properties.

⇒ deep learning is undetermined (often)
 AND so INITIALIZATION plays major role!

IN FACT, SGD plays a starring role in
MODERN theory of BIAS VARIANCE

Belkin et. al 2018 "DOUBLE DESCENT"



MODEL complexity
(ABOVE degree)

loss "INTERPOLATING"

first OBSERVED (widely) for DEEP NETS,
BUT ALSO TRUE for CNNs.

SGD Regularizes by picking min norm solution!

MEMORIZATION AND GENERALIZATION!

Other methods of bias & variance

- + DATA Augmentation (SEE STANON L1 Blog on JAIL)
- + Dropout "DATA ADAPTIVE"
- + OPTIMIZATION Algorithms (PROXIMAL POINT METHODS)

LOTS MORE TO GO!

Picking Hyper Parameters

THREE SETS OF LABELED DATA

TRAIN - FIT PARAMETERS

DEV - "FIT" "HYPERPARAMETERS" e.g. λ

TEST (BLIND)

Our first example

for degree $d \in \{0, 1, \dots, k\}$

TRAIN model(d) on TRAIN SET
 $\rightarrow h_d$

SCORE h_d on DEV SET

Pick BEST SCORE, hope for best on TEST

If WE HAVE INFINITELY many models
WE CAN grid SEARCH e.g.

for each $\lambda \in \{0, 10^{-4}, 10^{-3}, 10^{-2}, \dots\}$

\rightarrow SAME PROCESS

Why do we score on DEV, NOT TRAIN?

IMPROVEMENTS

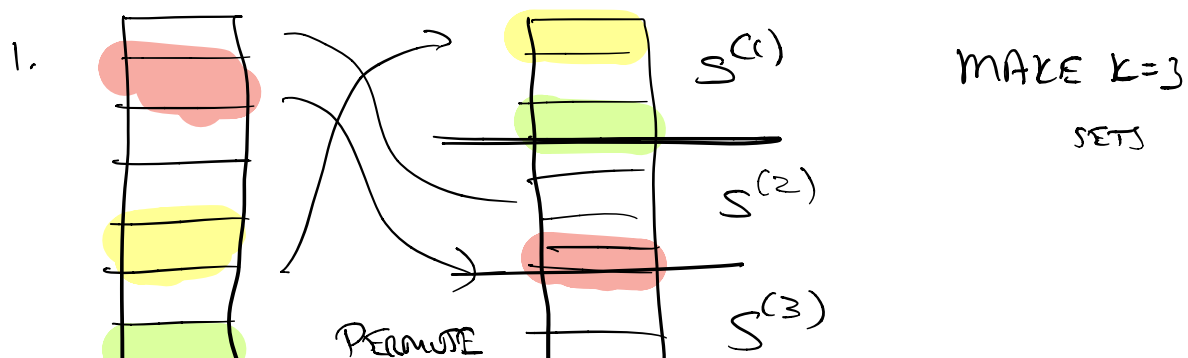
→ Data Efficiency: MAKE BEST USE
of DATA IN TRAIN/DEV
"CLASSICAL STATS"

→ Compute Efficiency: MANY RELATED
HYPER PARAMETERS → MANY MODELS

"MODERN ML SITUATION"
COMBINATORIAL EXPLOSION.

DATA K-fold CROSS VALIDATION

$k=3$ but 5, 10, ... typical





2. TRAIN
 $S^{(1)}, S^{(2)}$
 $S^{(1)}, S^{(3)}$
 $S^{(2)}, S^{(3)}$
- SCORE
 $S^{(3)}$
 $S^{(2)}$
 $S^{(1)}$
3. COMBINE SCORES
 (AVERAGE)
 → USE THIS
 TO PICK BEST.

Computational

Motivation: Regularizer, dropout rate, STEPSIZE, DIMENSIONS - MANY LAYERS!

Practical trick:

1. TUNE 1 PARAMETER AT A TIME
2. SWEEP OVER ALL PARAMETERS

CHOICES

$$2(5 + 6 + 7) < 5 \cdot 6 \cdot 7$$

MORE ADVANCED Hyperband (Jamieson 15)

Run all 5.6.7 models

→ But for just a few steps

Pick top half, run for a 2x steps

⇒ REPEAT ---

EACH ROUND, WE USE SAME NUMBER OF RESOURCES → but with fewer models

Run $\log_2(5-6-7)$ ~~times~~ rounds

lots MORE TO DO HERE (learn across runs?)

RECAP:

BIAS \updownarrow VARIANCE $\begin{matrix} \curvearrowright \text{classical} \\ \rightarrow \text{modern} \end{matrix}$

REGULARIZATION Explicit AND Implicit

Tuning Cross Validation $\&$ Hyperparameter Search