

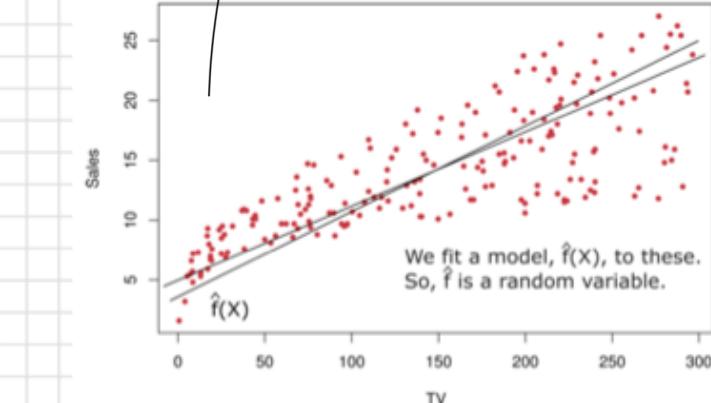
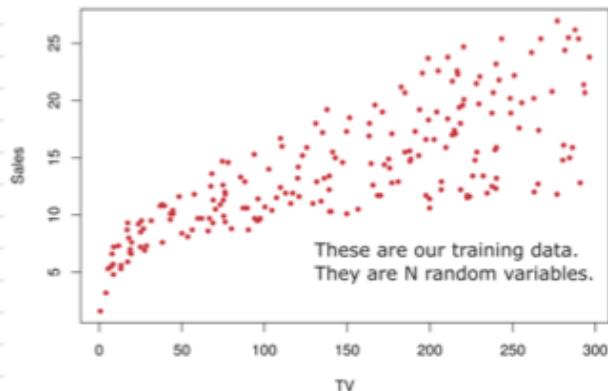
# LECTURE 3 - 03/09/24

- Quan notions

THE  $N$  OBSERVATIONS WE REFER TO AS **TRAINING DATA**

THE ALGORITHM USED FOR FITTING THE FUNCTION IS CALLED **LEARNER**

APPLYING THIS **LEARNER** TO THE **TRAINING DATA**, WE CAN **PREDICT**



- WHAT IS OUR GOAL?

OUR GOAL IS TO **PREDICT UNSEEN DATA POINTS:  $(x_0, y_0)$**

IMAGINE WE HAVE 2 MODELS FITTED TO THE DATA:

$$f_1(x)$$

$$f_2(x)$$

WE NOW **GIVE THEM AN UNSEEN DATA POINT  $(x_0, y_0)$**

HOW DO WE KNOW **WHICH OF THE TWO IS BEST?**

b) IMAGINE NOW WE GIVE THEM  $n$  UNSEEN DATA POINTS  $(x_i, y_i)$  FOR  $i=1\dots n$

**WHICH OF THE TWO IS BEST?**

- THE ANSWER TO PREVIOUS QUESTION...

**MEAN SQUARED ERROR**

RELIES ON THE **EXPECTED TEST MSE**:

$$\mathbb{E} \left[ (y_0 - \hat{f}(x_0))^2 \right]$$

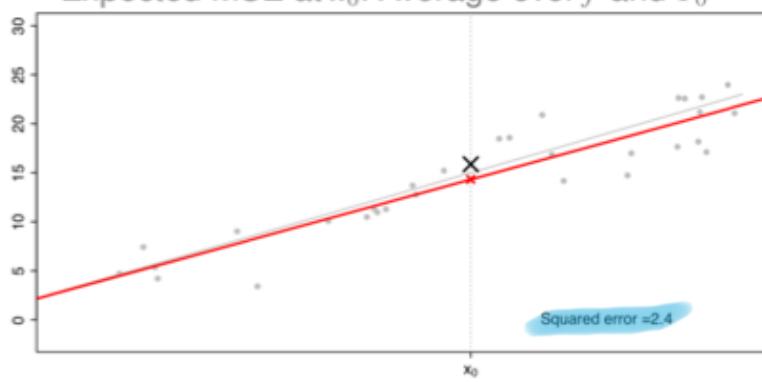
WE WANT THIS TO BE AS SMALL AS POSSIBLE.

THIS IS BASICALLY THE **POSITIVE DIFFERENCE** BETWEEN THE **REAL  $y_0$**  & OUR **PREDICTED  $y$**  COMPUTED BY OUR MODEL  $f(x_0)$

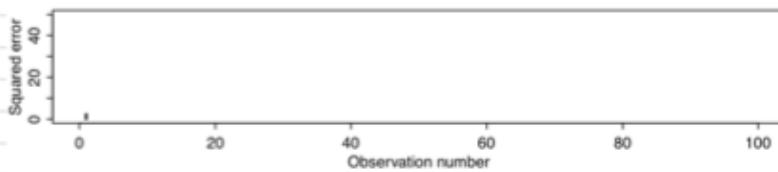
**REMEMBER A**

**THE UNSEEN DATA IS WHAT WE CARE ABOUT**

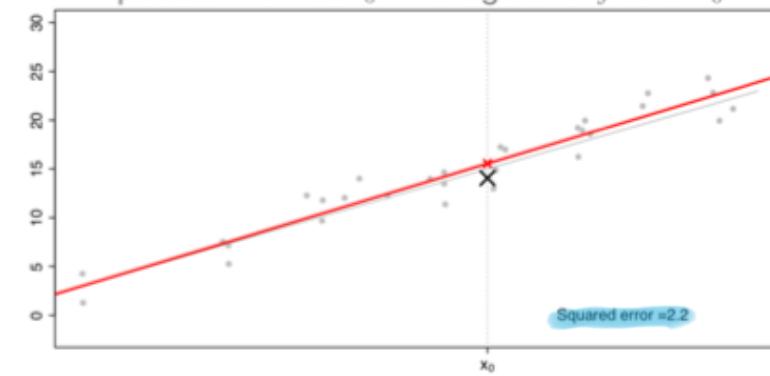
### Expected MSE at $x_0$ : Average over $\hat{f}$ and $Y_0$



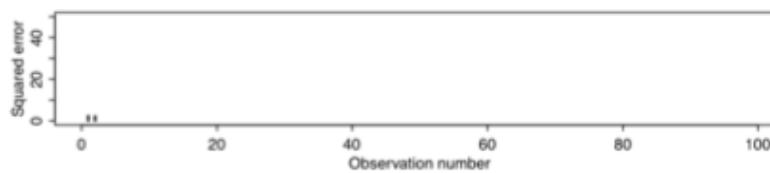
- In this case our  $f(x)$  predicted  $\hat{y}$  with a 1.2 unit difference from the real  $y$
- ↳ This leads to our squared error = 2.4



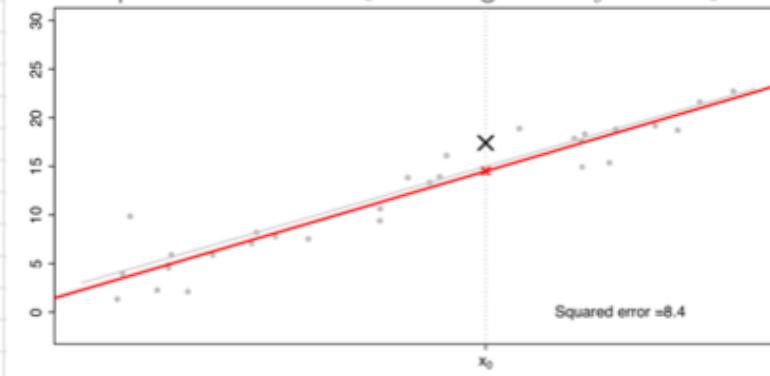
### Expected MSE at $x_0$ : Average over $\hat{f}$ and $Y_0$



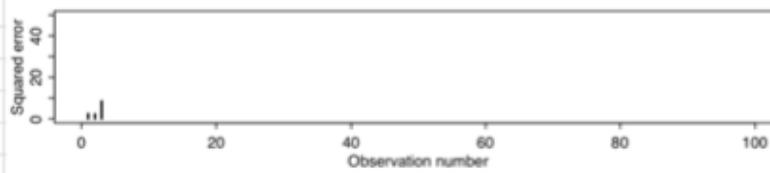
- In this case our  $f(x)$  predicted  $\hat{y}$  with a 1.1 unit difference from the real  $y$
- ↳ This leads to our squared error = 2.2



### Expected MSE at $x_0$ : Average over $\hat{f}$ and $Y_0$

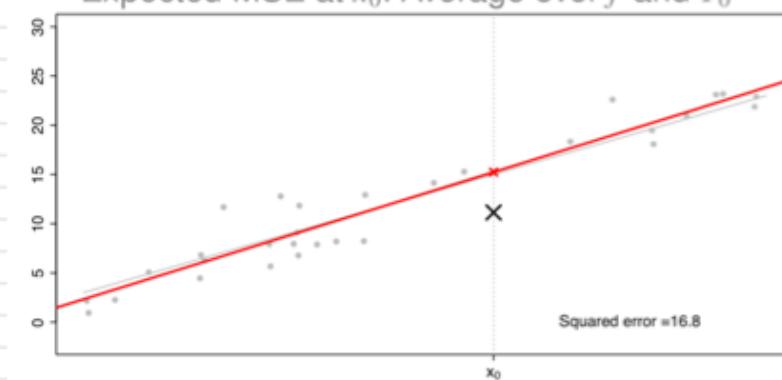


- In this case our  $f(x)$  predicted  $\hat{y}$  with a 4.2 unit difference from the real  $y$
- ↳ This leads to our squared error = 8.4

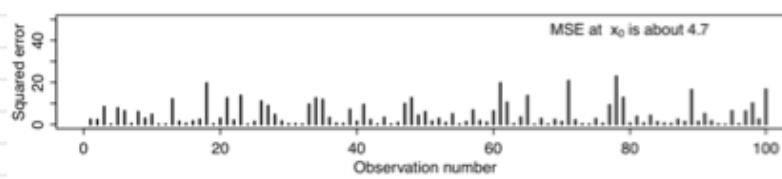


----- We make lots of them (watch slides)

Expected MSE at  $x_0$ : Average over  $\hat{f}$  and  $Y_0$



WE END UP WITH THIS.



- WE CAN NOW THINK ABOUT THE FOLLOWING QUESTION:

What is the expected squared error that we get from the entire process of first training the learner on some training data and then using it for predicting the response  $Y_0$  for a new observation  $X_0$ ?

Remember:

- The function  $\hat{f}$  is a random variable, because it is obtained from applying a learner to a dataset of  $N$  random variables  $(X, Y)$ .
- A new observation  $(X_0, Y_0)$  is a random variable.

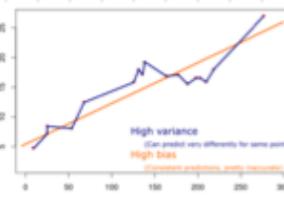
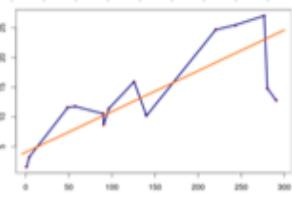
### Bias-Variance decomposition

- THE EXPECTED TEST MSE AT  $X_0$  CAN BE DECOMPOSED AS:

$$\mathbb{E} \left( Y_0 - \hat{f}(x_0) \right)^2 = \underbrace{\mathbb{E} \left( \hat{f}(x_0) - \mathbb{E} \hat{f}(x_0) \right)^2}_{\text{Variance of } \hat{f}(x_0)} + \underbrace{\left( \mathbb{E} \hat{f}(x_0) - f(x_0) \right)^2}_{\text{Bias of } \hat{f}(x_0)} + \text{Var}(\epsilon)$$

First term is the variance

RECALL: MEASURE OF SPREAD IN DATA FROM IT'S MEAN.  
AMOUNT OF CHANGE DUE TO NEW SUBSETS OF DATA



Second term is the expected deviation of our model prediction from the true value (bias)!

RECALL: DIFF BETWEEN MODEL PREDICTED VALUE & REAL VALUE.

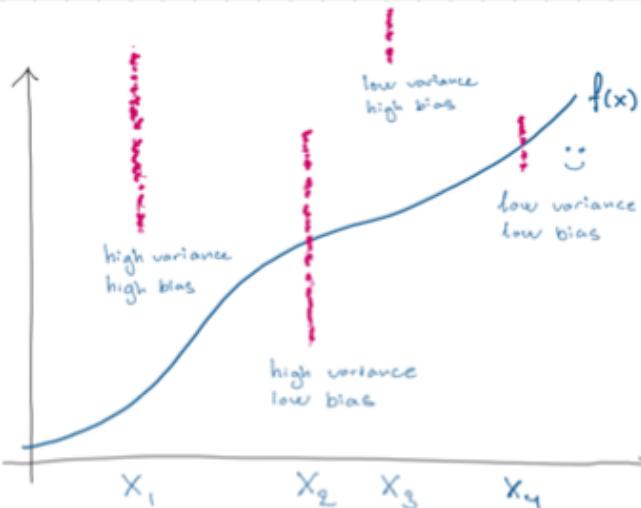
→ PROBLEMS:

OVERTHINKING  
UNDERFITTING

From ISL:

"Variance refers to the amount by which  $\hat{f}$  would change if we estimated it using a different training data set."

"Bias refers to the error that is introduced by approximating a real-life problem, which may be extremely complicated, by a much simpler model."



(In this illustration, red dots are predictions made by different instances of  $\hat{f}$ )

As said, **PREDICTING BOTH LOW & HIGH VARIANCE IS VERY DIFFICULT THAT'S WHY WE TALK ABOUT BIAS - VARIANCE TRADEOFF**

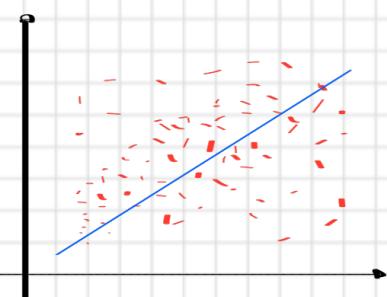
$$\mathbb{E} \left( Y_0 - \hat{f}(x_0) \right)^2 = \underbrace{\mathbb{E} \left( \hat{f}(x_0) - \mathbb{E} \hat{f}(x_0) \right)^2}_{\text{Reducible error}} + \underbrace{\left( \mathbb{E} \hat{f}(x_0) - f(x_0) \right)^2}_{\text{Irreducible error}} + \underbrace{\text{Var}(\epsilon)}$$

All three terms are non-negative, so if any is large, the MSE is large.

- The **Reducible Error** can be lowered by using an estimator  $\hat{f}$  that has both **LOW VARIANCE & LOW BIAS**
- The **Irreducible Error** is a **lower bound** on the accuracy of our prediction  $\hat{Y}$

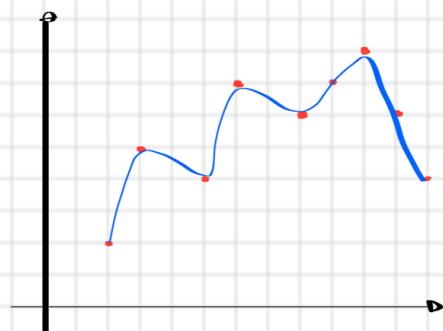
We can think about two examples:

Fit a **constant line**



Low Variance  
High Bias

Fit a **wave that interpolates all points**



High Variance  
Low Bias

\* Note

It is very easy to find such an  $\hat{f}$  with either very low bias or very low variance. But less straightforward to find one with **BOTH**

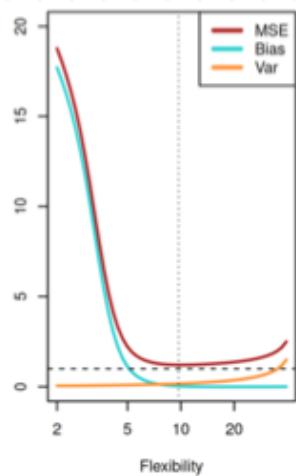
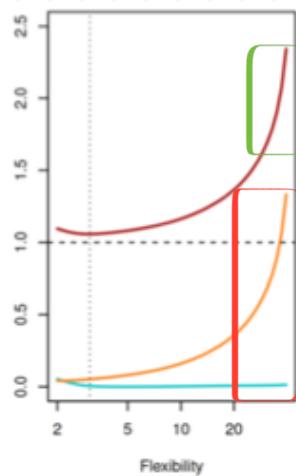
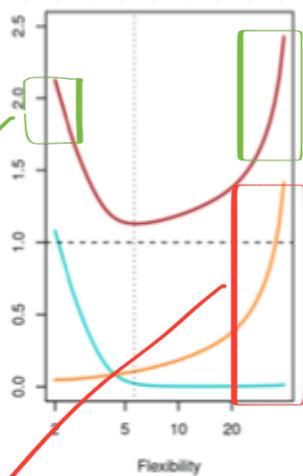
• RECALL PROFESSOR EXAMPLE ABOUT THE DIFFERENT SCENARIO

• IT IS VERY DIFFICULT TO GET **LOW BIAS & LOW VARIANCE**

• BY PREDICTING LOW BIAS & LOW VARIANCE WE PREVENT THE TWO PROBLEMS:

- OVERRFITTING** (**SHOES TOO MUCH TO TRAIN**)
- UNDERFITTING** (**NOT ENOUGH TRAIN**)

- We can further analyze Bias, Variance & MSE with respect to complexity



"Flexibility" is model complexity.

Panels are 3 different data sets.

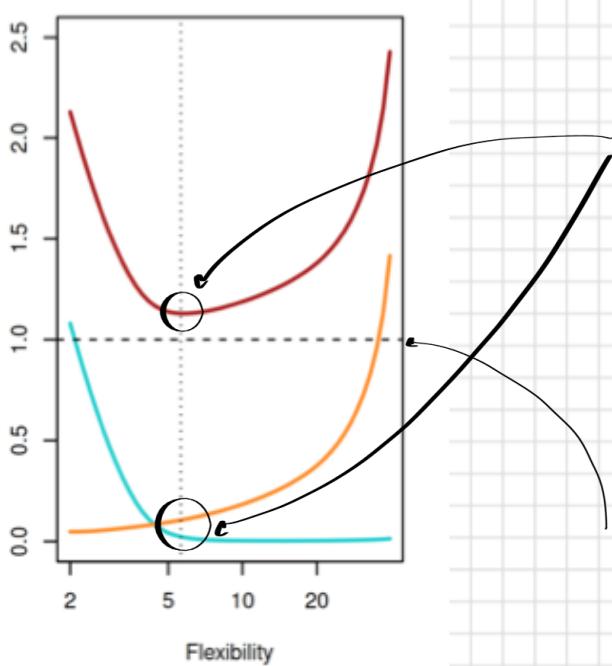
Q: What do these 3 panels have in common?

- From the three graphs we see multiple things:

MSE is very HIGH when model is too simple or too complex  
 Remember!!! → UNDERFITTING & OVERFITTING

Variance grows as complexity arises if it is the inverse  
 For Bias that decreases as model gets more complex.  
 Remember!!! → Easy to find inverse match, very difficult to find both now

- So where do we want our MSE to be?

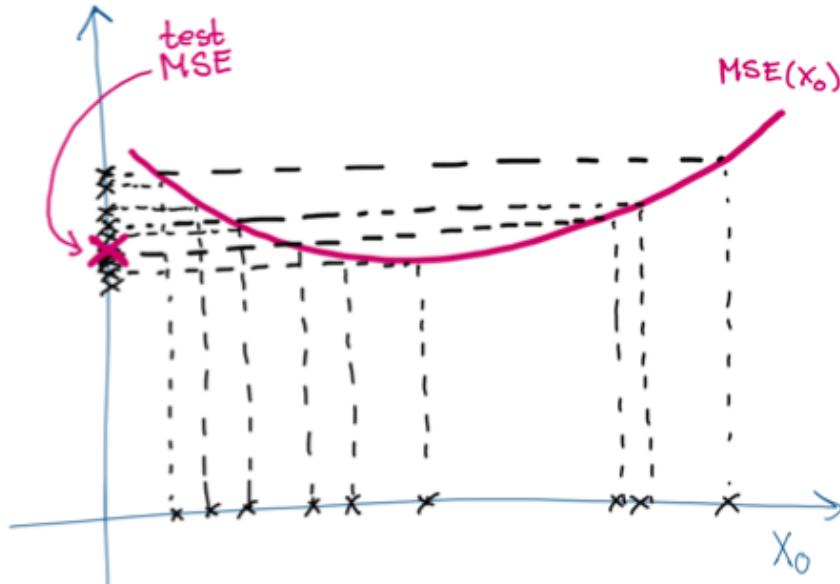


Sweet Spot

This is called SWEET SPOT, we see that we have a GOOD TRADEOFF between VARIANCE & BIAS & THE MSE IS AT IT'S MINIMUM.

This is the irreducible variance, an error that we can't minimize

- So far we have considered the test MSE at an individual point  $x_0$
- To have made a great model fit, we want to do well in **ALL** points



### Expected test MSE

- What we want to do is **take expectation** over new values  $x_0, y_0, \hat{f}$
- This can be done sequentially by:
  - ① First finding the MSE at  $x_0$ , which is then a function of  $x_0$
  - ② After **FINDING** the expectation of the MSE at  $x_0$  when it varies:

$$\mathbb{E}(MSE(\hat{f}, Y_0, X_0)) = \mathbb{E}_{X_0} \left( \mathbb{E}_{\hat{f}, Y_0}(MSE(\hat{f}, Y_0, X_0) | X_0 = x_0) \right)$$

- By doing this we take **expectations** over  $X_0$  in the bias-variance decomposition that result into:

$$MSE = \underbrace{\mathbb{E}(\text{Var}(\hat{f}(X_0)))}_{\text{"Average" variance of } \hat{f}(x_0) \text{ across the range of } X_0.} + \underbrace{\mathbb{E}(\text{Bias}(\hat{f}(X_0)))^2}_{\text{"Average" squared bias of } \hat{f}(x_0) \text{ across the range of } X_0.} + \underbrace{\text{Var}(\epsilon)}_{\text{Residual variance}}$$

### Test MSE vs Training MSE

- The MSE is a theoretical quantity that we can **estimate** from data
- What we do is we **use** data points that were **NOT** used during **TRAINING** (train set) so this is always due to the **UNDERFITTING** & **OVERFITTING** problems

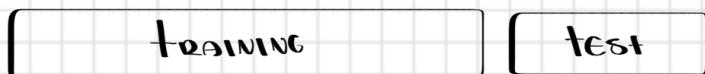
Always **KEEP IN MIND** (⌚) this distinction

**Model Selection**: Estimating the prediction error of different models aiming to choose best one.

**Model Assessment**: Having chosen a final model, estimating its prediction error.

## HAVING OUR DATA FOR TESTING

- In linear regression model we used a single dataset to train the models & selected between models using p-values or AIC.



- This is very good because we lefted out data that we can now give to our model & test it has never seen.
- Sometimes we have Hyperparameters to specify in our model:  
i.e.  $\lambda \in [0, 1]$  where a specific val gives best-performing model.

### Question

How can we train model for different  $\lambda$ 's & choose the best?

- The answer is to change the dataset division into:



From here we can retrieve  $\lambda$ .

- The process now sees TRAINING, TUNING hyperparams, TEST.
- A typical split is:

TRAINING : 60%

VALIDATION : 20%

TEST : 20%