

# Bias-Variance Tradeoff

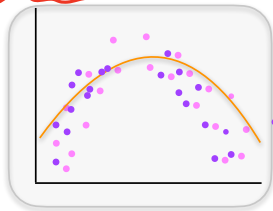
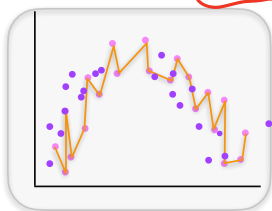
complexity  $\uparrow \Rightarrow$  bias  $\downarrow$  variance  $\uparrow$   
 complexity  $\downarrow \Rightarrow$  bias  $\uparrow$  variance  $\downarrow$

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

irreducible error

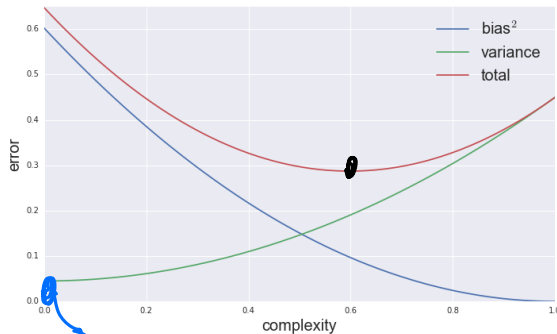
learning err: 
$$+(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 + \underbrace{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}}$$

bias squared variance



If we re-drew our data, what the LS training error estimator look like for generalized linear functions in small p/large p dimensions?

Q: what is variance of a constant predictor  $f_{\eta}(x) = c$ ,  $c$  is independent of data



## Example: Linear LS

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

$n \times d$        $n$

$$w \in \mathbb{R}^d$$

$$Y = Xw + \epsilon$$

$\epsilon \in \mathbb{R}^n$

Assumption

if  $y_i = x_i^T w + \epsilon_i$  and  $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$x \in \mathbb{R}^d \quad \eta(x) = \mathbb{E}_{Y|X} [Y | X=x] = \mathbb{E}_{Y|X} [x^T w + \epsilon | X=x] = w^T x$$

MLE

$$\hat{w} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (Xw + \epsilon) = w + (X^T X)^{-1} X^T \epsilon$$

new point  $x_{\text{new}}$   $\hat{f}_D(x_{\text{new}}) = x_{\text{new}}^T \hat{w} = x_{\text{new}}^T w + x_{\text{new}}^T (X^T X)^{-1} X^T \epsilon$

irreducible error

$$\begin{aligned} \mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X=x] &= \mathbb{E}_{Y|X} [(w^T x + \epsilon - w^T x)^2 | X=x] \\ &= \mathbb{E} [\epsilon^2] = \sigma^2 \end{aligned}$$

## Example: Linear LS: compute bias

$$\mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\mathcal{D} = \{ (x_i, y_i) \}_{i=1}^n$$

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

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

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

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

bias squared

$$\begin{aligned} x_{\text{new}} \text{ is fixed } \quad \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x_{\text{new}})] &= \mathbb{E}_{\mathcal{D}}[x_{\text{new}}^T w + x_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \\ &= x_{\text{new}}^T w + \mathbb{E}_{\mathcal{X}}[x_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underbrace{\mathbb{E}_{\mathcal{Y}}[\epsilon]}_{=0}] \\ &= x_{\text{new}}^T w = \eta(x_{\text{new}}) \end{aligned}$$

$\Rightarrow$  unbiased

## Example: Linear LS: compute variance

$$\text{tr}(\mathbf{I}_d) = \text{tr}\begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix} = d.$$

$$\mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

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

$$\hat{f}_D(x) = \hat{w}^T x = w^T x + \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x$$

$$\mathcal{X} = \mathcal{X}_{\text{new}}$$

$$\mathbb{E}[\epsilon \epsilon^T] = \sigma^2 \mathbf{I}$$

$$\mathbb{E}[\epsilon \epsilon^T]_{ij} = \mathbb{E}[\epsilon_i \epsilon_j]$$

$$\mathbf{X}^T \mathbf{X} = \sum_{i=1}^n x_i x_i^T \xrightarrow{n \rightarrow \infty} n \Sigma$$

central limit theorem  
 $\Sigma = \mathbb{E}_X[\mathbf{X} \mathbf{X}^T]$

$$\begin{aligned} \text{variance} \quad \mathbb{E}_D[(\mathbb{E}_D[\hat{f}_D(x)] - \hat{f}_D(x))^2] &= \mathbb{E}_D[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x] \\ &= \mathbb{E}_X[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}_Y[\epsilon \epsilon^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x] \\ &= \sigma^2 \mathbb{E}_X[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x] \\ &= \sigma^2 \mathbb{E}_X[x^T (\mathbf{X}^T \mathbf{X})^{-1} x] \end{aligned}$$

$$\xrightarrow{n \rightarrow \infty} \mathbb{E}_X[\mathbb{E}_D[(\mathbb{E}_D[\hat{f}_D(x)] - \hat{f}_D(x))^2]] = \mathbb{E}_X[\sigma^2 x^T (n \Sigma)^{-1} x]$$

Trace formula  

$$\begin{aligned} &= \mathbb{E}_X[\sigma^2 \text{tr}((n \Sigma)^{-1} \mathbf{X} \mathbf{X}^T)] \\ &= \sigma^2 \text{tr}((n \Sigma)^{-1} \Sigma) \\ &= \frac{\sigma^2}{n} \text{tr}(\mathbf{I}) \\ &= \frac{\sigma^2 d}{n} \end{aligned}$$

## Example: Linear LS

$$\mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

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

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

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

$$\underbrace{\mathbb{E}_{XY}[(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}} = \sigma^2 \quad \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{bias squared}} = 0$$

$$\mathbb{E}_{X=x} \underbrace{[\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]]}_{\text{variance}} = \frac{\sigma^2}{n}$$

as  $d \uparrow \Rightarrow \uparrow \text{variance}$

# Overfitting

---



# Bias-Variance Tradeoff

---

> Choice of hypothesis class introduces learning bias

- More complex class  $\rightarrow$  less bias  $\exists f \in \mathcal{F} \quad \eta \approx f$
- More complex class  $\rightarrow$  more variance

> But in practice??

# Bias-Variance Tradeoff

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

linear   quadratic   cubic   . . .

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

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

Complexity grows as k grows



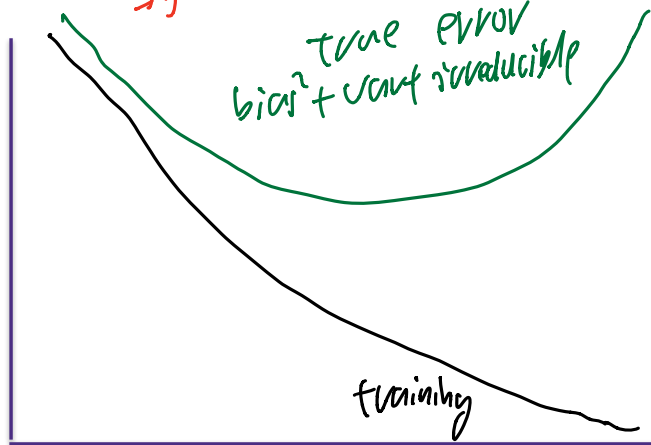
# Training set error as a function of model complexity

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

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

*Handwritten red note:  $\hat{f}$*

Error



Complexity (k)

**TRAIN error:**

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

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

**TRUE error:**

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

## Training set error as a function of model complexity

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

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

all data  $\mathcal{D} \cup \mathcal{T} \sim P_{XY}$   
↑  
test set

**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:**

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

$$O\left(\frac{1}{\sqrt{|\mathcal{T}|}}\right)$$

**TEST error:**

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

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

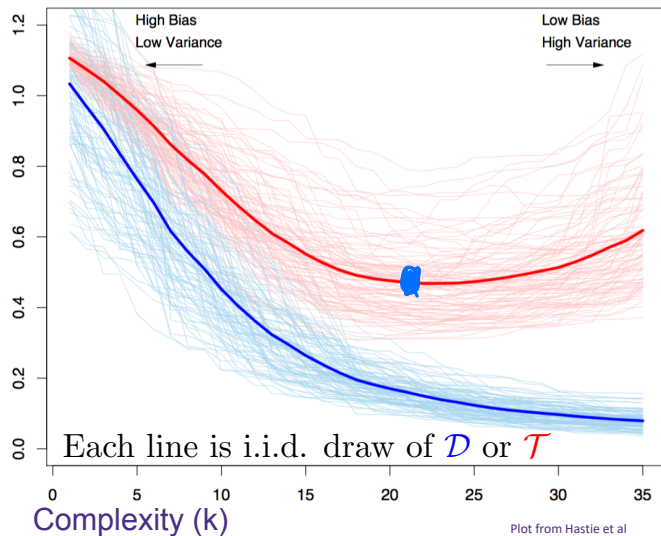
Important:  $\mathcal{D} \cap \mathcal{T} = \emptyset$

# Training set error as a function of model complexity

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

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

error



**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:**

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

plug in

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

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

# 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  $\mathcal{T}$  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:**

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

$\mathbb{E}_{\mathcal{T}}[\text{test error}] = \text{True error}$

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

# 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

## > 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)