

BIAS \neq VARIANCE: CLASSICAL \neq MODERN Elements

+ CLASSICAL Theory

- Regularization

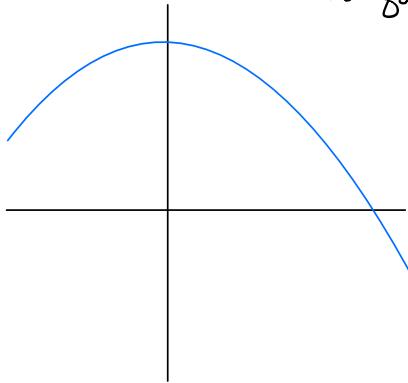
Compute Efficient : Successively

- PARAMETER SELECTION

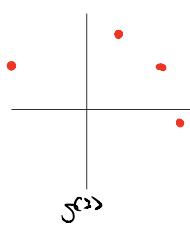
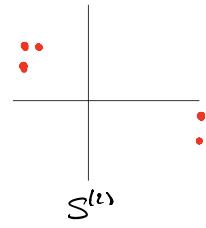
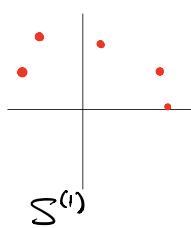
DATA * : K-Fold

+ MODERN Theory (Bonus)

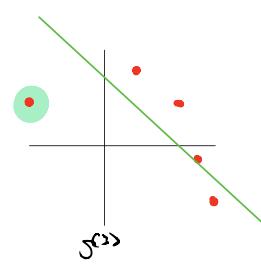
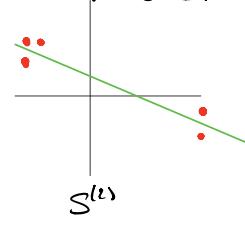
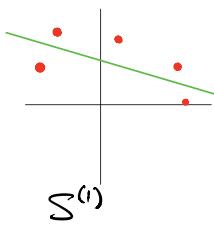
TRUE function $f_0(x) = \theta_2 x^2 + \theta_1 x + \theta_0$



WE DON'T get to SEE h directly
only samples from it!

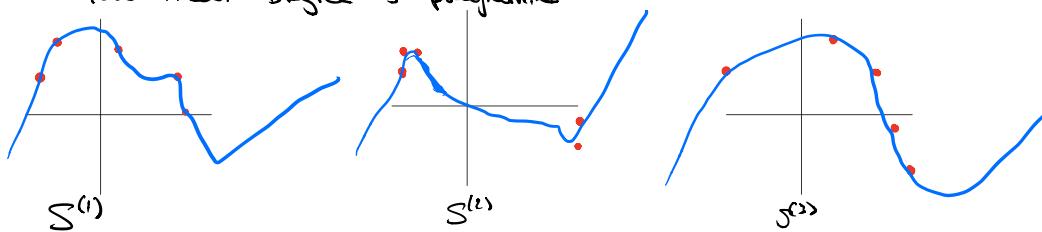


WHAT IF WE FIT A LINE?



Informally, we call this **UNDERFIT**. The error is pretty big (**HIGH bias**)

How about degree \leq polynomial



OVERFIT

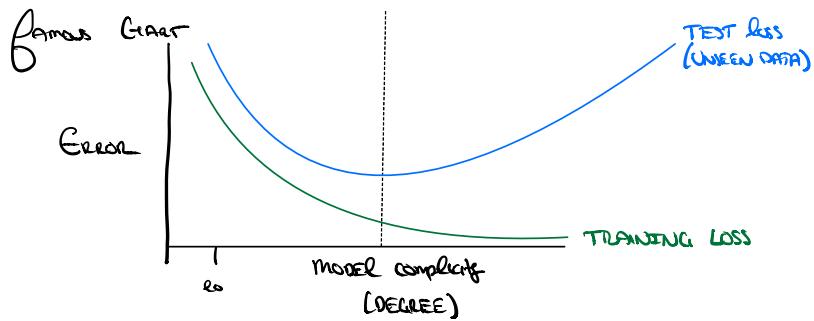
fits the data well but changes on each sample.

(High variance)

Recall h is quadratic, what if we pick that model class?

→ low bias & low variance

→ IT FITS!



UNDERFIT: TRAIN loss \approx TEST loss . But Error is High

OVERFIT : TRAIN loss < TEST loss . But Error maybe low

this is classical BIAS-VARIANCE

- + helpful to understand ML IDEAS
- + INCOMPLETE for modern models (more later)

MORE Formal Bias-Variance TRADEOFF

$$y = h_{\theta}(x) + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

FEATURES (DATA) GAUSSIAN NOISE

PARAMETERS e.g. $h_{\theta}(x) = \theta \cdot x$ (linear model)

$y, \epsilon \in \mathbb{R} \quad \theta, x \in \mathbb{R}^d$

PROCEDURE

Pick an $x \in \mathbb{R}^d$, A TEST POINT (REASON ABOUT THIS LATER)

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

NB: $y^{(i)} = h_{\theta}(x^{(i)}) + \epsilon^{(i)}$ AS ABOVE.

2. TRAIN A MODEL ON S CALL IT $\hat{h}_S: \mathbb{R}^d \rightarrow \mathbb{R}$

3. DRAW A TEST SAMPLE y

$$y = h_{\theta}(x) + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

4. MEASURE $(\hat{h}_S(x) - y)^2$ RISK

WE EXAMINE $\mathbb{E}_{S, \epsilon} [(\hat{h}_S(x) - y)^2]$ NB: TWO INDEPENDENT SOURCES OF RANDOMNESS

Goal: Decompose this Error.

$$\begin{aligned} & \mathbb{E}_{S, \epsilon} [(\hat{h}_S(x) - h_{\theta}(x) + \epsilon)^2] \\ &= \mathbb{E}_{\epsilon, S} [\epsilon^2] + 2 \mathbb{E}_{\epsilon, S} [\epsilon (\hat{h}_S(x) - h_{\theta}(x))] + \mathbb{E}_{\epsilon, S} [(\hat{h}_S(x) - h_{\theta}(x))^2] \\ & \quad \downarrow \qquad \qquad \qquad \text{NB: } \mathbb{E}[\epsilon] = 0 \quad (\epsilon, S \text{ indep}) \\ &= \sigma^2 + 0 + \mathbb{E}_S [(\hat{h}_S(x) - h_{\theta}(x))^2] \end{aligned}$$

UNAVOIDABLE ERROR TERM

DEFINE $\hat{h}_{\text{avg}}(x) \triangleq \mathbb{E}_S [\hat{h}_S(x)]$ "long run avg of S "

ONCE: SELECT AN S , TRAIN \hat{h}_S , EVALUATE ON x

this IS RANDOM VARIABLE, AND HAS EXPECTATION.

$$\mathbb{E}_S \left[(h_S(x) - h_\theta(x))^2 \right] = \mathbb{E}_S \left[(h_S(x) - h_{\text{avg}}(x) + h_{\text{avg}}(x) - h_\theta(x))^2 \right]$$

$$= \mathbb{E}_S \left[(h_S(x) - h_{\text{avg}}(x))^2 \right] + (h_{\text{avg}}(x) - h_\theta(x))^2 + 0$$

VARIANCE OF TRAINING Procedure

BIAS² DOES NOT

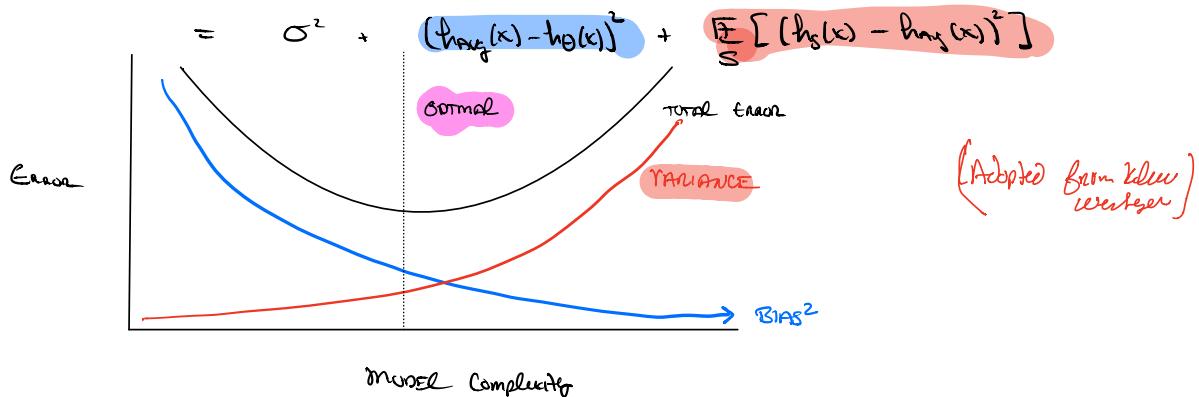
VAR_S(h)

DEPEND ON S \Rightarrow JUST THE MODEL CLASS!

VARIANCE

BIAS

Summary $\mathbb{E}_{S,S} \left[(h_S(x) - h_\theta(x))^2 \right] = \text{UNFAVORABLE ERROR} + \text{BIAS}^2 + \text{VARIANCE}$



NOTE: + Using distinct Dev / Train sets lets us

Assess VARIANCE AND STABILITY

+ If we use EXPRESSIVE model class need to "trust points less"

→ Regularization

REGULARIZATION IS AT THE HEART of Classical & modern Theory!

Regularization

REDUCE VARIANCE TO OBTAIN MORE ROBUST MODEL (TO TRAINING SET VARIATION)

→ EXPLICIT (CHANGE LOSS WE ARE OPTIMIZING)

→ IMPLICIT (BY PRODUCT OF ALGORITHM)

CLASSICAL SETTING

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

↑ λ → **REGULARIZATION PARAMETER**
↑ $\|\theta\|_2^2$ → **PENALTY FOR MODEL COMPLEXITY.**

$\lambda = 0 \Rightarrow$ ORDINARY LEAST SQUARES

$\lambda = 1000 \Rightarrow \theta \approx 0$ PROBABLY LOOKS GOOD!

SET λ TO BALANCE TRADEOFF! (HOW WELL SEE YOU)

Solution? Fix $\lambda > 0$. Then:

$$\nabla_{\theta} \left(\frac{1}{2} (x\theta - y)^T (x\theta - y) + \frac{\lambda}{2} \theta^T \theta \right) = 0$$

$$\Leftrightarrow x^T(x\theta - y) - \lambda\theta = 0 \Leftrightarrow (x^T x + \lambda I)\theta = x^T y \quad (\text{NORMAL EQUATIONS})$$

UNDETERMINED CASE (MODERN PERS.)

$\text{rank}(x^T x) < d$ if $\lambda = 0$ there is NOT A UNIQUE SOLUTION!

↪ ∃ v. s.t. $v \neq 0 \notin (x^T x)v = 0$ hence $(x^T x)/v \neq x^T y$.

If $\lambda = 0$, then it does have unique solution!

Why? $x^T x$ has eigenvalues $\sigma_1^2 \geq \dots \geq \sigma_n^2 \geq 0$ (positive semidefinite)

$$x^T x + \lambda I \quad \sigma_i^2 + \lambda \geq \dots \geq \sigma_n^2 + \lambda \geq 0 \quad \text{Positive Definite}$$

$$\text{so, } \theta_{\lambda} = (x^T x + \lambda I)^{-1} x^T y.$$

BACK TO VARIANCE, DID WE REDUCE IT?

$$\text{vars}(\hat{y}) = \mathbb{E}_s [(\hat{y}_{\text{OLS}}(x) - \hat{y}_{\text{REG}}(x))^2]$$

WE CONSIDERED FIXED DESIGN SETTING (POINTS x FIXED, ONLY LABEL NOISE)

θ_{λ} depends only on choice of $x^T - y$.

$$\begin{aligned}
 &= \mathbb{E}_{\theta \sim P} [((\theta_x(y) - \theta_x(\mathbb{E}(y))) \cdot x)] \leq \mathbb{E}_y [\|\theta_x(y) - \theta_x(\mathbb{E}(y))\|^2] \|x\|^2 \\
 &\quad \text{How much the model changes} \quad (\text{Cauchy-Schwarz}) \\
 \text{Recall } y &= X \cdot \theta + v \quad v \sim N(0, \sigma^2 I) \quad v \in \mathbb{R}^n \quad c \in \mathbb{R} \\
 &\quad \in \mathbb{R}^{n \times d} \\
 &= \mathbb{E}_A [\|A^T x + \lambda I\|^{-1} x^T (y - \mathbb{E}(y))\|^2] \cdot \|x\|^2 \\
 &= \mathbb{E}_A [\|Av\|^2] \quad Av \sim N(0, \sigma^2 A A^T) \\
 &\leq \frac{\sigma^2}{(\sigma^2 + \lambda)^2} \|x\|^2
 \end{aligned}$$

As λ increases, variance goes down.

Bonus Observation Implicitly Regularize as well Important in model theory

Thought Experiment Run gradient descent with $\lambda = 0$

If we initialize at $\theta^{(0)}$ note that

$$\theta_{00} = \underbrace{P_{\text{Model}}(\theta_{00})}_{\text{Claim}} + P_{\text{Norm}}(\theta_{00})$$

$$= P_{\text{Model}}(\theta^{(0)}) \quad \text{why?}$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha x^T (x^T \theta - y)$$

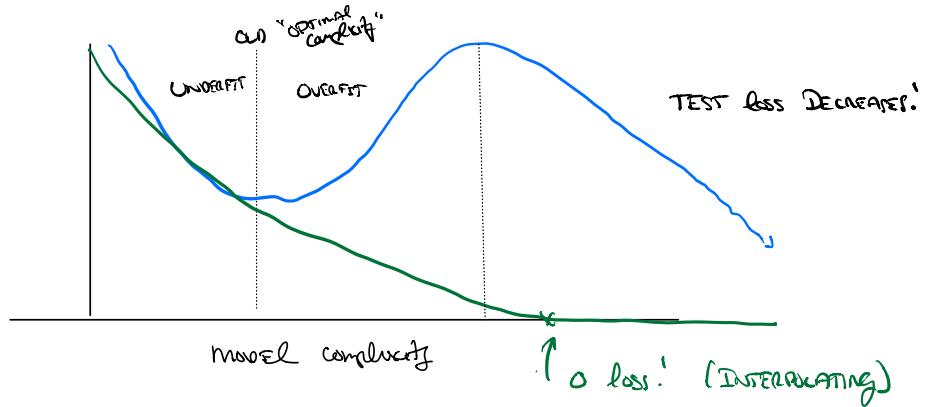
↑ This update is always in $\text{span } x^T$

Observation If we set $\theta^{(0)} = 0$, we get minimum norm solution!

(without regularization!)

Deep learning Underdetermined so initialization plays a major role.

EXTRN Belkin's 2018 "DOUBLE DESCENT"



↑ DESERVED FOR DEEP NETS, ALSO TRUE FOR LEARNERS

→ SGD Regularizes by picking minimum norm solution
⇒ OVERDETERMINED Regn.

Memorization \neq Generalization (CONCRETE SHIFT TO FIELD)

Other methods AMENABLE TO THIS TREATMENT

- + Augmentation
- + Dropout (Adaptive Regularization) MANY MORE!
- + Optimization Algorithms

Picking hyperparameters

THREE SETS of LABELED DATA

TRAIN — FIT PARAMETERS of the model

DEV — "FIT" hyperparameters eg. λ

TEST — Bling for evaluation

Example Pick Degree

For Degree $\in \{0, 1, 2, \dots, k\}$

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

Score Each h_d on DEV SET

Why do we score

Pick BEST score, As h

ON DEV NOT TRAIN?

Hope for best on test set

If we have infinitely many models eg λ , grid

$\lambda \in \{0, 10^{-4}, 10^{-3}, \dots\} \nsubseteq$ SAME ACCES

IMPROVEMENTS TO BASIC SCHEME

DATA Efficiency MAKE BEST USE OF DATA IN TRAIN/DEV/TEST

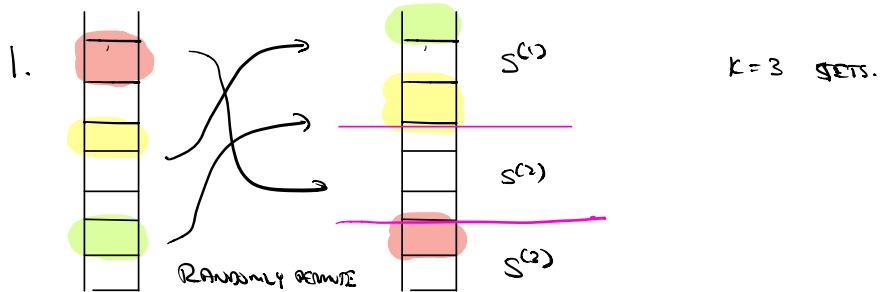
CLASSICAL STATS : K-FOLD CROSSVALIDATION (CV)

COMPUTE Efficiency MANY RELATED HYPERPARAMETERS

MODERN ML : Successive Halving

DATA K-fold CV

K=3 (but 5, 10, 20 are typical)



2.	TRAIN	SCORE	3. COMBINE SCORES (Average, ...)
	$S^{(1)}, S^{(2)}$	$S^{(3)}$	
	$S^{(1)}, S^{(3)}$	$S^{(2)}$	\Rightarrow USE THIS SCORE TO PICK THE BEST
	$S^{(2)}, S^{(3)}$	$S^{(1)}$	

Computational MANY TECHNIQUES

MOTIVATION EXPLOSION OF PARAMETERS WE DON'T KNOW HOW TO SET

REGULARIZER, LEARNING RATES, STEPSIZE, LAYER SIZE...

PRACTICAL TRICK 1

1. TUNE 1 PARAMETER AT A TIME (ω) GRID SEARCH
2. SWEEP OVER RANGE.
 $2(5 + 6 + 7)$ vs. $5 \cdot 6 \cdot 7$
GRID SEARCHES GRID SEARCHES (NARROW)

MORE ADVANCED HYPERSPACE JAMESON '15

- Let $M =$ ALL $5 \cdot 6 \cdot 7$ MODELS $T=1$ (small #)
1. RUN ALL MODELS IN M FOR T STEPS
 2. SCORE ALL MODELS IN M
 3. SET $M =$ TOP $\frac{1}{2}$ OF BEST MODELS
 $T = 2T$

THIS USES COMPUTATIONAL WORK PER ROUND $M \cdot T$ IS CONSTANT

RUNS FOR $\lceil \log_2 M \rceil$ ROUNDS SOME GUARANTEES
... lots more to do here...

RECAP

