

Bias-Variance Tradeoff

UNIVERSITY *of* WASHINGTON

W

Optimal Prediction

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

Find function η that minimizes

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X}[(Y - c_x)^2 | X = x]$)

Optimal Prediction

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

Find function η that minimizes

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X}[(Y - c_x)^2 | X = x]$)

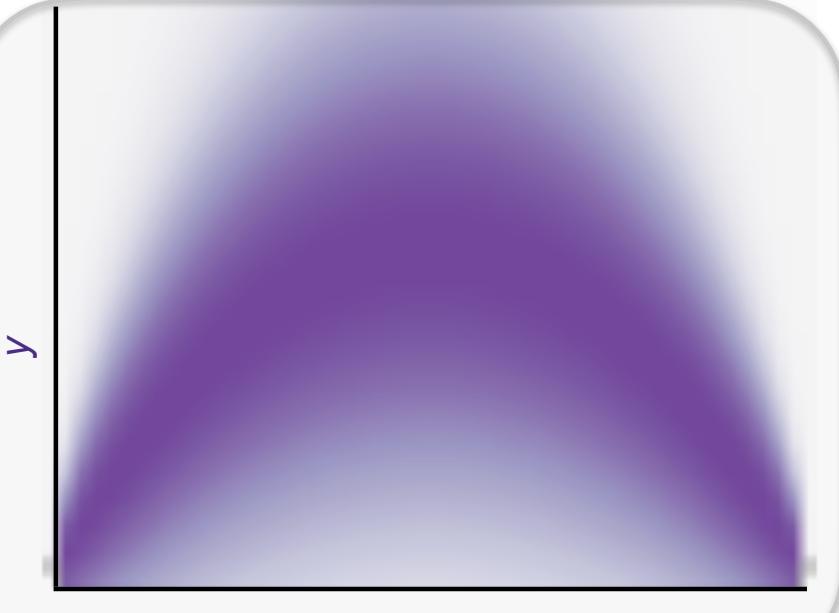
$$\begin{aligned} 0 &= \frac{d}{dc_x} \mathbb{E}_{Y|X}[(Y - c_x)^2 | X = x] \\ &= \mathbb{E}_{Y|X} \left[\frac{d}{dc_x} (Y - c_x)^2 | X = x \right] \\ &= \mathbb{E}_{Y|X}[-2(Y - c_x) | X = x] = -2\mathbb{E}_{Y|X}[Y | X = x] + 2c_x \end{aligned}$$

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

Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

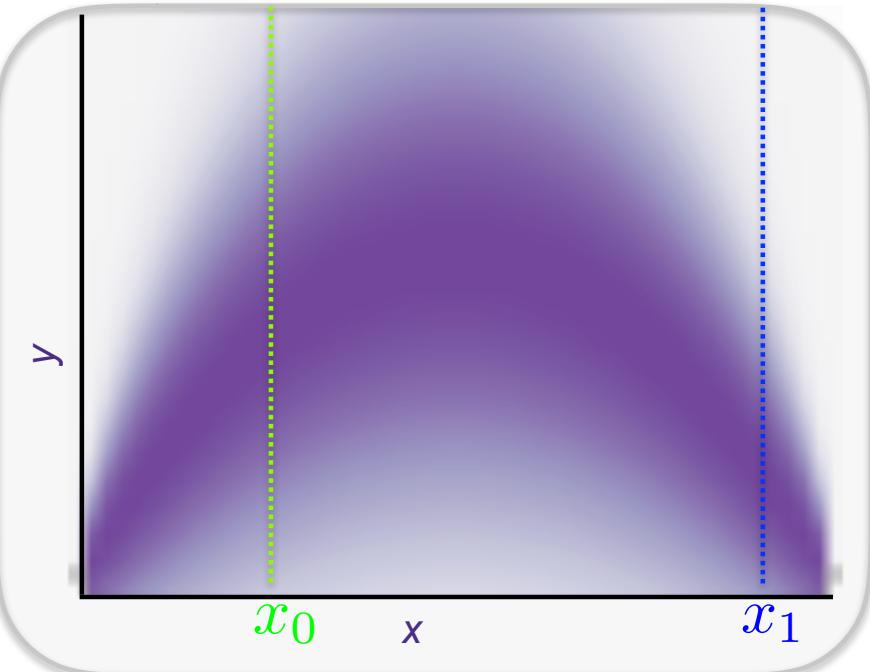
$$P_{XY}(X = x, Y = y)$$



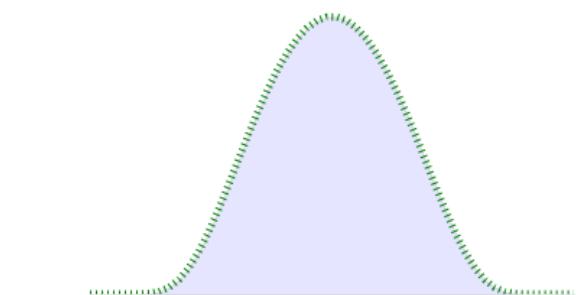
Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

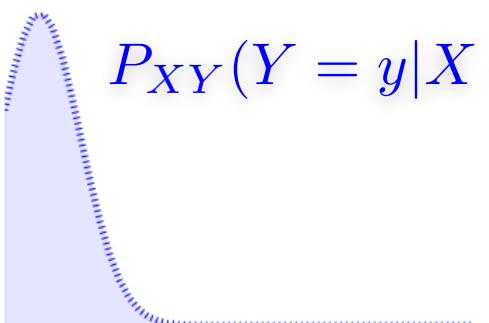
$$P_{XY}(X = x, Y = y)$$



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



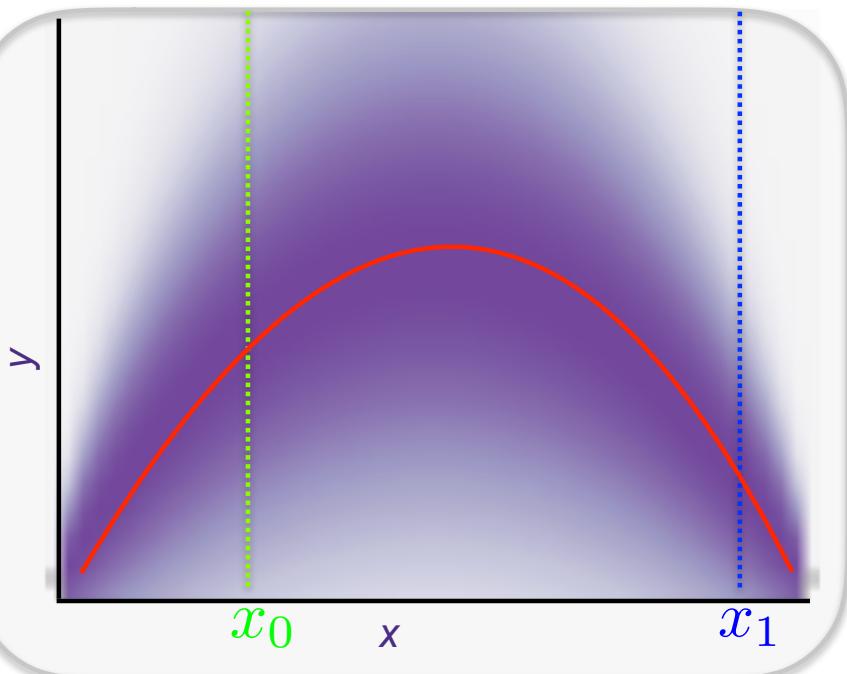
$$P_{XY}(Y = y|X = x_1)$$



Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

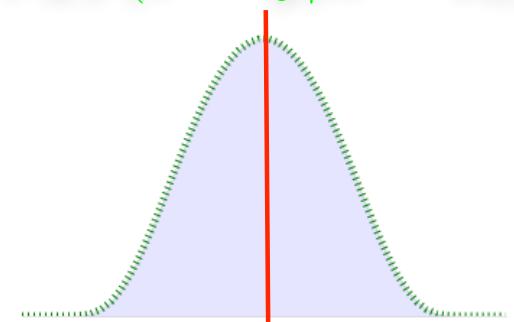
$$P_{XY}(X = x, Y = y)$$



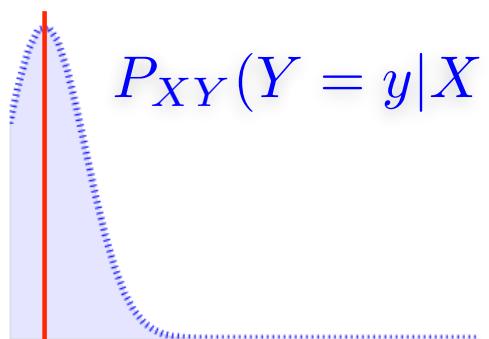
Ideally, we want to find:

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

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

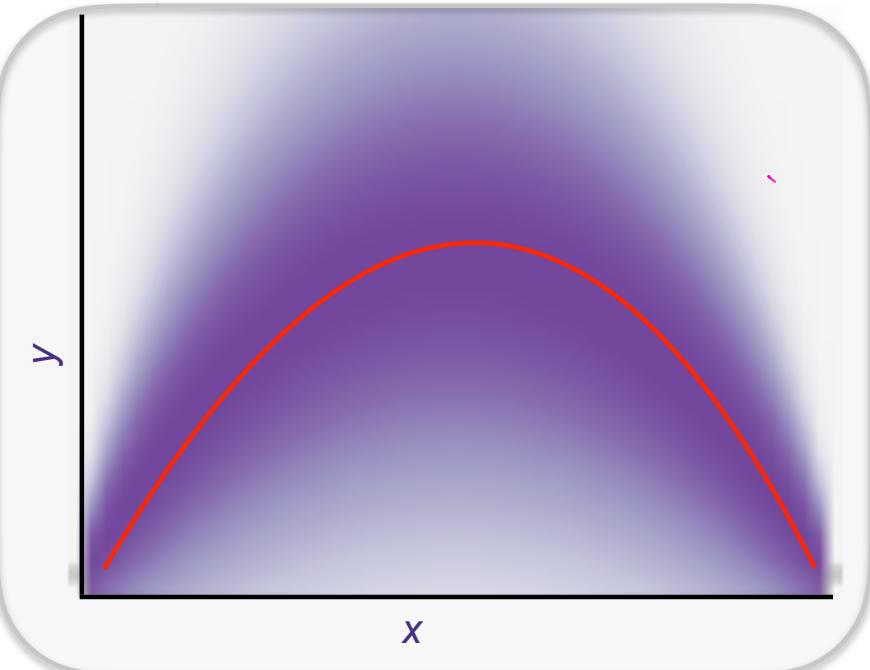


$$P_{XY}(Y = y|X = x_1)$$



Statistical Learning

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

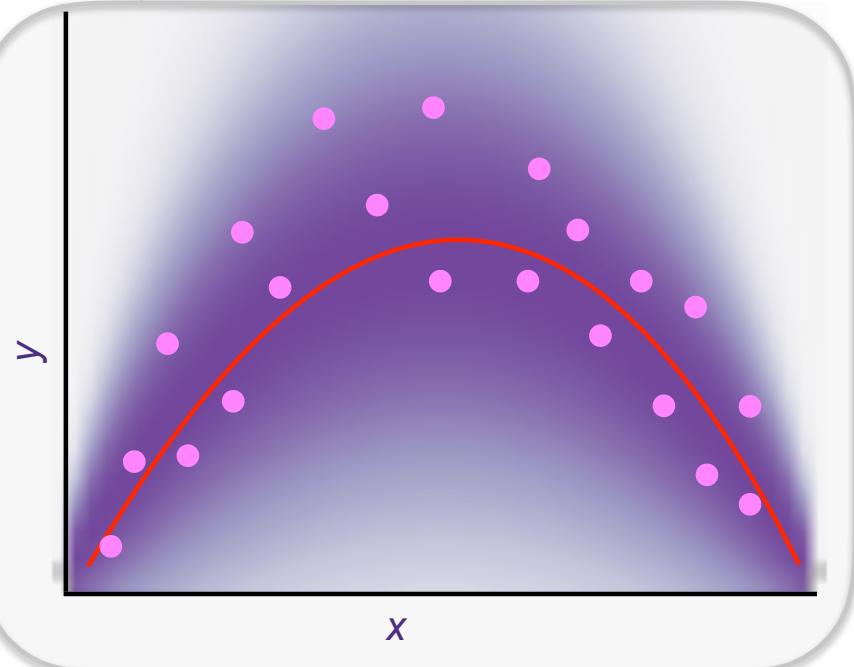


Ideally, we want to find:

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

Statistical Learning

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



Ideally, we want to find:

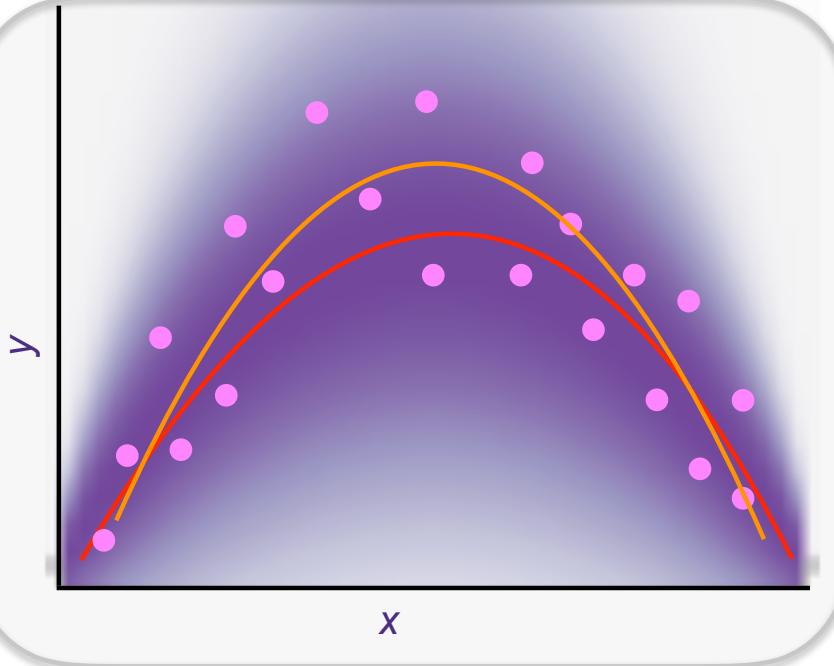
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

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:

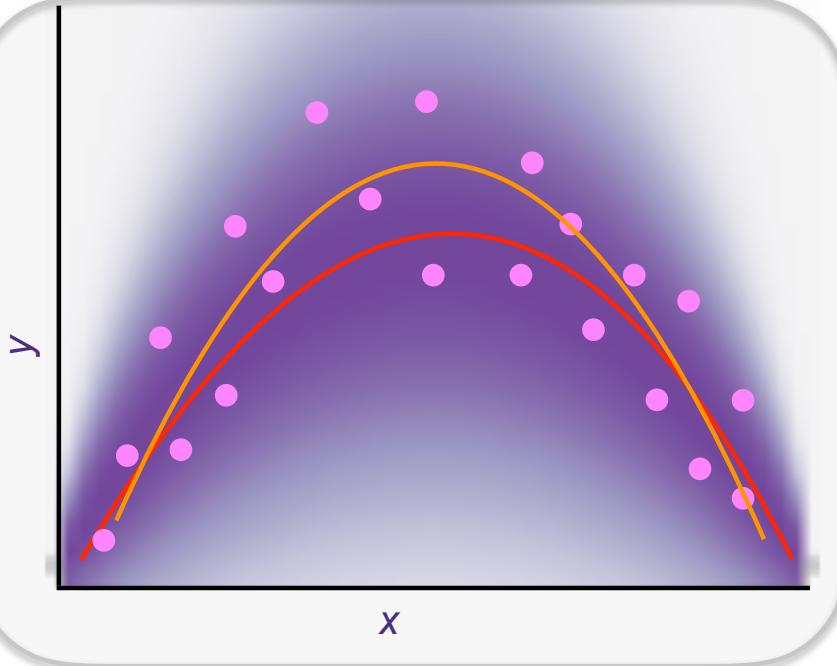
$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear)
so we compute:

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

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:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear)
so we compute:

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

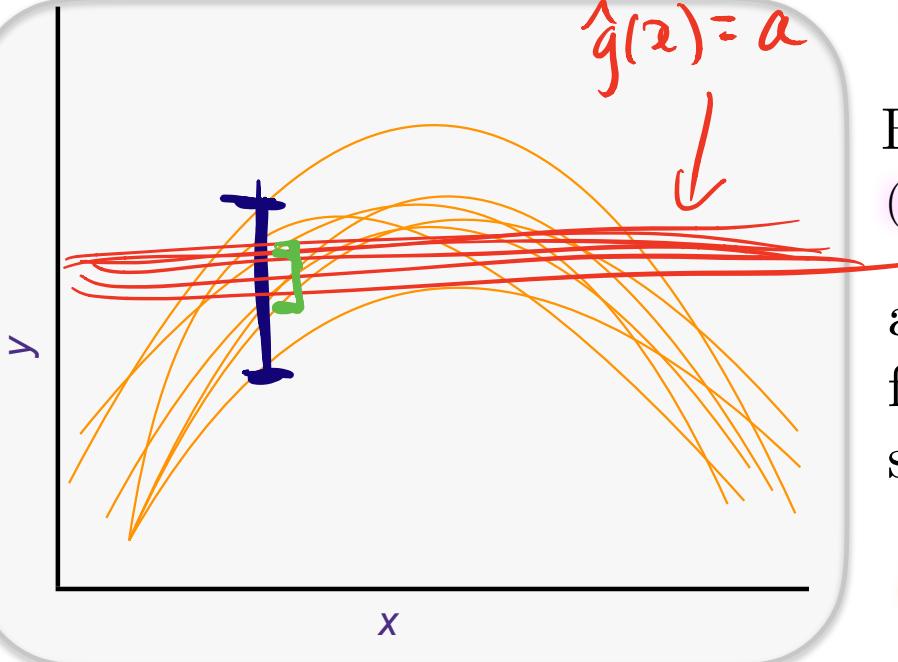
We care about future predictions: $\mathbb{E}_{XY}[(Y - \hat{f}(X))^2]$

Statistical Learning

$$\hat{f}(x) = a + bx + cx^2$$

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

$$\hat{g}(x) = a$$



Ideally, we want to find:

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

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear) so we compute:

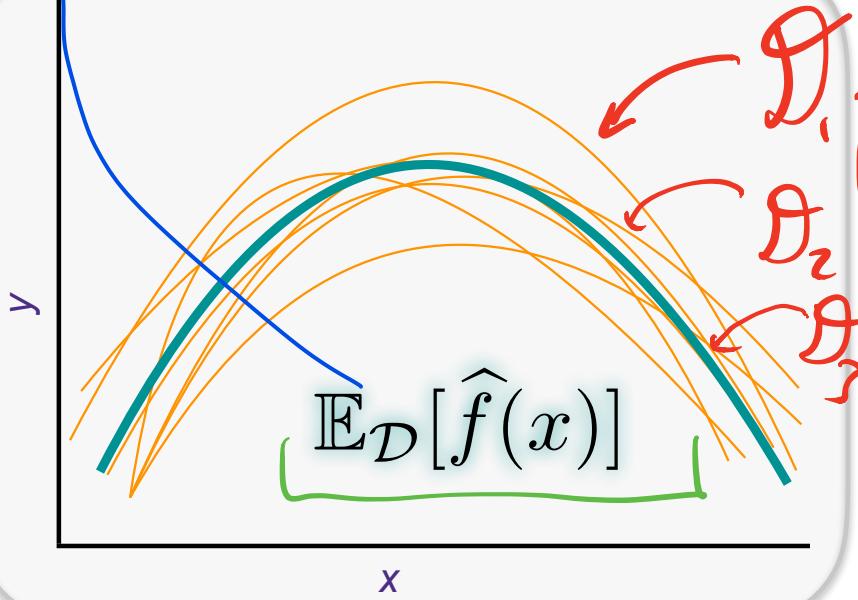
$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

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

Statistical Learning

farkli veri setleri üzerinde estimator ortalamasi (expected value)

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



Ideally, we want to find:

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

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear) so we compute:

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

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

Bias-Variance Tradeoff

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

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

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

Bias-Variance Tradeoff

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

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

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

irreducible error

Caused by stochastic
label noise

learning error

Caused by either using too
“simple” of a model or not
enough data to learn the model
accurately

Bias-Variance Tradeoff

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

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

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

Bias-Variance Tradeoff

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

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

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

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

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

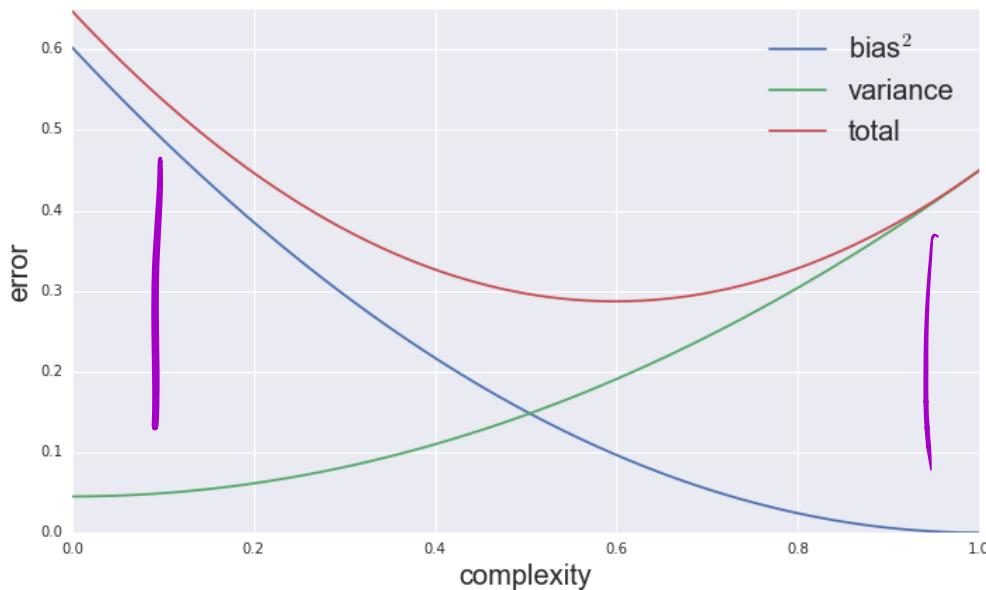
biased squared

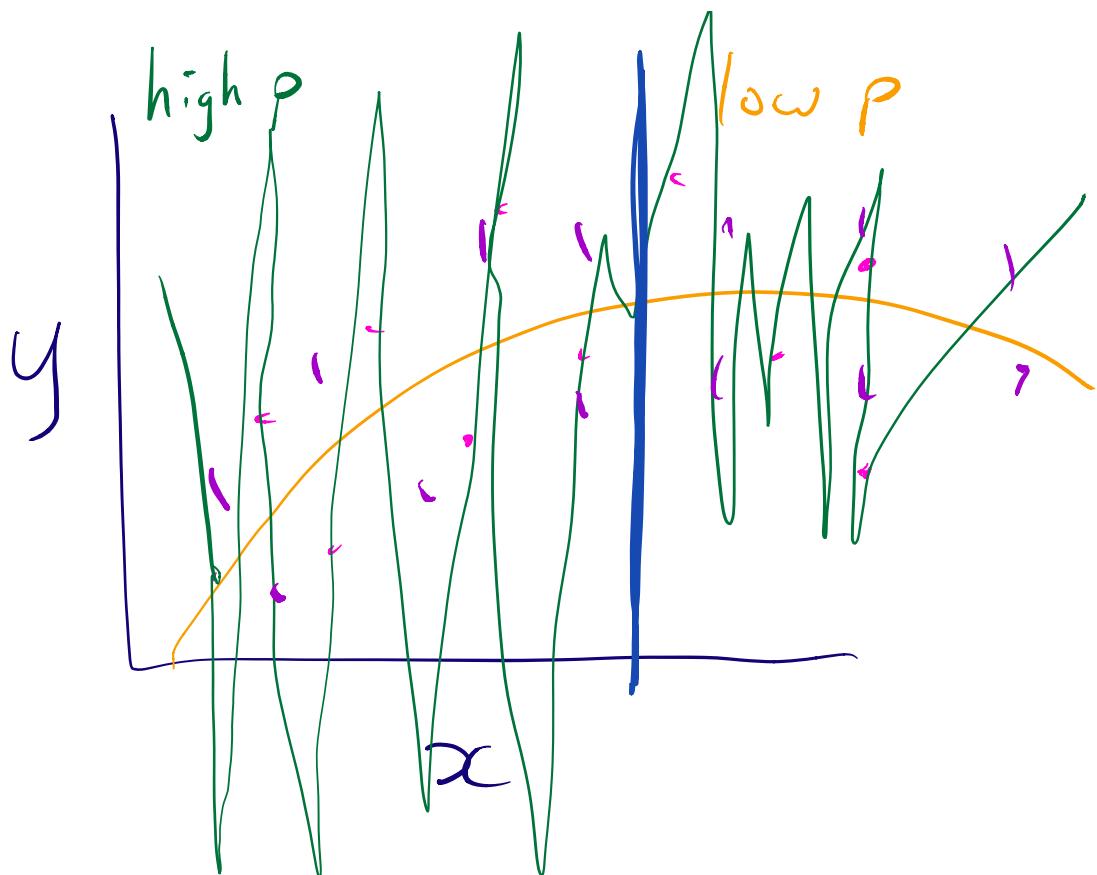
variance

Bias-Variance Tradeoff

$\gamma(x)$

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2] | X = x] = \underbrace{\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}} + \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 + \mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{biased squared variance}}$$





$$\hat{f}(x) = \sum_{n=0}^{\rho} a_n x^n$$

MORE PARAMETERS TO FIT, MORE VARIANCE



$$X = \begin{bmatrix} -x_1^T & - \\ -x_2^T & - \\ \vdots & \\ -x_n^T & - \end{bmatrix}$$

"Standardize data"

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{column vector in } \mathbb{R}^d)$$

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_{i,j} \leftarrow \begin{array}{l} \text{jth component} \\ \text{of the ith example} \end{array}$$

$$\sigma_j^2 = \frac{1}{n-1} \sum_{i=1}^n (x_{i,j} - \mu_j)^2$$

Overfitting

$\forall i, j$ set $\tilde{x}_{i,j} = (x_{i,j} - \mu_j) / \sigma_j$

$$\tilde{\mu}_j = \frac{1}{n} \sum_{i=1}^n \tilde{x}_{i,j} = 0$$

$$\tilde{\sigma}_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\tilde{x}_{i,j} - \tilde{\mu}_j)^2$$

$$= 1$$

W

$X^T \mathbf{1} = 0$ çünkü X mean 0.

Your data $\{(x_i, y_i)\}_{i=1}^n$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Want to fit $w \in \mathbb{R}^d$, $b \in \mathbb{R}$ s.t. $y_i \approx x_i^T w + b$

$$\hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (w^T x_i + b))^2 = \underset{w, b}{\operatorname{argmin}} \|y - (Xw + b\mathbf{1})\|_2^2$$

$$\Rightarrow \begin{cases} \tilde{X}^T \tilde{X} \tilde{w} + \tilde{b} \tilde{\mathbf{1}}^T \tilde{\mathbf{1}} = \tilde{X}^T \tilde{y} \\ \tilde{y}^T \tilde{X} \tilde{b} + \tilde{b} \underbrace{\mathbf{1}^T \mathbf{1}}_n = \mathbf{1}^T \tilde{y} \end{cases}$$

$$\tilde{w} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{y}$$

$$\tilde{b} = \frac{1}{n} \mathbf{1}^T \tilde{y} = 0$$

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_{i,j}$$

$$\tilde{x}_{i,j} = x_{i,j} - \mu_j$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$$

$$\tilde{y}_i = y_i - \bar{y}$$

$\Rightarrow \tilde{x}, \tilde{y}$ are both mean 0.

$$\tilde{y}_i \approx \tilde{x}_i^T \tilde{w}$$

$$\tilde{y}_i = y_i - \bar{y} \approx (x_i - \mu)^T \tilde{w}$$

$$\begin{aligned} \Rightarrow y_i &\approx (x_i - \mu)^T \tilde{w} + \bar{y} \\ &= x_i^T \tilde{w} - \mu^T \tilde{w} + \bar{y} \\ &= \underline{x_i^T \tilde{w}} + (\bar{y} - \mu^T \tilde{w}) \end{aligned}$$

$$\text{So let } \hat{w} = \tilde{w}, \quad \hat{b} = \bar{y} - \mu^T \tilde{w}$$

$$y_i \approx x_i^T \hat{w} + \hat{b}$$

$$\hat{\omega}, \hat{b} = \underset{\omega, b}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\omega^T x_i + b))^2$$

$$= \underset{\omega, b}{\operatorname{argmin}} \|y - (X\omega + b\mathbb{1})\|_2^2$$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} -x_1^T - \\ \vdots \\ -x_n^T - \end{bmatrix}$$

$$\nabla_{\omega}(\cdot) = 2 X^T (y - (X\omega + b\mathbb{1}))$$

$$= 2(X^T y - X^T X\omega - b X^T \mathbb{1}) = 0$$

$$X^T X\omega + b X^T \mathbb{1} = X^T y$$

$$\nabla_b(\cdot) = \rightarrow$$

$$\{(x_i, y_i)\}_{i=1}^n, \quad \text{Fit } \omega, b \text{ s.t. } y_i \approx x_i^T \omega + b$$

$$x_i^T \omega + b = [x_{i,1}, \dots, x_{i,d}]^T \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_d \end{bmatrix} + 1 \cdot b$$

$$\vec{x} = \begin{bmatrix} x, 1 \end{bmatrix} = \underbrace{[x_{i,1}, \dots, x_{i,d}, 1]}_{\vec{x}_i} \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_d \\ b \end{bmatrix}$$

$$\|y - \vec{x}\vec{\omega}\|_2^2 \equiv \|y - (X\omega + b\mathbb{1})\|_2^2 =: \vec{x}_i^T \vec{\omega}$$

$x_i \rightarrow \vec{x}_i$ by appending a 1 to the end

$\omega \rightarrow \vec{\omega}$ by inflating dim by 1 $\sum_{i=1}^n (y_i - \vec{\omega}^T \vec{x}_i)^2$

Bias-Variance Tradeoff

- > Choice of hypothesis class introduces learning bias
 - More complex class → less bias
 - More complex class → more variance
- > But in practice??

Bias-Variance Tradeoff

- > Choice of hypothesis class introduces learning bias
 - More complex class → less bias
 - More complex class → more variance
- > But in practice??
- > Before we saw how increasing the feature space can increase the complexity of the learned estimator:

$$\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \subset \dots$$

degree 1 poly *degree 2* *3* *4.* *---*

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

Complexity grows as k grows

Training set error as a function of model complexity

$$\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \subset \dots$$

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

TRAIN error:

$$\mathcal{D} \stackrel{i.i.d.}{\sim} P_{XY}$$
$$\frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - \widehat{f}_{\mathcal{D}}^{(k)}(x_i))^2$$

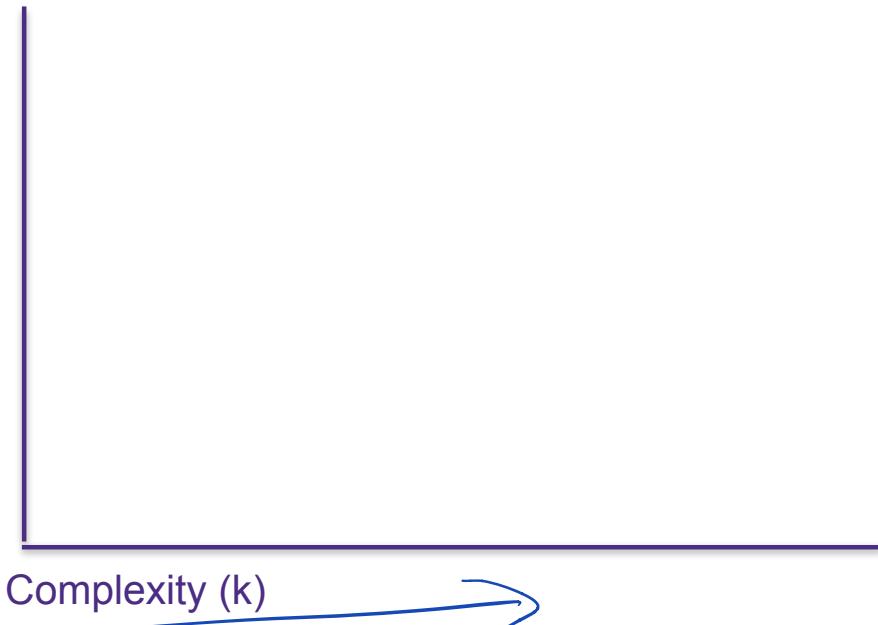
TRUE error:

$$\mathbb{E}_{XY}[(Y - \widehat{f}_{\mathcal{D}}^{(k)}(X))^2]$$

Training set error as a function of model complexity

$$\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \subset \dots$$

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



TRAIN error:

$$\mathcal{D} \stackrel{i.i.d.}{\sim} P_{XY}$$

$$\frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - \widehat{f}_{\mathcal{D}}^{(k)}(x_i))^2$$

TRUE error:

$$\mathbb{E}_{XY}[(Y - \widehat{f}_{\mathcal{D}}^{(k)}(X))^2]$$

TEST error:

$$\mathcal{T} \stackrel{i.i.d.}{\sim} P_{XY}$$

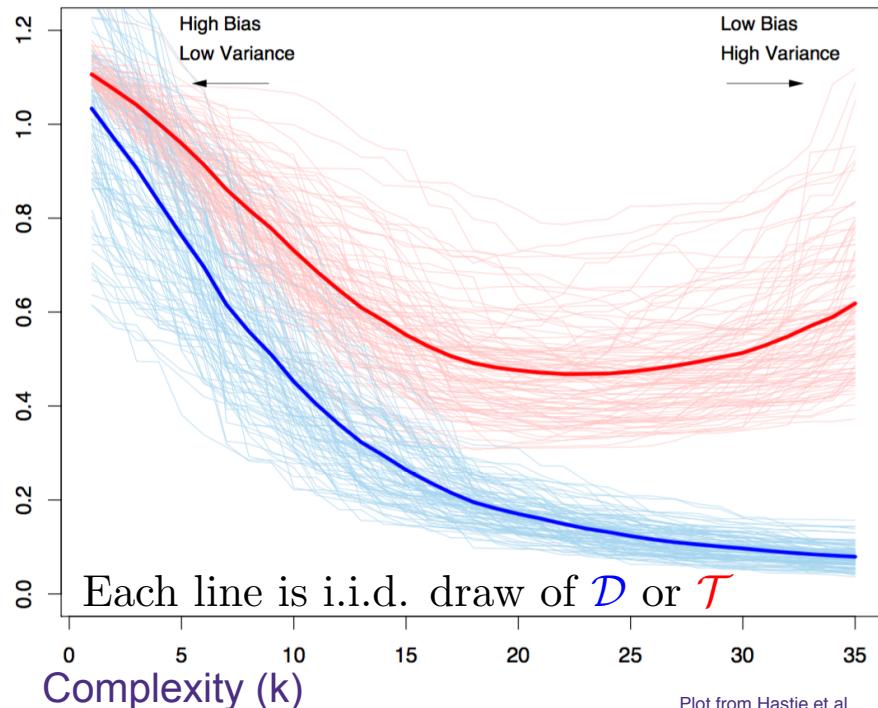
$$\left| \frac{1}{|\mathcal{T}|} \sum_{(x_i, y_i) \in \mathcal{T}} (y_i - \widehat{f}_{\mathcal{D}}^{(k)}(x_i))^2 \right|$$

Important: $\mathcal{D} \cap \mathcal{T} = \emptyset$

Training set error as a function of model complexity

$$\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \subset \dots$$

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



TRAIN error:

$$\mathcal{D} \stackrel{i.i.d.}{\sim} P_{XY}$$

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

TRUE error:

$$\mathbb{E}_{XY}[(Y - \hat{f}_{\mathcal{D}}^{(k)}(X))^2]$$

= int. err. + bias² + variance

TEST error:

$$\mathcal{T} \stackrel{i.i.d.}{\sim} P_{XY}$$

$$\mathbb{E}_{\mathcal{T}} \left[\frac{1}{|\mathcal{T}|} \sum_{(x_i, y_i) \in \mathcal{T}} (y_i - \hat{f}_{\mathcal{D}}^{(k)}(x_i))^2 \right]$$

Important: $\mathcal{D} \cap \mathcal{T} = \emptyset$

Training set error as a function of model complexity

$$\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \subset \dots$$

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

TRAIN error is optimistically biased because it is evaluated on the data it trained on. TEST error is unbiased only if T is never used to train the model or even pick the complexity k.



TRAIN error:

$$\mathcal{D} \stackrel{i.i.d.}{\sim} P_{XY}$$

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

TRUE error:

$$\mathbb{E}_{XY}[(Y - \hat{f}_{\mathcal{D}}^{(k)}(X))^2]$$

TEST error:

$$\mathcal{T} \stackrel{i.i.d.}{\sim} P_{XY}$$

$$\frac{1}{|\mathcal{T}|} \sum_{(x_i, y_i) \in \mathcal{T}} (y_i - \hat{f}_{\mathcal{D}}^{(k)}(x_i))^2$$

Important: $\mathcal{D} \cap \mathcal{T} = \emptyset$

Test set error

> Given a dataset, randomly split it into two parts:

- Training data: \mathcal{D}
- Test data: \mathcal{T}

Important: $\mathcal{D} \cap \mathcal{T} = \emptyset$

> Use training data to learn predictor

- e.g., $\frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - \hat{f}_{\mathcal{D}}^{(k)}(x_i))^2$
- use training data to pick complexity k

> Use test data to report predicted performance

$$\frac{1}{|\mathcal{T}|} \sum_{(x_i, y_i) \in \mathcal{T}} (y_i - \hat{f}_{\mathcal{D}}^{(k)}(x_i))^2$$

How many points do I use for training/testing?

- > Very hard question to answer!
 - Too few training points, learned model is bad
 - Too few test points, you never know if you reached a good solution
- > Bounds, such as Hoeffding's inequality can help:

$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$

- > More on this later the quarter, but still hard to answer
- > Typically:

- If you have a reasonable amount of data 90/10 splits are common
- If you have little data, then you need to get fancy (e.g., bootstrapping)