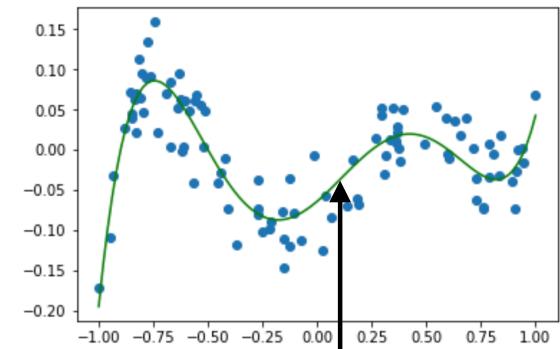


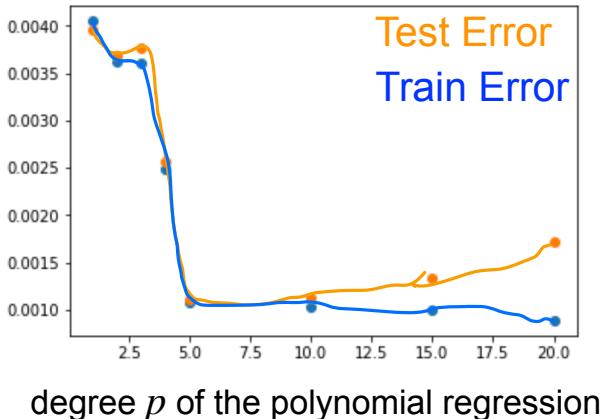
Lecture 6: Bias-Variance Tradeoff (continued)

W

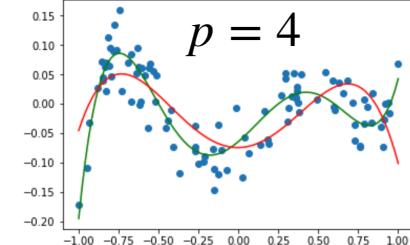
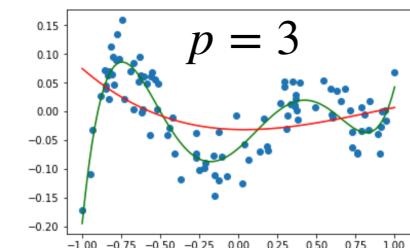
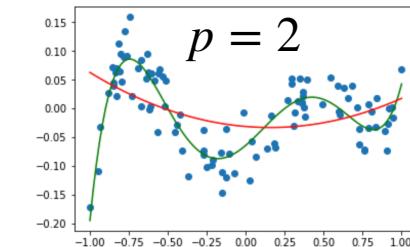
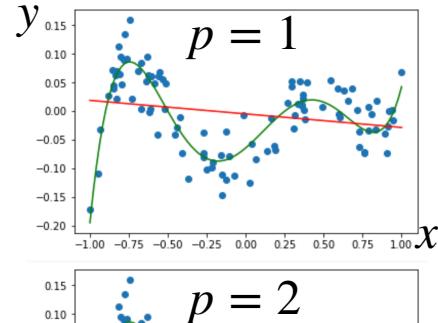
Test error vs. model complexity



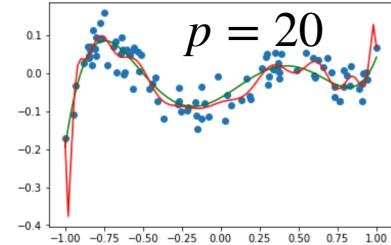
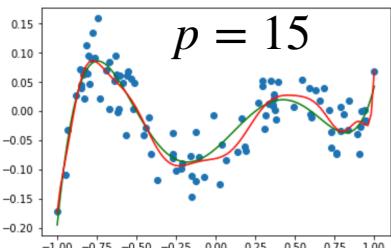
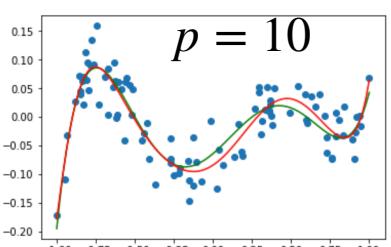
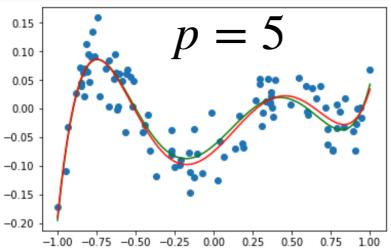
Error



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

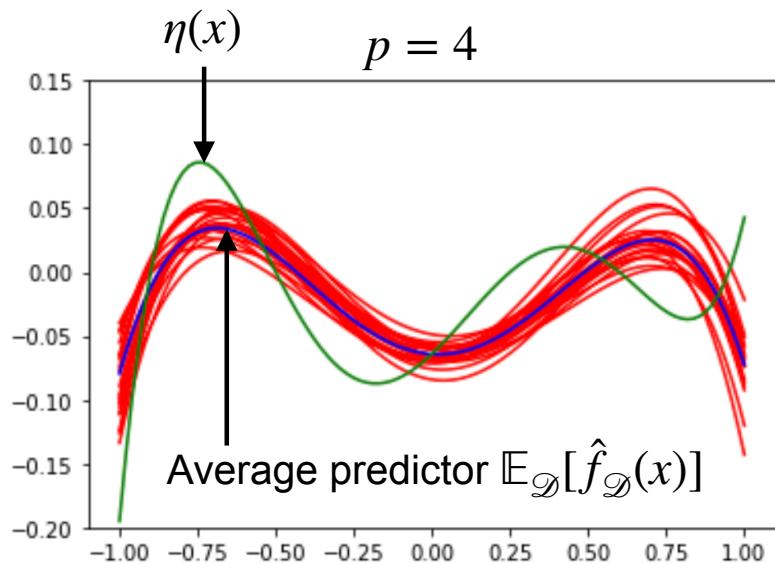
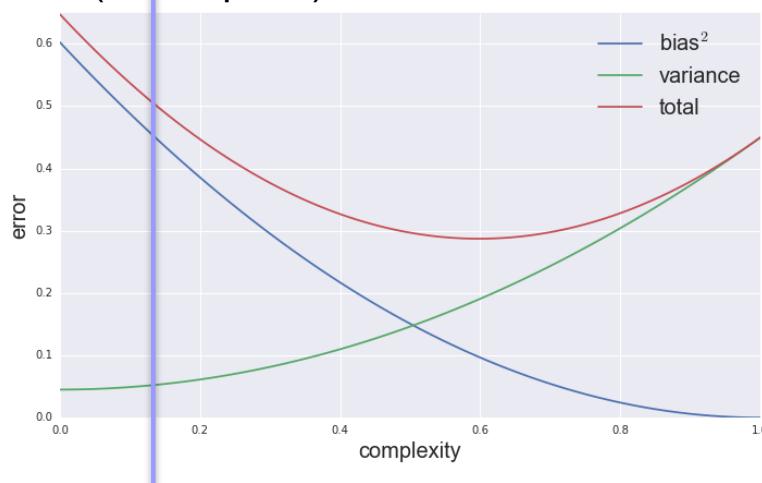


Complex model:



Recap: Bias-variance tradeoff with simple model

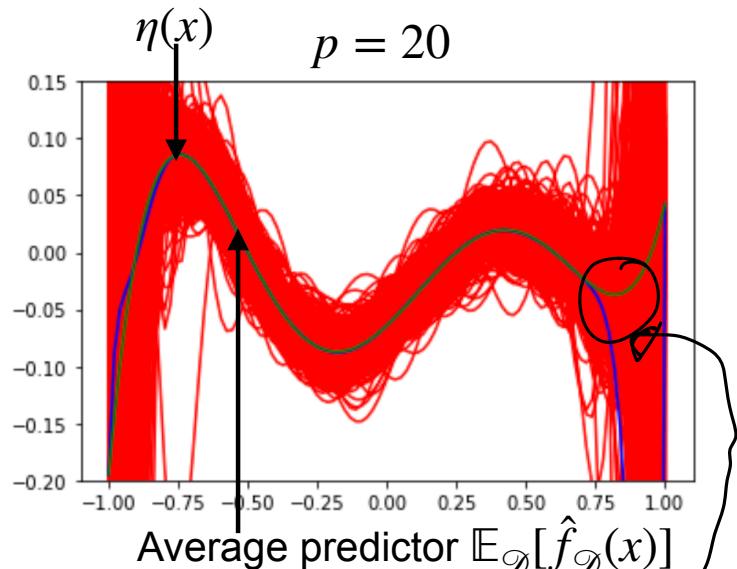
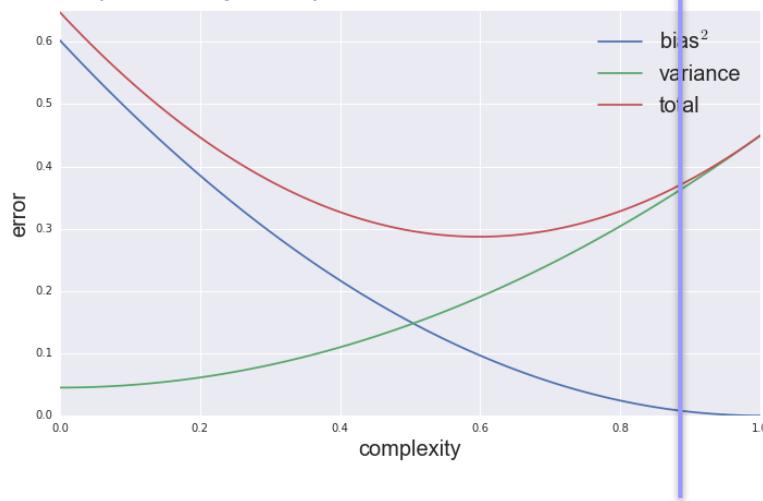
(Conceptual) bias variance tradeoff



- When model **complexity is low** (lower than the optimal predictor $\eta(x)$)
 - Bias² 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* ← circled Q
 - Variance *goes down*
 - Because Variance is already small, overall test error *does not change much*

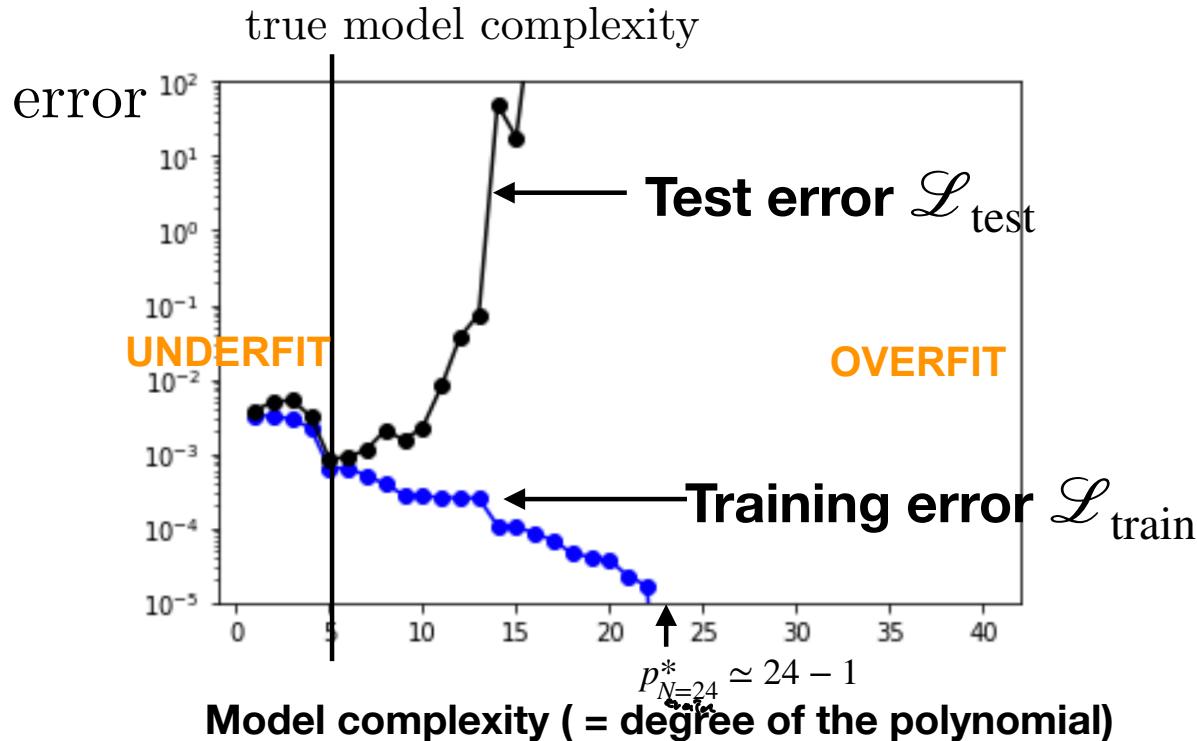
Recap: Bias-variance tradeoff with simple model

(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 *decreases*
 - Because Variance is dominating, overall test error *decreases*

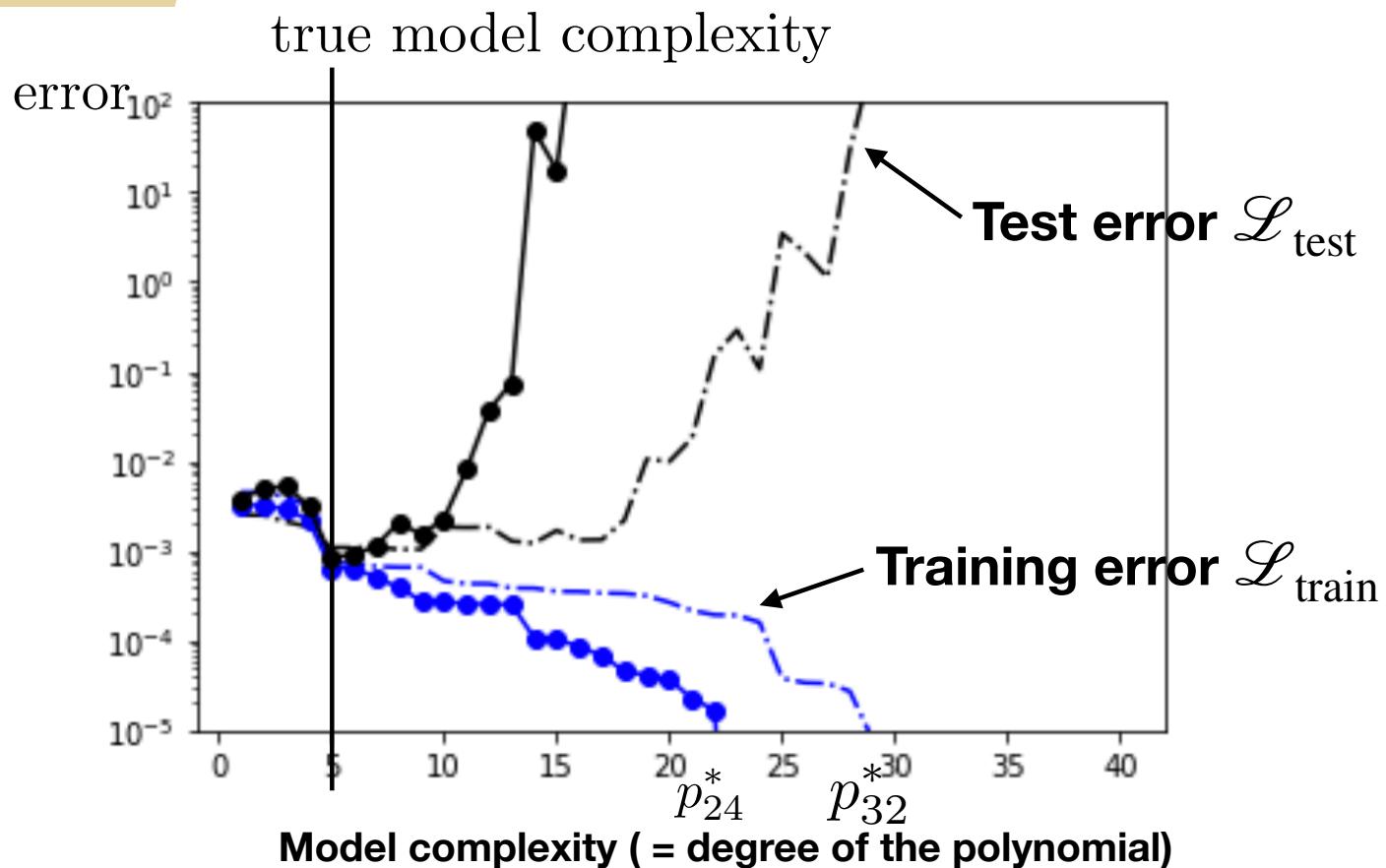
- let us first fix sample size $N=30$, collect one dataset of size N i.i.d. from a distribution, and fix one training set S_{train} and test set S_{test} via 80/20 split
- then we run multiple validations and plot the computed MSEs for all values of p that we are interested in



- Given sample size N there is a threshold, $p_{N_{\text{train}}}^*$, 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

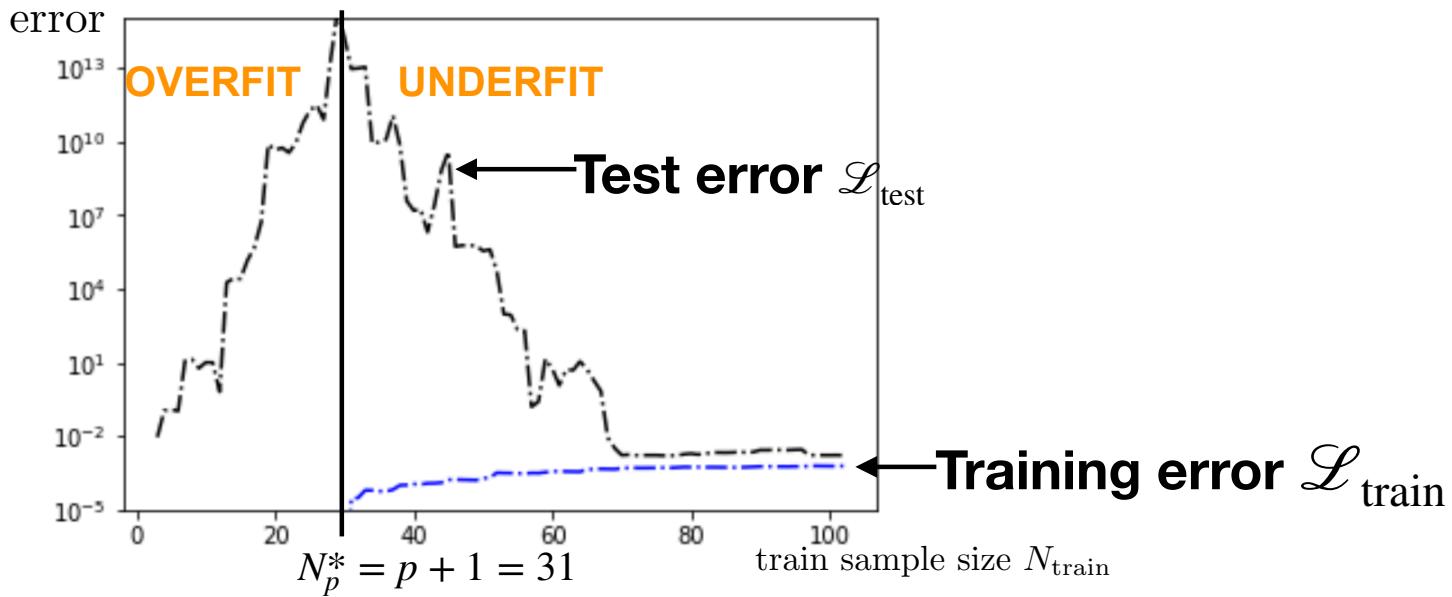
(Q) what happens if
 $N \uparrow$

- 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_{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 because they are R.V.

$$\begin{aligned} \mathbf{y} &= \mathbf{X}w^* + \boldsymbol{\epsilon} && \text{Bold because they are vector } \mathbf{y} \in \mathbb{R}^n \text{ and matrix } \mathbf{X} \in \mathbb{R}^{n \times d} \\ \hat{\mathbf{w}}_{\text{MLE}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \boldsymbol{\epsilon}) \\ &= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \end{aligned}$$

$$\eta(x) = \mathbb{E}_{Y|X} [Y | \underset{\mathbb{R}^d}{X} = x] = x^T w^* \quad \leftarrow Y_i = X_i^T w^* + \epsilon_i$$

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

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: $\mathbb{E}_{X,Y}[(Y - \eta(x))^2 | X = x] = \mathbb{E}[(\mathbf{x}^T w^* + \boldsymbol{\epsilon} - \mathbf{x}^T w^*)^2] = \mathbb{O}^2$
- Bias squared: $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 = x^T w^* - \mathbb{E}[x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}]$
(is independent of the sample size!) ↑ zero mean

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: $\mathbb{E}_{\mathcal{D}} \left[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 \right] = \mathbb{E} \left[x^T (x^T x)^{-1} x \underbrace{\sum \epsilon \epsilon^T}_{\mathbb{E}[\epsilon \epsilon^T] = \sigma^2 I} (x^T x)^{-1} x \right]$

$$= \sigma^2 x^T \mathbb{E}[(x^T x)^{-1}] x$$
$$= \sigma^2 x^T \mathbb{E}[(x^T x)^{-1}] x$$

Bias-variance tradeoff for linear models

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

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

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

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

- Variance:
$$\begin{aligned}\mathbb{E}_{\mathcal{D}} \left[\left(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2 \right] &= \mathbb{E}_{\mathcal{D}}[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x] \\ &= \sigma^2 \mathbb{E}_{\mathcal{D}}[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x] \\ &= \sigma^2 x^T \mathbb{E}_{\mathcal{D}}[(\mathbf{X}^T \mathbf{X})^{-1}] x\end{aligned}$$
- 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 $\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1}] \simeq \frac{1}{n} \mathbf{I}$, then
 - Variance is $\frac{\sigma^2 x^T x}{n}$, and decreases with increasing sample size n

Questions?
