

Bias-Variance Trade-off

Reading: - Murphy 4.7.6.  
- MIT Notes 2.8

Code Demo:

bias-variance\_demo.ipynb

- Let's start by recalling that for:

Optimal Prediction

, the goal is:

Predict  $Y \in \mathbb{R}$  given  $X \in \mathbb{R}^d$  if  $(X, Y) \sim P_{XY}$

← true distribution out of which dataset  $D$  is sampled from. We do not have access to it (do not know it)

- We want to find function  $\eta$  that minimizes:

$$E_{XY} [(Y - \eta(X))^2] = E_X [E_{Y|X} [(Y - \eta(X))^2 | X=x]]$$

also known/called the risk  $R(\eta)$

↑ we condition on X trick

(See Mostafa, pp.61)

Outer expectation is just an average over distribution  $X \Rightarrow$  to minimize  $R(\eta)$ , it suffices to minimize the inner conditional expectation for each fixed value  $X=x$ .

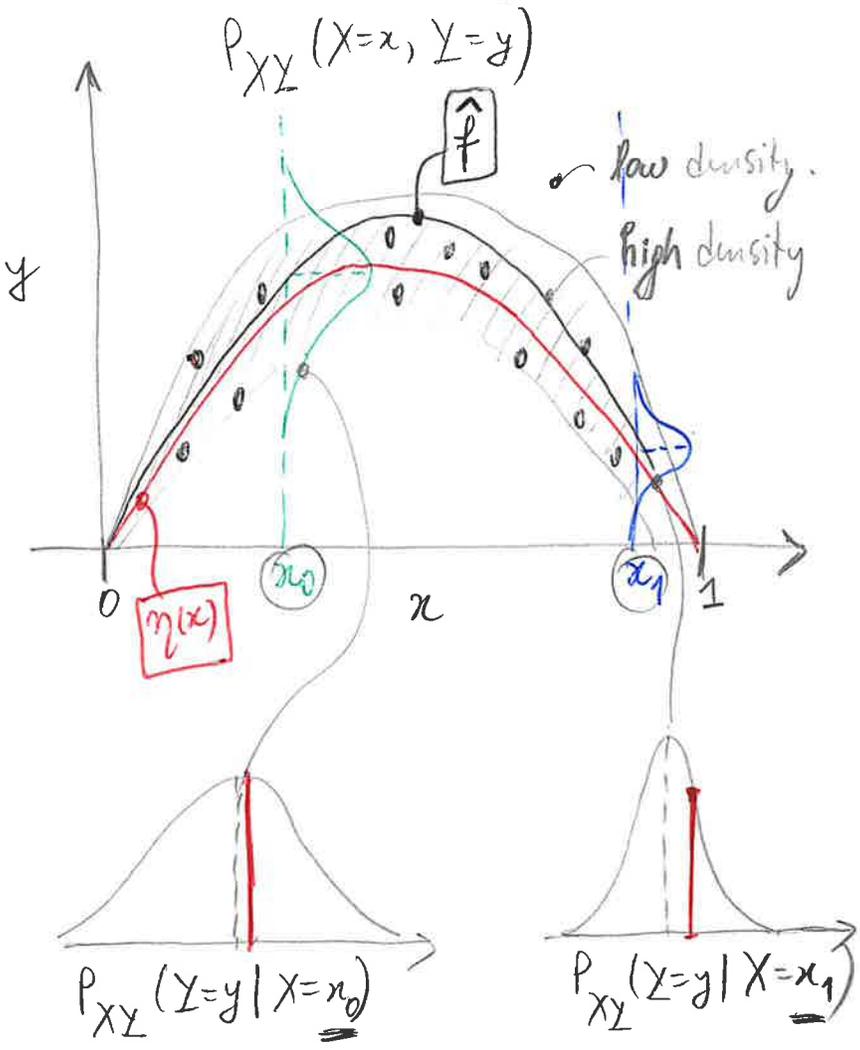
For any fixed value  $X=x$ ,  $\eta(x) = c_x$ : a simple number because once  $x$  is fixed dependence on  $\eta$  is only through the scalar value  $\eta(x)$ .

$c_x$  minimizes  $E_{Y|X} [(Y - c_x)^2 | X=x]$ , so, to get it, take derivative and set to 0:

$$\begin{aligned} 0 &= \frac{d}{dc_x} E_{Y|X} [(Y - c_x)^2 | X=x] \\ &= E_{Y|X} \left[ \frac{d}{dc_x} (Y - c_x)^2 | X=x \right] \\ &= E_{Y|X} [-2(Y - c_x) | X=x] \\ &= -2 E_{Y|X} [Y | X=x] + 2 c_x \end{aligned}$$

$\Rightarrow$  Squared Error Optimal Predictor:  $c_x = \eta(x) = E_{Y|X} [Y | X=x]$

# Statistical Learning



- Find  $\eta$  that minimizes:  
 $E_{X,Y} [(Y - \eta(X))^2]$

- Ideally, we want to find:  
 $\eta(x) = E_{Y|X} [Y | X=x]$

- But, we only have samples:  
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{X,Y}, i=1,2,3,\dots,n$

- And are restricted to a function class (e.g. linear)

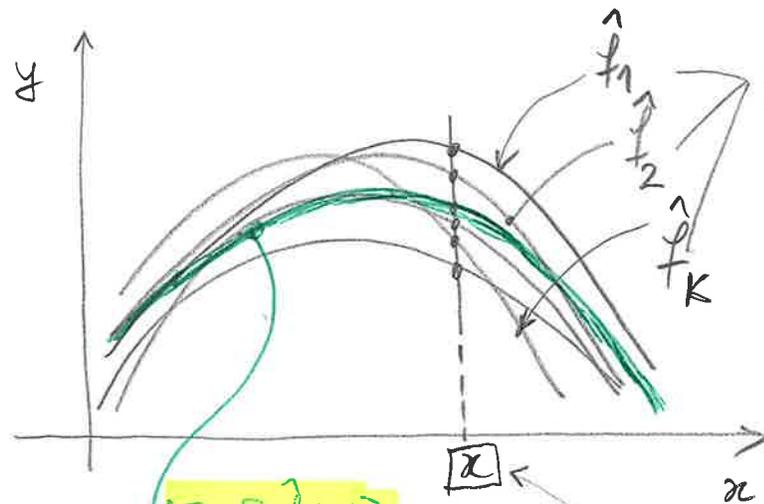
- so, we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$\hat{f} \neq \eta(x)$ , but this approach is the best we can do.

We care about future predictions:  $E_{X,Y} [(Y - \hat{f}(X))^2]$

Note: Each sample/draw  $\mathcal{D} = \{ (x_i, y_i) \}_{i=1}^n$  results in different  $\hat{f}$



obtained by applying the learning algorithm on each  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

(See Mostefa pp 63)  
(See Bishop pp. 149)

$E_{\mathcal{D}}[f(x)]$

Gives an "average function", which can be interpreted in the following operational way:

- Generate many datasets  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$  and apply learning algo to produce  $\hat{f}_1, \hat{f}_2, \dots, \hat{f}_K$

- Estimate "average function" for any  $x$  by:  $\hat{f}(x) \approx \frac{1}{K} \sum_{j=1}^K \hat{f}_j(x)$

- We are viewing  $\hat{f}(x)$  as a random variable; with the randomness coming from the randomness in the dataset.

-  $E_{\mathcal{D}}[\hat{f}(x)]$  is the expected value of this random variable (for a particular  $x$ )

Note: The "average function"  $\hat{f}(x)$  need not be in the model's hypothesis set, even though it is the average of functions that are.

### Empirical Reconstruction Error. (ERE)

- Given a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ , the ERE is the finite-sample average:

$$\text{ERE}(f; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \quad : \text{minimized in practice}$$

- This is what we actually compute from data.
- In Statistical Learning theory, it is called the "empirical risk"

### Expected (population) Reconstruction Error

- if instead we take the expectation under the true data-generating  $P_{XY}$ , we get the population risk:

$$R(f) = \mathbb{E}_{XY} [(Y - f(X))^2] \quad : \text{analyzed in theory}$$

- This is what we would minimize if we had access to the true distribution!!
- Its minimizer is the Bayes-optimal predictor  $\eta(x) = \mathbb{E}[Y|X=x]$

### Relationship between the two:

- $\text{ERE}(f, \mathcal{D})$  is a Monte Carlo approximation of  $R(f)$
- By the Law of Large Numbers, for fixed  $f$

$$\text{ERE}(f, \mathcal{D}) \xrightarrow[n \rightarrow \infty]{} R(f)$$

• So, minimizing ERE is the empirical counterpart of minimizing the true squared error loss/risk !!

**Proof Sketch:** • Start w/  $\text{ERE}(f; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$

• Take expectation wrt data-generating distribution:

$$\begin{aligned} \mathbb{E}_{XY} [\text{ERE}(f; \mathcal{D})] &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(y_i - f(x_i))^2] = \\ & \quad \uparrow \text{by linearity of expectation} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(y - f(x))^2] = \mathbb{E} [(y - f(x))^2] \end{aligned}$$

by IID:  $\mathbb{E}(x_i), \forall i \equiv \mathbb{E}(x), \forall x$

# Bias-Variance Tradeoff Decomposition

- We wish to break down the squared error of model  $\hat{f}(\cdot)$  that makes a prediction of our target  $y$  given input data point  $x$ .

- True predictor  $E[Y|X]$  not available; unknown to us:

$$\eta(x) = E_{Y|X}[Y|X=x]$$

- But, for given dataset  $\mathcal{D}$ , we can apply learning algo to produce

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

And the idea is to use this as a substitute/surrogate for the ground true predictor.

- The **generalization error** we want to minimize is the usual expected squared error:

$$E_{X,Y}[(Y - \eta(x))^2] \approx E_{X,Y}[(Y - \hat{f}(x))^2] \stackrel{\text{use tower trick}}{=} \text{based on what we discussed on paper 2,3,3}$$

(see also MIT Notes pp.23)

$$= E_{Y|X} [ E_{\mathcal{D}} [(Y - \hat{f}_{\mathcal{D}}(x))^2] | X=x ]$$

to acknowledge dependence of our fitted model  $\hat{f}$  on the dataset  $\mathcal{D}$ .

$$= E_{Y|X} [ E_{\mathcal{D}} [ \underbrace{(Y - \eta(x))}_a + \underbrace{\eta(x) - \hat{f}_{\mathcal{D}}(x)}_b ]^2 | X=x ]$$

$$= E_{Y|X} [ E_{\mathcal{D}} [ \underbrace{(Y - \eta(x))^2}_{a^2} + \underbrace{2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x))}_{2ab \rightarrow 0 \text{ (see next page)}} + \underbrace{(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2}_{b^2} | X=x ]$$

$$= E_{Y|X} [(Y - \eta(x))^2 | X=x] + E_{\mathcal{D}} [ (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2 ]$$

**irreducible error**

Caused by (stochastic) label noise

**learning error**

we want to minimize, coz we can control!!!

Caused by:

- either using too 'simple' model
- or not enough data to learn model accurately

Show:  $2ab \rightarrow 0$

$$E_{Y|X} \left[ E_D \left[ \underbrace{2(Y - \eta(x))}_{\text{no } D} \underbrace{(\eta(x) - \hat{f}_D(x))}_{\text{no } Y} \right] \mid X=x \right]$$

$$E_{Y|X} \left[ \cancel{E_D} \left[ 2(Y - \eta(x)) \right] \mid X=x \right]$$

$$E_{Y|X} \left[ Y - E[Y|X=x] \mid X=x \right] \stackrel{\eta(x) = E_{Y|X}[Y|X=x]}{=} \dots$$

$$= E_{Y|X} [Y|X=x] - E_{Y|X} [Y|X=x]$$

$$= 0$$

$Y$  &  $D$  are independent of each other w.r.t. the sources of randomness we are considering.

---

Now, looking at the **Learning error** separately:

$$E_D [(\eta(x) - \hat{f}_D(x))^2] = E_D [(\eta(x) - \underbrace{E_D[\hat{f}_D(x)]}_a + \underbrace{E_D[\hat{f}_D(x)] - \hat{f}_D(x)}_b)^2] =$$

added and subtracted trick

$$= E_D [(\eta(x) - E_D[\hat{f}_D(x)])^2 + 2(\eta(x) - E_D[\hat{f}_D(x)]) \cdot (E_D[\hat{f}_D(x)] - \hat{f}_D(x)) + (E_D[\hat{f}_D(x)] - \hat{f}_D(x))^2]$$

$a^2$ 
 $ab \rightarrow$  because  $b \rightarrow 0$ 
 $b^2$

$$= (\eta(x) - E_D[\hat{f}_D(x)])^2 + E_D [(E_D[\hat{f}_D(x)] - \hat{f}_D(x))^2]$$

bias squared
Variance

(see MIT Notes pp. 23)  
(see Bishop, pp. 149)

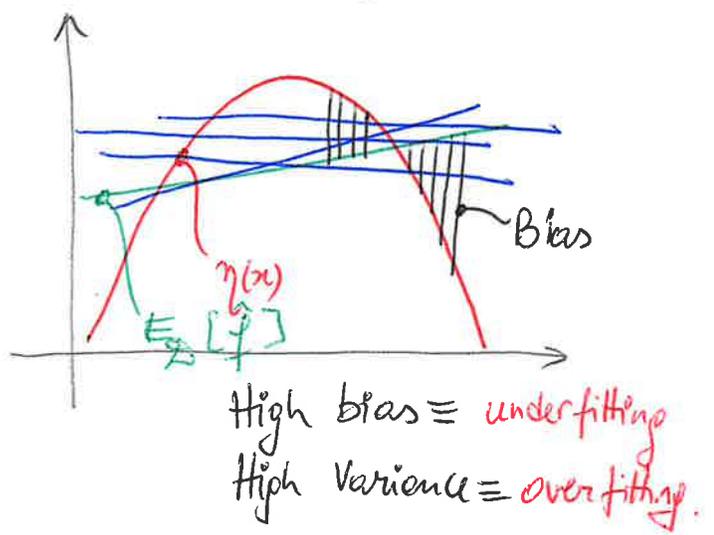
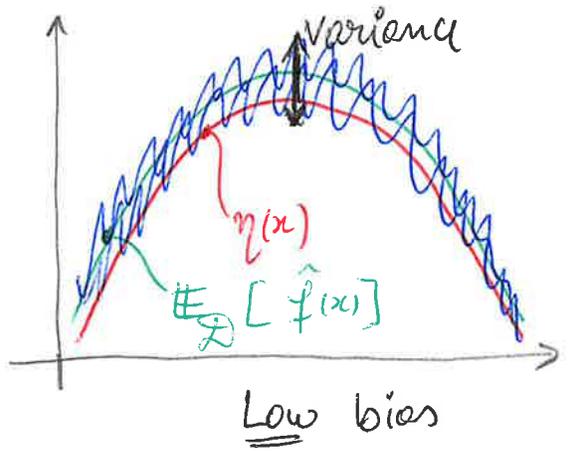
- we arrived at important result:

$$\text{Learning error} = \text{bias}^2 + \text{Variance}$$

**Bias** - Represent the extent to which the average prediction over **all** datasets differs from the ideal regression function/predictor  $\eta(x)$

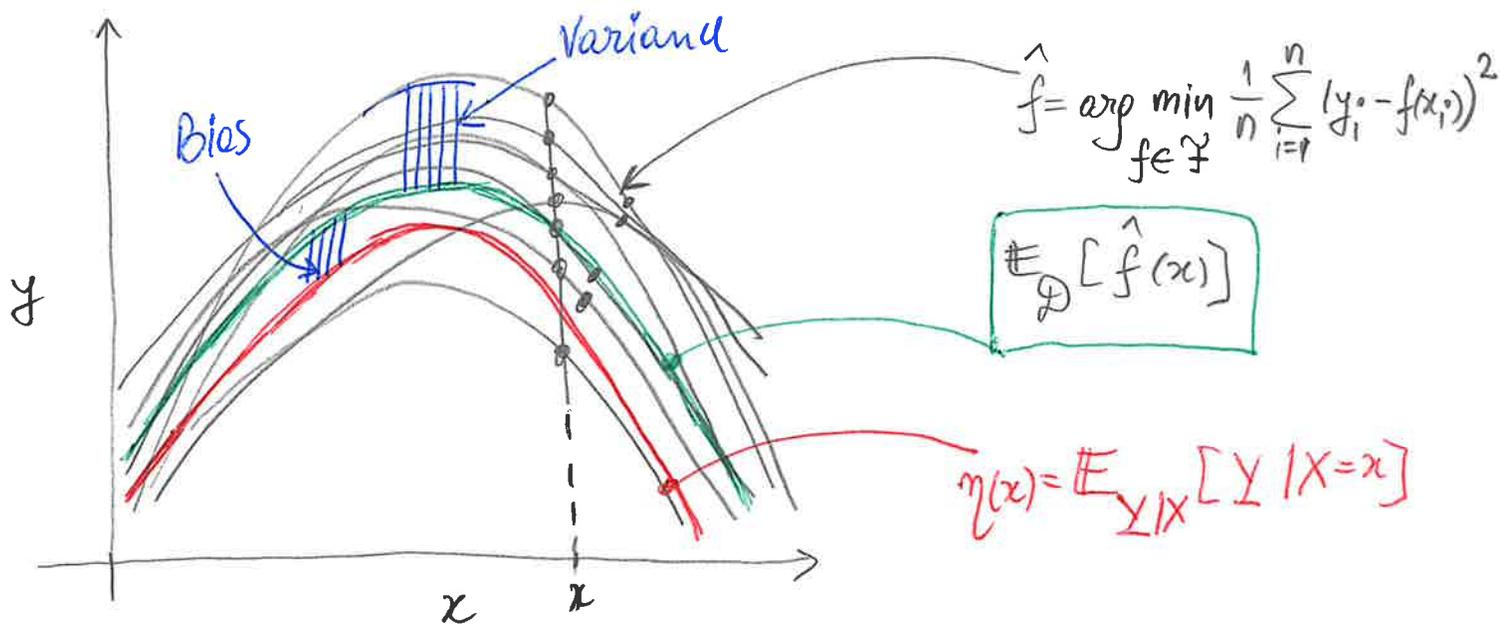
- if it is large  $\Rightarrow$  function class  $\hat{f}$  cannot describe data well.

**Variance** - for a particular sample dataset  $D \sim P_{X,Y}$ , indicates how different is  $\hat{f}_D$  from average  $\hat{f}$ ; if varying wildly  $\Rightarrow$  error!



$$\mathbb{E}_{\mathcal{D}} [\eta(x) - \hat{f}_{\mathcal{D}}(x)]^2 = (\eta(x) - \mathbb{E}_{\mathcal{D}} [\hat{f}_{\mathcal{D}}(x)])^2 + \mathbb{E}_{\mathcal{D}} [(\mathbb{E}_{\mathcal{D}} [\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$$

Learning error = Biased squared + Variance.



So, **True Generalization Error:** (E notation in Bishop, pp. 149)

$$\mathbb{E}_{Y|X} [\mathbb{E}_{\mathcal{D}} [(Y - \hat{f}_{\mathcal{D}}(x))^2] | X=x] = \mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X=x] +$$

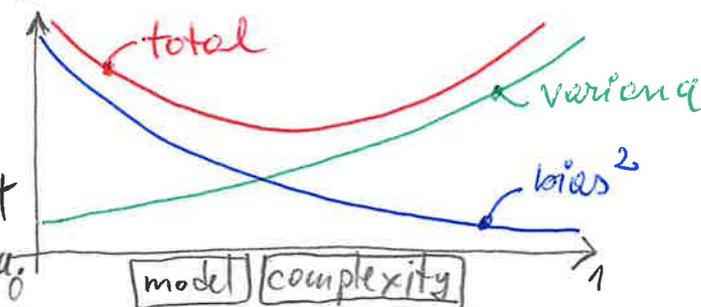
irreducible error

Learning error

$$+ (\eta(x) - \mathbb{E}_{\mathcal{D}} [\hat{f}_{\mathcal{D}}(x)])^2 + \mathbb{E}_{\mathcal{D}} [(\mathbb{E}_{\mathcal{D}} [\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$$

Biased squared      Variance

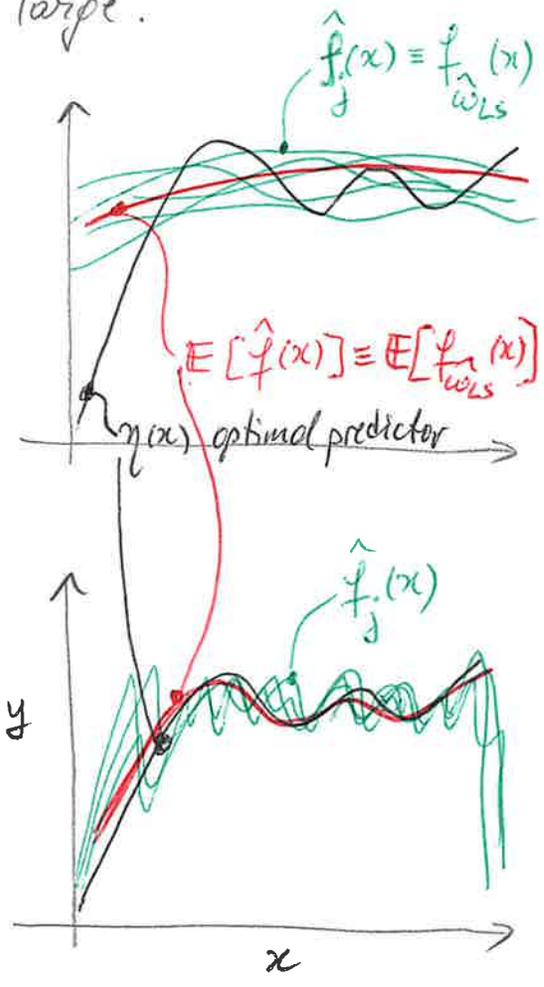
- Main takeaway:
- Model w/ optimal **error** predictive ability is the one that leads to best balance between bias & variance.



(\*)  
 "Bias" goes wrong in Training  
 "Variance" goes wrong in Testing.

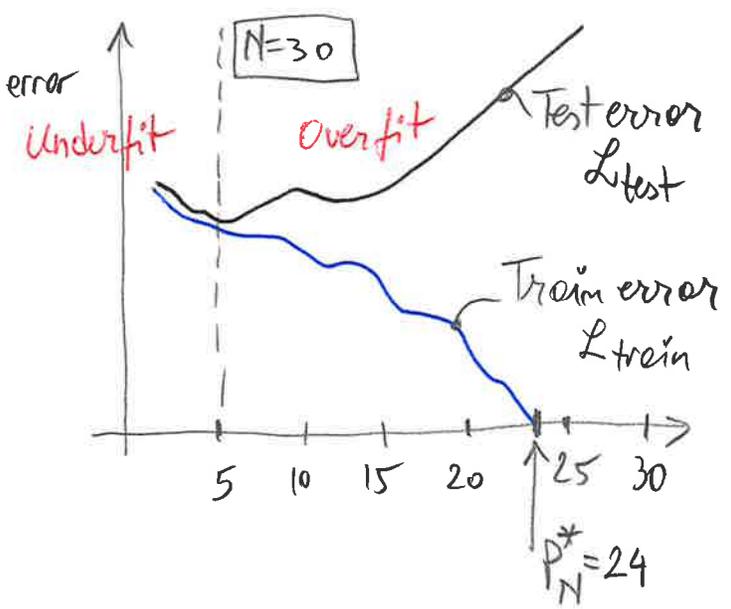
### Summary

- With **simple** model (model complexity is low)
  - Bias<sup>2</sup> of model of our predictor is large.
  - variance of predictor is small.
  - if more samples (larger  $n$ ):
    - Bias remains the same
    - Variance goes down.
    - overall test-error the same.



- With **complex** model (higher than that of optimal predictor  $\eta(x)$ )
  - Bias<sup>2</sup> of our predictor is small
  - Variance is large
  - if more samples (larger  $n$ ):
    - Bias remains the same
    - Variance goes down.
    - Overall test error goes down.

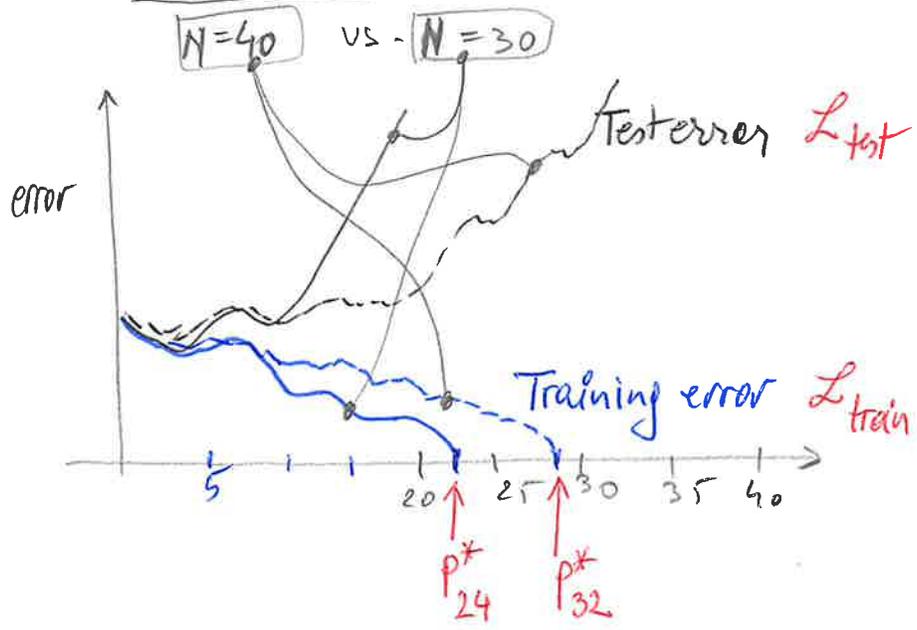
→ **Optimal model complexity** depends on data size (i.e., degree of polynomial in our code demo)



- Given sample size  $N$ , there is a threshold  $P_N^*$  where train error = 0
- Training error is always monotonically non-increasing
- Test error goes down, then up and fluctuates.

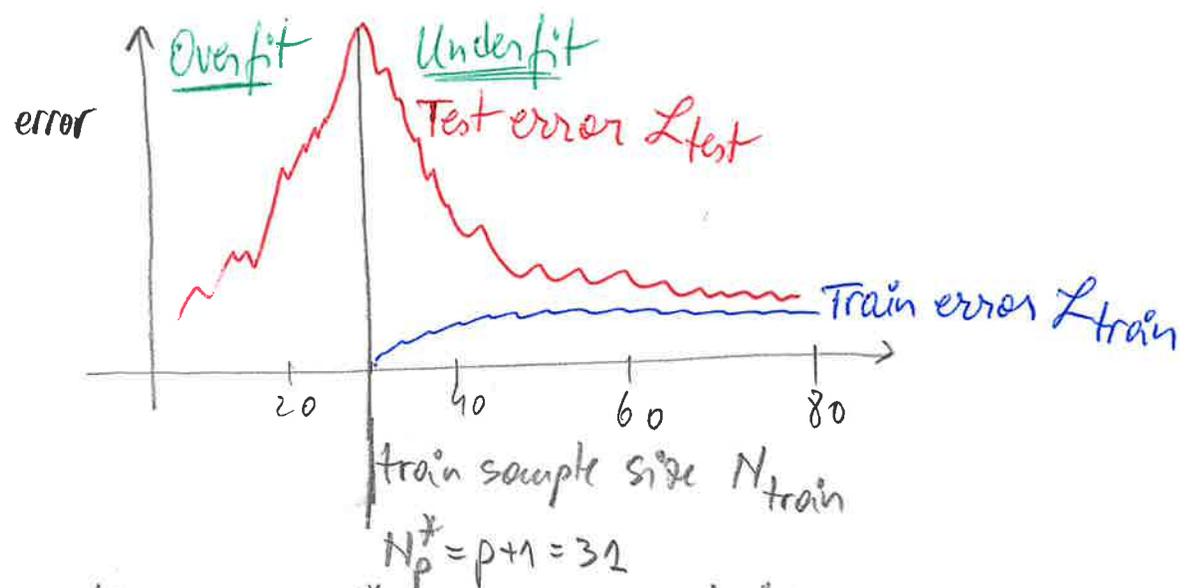
# Variance decreases with more data

→ lets us fit more complex models



- Threshold  $p_N^*$  moves to the right, as dataset size increases.
- Training error tends to increase, because more points need to fit.
- Test error tends to decrease, because variance decreases.

Choose model complexity  $p=30$ , vary data size  $n$



- Threshold  $N_p^*$  below which training error is zero (extreme overfit)
- Above threshold, test error tends to decrease.
- Training error tends to increase (harder to fit so much data w/ simple model!)