

Lectures 9

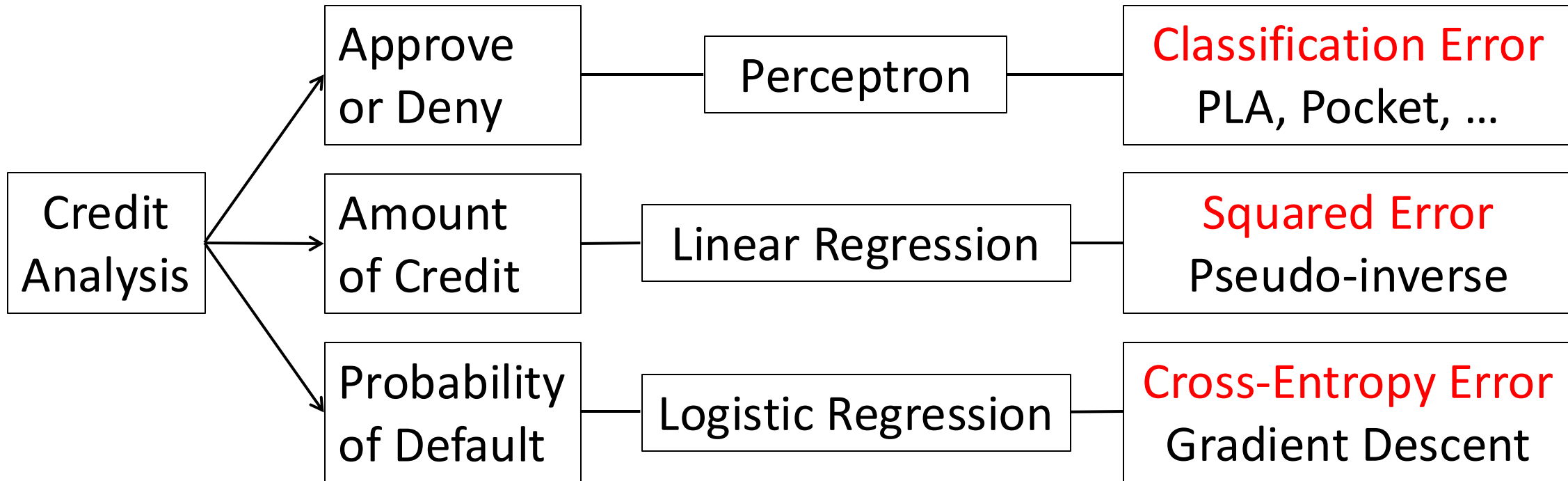
CS436/536: Introduction to Machine Learning

Zhaohan Xi

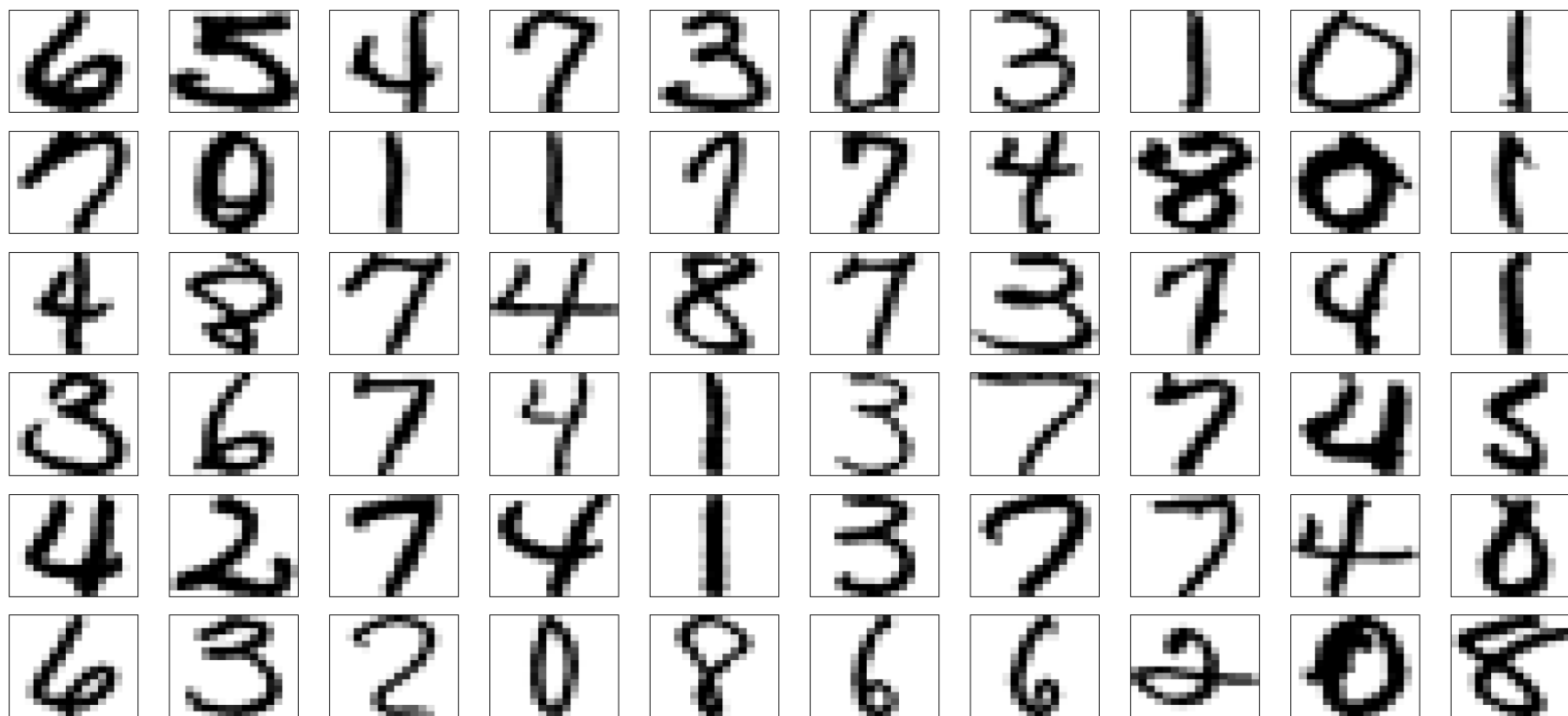
Binghamton University

zxi1@binghamton.edu

Recap: Linear Model for Three Learning Problems



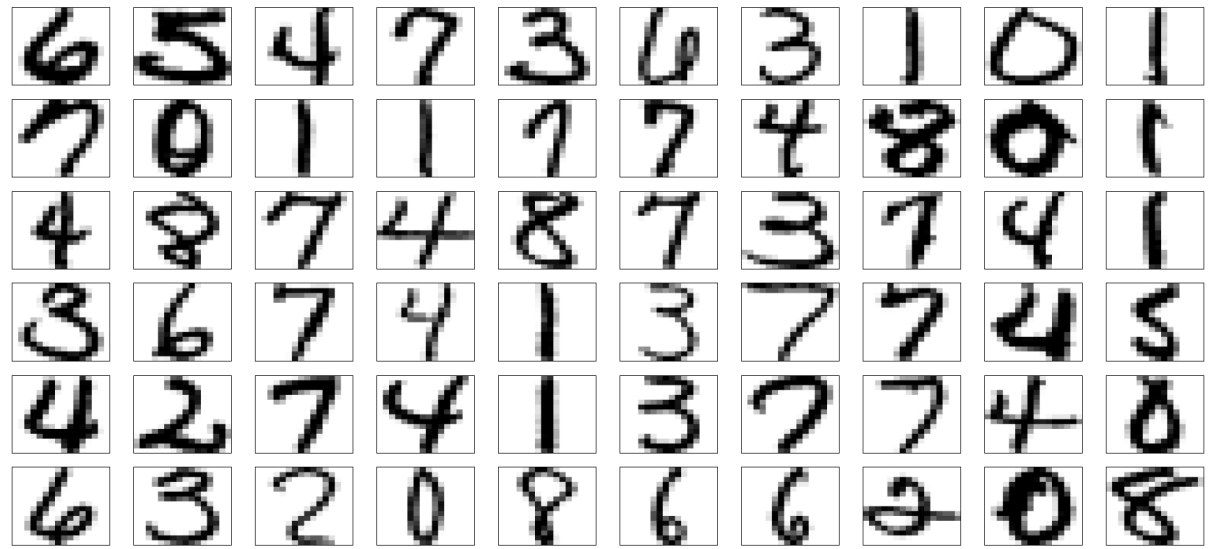
Digits Data



Each example of a digit is a 16×16 image

<http://yann.lecun.com/exdb/mnist/>

Digits Data



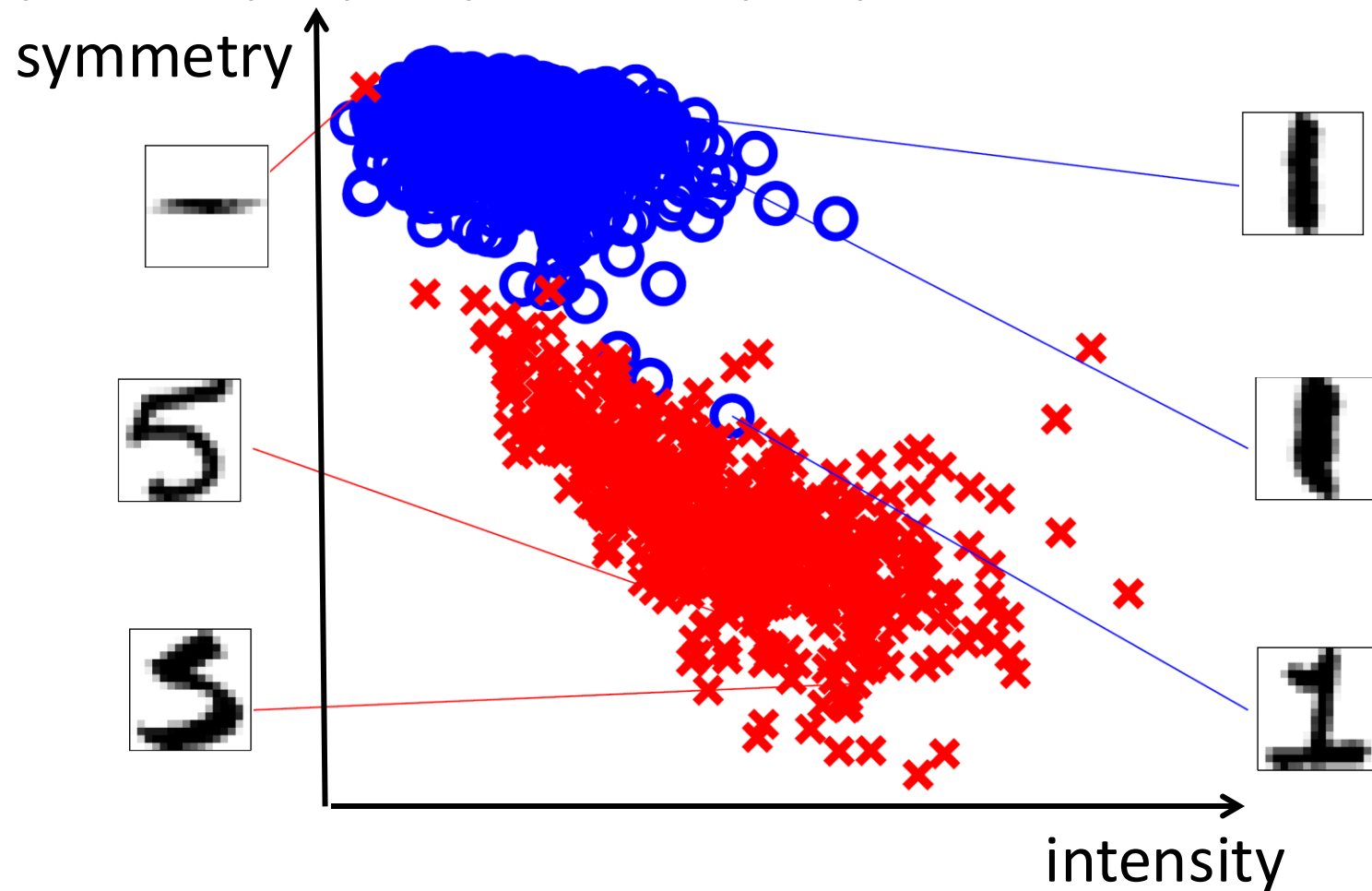
```
[-1 -1 -1 -1 -1 -1 -1 -0.63 0.86 -0.17 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -0.99 0.3 1 0.31 -1 -1 -1 -1 -1 -1
-1 -1 -1 -1 -1 -1 -1 -0.41 1 0.99 -0.57 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -0.68 0.83 1 0.56 -1 -1 -1 -1 -1 -1 -1
-1 -1 -1 -1 -1 -1 -1 -0.94 0.54 1 0.78 -0.72 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 0.1 1 0.92 -0.44 -1 -1 -1 -1 -1 -1 -1
1 -1 -1 -1 -1 -1 -1 -0.26 0.95 1 -0.16 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -0.8 0.91 1 0.3 -0.96 -1 -1 -1 -1 -1 -1 -1
0.49 1 0.88 0.09 -1 -1 -1 -1 0.28 1 0.88 -0.8 -1 -0.9 0.14 0.97 1 1 1 0.99 -0.74 -1 -1 -0.95 0.84 1
0.32 -1 -1 0.35 1 0.65 -0.10 -0.18 1 0.98 -0.72 -1 -1 -0.63 1 1 0.07 -0.92 0.11 0.96 0.30 -0.88 -1 -
0.07 1 0.64 -0.99 -1 -1 -0.67 1 1 0.75 0.34 1 0.70 -0.94 -1 -1 0.54 1 0.02 -1 -1 -1 -0.90 0.79 1 1 1 1
0.53 0.18 0.81 0.83 0.97 0.86 -0.63 -1 -1 -1 -1 -0.45 0.82 1 1 1 1 1 1 1 0.13 -1 -1 -1 -1 -1 -1 -1 -0.48
0.81 1 1 1 1 1 1 0.21 -0.94 -1 -1 -1 -1 -1 -1 -1 -1 -0.97 -0.42 0.30 0.82 1 0.48 -0.47 -0.99 -1 -1 -1 -1 -1]
```

$$\left. \begin{array}{l} x = (1, x_1, \dots, x_{256}) \\ w = (w_0, w_1, \dots, w_{256}) \end{array} \right\} d_{VC} = 257$$

Features: Intensity and Symmetry

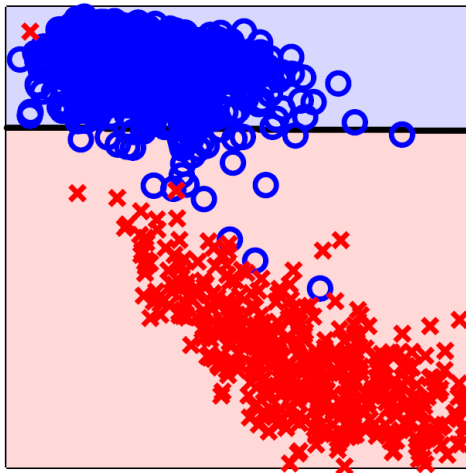
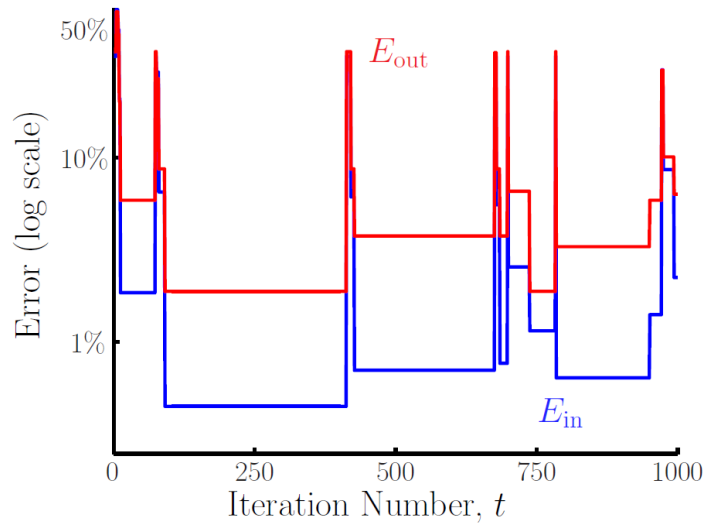
Human-like approach: Summarize image by a few features

feature: an important property of the input you think is useful for classification

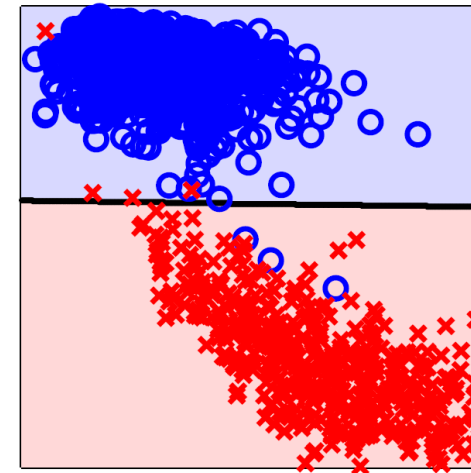
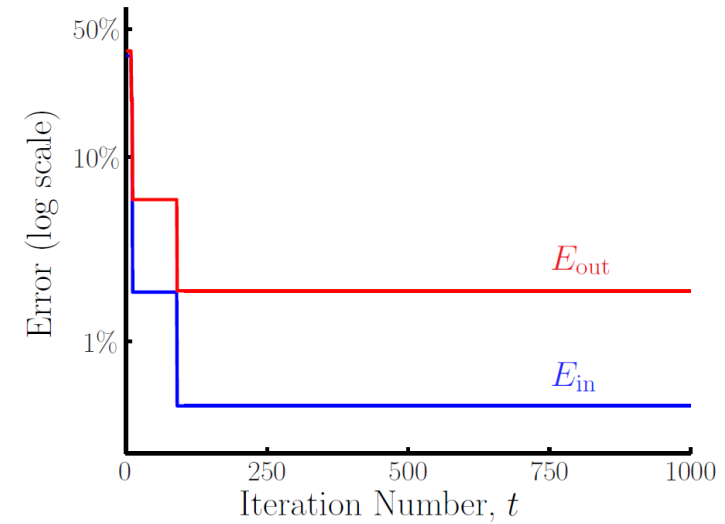


Perceptron Model on Digits Data

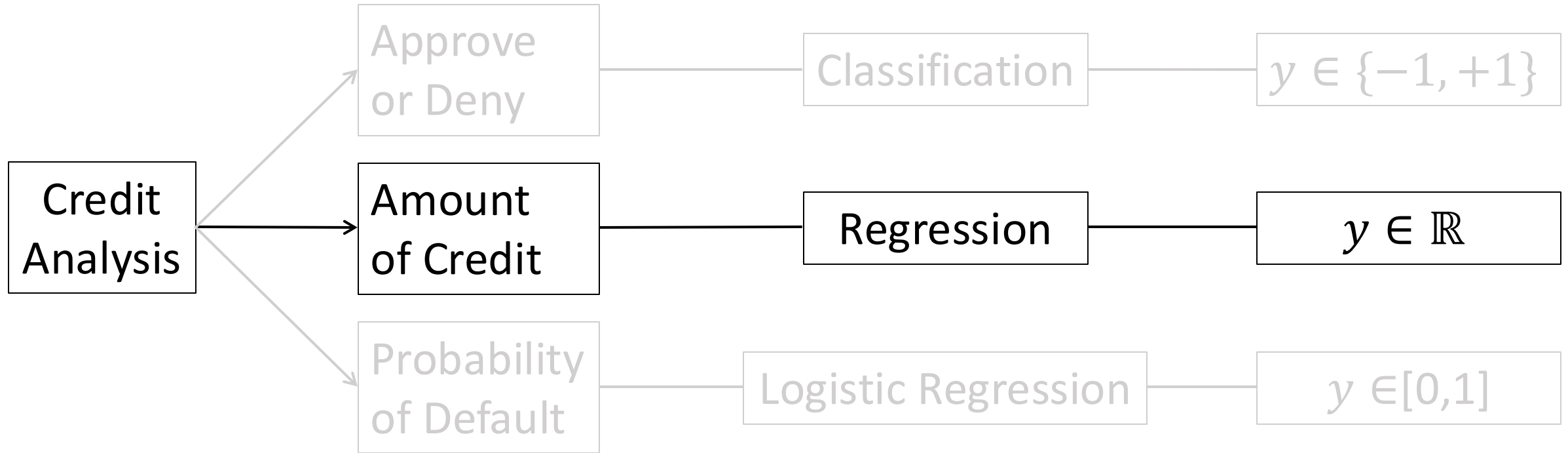
PLA



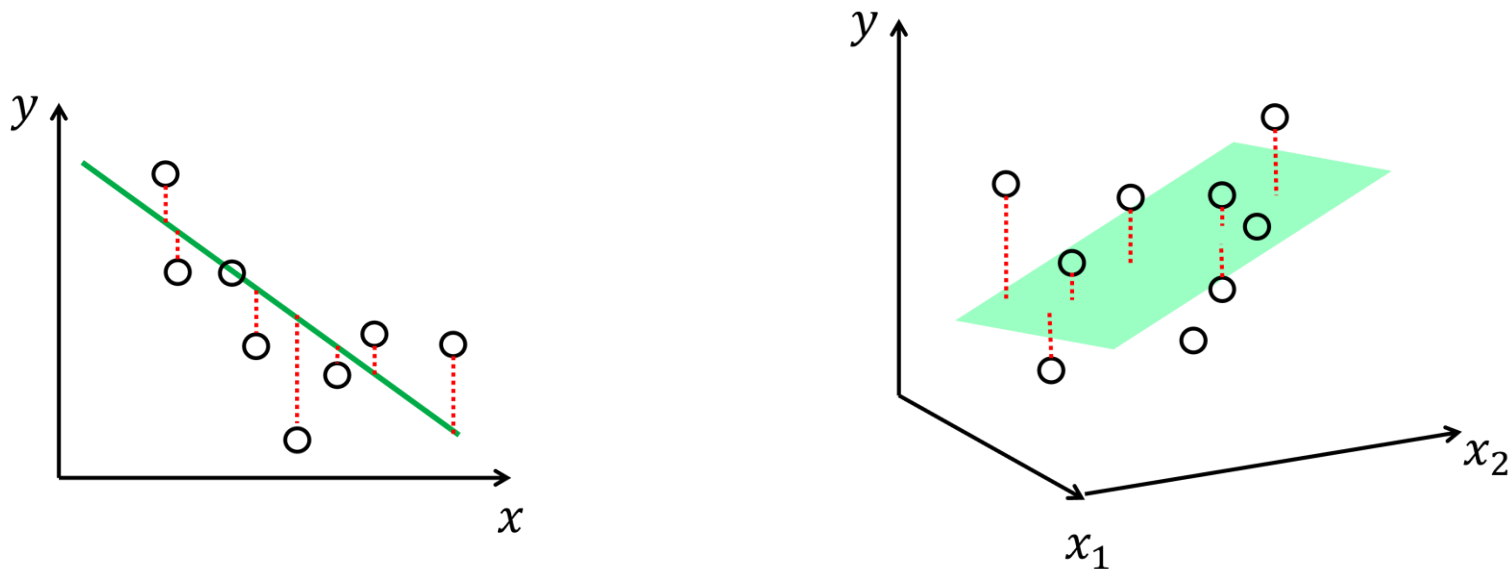
Pocket Algorithm



Recap: Linear Models for Credit Analysis



Least Squares Linear Regression

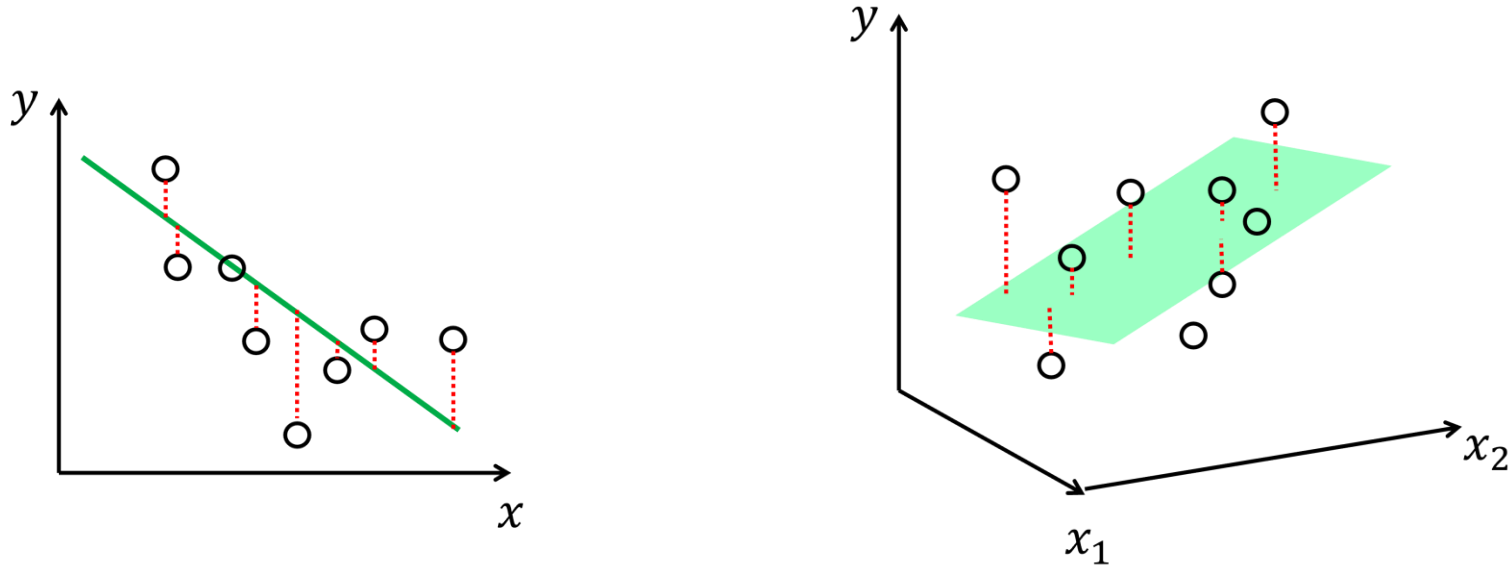


$$\text{error}(h(\mathbf{x}), f(\mathbf{x})) = (h(\mathbf{x}) - f(\mathbf{x}))^2$$

prediction

actual

Least Squares Linear Regression



$$y = f(\mathbf{x}) + \epsilon$$

$$\sim P(y|\mathbf{x})$$

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}) - f(\mathbf{x}))^2$$

$$E_{out}(h) = \mathbb{E}_{\mathbf{x}} \left[(h(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

Recap: Ordinary Least Squares: Minimizing E_{in}

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

Differentiable

Let $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ with dimensions $(d + 1) \times (d + 1)$

Let $\mathbf{b} = \mathbf{X}^T \mathbf{y}$ with dimensions $(d + 1) \times 1$

Let $c = \mathbf{y}^T \mathbf{y}$

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{w}^T \mathbf{A} \mathbf{w} - 2\mathbf{w}^T \mathbf{b} + c)$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{1}{N} [(\mathbf{A} + \mathbf{A}^T) \mathbf{w} - 2\mathbf{b}]$$

Useful gradient identities:

- $\nabla_{\mathbf{z}}(\mathbf{z}^T \mathbf{A} \mathbf{z}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{z}$
- $\nabla_{\mathbf{z}}(\mathbf{z}^T \mathbf{b}) = \mathbf{b}$

Recap: Ordinary Least Squares: Minimizing E_{in}

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

$$\text{Let } \mathbf{A} = \mathbf{X}^T \mathbf{X}, \mathbf{b} = \mathbf{X}^T \mathbf{y}, c = \mathbf{y}^T \mathbf{y}$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{1}{N} [(\mathbf{A} + \mathbf{A}^T) \mathbf{w} - 2 \mathbf{b}]$$

$$(\mathbf{X}^T \mathbf{X})^T = (\mathbf{X})^T (\mathbf{X}^T)^T = \mathbf{X}^T \mathbf{X}$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} [\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}]$$

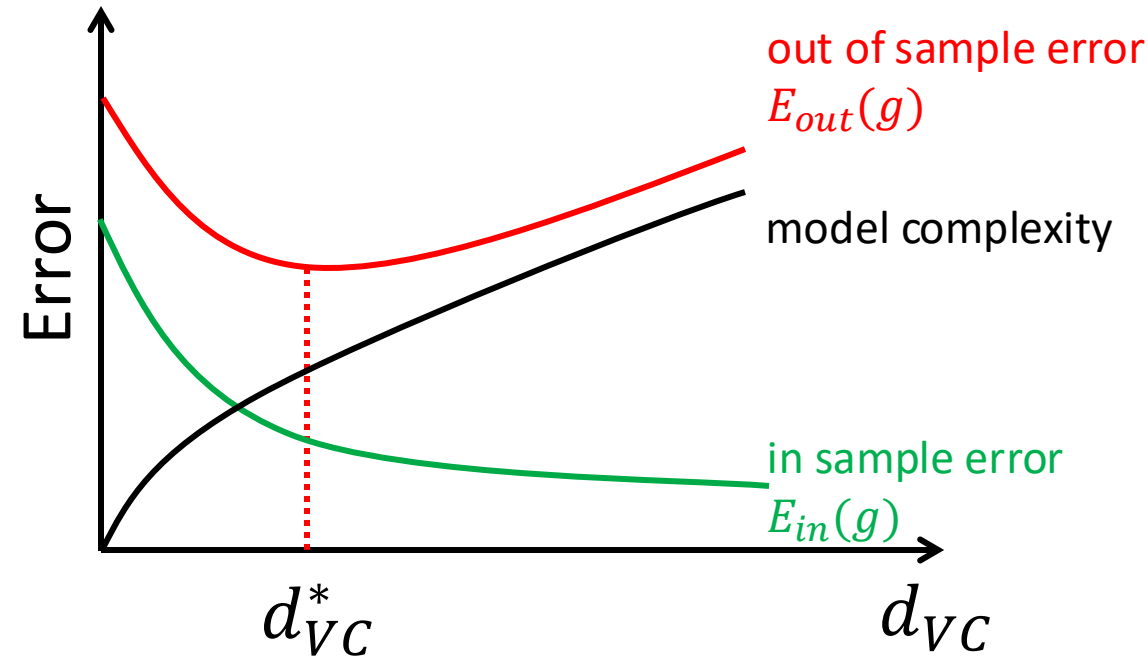
Solve for \mathbf{w} :

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}_{lin} = \mathbf{X}^\dagger \mathbf{y}$$

where $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

Flashback: VC Analysis, Approximation vs. Generalization



$d_{VC} \uparrow \Rightarrow E_{in} \approx 0$: Higher chance of approximating target function on data

$d_{VC} \downarrow \Rightarrow E_{in} \approx E_{out}$: Higher chance of generalizing to out of data

Depends only on \mathcal{H} ; Independent of $f, P(\mathbf{x}), \mathcal{A}$ (the learning algorithm) ¹²

Bias – Variance Analysis

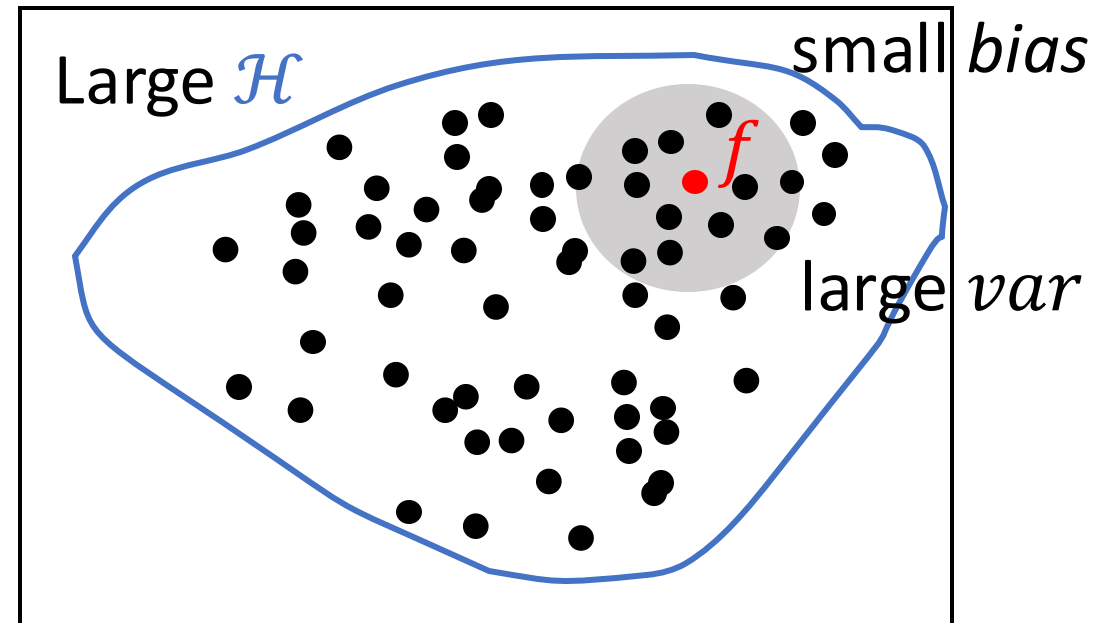
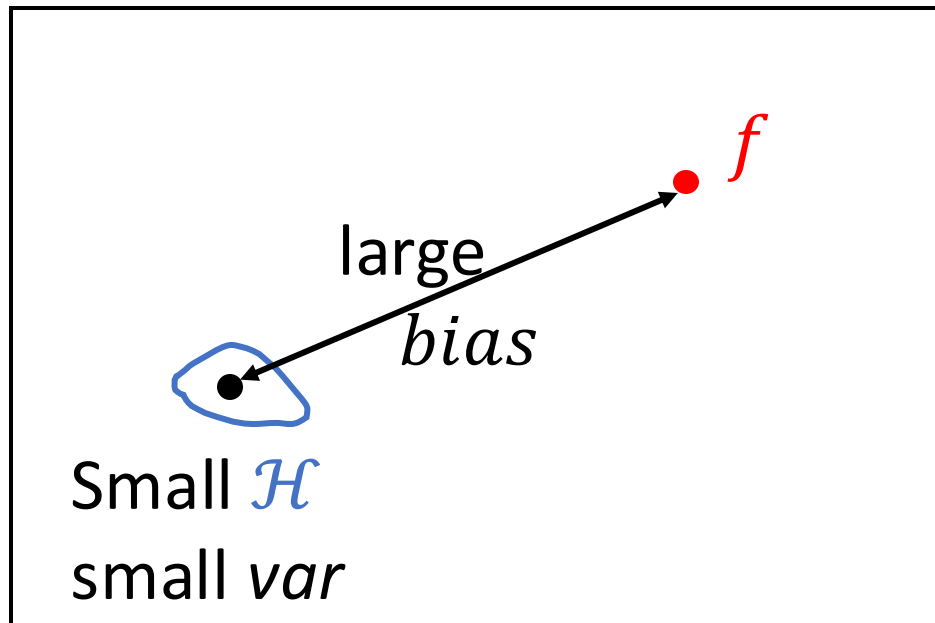
- An alternate view of the approximation-generalization tradeoff for squared error measures (e.g. regression)

1. How well ***can*** a hypothesis in \mathcal{H} approximate f ?
2. How close can we get to this using a finite dataset?

Bias-Variance Tradeoff

bias: How well can \mathcal{H} fit target f ?

var: How often can we find a good approximation?

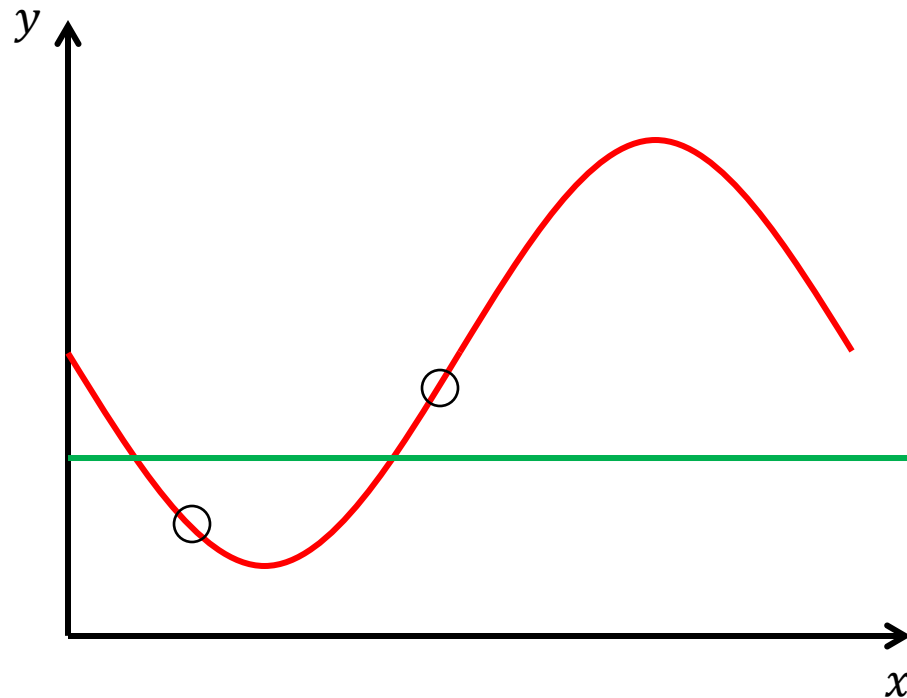


A Simple Learning Problem $f(x) = \sin \pi x$

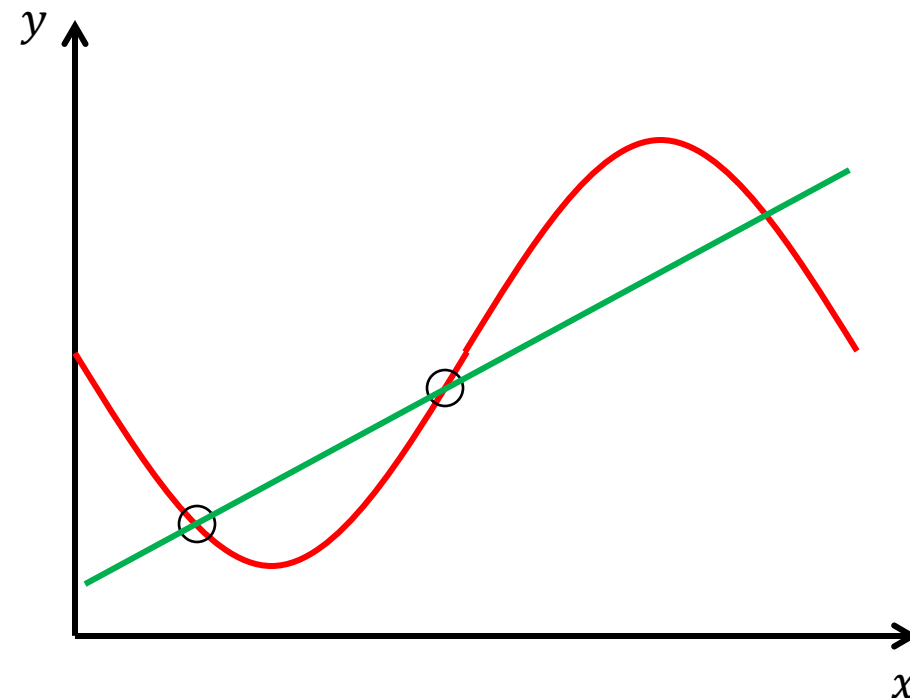
2 Data points

2 Hypothesis sets:

Flat lines $\mathcal{H}_0: h(x) = b$



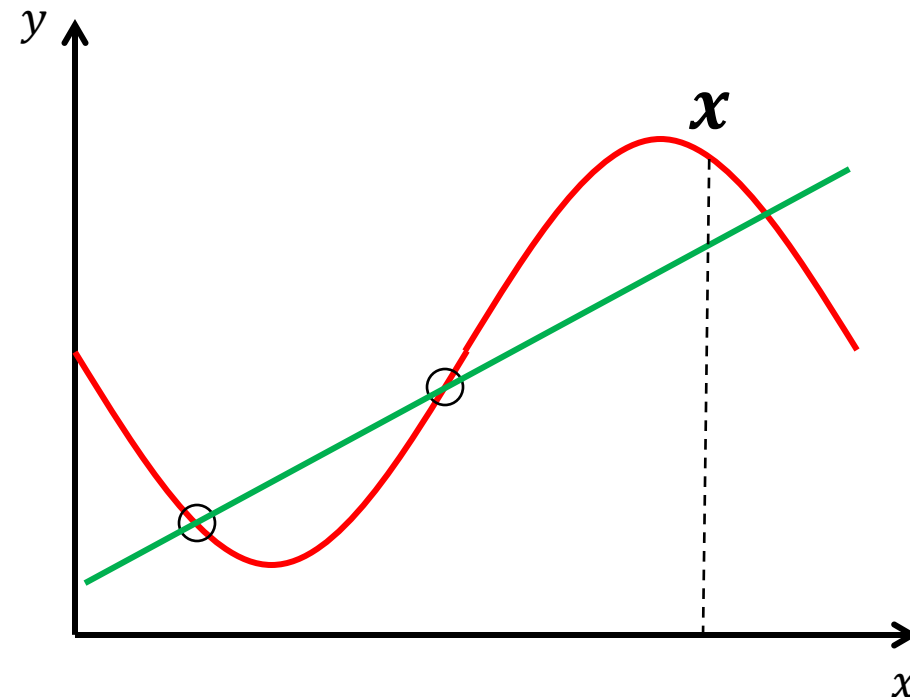
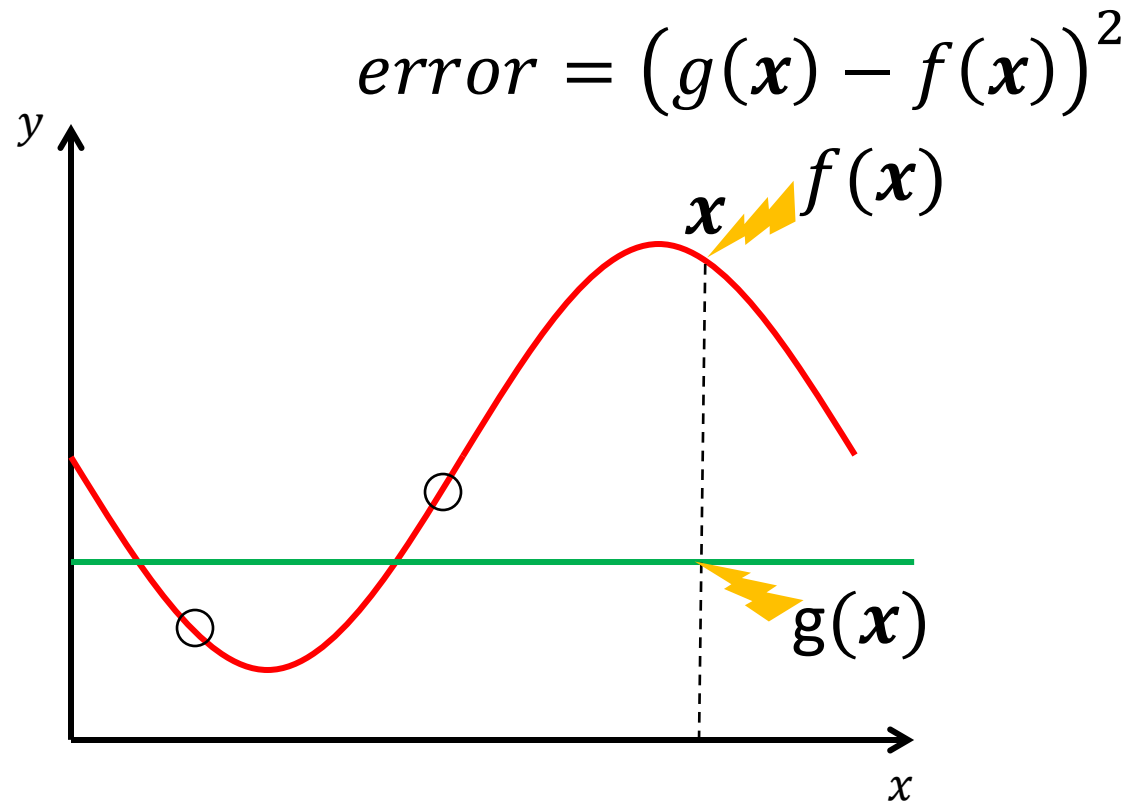
All lines $\mathcal{H}_1: h(x) = ax + b$



A Simple Learning Problem $f(x) = \sin \pi x$

Flat lines $\mathcal{H}_0: h(x) = b$

All lines $\mathcal{H}_1: h(x) = ax + b$

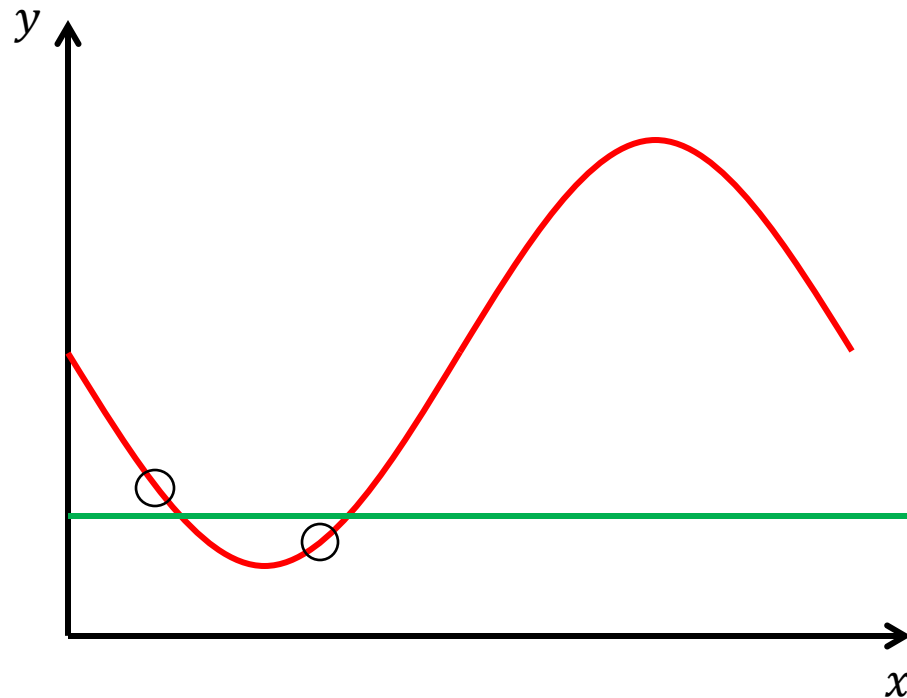


A Simple Learning Problem $f(x) = \sin \pi x$

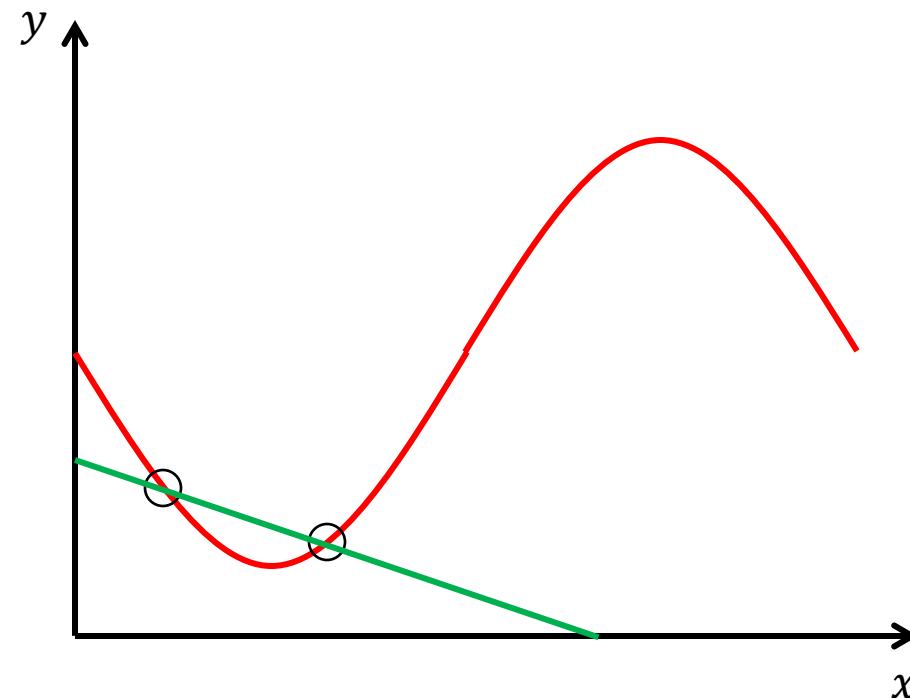
2 Data points

2 Hypothesis sets:

Flat lines $\mathcal{H}_0: h(x) = b$



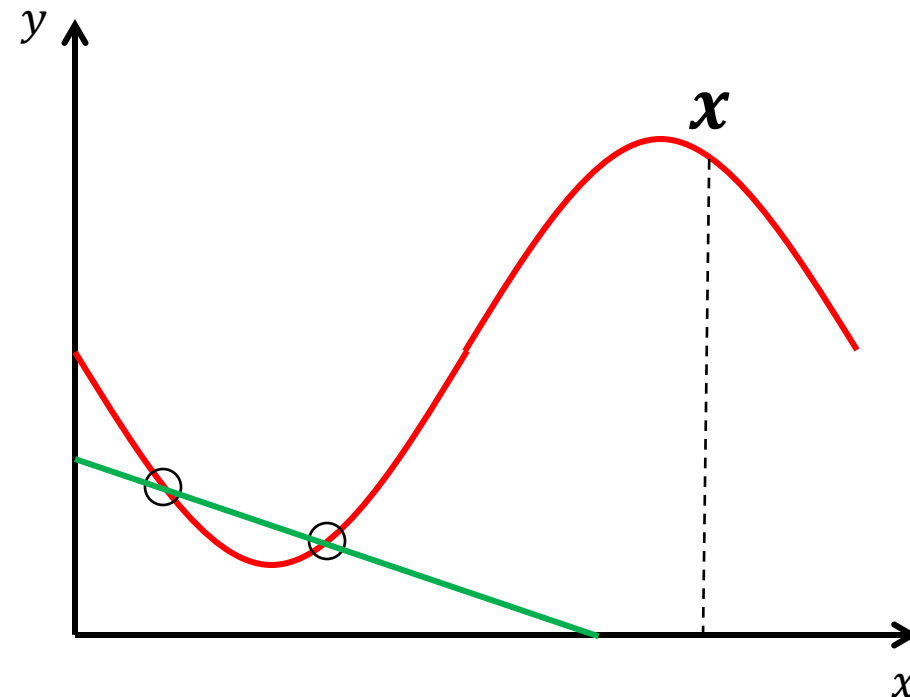
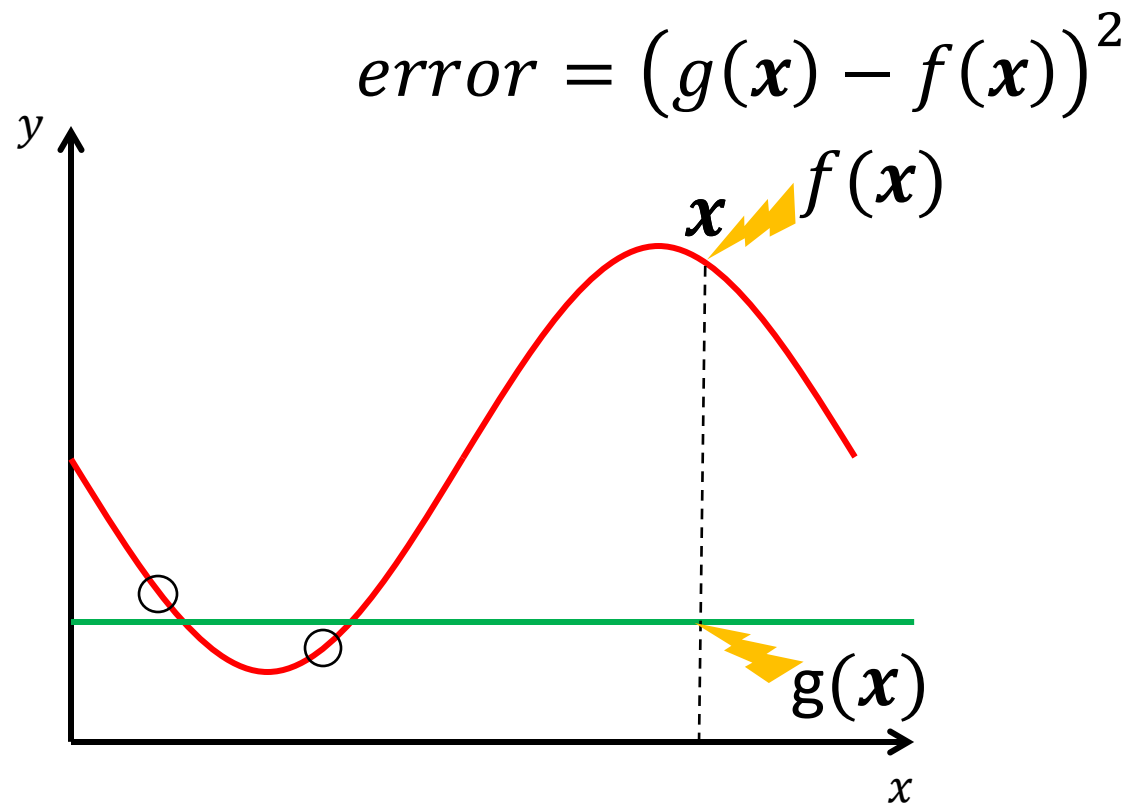
All lines $\mathcal{H}_1: h(x) = ax + b$



A Simple Learning Problem $f(x) = \sin \pi x$

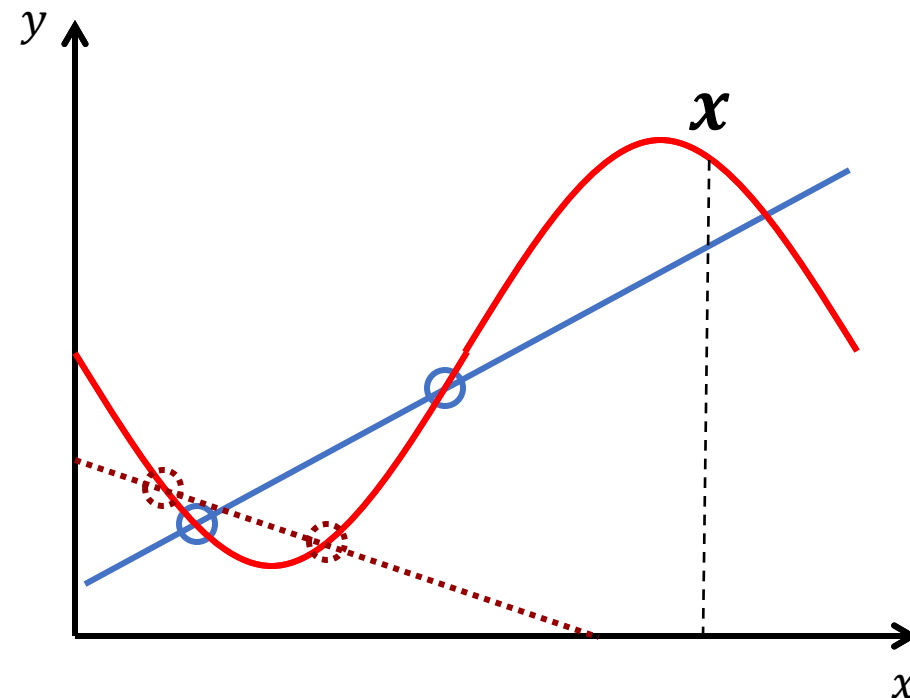
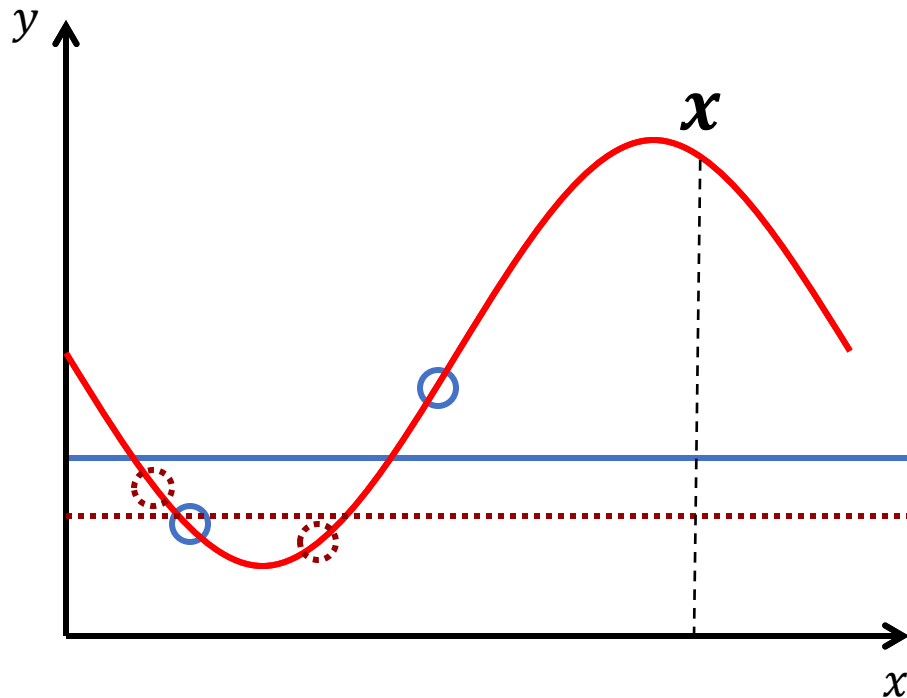
Flat lines $\mathcal{H}_0: h(x) = b$

All lines $\mathcal{H}_1: h(x) = ax + b$



Different Dataset \mathcal{D} , Different Output $g^{\mathcal{D}}$

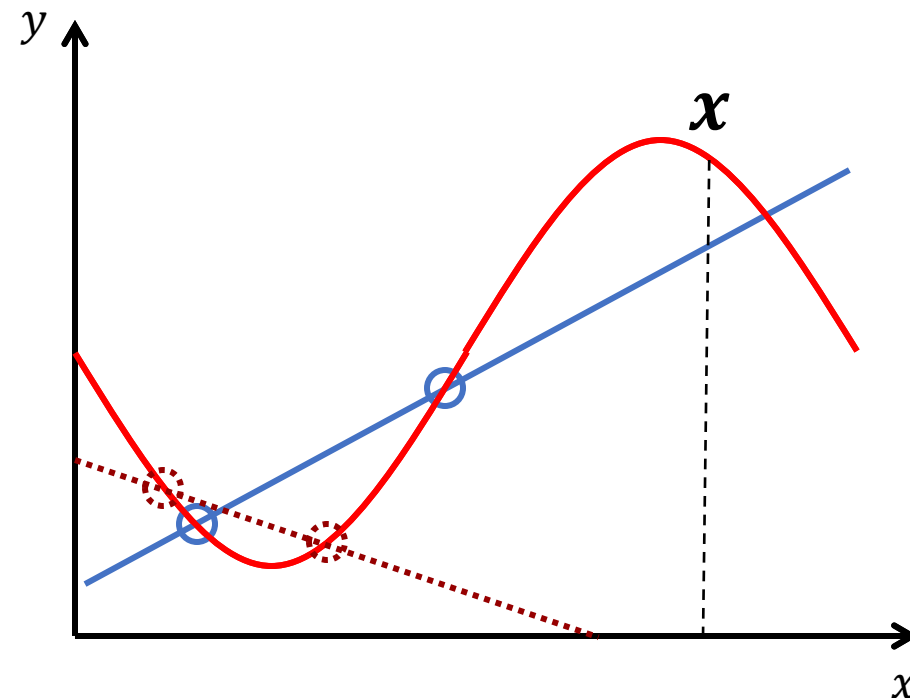
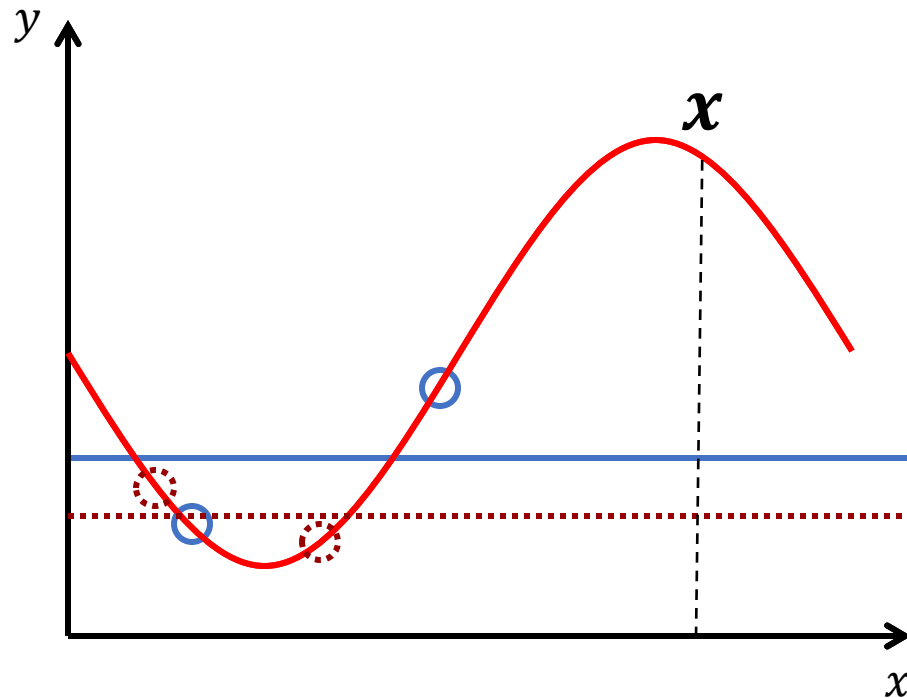
- For a fixed test point \mathbf{x} ,
 $g^{\mathcal{D}}(\mathbf{x})$ is a random value that depends on \mathcal{D}



We can't pick our data set!

We must analyze the entire process:

- Sample a data set
- Fit it (pick g from \mathcal{H} using \mathcal{A})
- Measure E_{in}



Expected Behavior w.r.t. Data

Dataset		Fit		Test
\mathcal{D}_1	\rightarrow	$g^{\mathcal{D}_1}$	\rightarrow	$g^{\mathcal{D}_1}(\mathbf{x})$
\mathcal{D}_2	\rightarrow	$g^{\mathcal{D}_2}$	\rightarrow	$g^{\mathcal{D}_2}(\mathbf{x})$
...				
\mathcal{D}_K	\rightarrow	$g^{\mathcal{D}_K}$	\rightarrow	$g^{\mathcal{D}_K}(\mathbf{x})$

Average (expected) prediction:

$$\mathbb{E}_{\mathcal{D}}[g^{\mathcal{D}}(\mathbf{x})] = \bar{g}(\mathbf{x}) \approx \frac{1}{K} \left(g^{\mathcal{D}_1}(\mathbf{x}) + g^{\mathcal{D}_2}(\mathbf{x}) + \dots + g^{\mathcal{D}_K}(\mathbf{x}) \right)$$

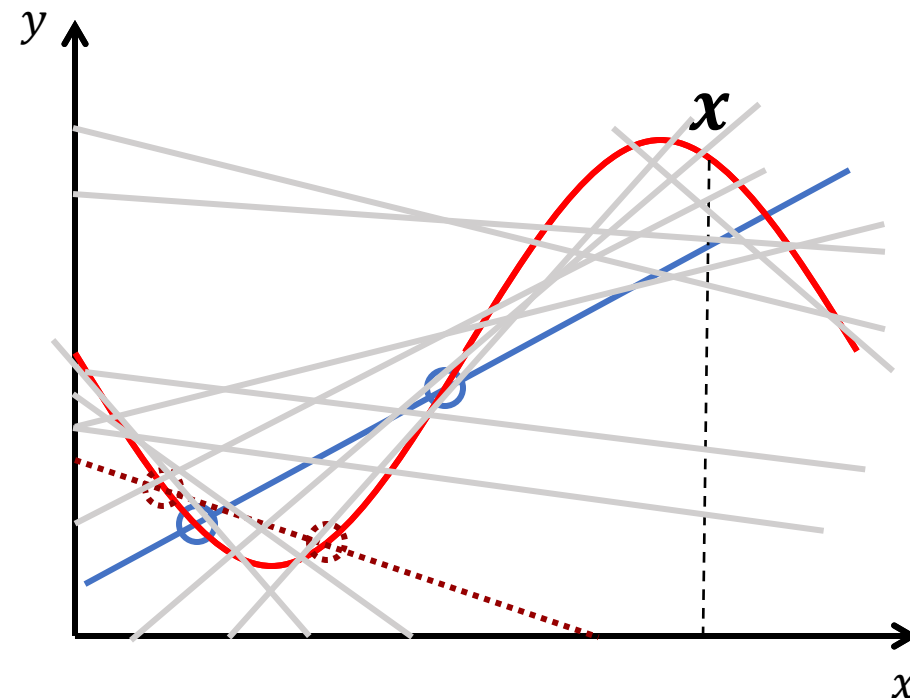
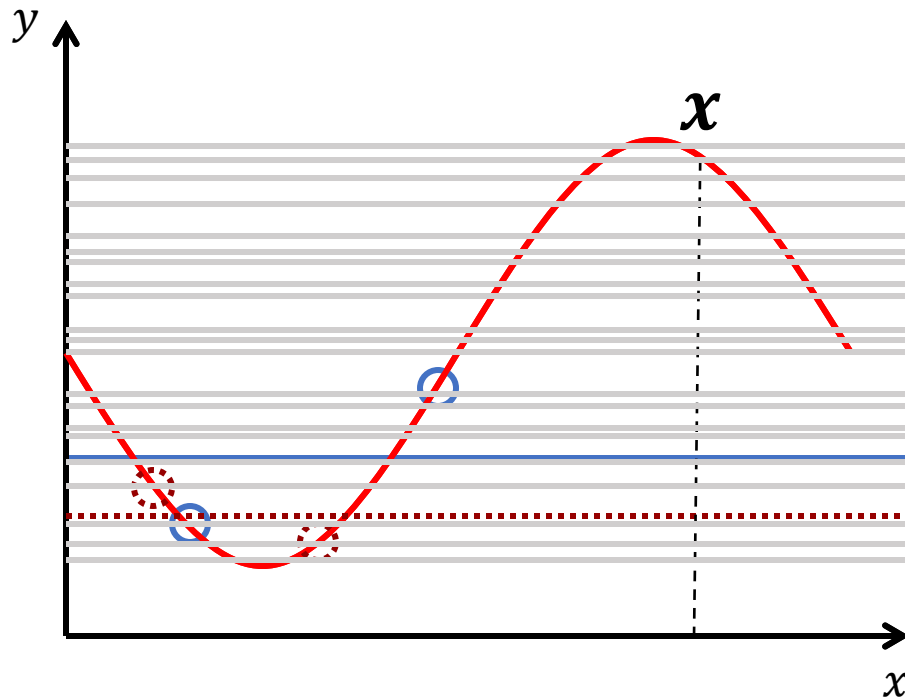
Variance of prediction: $var(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right]$

Different Dataset \mathcal{D} , Different Output $g^{\mathcal{D}}$

$g^{\mathcal{D}}(\mathbf{x})$ is a random value depending on \mathcal{D} (randomly generated data)

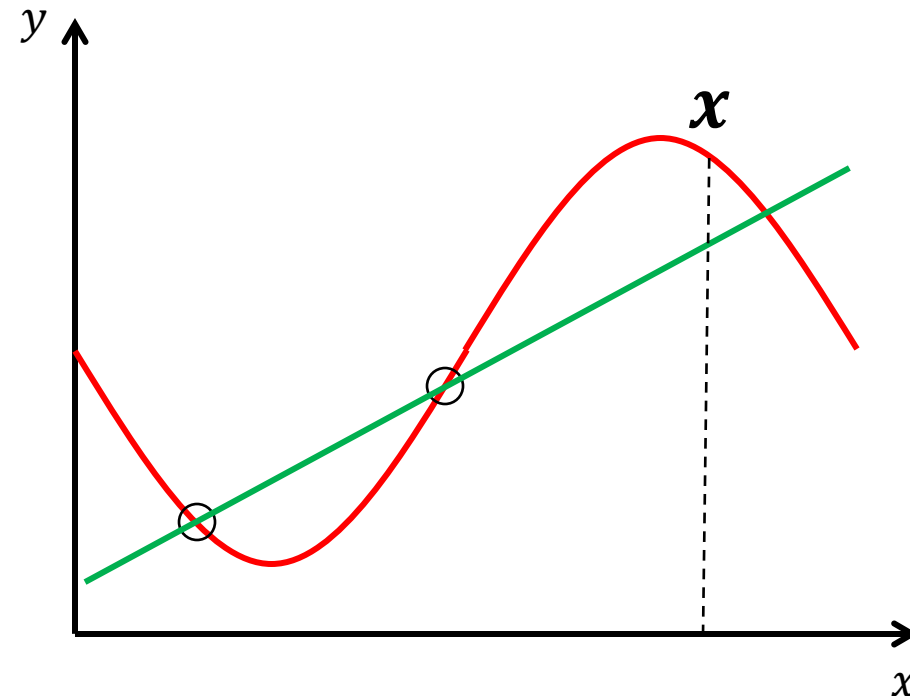
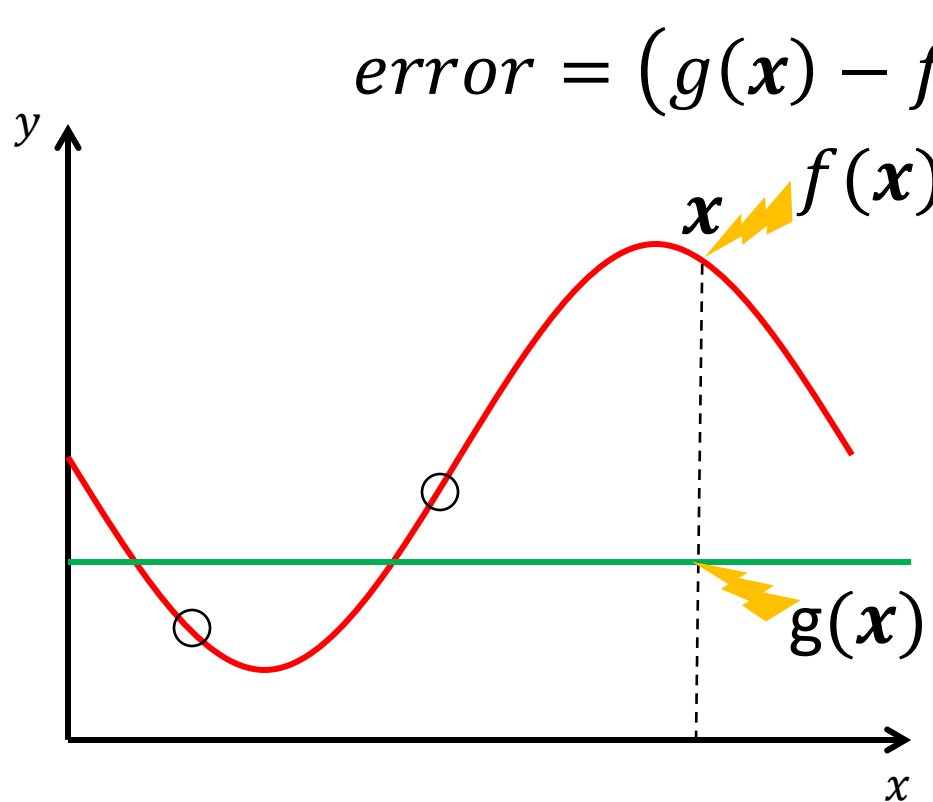
$\mathbb{E}_{\mathcal{D}}[g^{\mathcal{D}}(\mathbf{x})] = \bar{g}(\mathbf{x}) \approx \frac{1}{K} \left(g^{\mathcal{D}_1}(\mathbf{x}) + g^{\mathcal{D}_2}(\mathbf{x}) + \dots + g^{\mathcal{D}_K}(\mathbf{x}) \right)$, average prediction

$var(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right]$, variance of prediction



Out of Sample Error: E_{out} on test point \mathbf{x} for data \mathcal{D}

- $E_{out}^{\mathcal{D}}(\mathbf{x}) = \left(g^{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x})\right)^2$, squared error depending on random \mathcal{D}
- E_{out} before seeing the data: $E_{out}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[E_{out}^{\mathcal{D}}(\mathbf{x})]$



Bias-Variance Decomposition: Expected Error on Test Point

$$\begin{aligned} E_{out}(\mathbf{x}) &= \mathbb{E}_{\mathcal{D}}[E_{out}^{\mathcal{D}}(\mathbf{x})] &= \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} [g^{\mathcal{D}}(\mathbf{x})^2 - 2g^{\mathcal{D}}(\mathbf{x})f(\mathbf{x}) + f(\mathbf{x})^2] \\ [\mathbb{E}_{\mathcal{D}}[g^{\mathcal{D}}(\mathbf{x})] = \bar{g}(\mathbf{x})] &= \mathbb{E}_{\mathcal{D}} [g^{\mathcal{D}}(\mathbf{x})^2] - 2\bar{g}(\mathbf{x})f(\mathbf{x}) + f(\mathbf{x})^2 \\ &= \mathbb{E}_{\mathcal{D}} [g^{\mathcal{D}}(\mathbf{x})^2] + \bar{g}(\mathbf{x})^2 - \bar{g}(\mathbf{x})^2 - 2\bar{g}(\mathbf{x})f(\mathbf{x}) + f(\mathbf{x})^2 \\ &= \mathbb{E}_{\mathcal{D}} [g^{\mathcal{D}}(\mathbf{x})^2] - \bar{g}(\mathbf{x})^2 + (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2 \\ \mathbb{E}_{\mathcal{D}} [-2g^{\mathcal{D}}(\mathbf{x})\bar{g}(\mathbf{x}) + \bar{g}(\mathbf{x})^2] = -\bar{g}(\mathbf{x})^2 &= \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right] + (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2 \\ &= \text{var}(\mathbf{x}) + \text{bias}(\mathbf{x}) \end{aligned}$$

The Bias-Variance Decomposition

- Fact: For any random variable X , (and also for any $X = \hat{\theta} - \theta$),
 $\mathbb{E}[X^2] = (\mathbb{E}[X])^2 + \text{var}(X)$

 estimator  unknown, fixed
parameter

Theorem: $\mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] = \left(\bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 + \text{var}(\mathbf{x})$
 $= \text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})$
 $= \text{expected out of sample error at test point } \mathbf{x}$

Here, $\text{var}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right]$

E_{out} : Average over \mathbf{x}

$$\begin{aligned} E_{out} &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}} \left[\left(g^{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \right] = \mathbb{E}_{\mathbf{x}} \left[\left(\bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 + var(\mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{x}} [bias(\mathbf{x}) + var(\mathbf{x})] \\ &= \quad \quad \quad bias \quad + \quad var \end{aligned}$$



generalization

How close is average
learned hypothesis to
target function?



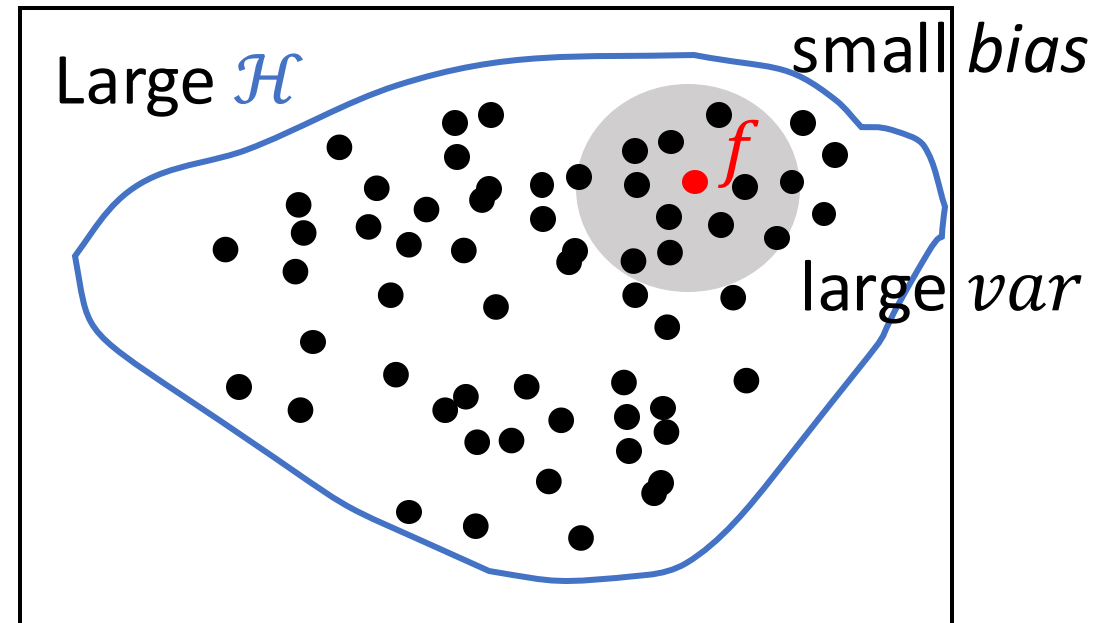
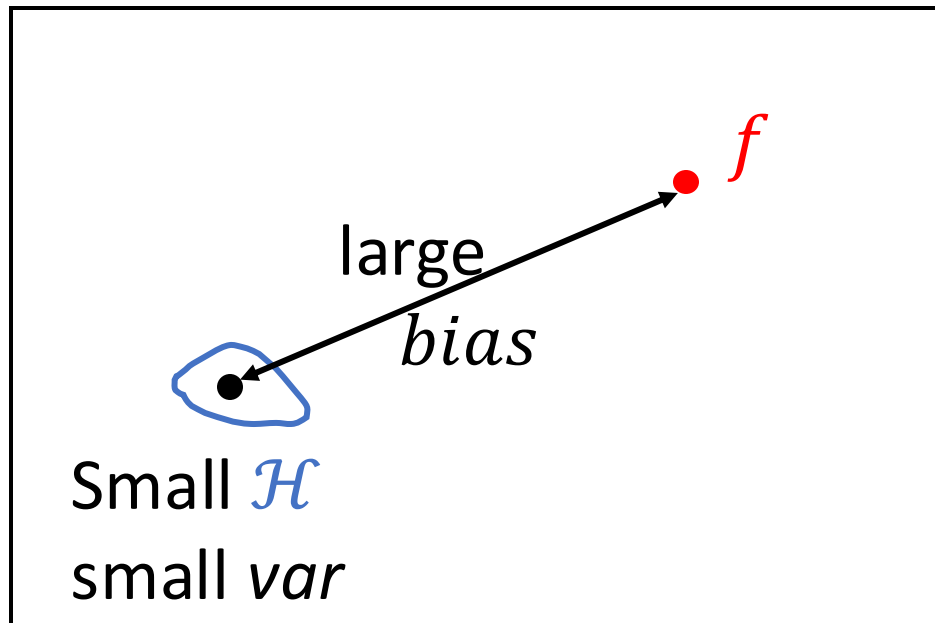
approximation

How often can we
find a good
hypothesis?

Bias-Variance Tradeoff

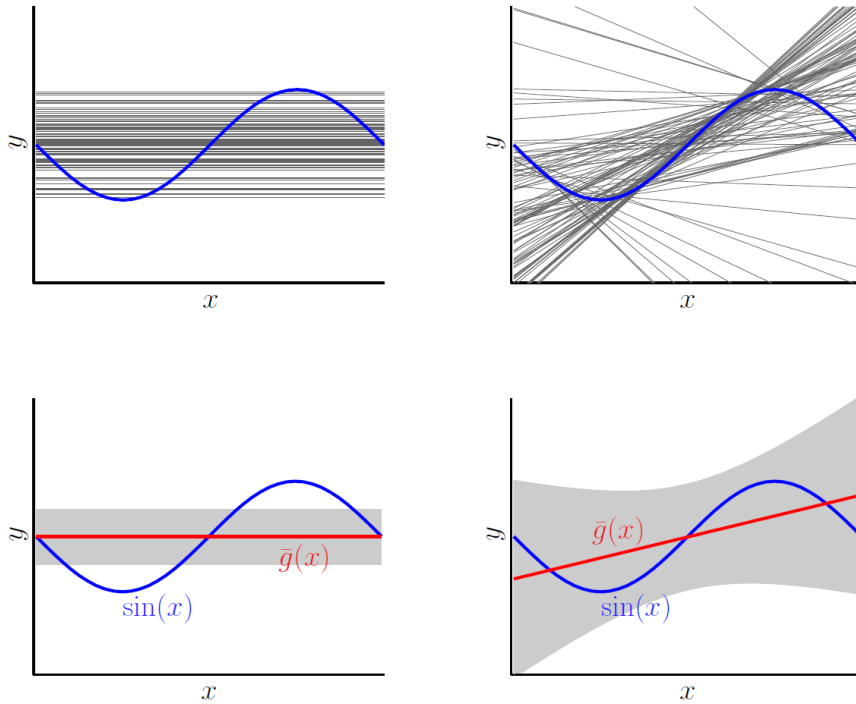
bias: How well can \mathcal{H} fit target f ?

var: How often can we find a good approximation?



Match Learning Power to Data, Not the Target f

2 Data points



\mathcal{H}_0

$$bias = 0.5$$

$$var = 0.25$$

$$\checkmark E_{out} = 0.75$$

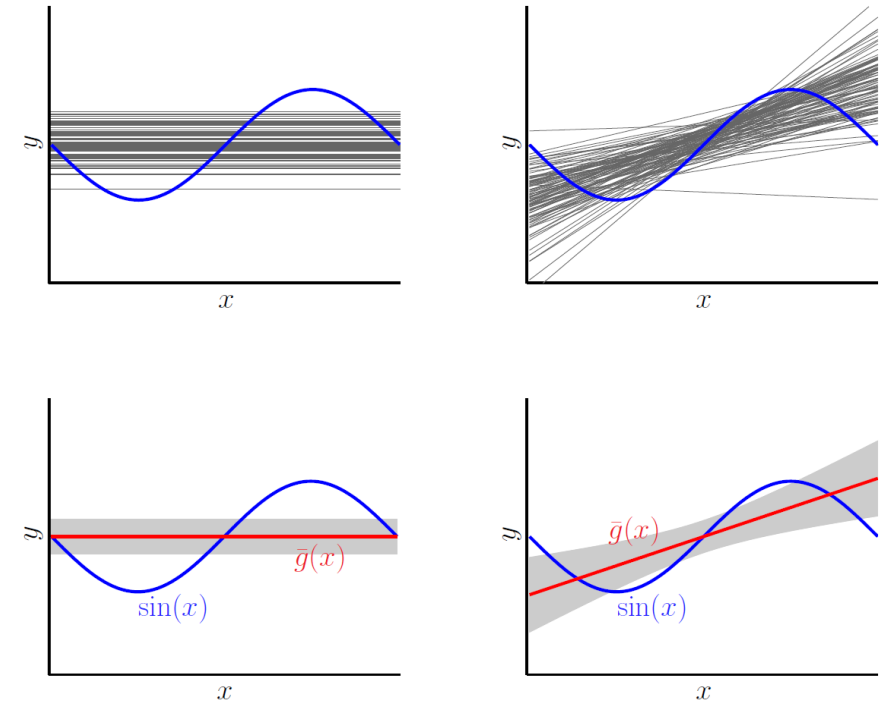
\mathcal{H}_1

$$bias = 0.21$$

$$var = 1.69$$

$$E_{out} = 1.90$$

5 Data points



\mathcal{H}_0

$$bias = 0.5$$

$$var = 0.1$$

$$E_{out} = 0.6$$

\mathcal{H}_1

$$bias = 0.21$$

$$var = 0.21$$

$$\checkmark E_{out} = 0.42$$

Bias-Variance Analysis: A Useful Conceptual Tool

- Depends on $f, P(\mathbf{x})$ – both **unknown**
- Depends on \mathcal{A}

The objective of \mathcal{A} is to minimize squared error

But for Bias-Variance analysis, we will use the squared error of g selected by \mathcal{A}

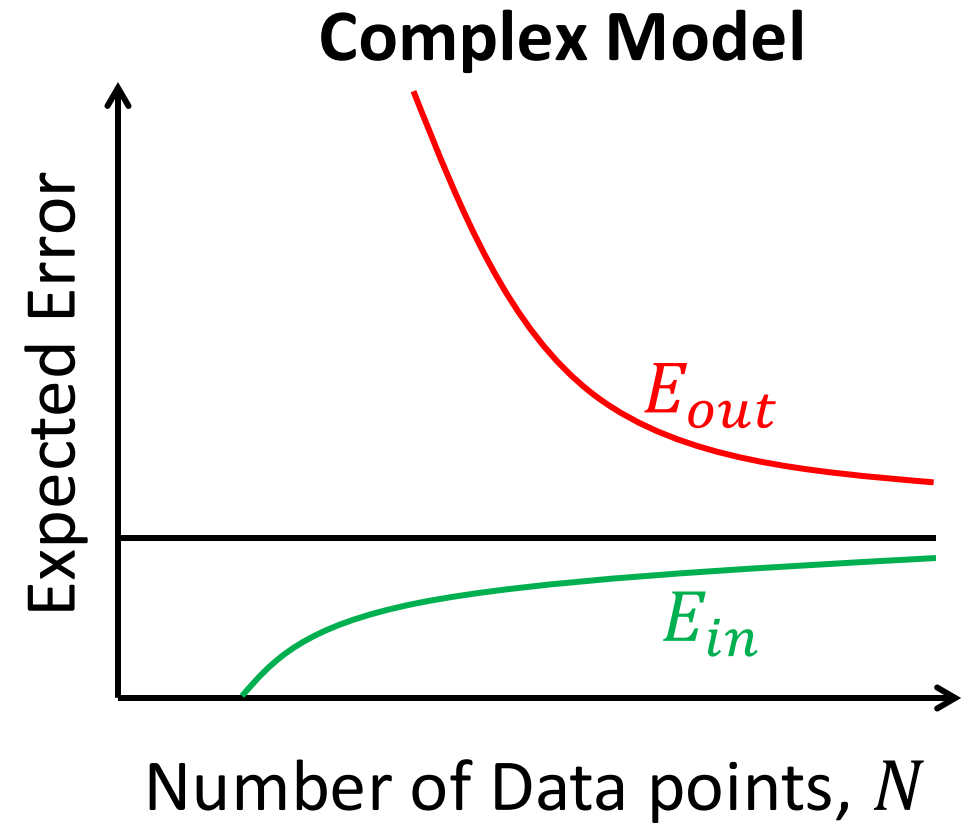
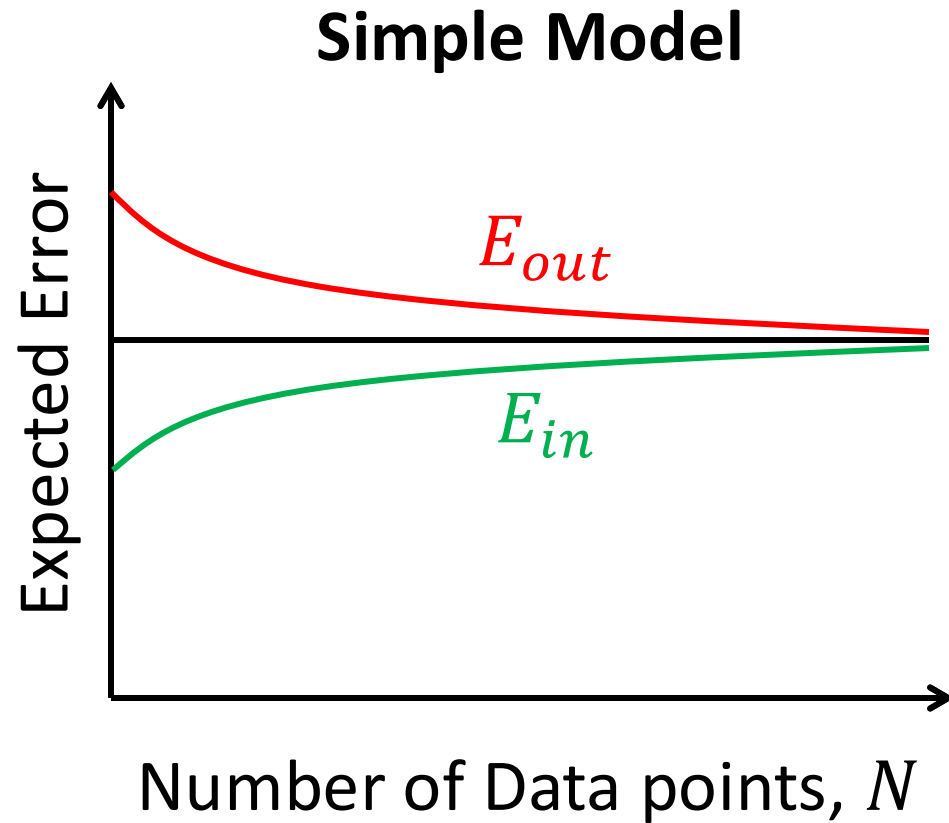
- Developing a model:
 - Lower variance without increasing bias
 - Lower bias without increasing variance
- Techniques discussed later in this course

Summary: The Learning Process

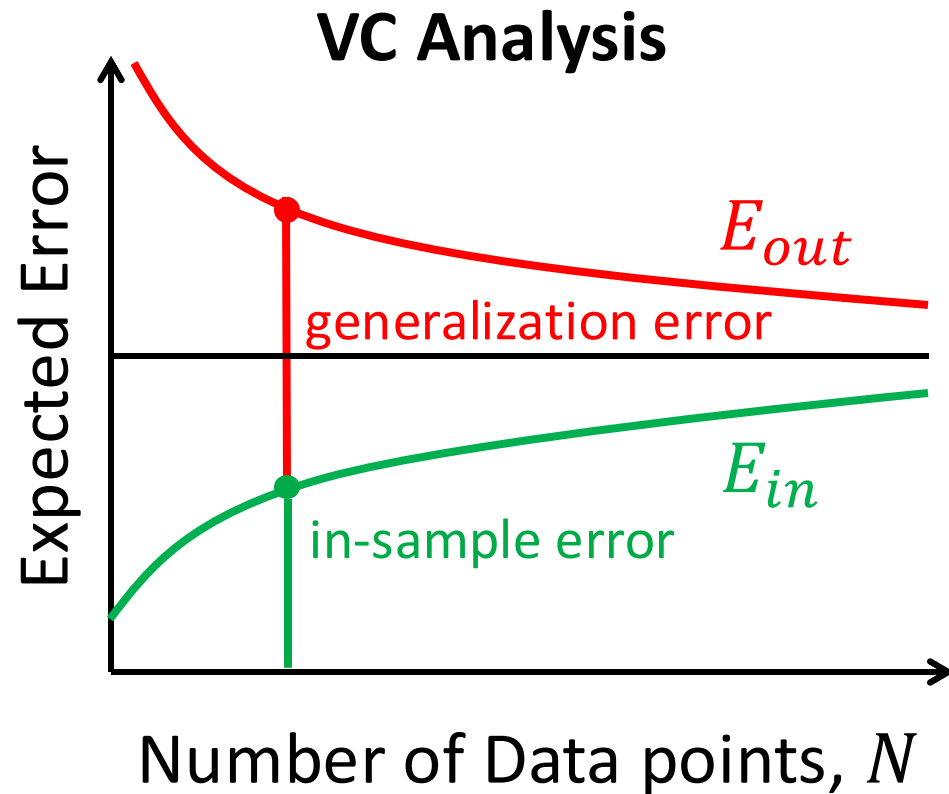
The Learning Process: For unknown $f(\mathbf{x})$ [or $P(y|\mathbf{x})$] and $P(\mathbf{x})$

- Fix \mathcal{H}
- Draw N data points \mathcal{D} from \mathcal{X} i.i.d. at random according to $P(\mathbf{x})$
- Pick $g^{\mathcal{D}}$ from \mathcal{H}
 - Which has $E_{in}(g^{\mathcal{D}})$ [measured on \mathcal{D}] and $E_{out}(g^{\mathcal{D}})$
- Expected error of learning process: Expectation over all \mathcal{D}
 - $\mathbb{E}_{\mathcal{D}}[E_{in}(g^{\mathcal{D}})]$
 - $\mathbb{E}_{\mathcal{D}}[E_{out}(g^{\mathcal{D}})]$

Summary: The Learning Curve

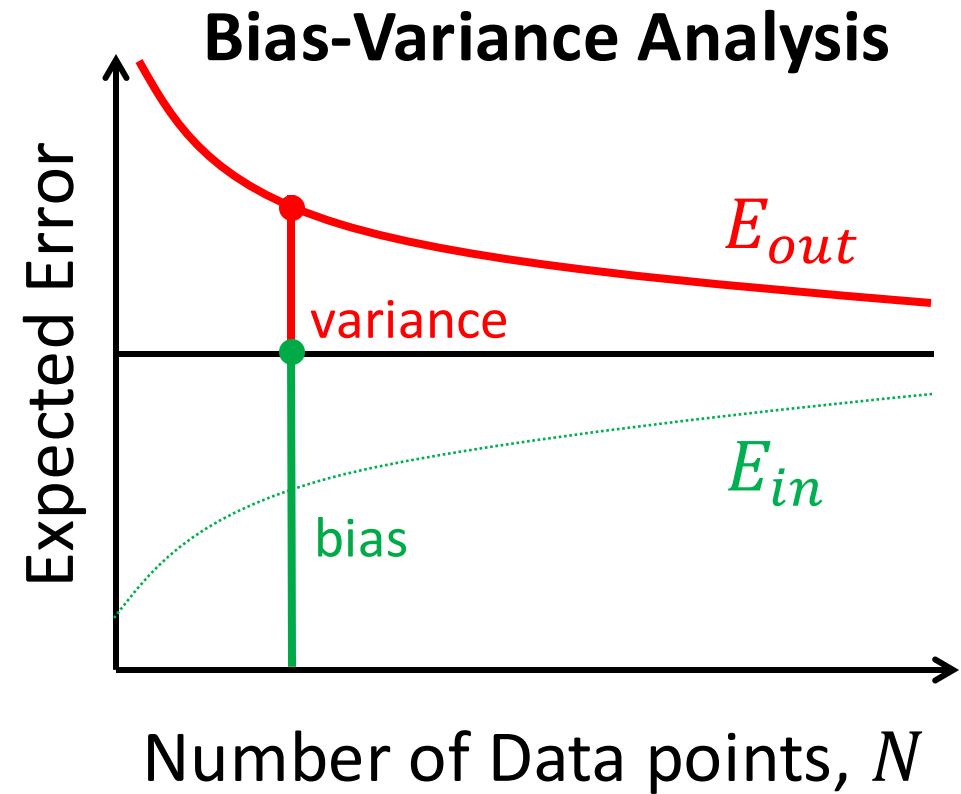


Decomposing the Learning Curve



Pick \mathcal{H} to:

1. Generalize well, i.e., ensure $E_{out} \approx E_{in}$
2. Fit \mathcal{D} , i.e., get small E_{in}



Pick \mathcal{H}, \mathcal{A} to:

1. Approximate f
2. Not vary too wildly with \mathcal{D}