

- HW1 out today due Jan 25th Tuesday midnight.
- Pennī

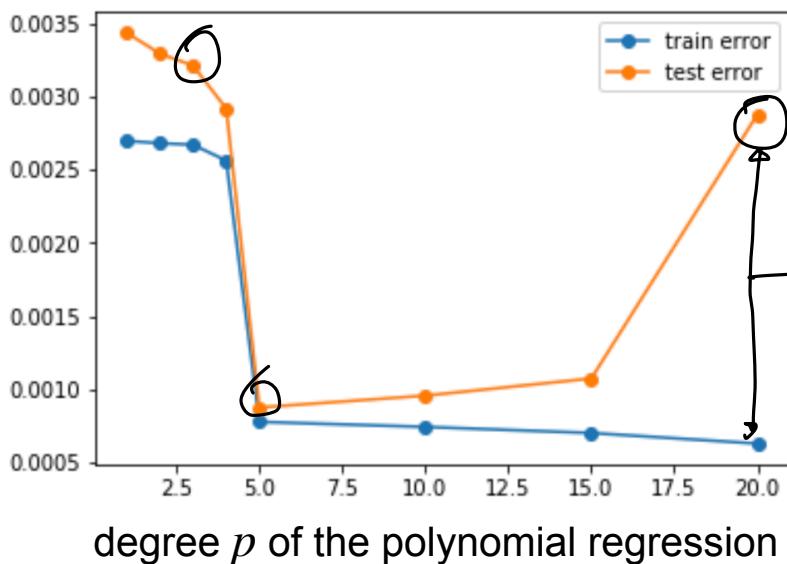
Lecture 5: Bias-Variance Tradeoff

- explaining test error using theoretical analysis

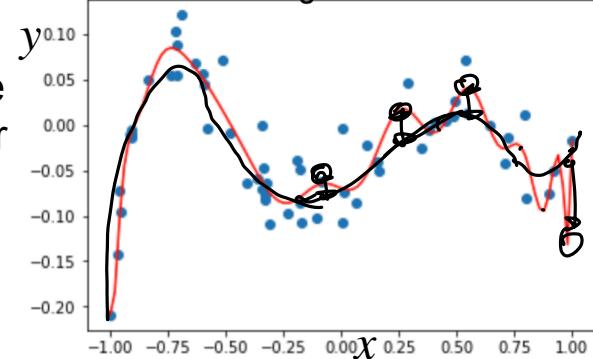
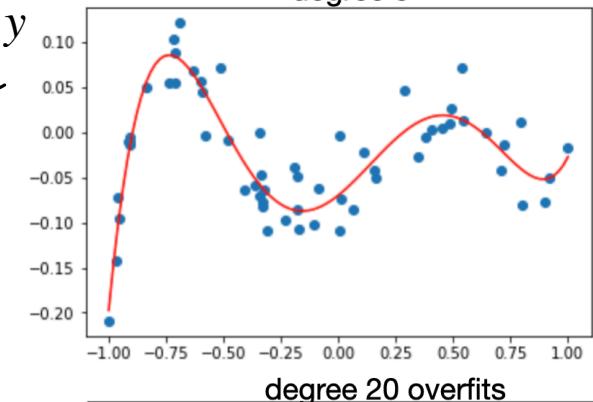
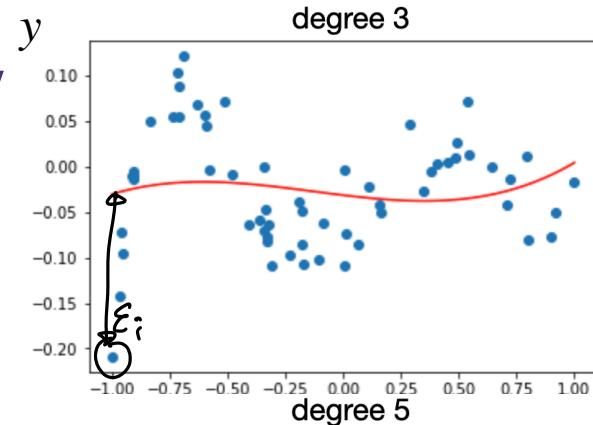
W

Train/test error vs. complexity

Error



degree p of the polynomial regression



- **Model complexity** e.g., degree p of the polynomial model, number of features used in diabetes example
 - Related to the dimension of the model parameter
- **Train error** monotonically decreases with model complexity
- **Test error** has a U shape

Statistical learning

Typical notation:

X denotes a random variable

\overline{x} denotes a deterministic instance

- Suppose data is generated from a statistical model $(X, Y) \sim P_{X,Y}$
 - and assume we know $P_{X,Y}$ (just for now to explain statistical learning)
 - Then **learning** is to find a predictor $\eta : \mathbb{R}^d \rightarrow \mathbb{R}$ that minimizes
 - the expected error $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2]$ = $\underbrace{\text{True Error}}_{\hookrightarrow \text{True dist } P_{X,Y}}$
 - think of this random (X, Y) as a new sample you will encounter when you deployed your learned model, and we care about its average performance
 - Since, we do not assume anything about the function $\eta(x)$, it can take any value for each $X = x$, hence the optimization can be broken into sum (or more precisely integral) of multiple objective functions, each involving a specific value $X = x$
 - $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2] = \mathbb{E}_{X \sim P_X} [\mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 | X = x]]$
 $= \int \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 | X = x] P_X(x) dx$
- Or for discrete X ,
- $$= \sum_x P_X(x) \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 | X = x]$$

Where we used the chain rule: $\mathbb{E}_{X,Y}[f(X, Y)] = \mathbb{E}_X [\mathbb{E}_{Y|X}[f(x, Y) | X = x]]$

Statistical learning

- We can solve the optimization for each $X = x$ separately

- $\eta(x) = \arg \min_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 | X = x]$

- The optimal solution is $\eta(x) = \mathbb{E}_{Y \sim P_{Y|X}}[Y | X = x]$,
which is the best prediction in ℓ_2 -loss/Mean Squared Error

- Claim: $\mathbb{E}_{Y \sim P_{Y|X}}[Y | X = x] = \arg \min_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 | X = x]$

- Proof:
$$\begin{aligned} \frac{\partial}{\partial a} \mathbb{E}_{Y|X} [Y^2 - 2aY + a^2] &= \frac{\partial}{\partial a} \left[\mathbb{E}[Y^2 | X=x] - 2a \mathbb{E}[Y | X=x] + a^2 \right] \\ &= -2 \mathbb{E}[Y | X=x] + 2a \Big|_{a=\eta(x)} = 0 \end{aligned}$$

Optimal Model / Predictor : $\eta(x) = \mathbb{E}[Y | X=x]$

- Note that this optimal statistical estimator $\eta(x) = \mathbb{E}[Y | X = x]$ cannot be implemented as we do not know $P_{X,Y}$ in practice
- This is only for the purpose of conceptual understanding

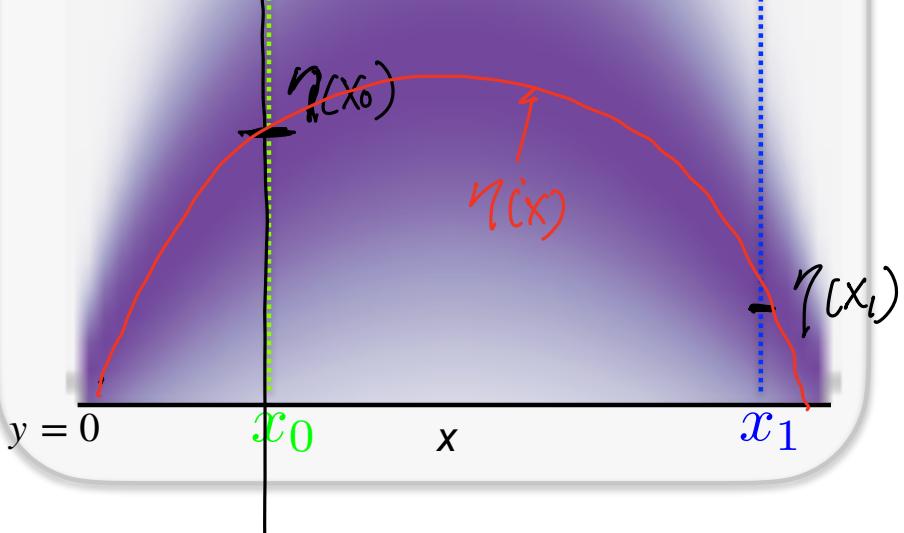
Statistical Learning

Ideally, we want to find:

Optimal Predictor $\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$

$$P_{XY}(X = x, Y = y)$$

$$y = 1$$



$$P_{XY}(Y = y|X = x_0)$$

$$y = 0$$

$$y = 1$$

$$\eta(x_0) = \mathbb{E}[Y|X = x_0]$$

$$P_{XY}(Y = y|X = x_1)$$

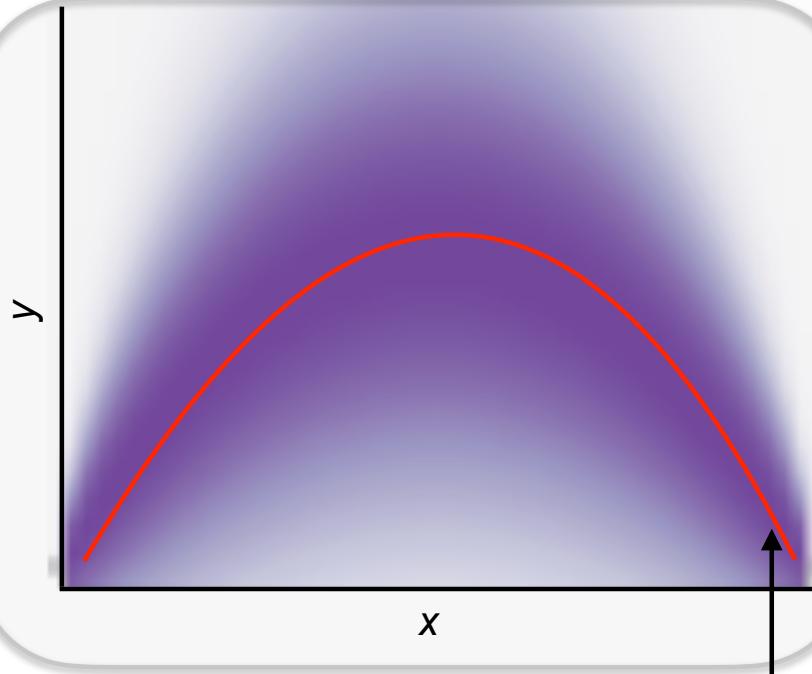
$$y = 0$$

$$y = 1$$

$$\eta(x_1) = \mathbb{E}[Y|X = x_1]$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

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

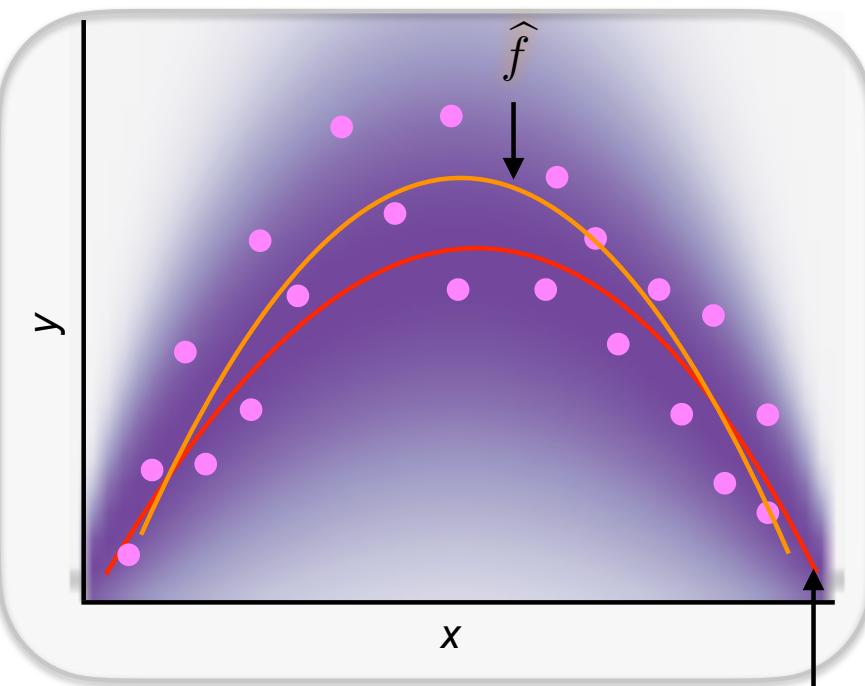
But we do not know $P_{X,Y}$

We only have samples.

Q. How is $\hat{f}(x)$ related to optimal predictor $\eta(x)$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

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

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

So we need to restrict our predictor to a function class (e.g., linear, degree- p polynomial) to avoid overfitting:

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

$$\mathbb{E}_{Y|X}[Y|X = x]$$

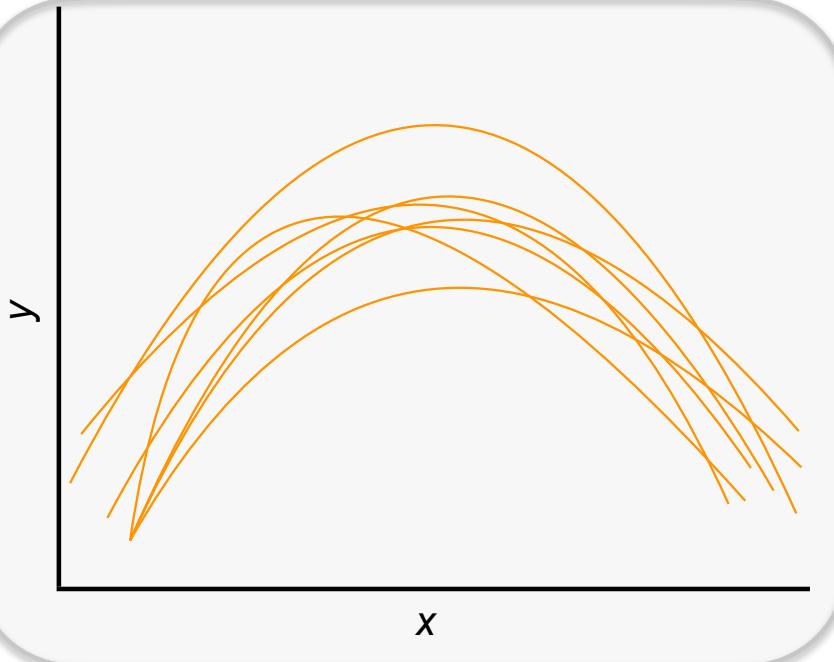
We care about how our predictor performs on future unseen data

True Error of \hat{f} : $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$

Future prediction error $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$ is random

because \hat{f} is random (whose randomness comes from training data \mathcal{D})

$$P_{XY}(X = x, Y = y)$$



True Error: $\mathbb{E}_{x,y}[(y - \hat{f}(x))^2]$

↑
Random

Average True Error :
$$\underline{\mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{x,y} [(y - \hat{f}(x))^2] \right]}$$

* Mayura

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

Bias-variance tradeoff

Notation:

I use predictor/model/estimate,
interchangeably

Ideal predictor

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

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- We are interested in the **True Error** of a (random) learned predictor:

$$\mathbb{E}_{X,Y}[(Y - \hat{f}_{\mathcal{D}}(X))^2] \quad \leftarrow \text{Random}$$

- But the analysis can be done for each $X = x$ separately, so we analyze the **conditional true error**:

$$\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathcal{D}}(x))^2 | X = x] \quad \leftarrow \text{Random}$$

- And we care about the **average conditional true error**, averaged over training data:

$$\mathbb{E}_{\mathcal{D}}[\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathcal{D}}(x))^2 | X = x]] \quad \leftarrow \text{deterministic}$$

written compactly as $= \mathbb{E}[(Y - \hat{f}_{\mathcal{D}}(x))^2]$

Bias-variance tradeoff

Ideal predictor

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

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- Average conditional true error:

$$\mathbb{E}_{\mathcal{D}, Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \underbrace{\eta(x) + \eta(x)}_0 - \hat{f}_{\mathcal{D}}(x))^2]$$

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

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

$Y|x$ and \mathcal{D} are independent.

$$= \frac{\mathbb{E}_{Y|x}[(Y - \eta(x))^2 | X=x]}{\text{Irreducible Error.}} + \underbrace{\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{Expected Learning Error}}$$

Bias-variance tradeoff

any function = Largest function class

*Matthew W.

*

Ideal predictor

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

best predictor you get if you have infinite samples

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

Average conditional true error:

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2] \\ &= \mathbb{E}_{\mathcal{D}, Y|x} \left[(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) + (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2 \right] \\ &= \mathbb{E}_{Y|x}[(Y - \eta(x))^2] + \underbrace{2\mathbb{E}_{\mathcal{D}, Y|x}[(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x))]}_{=0} + \mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] \end{aligned}$$

(this follows from independence of \mathcal{D} and (X, Y) and

$$\mathbb{E}_{Y|x}[Y - \eta(x)] = \mathbb{E}[Y|X = x] - \eta(x) = 0$$

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

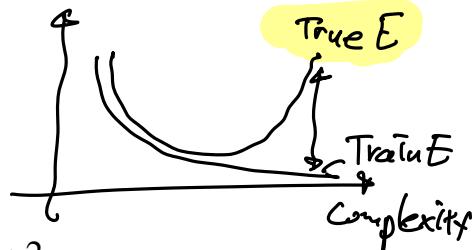
Irreducible error

- (a) Caused by stochastic label noise in $P_{Y|X=x}$
- (b) cannot be reduced

Average learning error

Caused by

- (a) either using too “simple” of a model or
- (b) not enough data to learn the model accurately



Bias-variance tradeoff

Ideal predictor

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

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

• Average learning error:

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] &= \mathbb{E}_{\mathcal{D}} \left[(\eta(x) - \underbrace{\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]}_{\approx 0} + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2 \right] \\ &= (\eta(x) - \bar{f}(x))^2 + 2 \cdot \mathbb{E} \left[(\eta(x) - \bar{f}(x))(\bar{f}(x) - \hat{f}_{\mathcal{D}}(x)) \right] + \mathbb{E} \left[(\bar{f}(x) - \hat{f}_{\mathcal{D}}(x))^2 \right] \\ &= \mathbb{E} \left[\bar{f}(x) - \hat{f}_{\mathcal{D}}(x) \right] \\ &= \bar{f}(x) - \mathbb{E}[\hat{f}_{\mathcal{D}}(x)] = 0 \\ &= \frac{(\eta(x) - \bar{f}(x))^2}{\text{Bias}^2} + \frac{\mathbb{E}[(\bar{f}(x) - \hat{f}_{\mathcal{D}}(x))^2]}{\text{Variance of } \hat{f}_{\mathcal{D}}(x)} = \bar{f}(x)\end{aligned}$$

Bias-variance tradeoff

Ideal predictor

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

Learned predictor

$$\hat{f}_{\mathcal{D}} = \arg \min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

$$(a-b+b-c)^2 = (a-b)^2 + 2(a-b)(b-c) + (b-c)^2$$

• Average learning error:

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

biased squared

variance

Bias-variance tradeoff

- Average conditional true error:

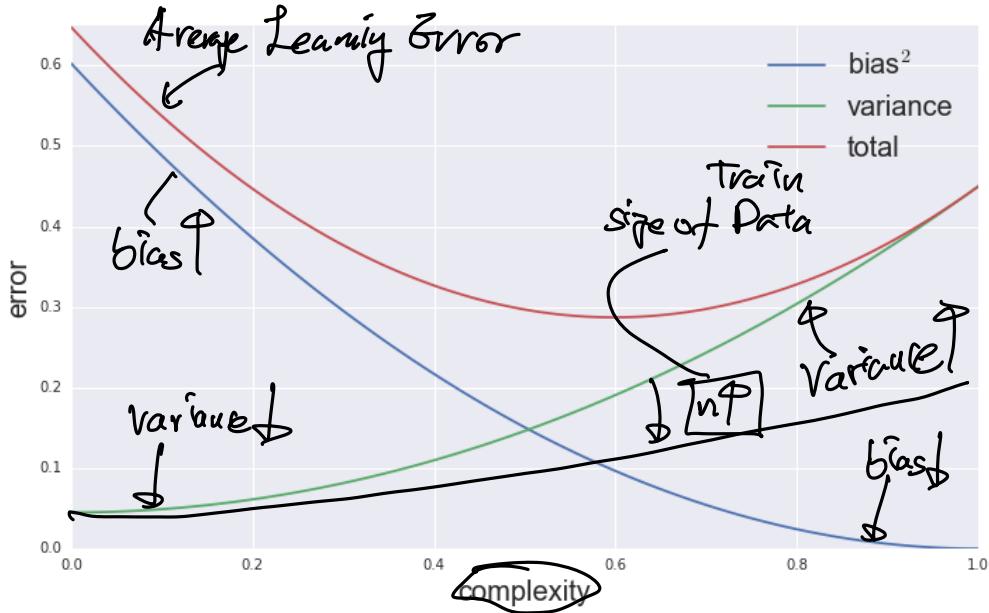
$$\mathbb{E}^2 \mathbb{E}_{\mathcal{D}, Y|x} [(Y - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{Y|x} [(Y - \eta(x))^2]$$

irreducible error

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

Bias squared:
measures how the predictor is mismatched with the best predictor in expectation

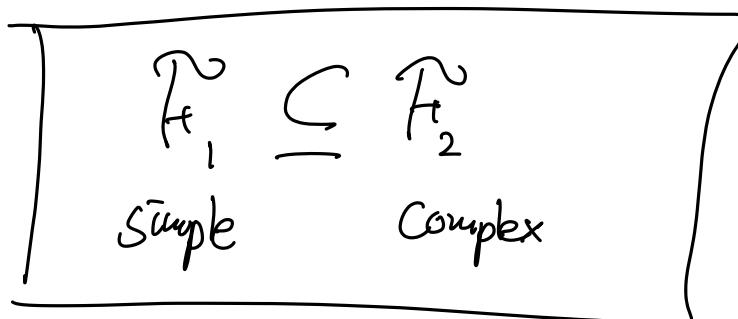
variance:
measures how the predictor varies each time with a new training datasets



Questions?

class F
Model Complexity

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



Constant \subseteq Linear \subseteq Quadratic $\subseteq \dots$