

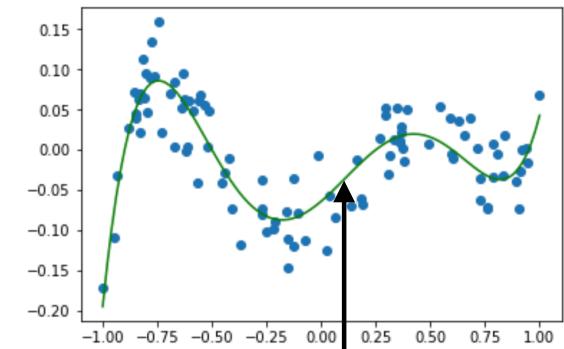
- \* Monday No class
- \* Tlai

# Lecture 6: Bias-Variance Tradeoff (continued)

---

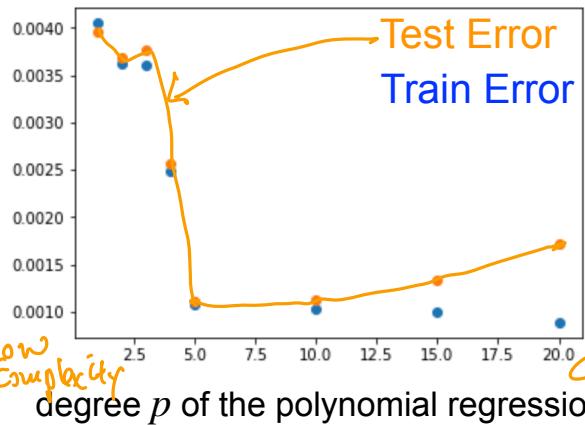
W

# Test error vs. model complexity

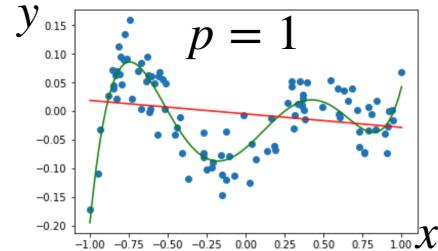


Optimal predictor  $\eta(x)$   
is degree-5 polynomial

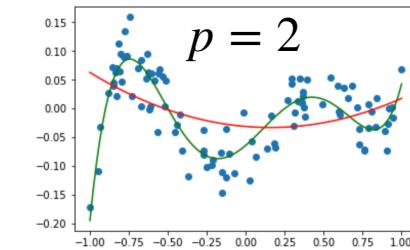
Error



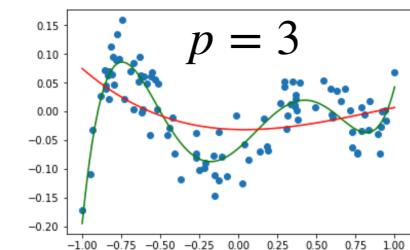
Simple model:  
Model complexity is below  
the complexity of  $\eta(x)$



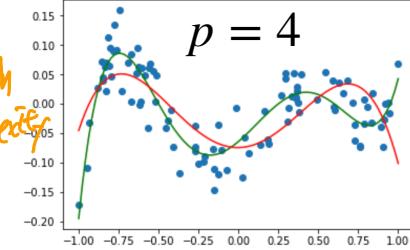
$p = 1$



$p = 2$

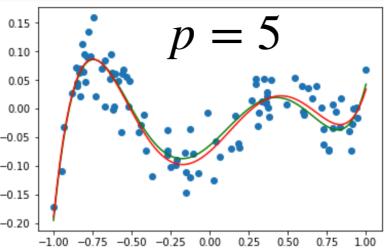


$p = 3$

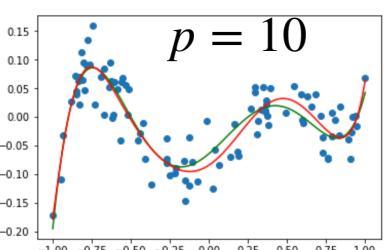


$p = 4$

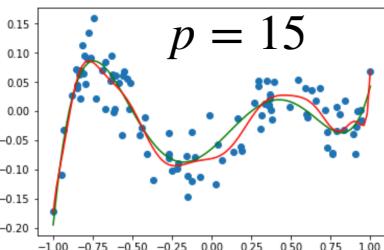
Complex model:



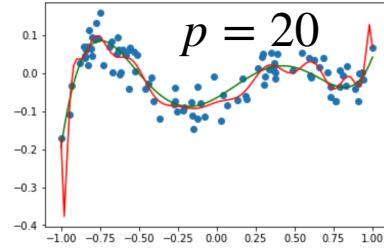
$p = 5$



$p = 10$

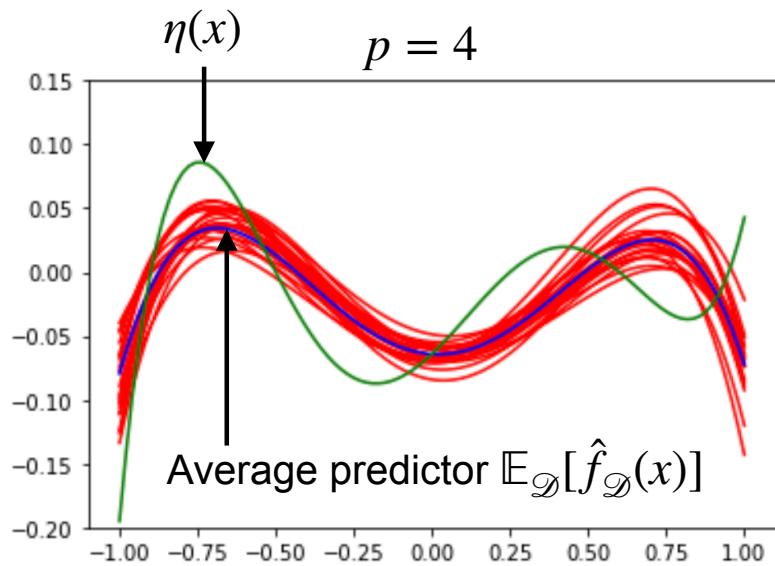
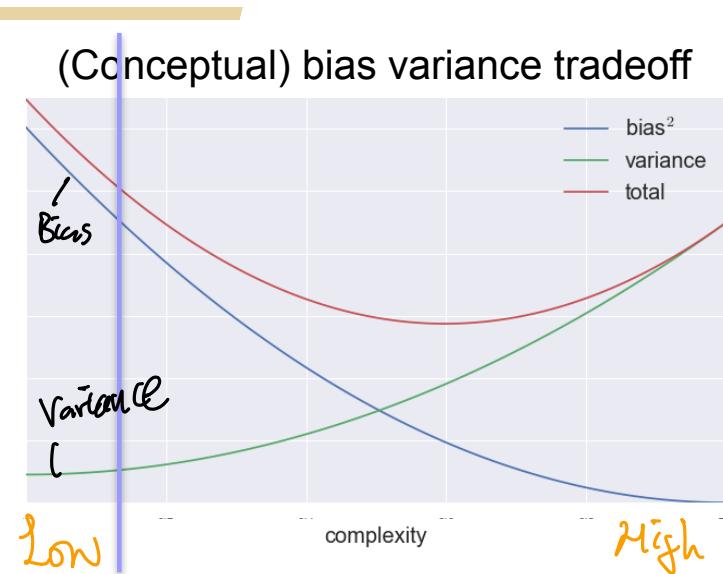


$p = 15$



$p = 20$

# Recap: Bias-variance tradeoff with simple model

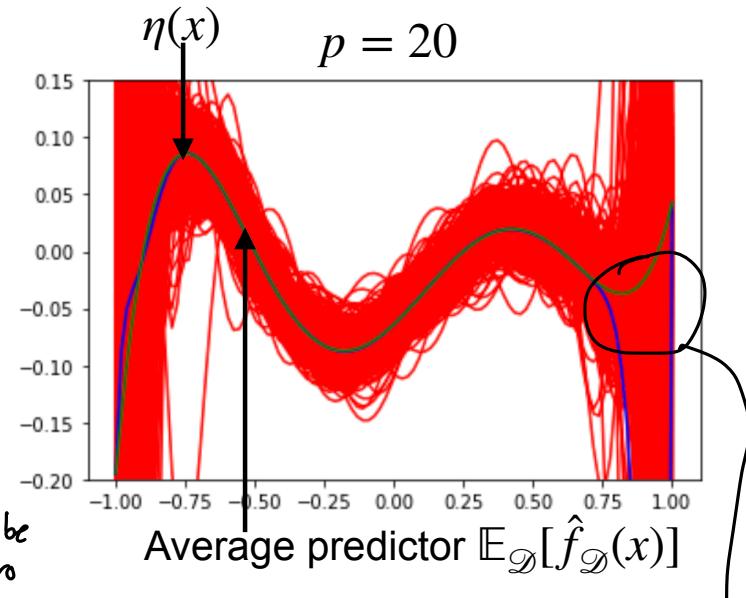
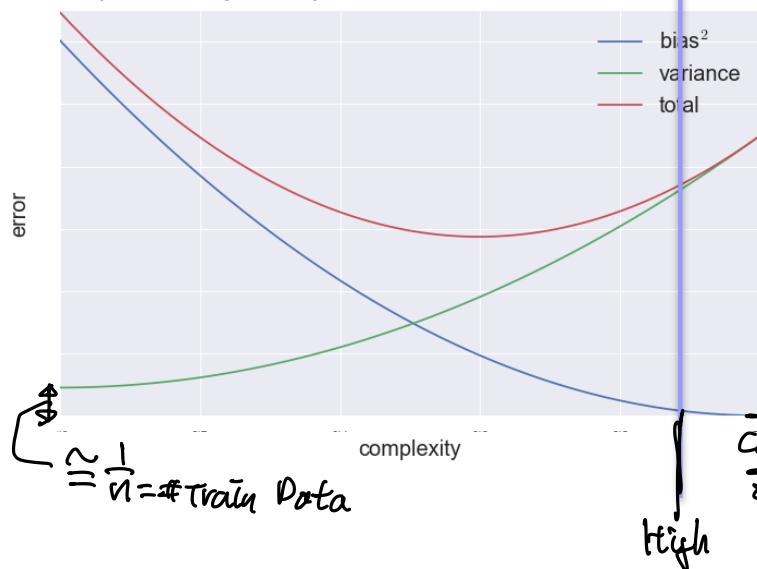


- When model **complexity is low** (lower than the optimal predictor  $\eta(x)$ )
  - Bias<sup>2</sup> of our predictor,  $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2$ , is large
  - Variance of our predictor,  $\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$ , is small
  - If we have more samples, then
    - Bias *does not change*
    - Variance *decreases*
    - Because Variance is already small, overall test error *does not change much*

# Recap: Bias-variance tradeoff with simple model

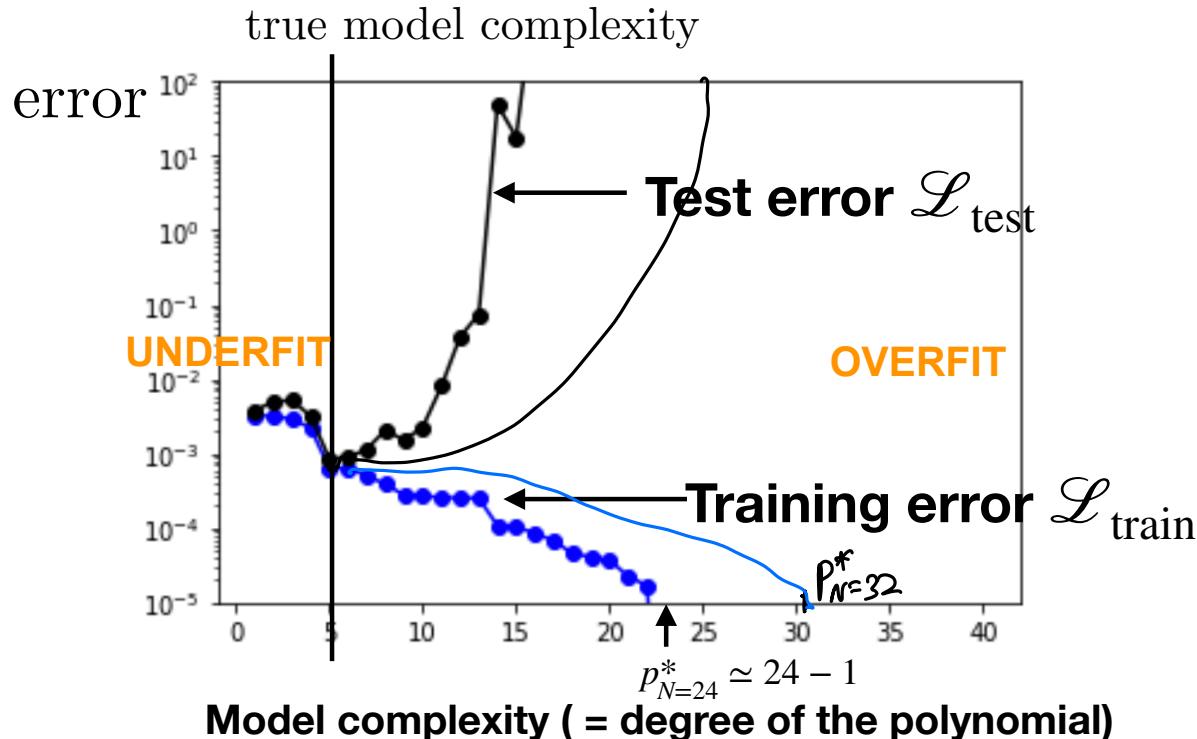
Complex

(Conceptual) bias variance tradeoff



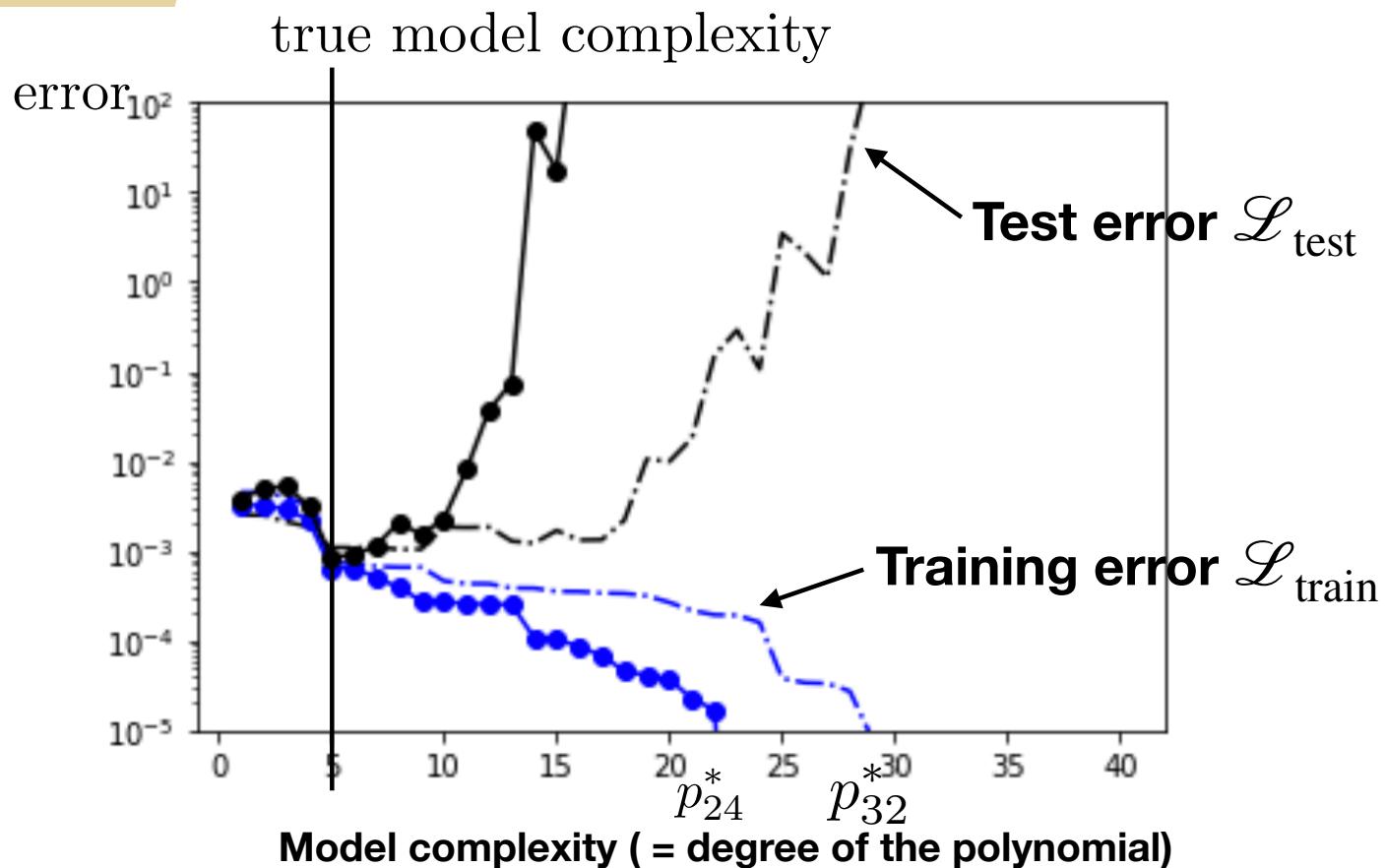
- When model complexity is high (higher than the optimal predictor  $\eta(x)$ )
  - Bias of our predictor,  $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2$ , is small
  - Variance of our predictor,  $\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$ , is large
  - If we have more samples, then
    - Bias *does not change*
    - Variance *goes down*
    - Because Variance is dominating, overall test error *goes down if more samples*

- let us first fix sample size  $N=30$ , collect one dataset of size  $N$  i.i.d. from a distribution, and fix one training set  $S_{\text{train}}$  and test set  $S_{\text{test}}$  via 80/20 split
- then we run multiple validations and plot the computed MSEs for all values of  $p$  that we are interested in



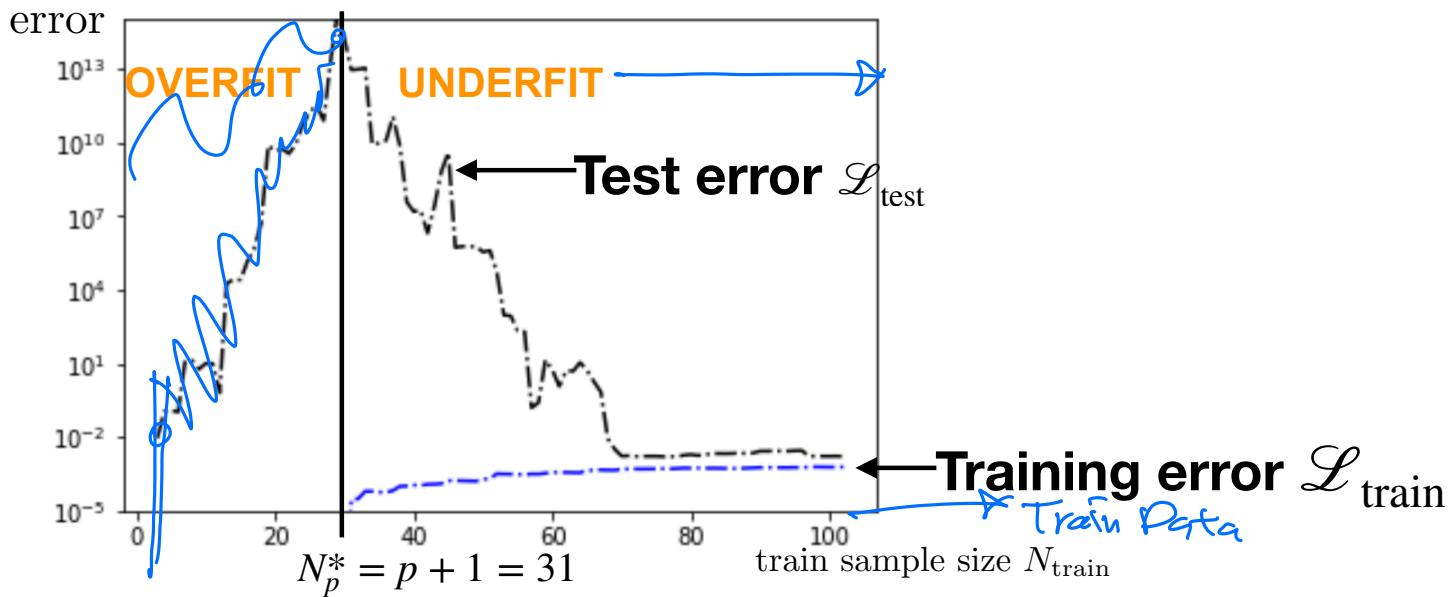
- Given sample size  $N$  there is a threshold,  $p_N^*$ , where training error is zero
- Training error is **always** monotonically non-increasing
- Test error has a trend of going down and then up, but fluctuates

- let us now repeat the process changing the sample size to **N=40** , and see how the curves change



- The threshold,  $p_N^*$ , moves right
- Training error tends to increase, because more points need to fit
- Test error tends to decrease, because Variance decreases

- let us now fix predictor model complexity  $p=30$ , collect multiple datasets by starting with 3 samples and adding one sample at a time to the training set, but keeping a large enough test set fixed
- then we plot the computed MSEs for all values of train sample size  $N_{\text{train}}$  that we are interested in



- There is a threshold,  $N_p^*$ , below which training error is zero (extreme overfit)
- Below this threshold, test error is meaningless, as we are overfitting and there are multiple predictors with zero training error some of which have very large test error
- Test error tends to decrease
- Training error tends to increase

# Bias-variance tradeoff for linear models

If  $Y_i = X_i^T w^* + \epsilon_i$  and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

↑  
↑ Capital Means Random

$$\mathbf{y} = \mathbf{X}w^* + \boldsymbol{\epsilon}$$

Bold means they are concatenated

$$\hat{w}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \boldsymbol{\epsilon})$$
$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}$$
$$= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}$$

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$
$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

$$\eta(x) \stackrel{\triangle}{=} \mathbb{E}_{Y|X}[Y | X = x] = \mathbb{E}[x^T w^* + \epsilon_i] = x^T w^*$$

def.

$$\hat{f}_{\mathcal{D}}(x) = \underbrace{x^T \hat{w}_{\text{MLE}}}_{=\eta(x)} + \underbrace{x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}}_{\text{error}}$$

# Bias-variance tradeoff for linear models

If  $Y_i = \mathbf{X}_i^T w^* + \epsilon_i$  and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \boldsymbol{\epsilon}$$

$$\begin{aligned}\widehat{w}_{\text{MLE}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \boldsymbol{\epsilon}) \\ &= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}\end{aligned}$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X=x] = x^T w^*$$

$$\hat{f}_{\mathcal{D}}(x) = x^T \widehat{w}_{\text{MLE}} = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}$$

- Irreducible error:  $\mathbb{E}_{X,Y}[(Y - \eta(x))^2 | X = x] = \mathbb{E}[(x^T w^* + \epsilon_i - x^T w^*)^2] = \mathbb{E}[\epsilon_i^2]$
- Bias squared:  $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 = x^T w^* - \mathbb{E}[x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}] = x^T w^* - x^T w^* = 0$  Remark

# Bias-variance tradeoff for linear models

If  $Y_i = X_i^T w^* + \epsilon_i$  and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\hat{w}_{\text{MLE}} = w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = x^T w^*$$

$$\hat{f}_{\mathcal{D}}(x) = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

• Variance:  $\mathbb{E}_{\mathcal{D}} \left[ (\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 \right] = \mathbb{E} \left[ (x^T w^* + x^T (X^T X)^{-1} X^T \epsilon - x^T w^*)^2 \right]$

$$= \mathbb{E} \left[ x^T (X^T X)^{-1} X^T \underbrace{\epsilon \epsilon^T}_\text{indep} X (X^T X)^{-1} x \right]$$
$$= \sigma^2 \mathbb{E} \left[ x^T (X^T X)^{-1} \cancel{X^T X} \cancel{(X^T X)^{-1}} x \right]$$
$$= \sigma^2 \underbrace{x^T \mathbb{E}[(X^T X)^{-1}] x}$$

# Bias-variance tradeoff for linear models

If  $Y_i = \mathbf{X}_i^T w^* + \epsilon_i$  and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\hat{w}_{\text{MLE}} = w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = x^T w^*$$

$$\hat{f}_{\mathcal{D}}(x) = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\begin{aligned}\mathbb{E} \epsilon &= 0 \\ \mathbb{E} \sum \epsilon^2 &= \sigma^2 \cdot \mathbf{I}_{d \times d} \\ \epsilon &= \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}\end{aligned}$$

- Variance:  $\mathbb{E}_{\mathcal{D}} \left[ (\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 \right] = \mathbb{E}_{\mathcal{D}}[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x]$   
 $= \sigma^2 \mathbb{E}_{\mathcal{D}}[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x]$   
 $= \sigma^2 x^T \mathbb{E}_{\mathcal{D}}[(\mathbf{X}^T \mathbf{X})^{-1}] x$   
• To analyze this, let's assume that  $X_i \sim \mathcal{N}(0, \mathbf{I})$  and number of samples,  $n$ , is large enough such that  $\mathbf{X}^T \mathbf{X} = n \mathbf{I}$  with high probability and  $\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1}] \approx \frac{1}{n} \mathbf{I}$ , then
  - Variance is  $\frac{\sigma^2 x^T x}{n}$ , and decreases with increasing sample size  $n$

$$(\mathbf{X}^T \mathbf{X})^{-1} \approx \frac{1}{n} \mathbf{I}_{d \times d}$$

# Questions?

---