# BIAS-VARIANCE

① How do we pick $d$ in $^{(1D)}$ polynomial regression?

One Idea:

$$\hat{y}^{(0)} = w_0$$

$$\hat{y}^{(1)} = w_0 + w_1 x^1$$

$x \in \mathbb{R}^1$

$$\hat{y}^{(d)} = w_0 + w_1 x^1 + \cdots + w_d x^d$$

$$\vec{w}^{(d)} = \arg\min_{w \in \mathbb{R}^d} \sum_{\substack{\text{training} \\ \text{data}}} \left( y_i - \hat{y}_i^{(d)} \right)^2$$

$$d = \arg\min_{d \in \{0, 1, \ldots, 100\}} \sum_{\substack{\text{training} \\ \text{data}}} \left( y_i - \hat{y}^{(d)} \right)^2$$

Won't work too well. Model classes are needed,
$\hat{d} = 100$ will always gives lowest training error.



all $10^{th}$ order polynomials

all $0^{th}$ order

all $9^{th}$ order

This is called the problem of "Model Selection"
→ How do I pick a model class to search in?

② Types of Error

→ Larger Model classes will always do better on training error, but that's not what we care about.

What we really care about → Expected Loss/Error

$$X, Y \sim P(X, Y) \quad [\text{Unknown}]$$

$$E_{P(X,Y)}\left[ L(y, \hat{y}(x; w)) \right]$$

all parameters to hyperparameters

What we ideally want to do

$$\min_{W} \int_{X}\int_{Y} L(y, g(x; w)) \, p(x, y) \, dx \, dy$$

Two problems
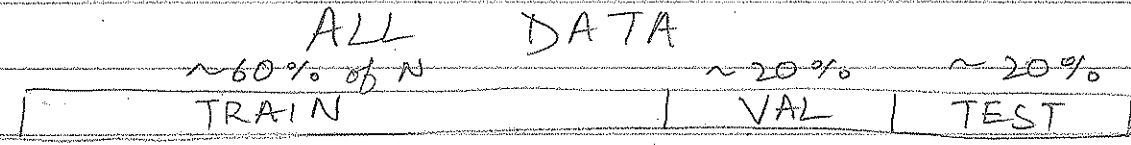→ Integral hard to compute
→ $p(x, y)$ unknown

So let's approximate integral with samples $(x_i, y_i)$
$\sim P(X, Y)$

$$E_{approx}(W) = \frac{1}{N} \sum_{i=1}^{N} L(y_i, g(x_i; w))$$

So here's how we optimize various quantities:

ALL DATA

| TRAIN | VAL | TEST |
|---|---|---|
| ~60% of N | ~20% | ~20% |

$E_{train}$ — Used to fit model parameters

$$\hat{w} = \operatorname*{argmin}_{w} E_{train}$$

$E_{val}$ — Used to choose model classes

$$\hat{d} = \operatorname*{argmin}_{d \in \{0, \ldots, 100\}} E_{val}$$

$E_{test}$ — Used to estimate Expected Error/Loss

No tweaking, No learning on this, otherwise becomes a biased estimate.

If not enough data to do: (train, val) split we do cross-validation.

③ Overfitting vs Under Fitting



$E_{val}$

$E_{train}$

Model Complexity (say d)

$E_{val} \approx E_{train}$ = high $\Rightarrow$ Underfitting
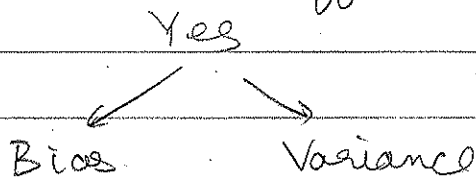
$E_{val} \gg E_{train}$ $\Rightarrow$ Overfitting

Overfitting ⟹ Model class too large, too expressive
Need more assumptions, Need smaller model class

Underfitting ⟹ Model class too small, too simple,
Need fewer assumptions

⟵——————————————⟶

④ Bias - Variance
Eval is high for both overfitting & underfitting
Can we differentiate?
Yes
Bias          Variance

Back to basics: Coin Toss.

$$D = \{x_1 \ldots, x_n\} \sim Ber(\theta^*) \quad (Say\ \theta^* = 0.5)$$
$$\underset{IID}{}$$

| $D$ | $\hat{\theta}_{MLE} = \frac{1}{N}\sum x_i$ | $\hat{\theta}_{lazy} = x_1$ | $\hat{\theta}_{silly} = 1$ |
|---|---|---|---|
| $\{H, H, T\}$ | 2/3 | 1 | 1 |
| $\{T, T, H\}$ | 1/3 | 0 | 1 |
| $\{H, T, H\}$ | 2/3 | 1 | 1 |
| ⋮ | ⋮ | ⋮ | ⋮ |

$$E[\hat{\theta}_{MLE}] = \frac{1}{N}\sum_{i=1}^{N} E[x_i] = \frac{1}{N}\sum_{i=1}^{N}[1 \cdot \theta^* + 0 \cdot (1-\theta^*)]$$

$$= \frac{1}{N}\sum_{i=1}^{N}\theta^* = \theta^*$$

$$\text{bias} = \underbrace{E[\hat{\theta}] - \theta^*}$$

Difference between what your estimator reports
on average to the Truth.

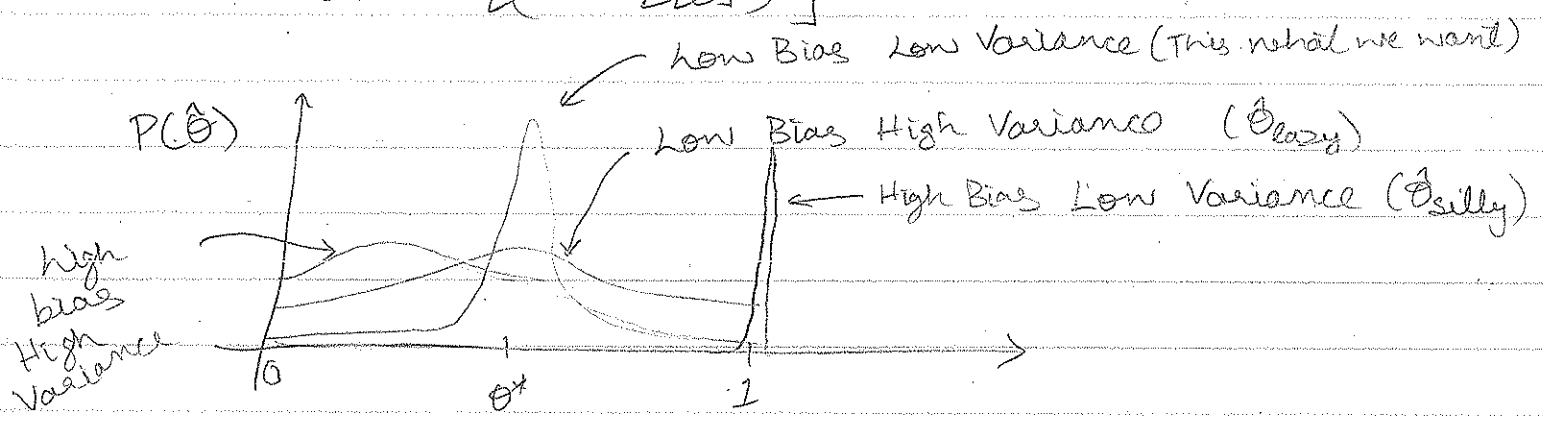$\hat{\theta}_{MLE}$ is UNBIASED $\because$ $E[\hat{\theta}_{MLE}] = \theta^*$

What about others?

$E[\theta_{lazy}] = E[X_1] = \theta^*$    (Also Unbiased)

$E[\theta_{silly}] = E[1] = \underline{1}$    (Biased if $\theta^* \neq 1$)

High
Variance     $\downarrow$ No variance

---

$$\text{Variance} = E\left[\left(\hat{\theta} - E[\hat{\theta}]\right)^2\right]$$

Low Bias Low Variance (This what we want)

Low Bias High Variance ($\hat{\theta}_{lazy}$)

High Bias Low Variance ($\hat{\theta}_{silly}$)

$P(\hat{\theta})$

High
bias
High
Variance

0       $\theta^*$       1

(5) Bias-Variance Decomposition for Squared-Loss

Say we want to estimate $\theta^*$

your estimator's avg value

Let $\bar{\theta} \equiv E[\hat{\theta}]$

Now $E[Loss] = E[(\hat{\theta} - \theta^*)^2]$

$$= E[\hat{\theta} - \bar{\theta} + \bar{\theta} - \theta^*{}^2]$$

$$= E[(\hat{\theta} - \bar{\theta})^2 + (\bar{\theta} - \theta^*)^2 - 2(\hat{\theta} - \bar{\theta})(\bar{\theta} - \theta^*)]$$

$$= E[(\hat{\theta} - \bar{\theta})^2] + E[(\bar{\theta} - \theta^*)^2] - 2\underbrace{(\bar{\theta} - \theta^*)}_{constants} E[(\hat{\theta} - \bar{\theta})]$$

$$= Var(\hat{\theta}) + E[bias^2] - 2(\bar{\theta} - \theta^*)\underbrace{(E[\hat{\theta}] - \bar{\theta})}_{=0}$$

$$= Var(\hat{\theta}) + bias^2$$

So $E[Loss] = bias^2 + Var(\hat{\theta})$

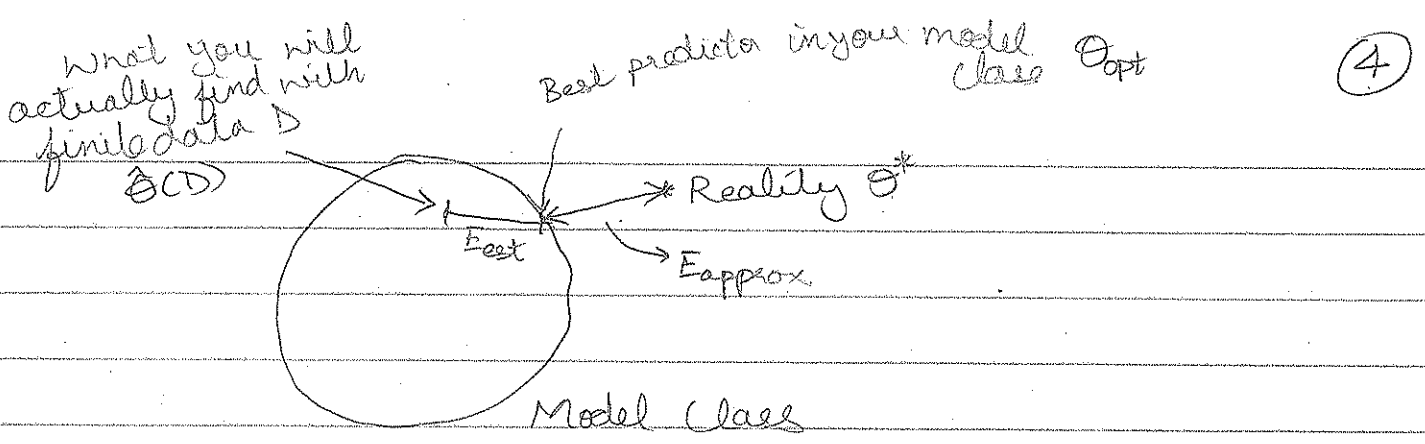2 predictors (one with high bias, one with high loss can have the same $L_2$ loss)

⟶

(6) General Error Decomposition

Expected Error = "Approximation Error" (or ≈ Bias)
→ you approximated reality with your model

+ "Estimation Error" (or ≈ Variance)
→ you estimated with finite data

+ "Optimization Error"
→ you didn't / couldn't maximize MLE exactly
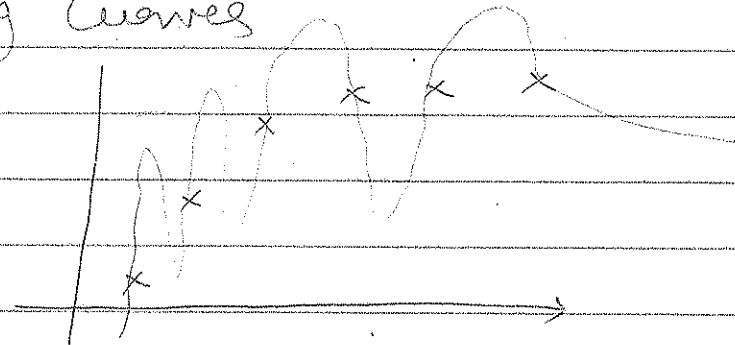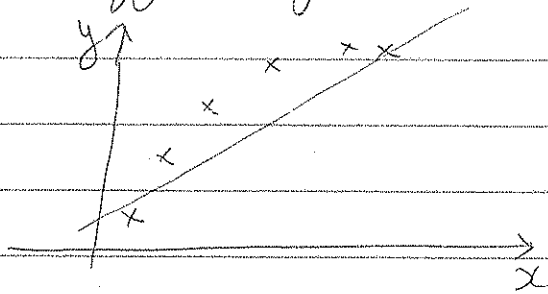
+ Bayes Error
→ lower bound on error

What you will actually find with finite data $D$
$\hat{\theta}(D)$

Best predictor in your model class $\theta_{opt}$

$\rightarrow$ * Reality $\theta^*$

$E_{est}$

$\rightarrow E_{approx}$

Model Class

$\longleftrightarrow$

⑦ Model Selection via Regularization
(Let's put in a preference for simple models)

$$\min_{W} \sum_{train} L(y_i, \hat{y}(x_i; W)) + \lambda \|W\|_2^2$$

Can also be viewed as a Gaussian prior on $W$
(MAP estimator)

⑧ Effect of Data / Learning Curves
                 More

$y$
×  ×  ×
    ×
  ×
×
$x$

Simple Model Class
$\Rightarrow$ More dat won't help

Eval
$E_{train}$

$N$
Data

×  ×  ×  ×  ×
  ×
×

Too Expressive Model Class
$\Rightarrow$ More data will help, but very slowly.

Eval

$E_{train}$

$N$
(Data)