



# Learning Theory: Why ML Works

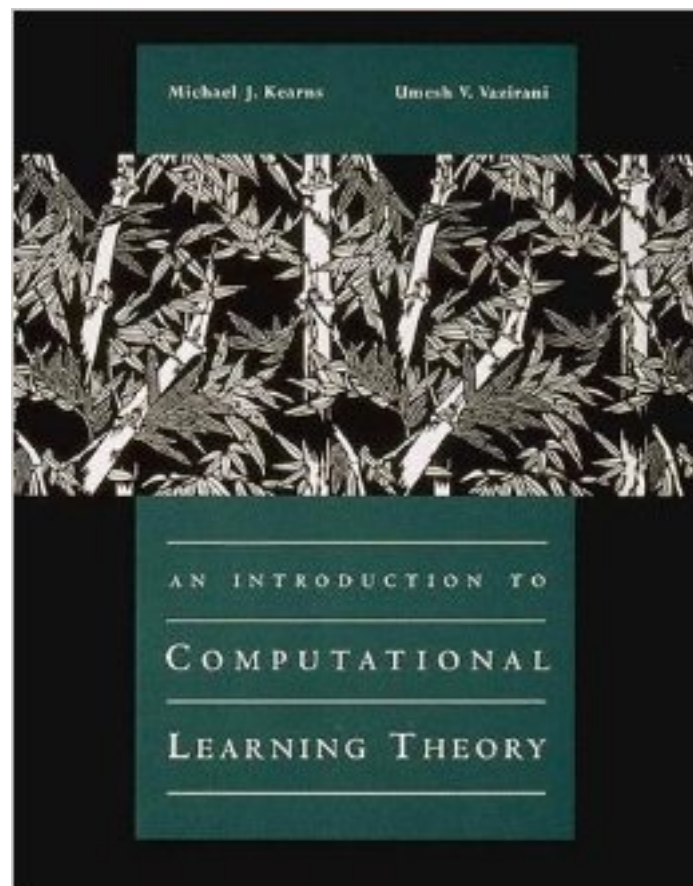
These slides were assembled by Byron Boots, with only minor modifications from Eric Eaton's slides and grateful acknowledgement to the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

# Computational Learning Theory

Entire subfield devoted to the mathematical analysis of machine learning algorithms

Has led to several practical methods:

- PAC (probably approximately correct) learning → boosting
- VC (Vapnik–Chervonenkis) theory → support vector machines



Annual conference: Conference on Learning Theory (COLT)

# Computational Learning Theory

Fundamental Question: What general laws constrain a system's ability to learn?

Seeks theory to relate:

- Probability of successful learning
- Number of training examples
- Complexity of hypothesis space
- Accuracy to which target function is approximated
- Manner in which training examples should be presented

# Sample Complexity

Assume that  $f : \mathcal{X} \mapsto \{0, 1\}$  is the target function

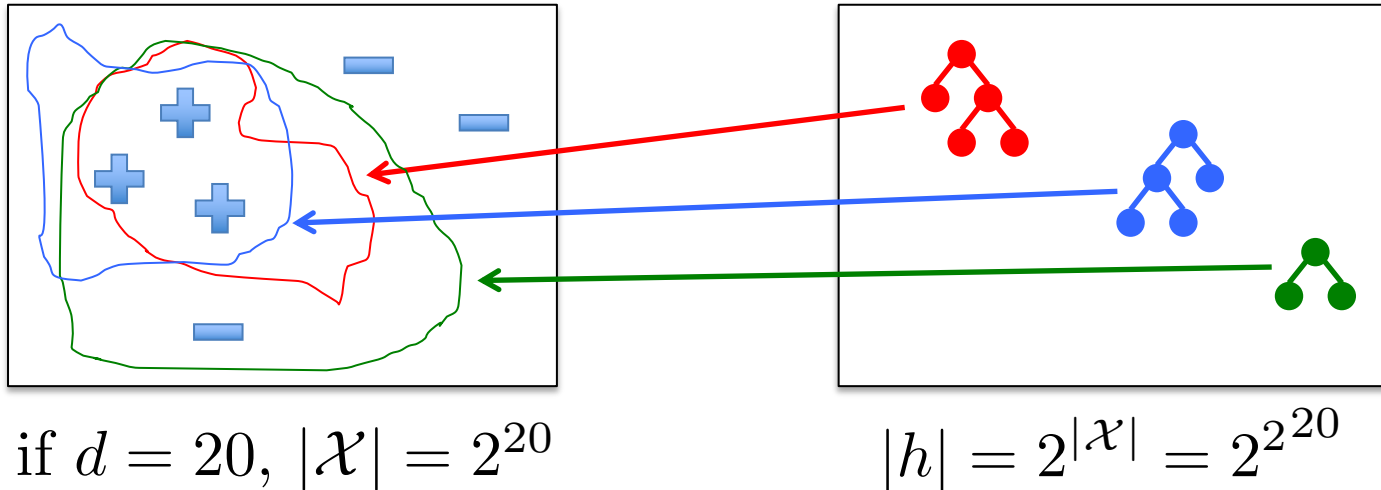
How many training examples are sufficient to learn the target function  $f$  ?

1. If learner proposed instances as queries to teacher
  - Learner proposes instance  $x$ , teacher provides  $f(x)$
2. If teacher (who knows  $f$ ) provides training examples
  - Teacher provides labeled examples in form  $\langle x, f(x) \rangle$
3. If some random process (e.g., nature) proposes instances
  - Instance  $x$  generated randomly, teacher provides  $f(x)$

# Function Approximation: The Big Picture

Instance Space  $\mathcal{X} = \{0, 1\}^d$   
 $\mathbf{x} = \langle x_1, x_2, \dots, x_d \rangle \in \mathcal{X}$

Hypothesis Space  
 $H = \{h \mid h : \mathcal{X} \mapsto \{0, 1\}\}$



- How many labeled instances are needed to determine which of the  $2^{2^{20}}$  hypotheses are correct?
  - All  $2^{20}$  instances in  $\mathcal{X}$  must be labeled!
- Generalizing beyond the training data (inductive inference) is impossible unless we add more assumptions (e.g., priors over  $H$ )

# Bias-Variance Decomposition of Squared Error

- Assume that  $y = f(\mathbf{x}) + \epsilon$ 
  - Noise  $\epsilon$  is sampled from a normal distribution with 0 mean and variance  $\sigma^2$ :  $\epsilon \sim N(0, \sigma^2)$
  - Noise lower-bounds the performance (error) we can achieve

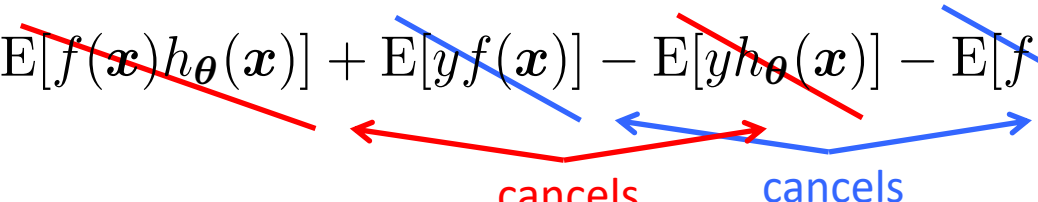
- Recall the following objective function:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left( y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^2$$

- We can view this as an approximation of the expected value of the squared error:  $E(y - h_{\boldsymbol{\theta}}(\mathbf{x}))^2$


# Bias-Variance Decomposition of Squared Error

$$\begin{aligned} E[(y - h_{\theta}(\mathbf{x}))^2] &= E[(y - f(\mathbf{x}) + f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &= E[(y - f(\mathbf{x}))^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &\quad + 2 E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))(y - f(\mathbf{x}))] \\ &= E[(y - f(\mathbf{x}))^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &\quad + 2 (E[f(\mathbf{x})h_{\theta}(\mathbf{x})] + E[yf(\mathbf{x})] - E[yh_{\theta}(\mathbf{x})] - E[f(\mathbf{x})^2]) \end{aligned}$$



Therefore,

$$\begin{aligned} E[(y - h_{\theta}(\mathbf{x}))^2] &= E[(y - f(\mathbf{x}))^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &= E[\epsilon^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \end{aligned}$$

 This is actually  $\text{var}(\epsilon)$ , since mean is 0

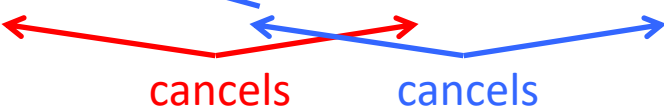
Aside:

Definition of Variance

$$\text{var}(z) = E[(z - E[z])^2]$$

# Bias-Variance Decomposition of Squared Error

$$\begin{aligned}
 \mathbb{E}[(y - h_{\theta}(\mathbf{x}))^2] &= \text{var}(\epsilon) + \mathbb{E}[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\
 &= \text{var}(\epsilon) + \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[h_{\theta}(\mathbf{x})] + \mathbb{E}[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2] \\
 &= \text{var}(\epsilon) + \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[h_{\theta}(\mathbf{x})])^2] + \mathbb{E}[(\mathbb{E}[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2] \\
 &\quad + 2\mathbb{E}[(\mathbb{E}[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))(f(\mathbf{x}) - \mathbb{E}[h_{\theta}(\mathbf{x})])] \\
 &= \text{var}(\epsilon) + \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[h_{\theta}(\mathbf{x})])^2] + \mathbb{E}[(\mathbb{E}[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2] \\
 &\quad + 2(\cancel{\mathbb{E}[f(\mathbf{x})\mathbb{E}[h_{\theta}(\mathbf{x})]]} - \cancel{\mathbb{E}[\mathbb{E}[h_{\theta}(\mathbf{x})]^2]} - \cancel{\mathbb{E}[f(\mathbf{x})h_{\theta}(\mathbf{x})]} + \cancel{\mathbb{E}[h_{\theta}(\mathbf{x})\mathbb{E}[h_{\theta}(\mathbf{x})]]})
 \end{aligned}$$



Therefore,

$$\mathbb{E}[(y - h_{\theta}(\mathbf{x}))^2] = \underbrace{\text{var}(\epsilon)}_{\text{noise}} + \underbrace{\mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[h_{\theta}(\mathbf{x})])^2]}_{\text{bias}} + \underbrace{\mathbb{E}[(\mathbb{E}[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2]}_{\text{variance}}$$

$$\mathbb{E}[(y - h_{\theta}(\mathbf{x}))^2] = \text{bias}(h_{\theta}(\mathbf{x}))^2 + \text{var}(h_{\theta}(\mathbf{x})) + \sigma^2$$



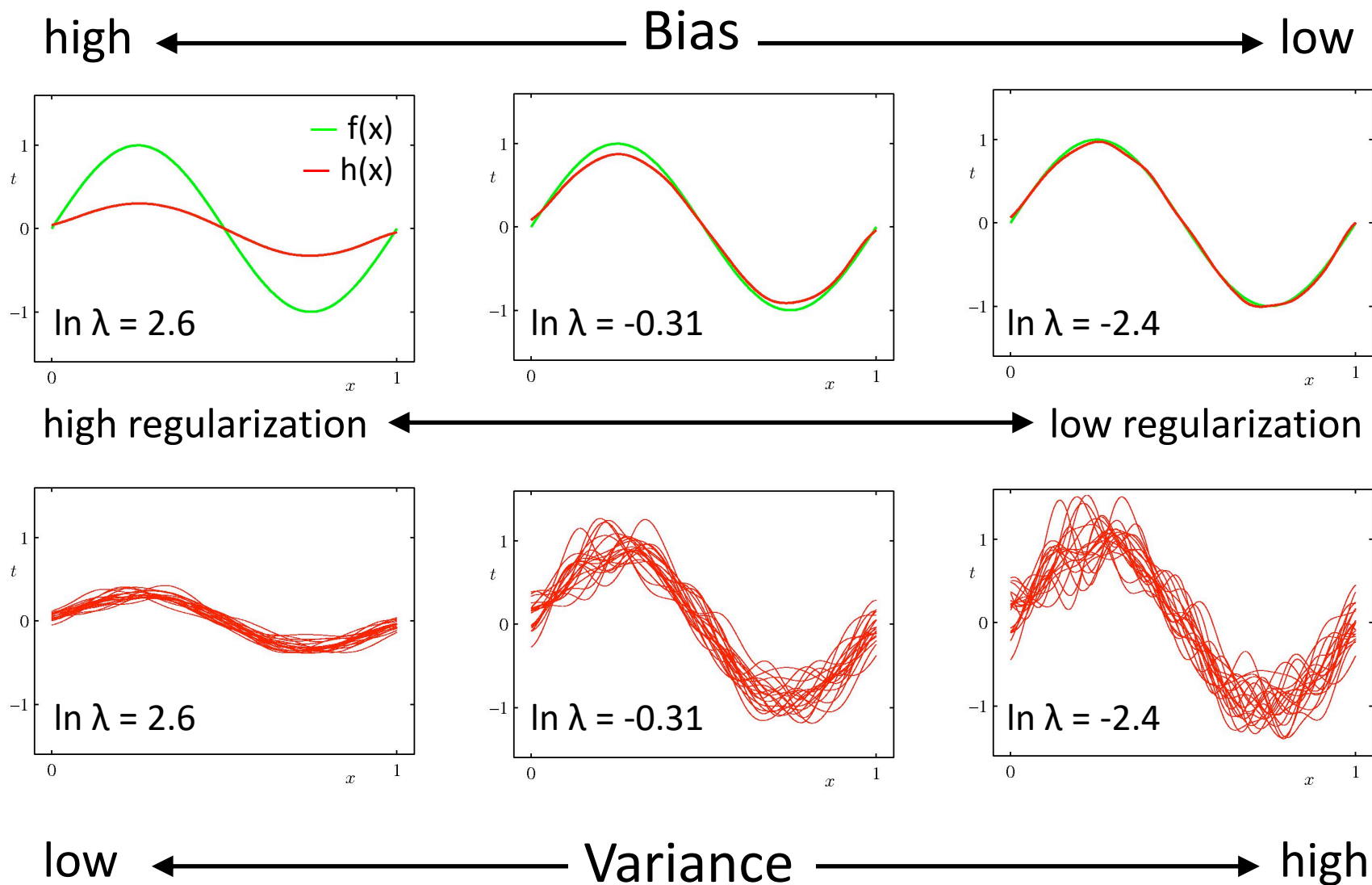
# Regularization

- Linear regression objective function

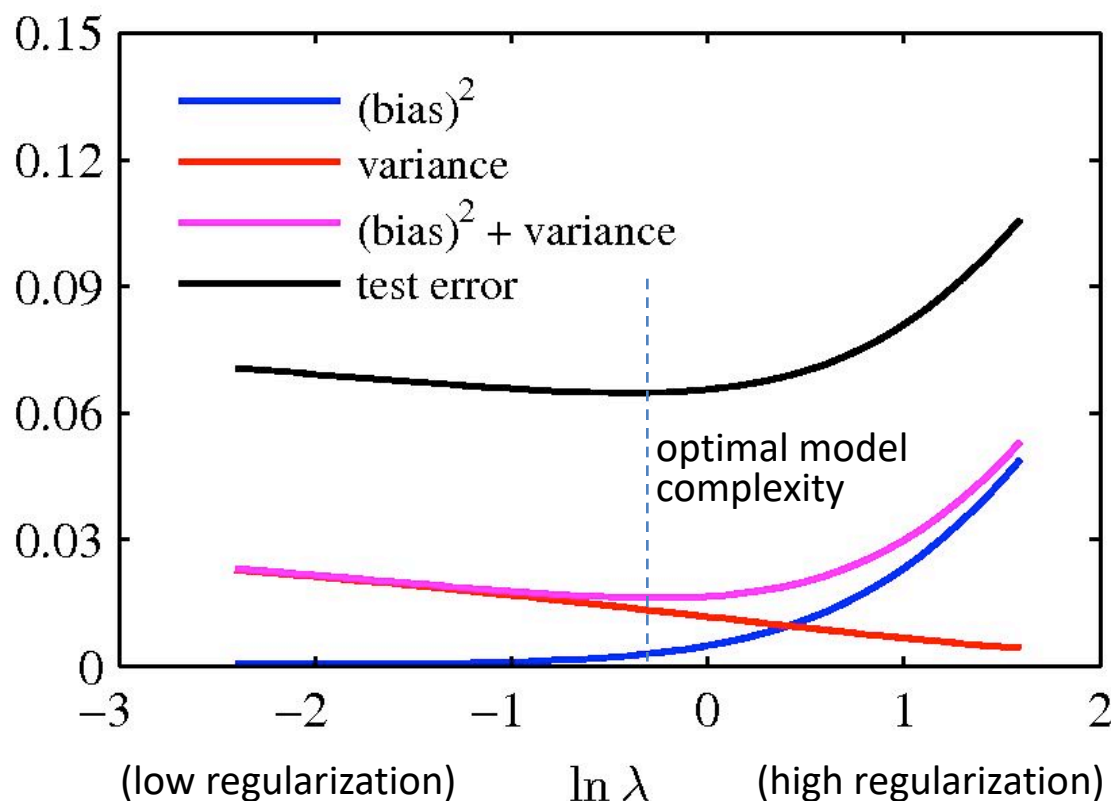
$$J(\boldsymbol{\theta}) = \underbrace{\frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}} \left( \mathbf{x}^{(i)} \right) - y^{(i)} \right)^2}_{\text{model fit to data}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d \theta_j^2}_{\text{regularization}}$$

–  $\lambda$  is the regularization parameter ( $\lambda \geq 0$ )

# Illustration of Bias-Variance



# Illustration of Bias-Variance



- Reducing training error drives down bias, but ignores variance