# Lecture 10

# CS436/536: Introduction to Machine Learning
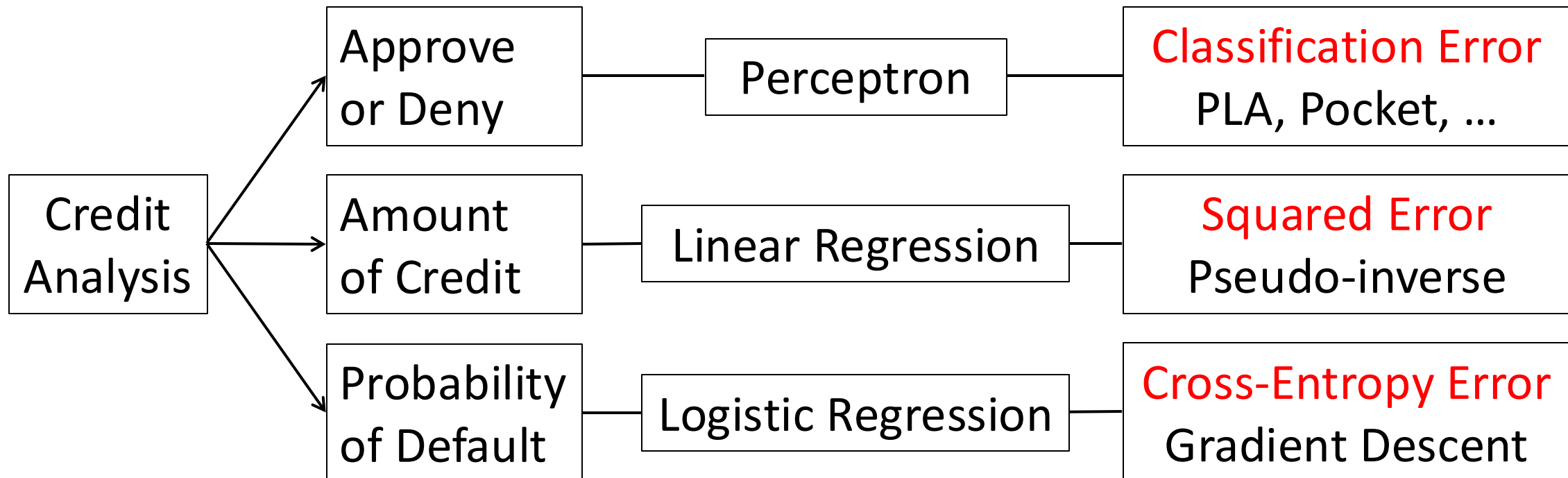
**Zhaohan Xi**
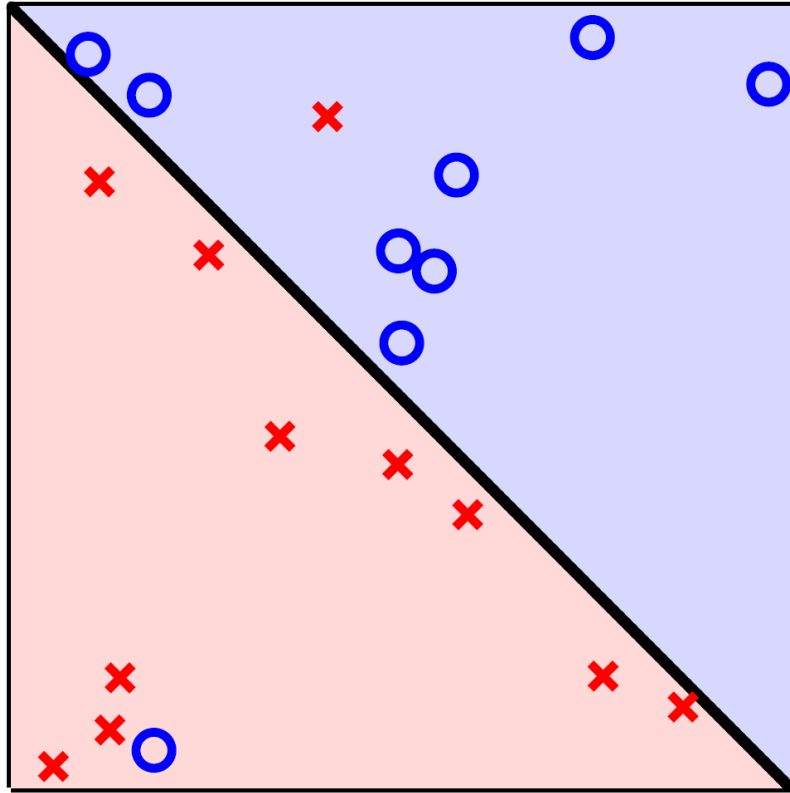**Binghamton University**

**zxi1@binghamton.edu**

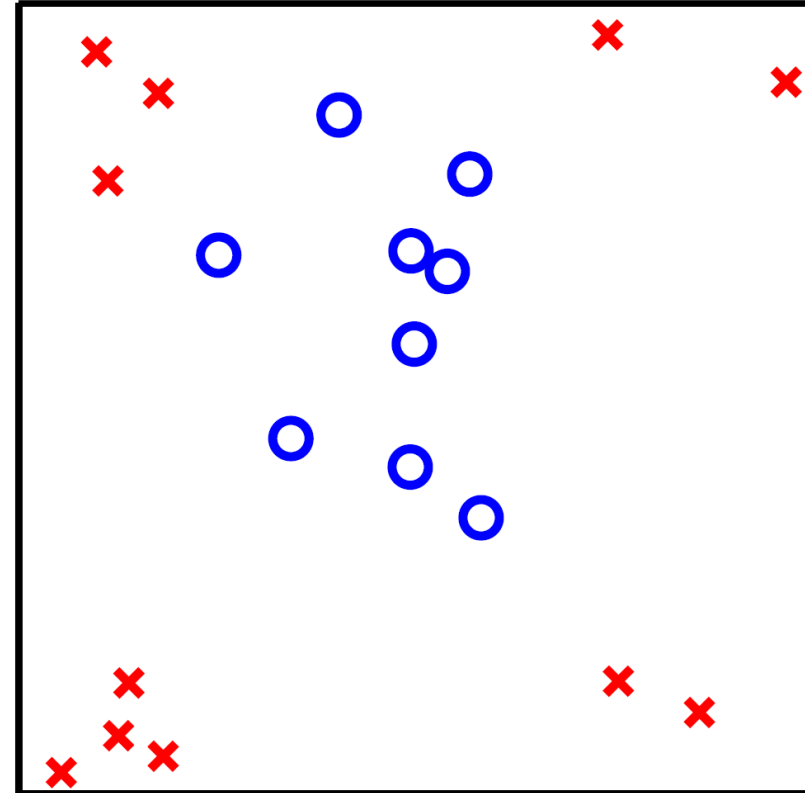# Recap: Linear Model for Three Learning Problems

$$h(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x}, \qquad \text{where } \boldsymbol{w} \in \mathbb{R}^d, \boldsymbol{x} \in 1 \times \mathbb{R}^d$$

| Credit Analysis | Approve or Deny | Perceptron | Classification Error<br>PLA, Pocket, … |
| --- | --- | --- | --- |
| | Amount of Credit | Linear Regression | Squared Error<br>Pseudo-inverse |
| | Probability of Default | Logistic Regression | Cross-Entropy Error<br>Gradient Descent |

# But: The Linear Model has Limitations



Linear with Outliers

Essentially Non-linear

Need something non-linear

Please remember these two scenarios for later

# Credit Approval

- Using salary, debt, years in residence, etc., approve for credit or not

| age | 33 years |
|---|---|
| salary | 50,000 |
| debt | 27,500 |
| years employed | 1 |
| years at residence | 2 |
| … | … |

**Approve for credit?**

# Use the Appropriate Feature



Income / Years in Residence ($Y$)

Stable High Income

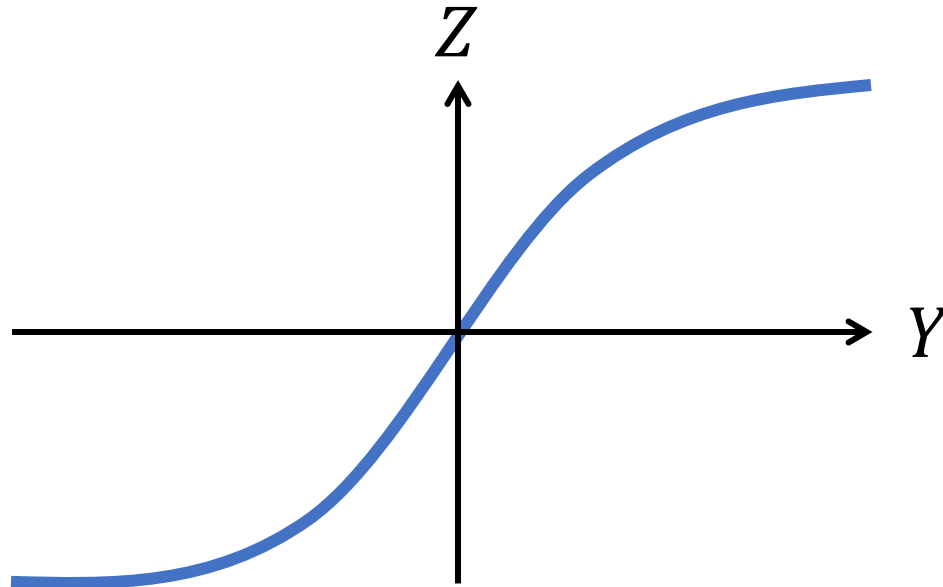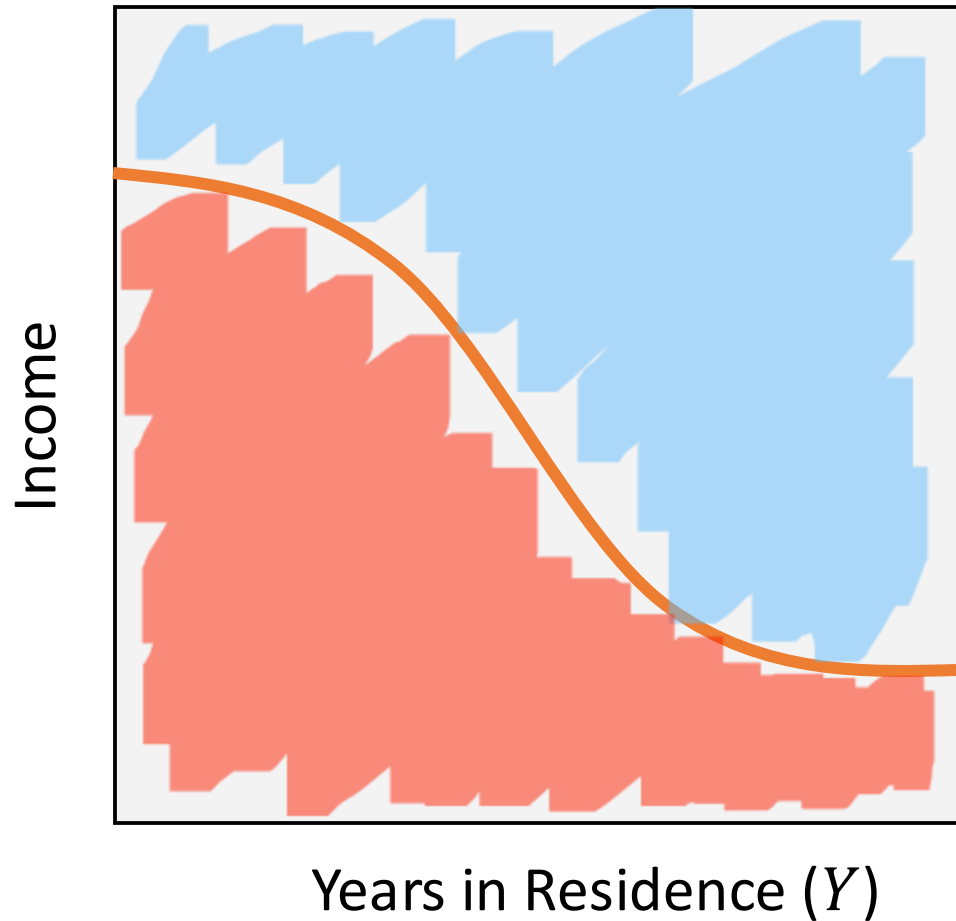Not Stable Low Income

$Y \gg 3$ years,    no additional effect beyond $Y = 3$
$Y \ll 0.3$ years,  no additional effect below $Y = 0.3$
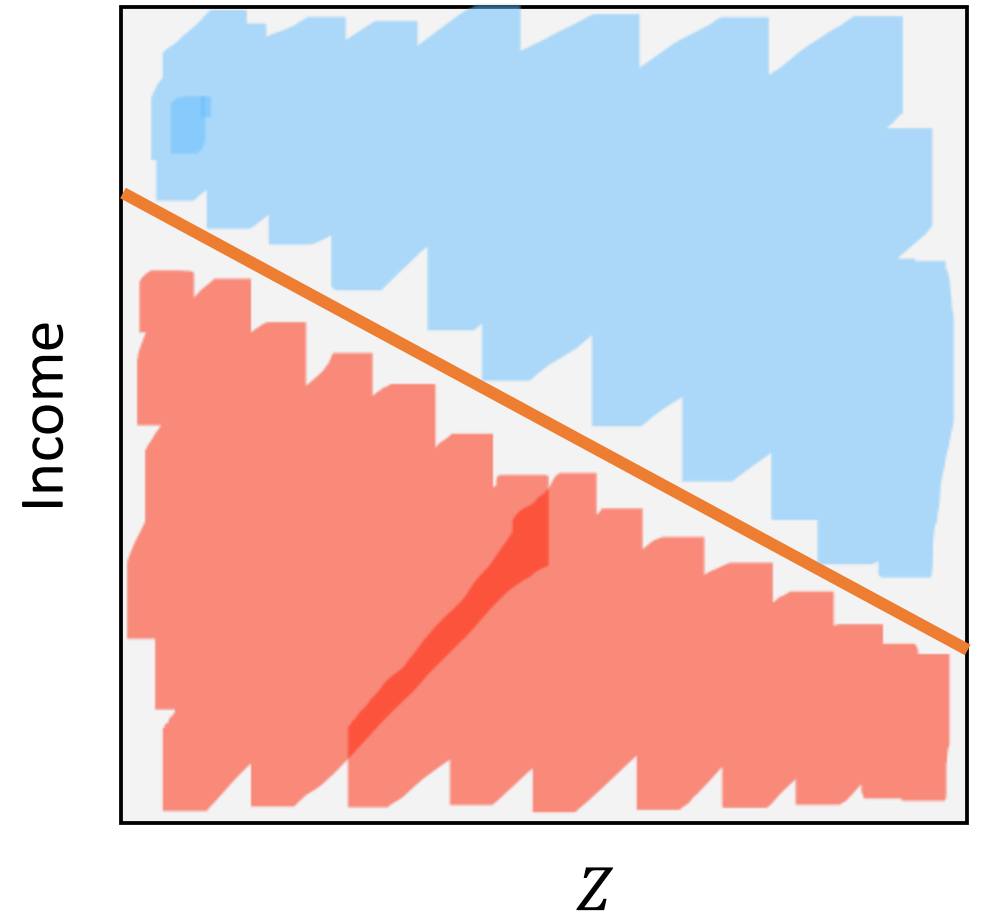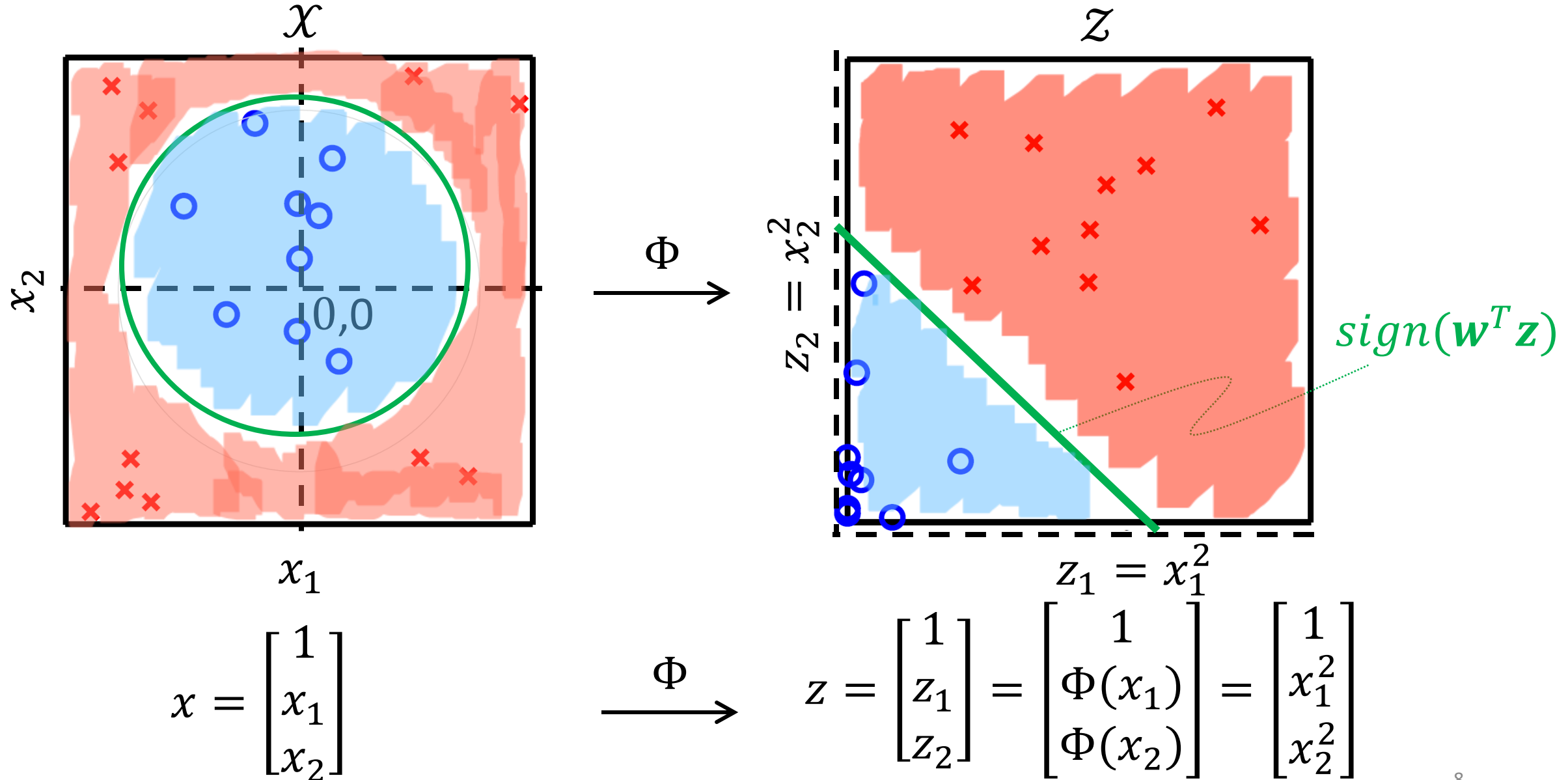
# Change the Feature Using a Transform

$$Y \overset{\Phi}{\to} Z$$

# Changing the Feature Using a Transform



$$Y \overset{\Phi}{\to} Z$$

Left plot: Income (vertical axis) vs Years in Residence ($Y$) (horizontal axis)

Right plot: Income (vertical axis) vs $Z$ (horizontal axis)

# Mechanics of Non-Linear Feature Transforms



$x$

$x_2$

$0,0$

$x_1$

$\Phi$

$z$

$z_2 = x_2^2$

$z_1 = x_1^2$

$sign(\boldsymbol{w}^T \boldsymbol{z})$

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$\Phi$

$$z = \begin{bmatrix} 1 \\ z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} 1 \\ \Phi(x_1) \\ \Phi(x_2) \end{bmatrix} = \begin{bmatrix} 1 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

8

# Feature Transforms In General

$\mathcal{X}$-space is $\mathbb{R}^d$

$\mathcal{Z}$-space is $\mathbb{R}^{\tilde{d}}$

$$x = \begin{bmatrix} 1 \\ x_1 \\ \dots \\ x_d \end{bmatrix}$$

$$z = \Phi(x) = \begin{bmatrix} 1 \\ \Phi_1(x) \\ \dots \\ \Phi_{\tilde{d}}(x) \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ \dots \\ z_{\tilde{d}} \end{bmatrix}$$

$(x_1, y_1), \dots, (x_N, y_N)$

$(z_1, y_1), \dots, (z_N, y_N)$

$g(x) = sign(\widetilde{w}^T \Phi(x))$

$$\widetilde{w} = \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_{\tilde{d}} \end{bmatrix}$$

Low $E_{in}$

# What about "Real Learning"? Generalization

- Want $g \approx f$ over all of $\mathcal{X}$,
i.e., $E_{out}(g) \approx 0$

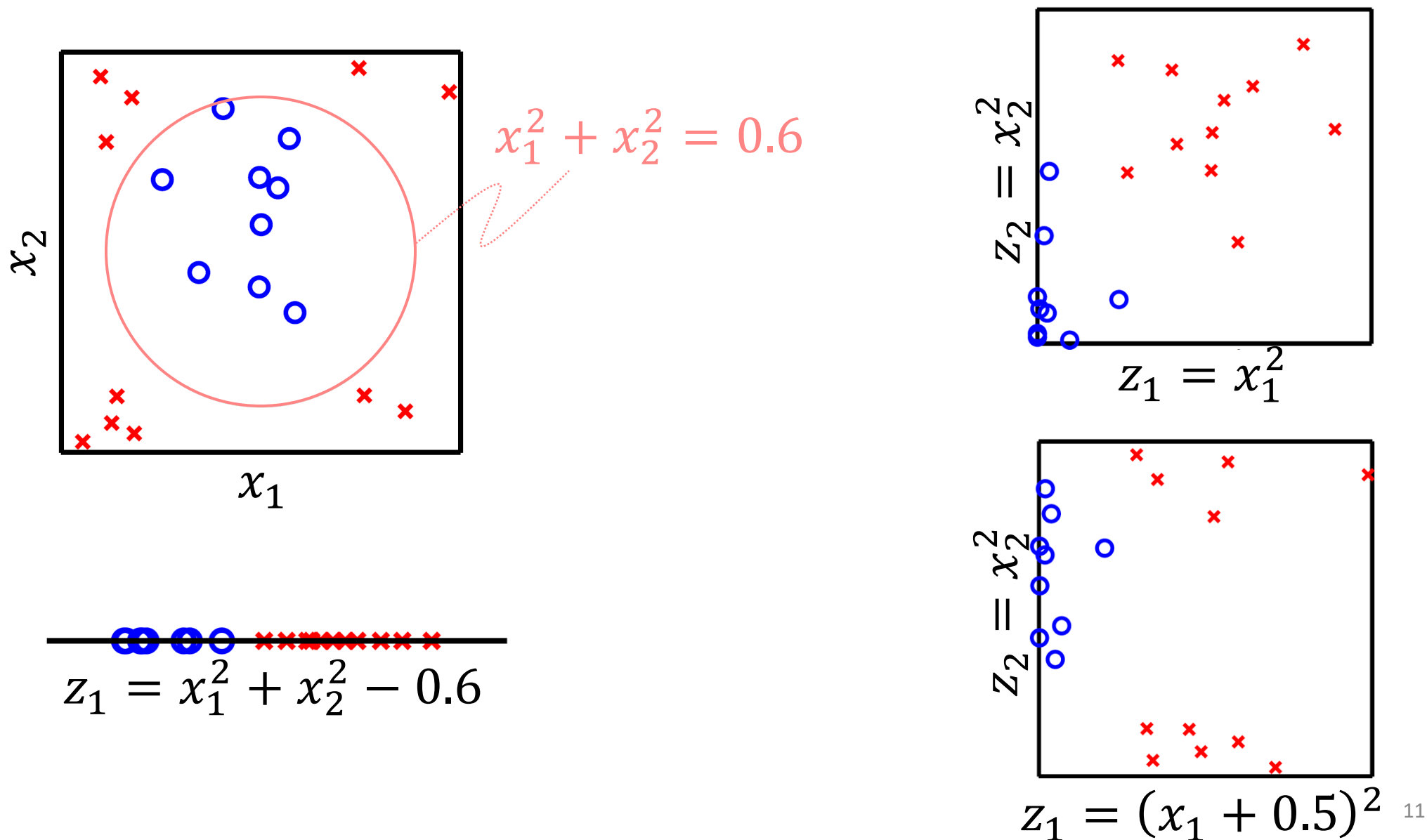[Approximation]    small $E_{in}(g)$

[Generalization]    $E_{in}(g) \approx E_{out}(g)$

$$
\begin{array}{ccc}
d_{VC} & \longrightarrow & \tilde{d}_{VC} \\
d + 1 & & \tilde{d} + 1
\end{array}
$$

Select transform $\Phi$ with small $\tilde{d}$

# Possibly Many 'Good' Nonlinear Transforms



$$x_1^2 + x_2^2 = 0.6$$

$$z_1 = x_1^2 + x_2^2 - 0.6$$

$$z_2 = x_2^2$$

$$z_1 = x_1^2$$

$$z_2 = x_2^2$$

$$z_1 = (x_1 + 0.5)^2$$

# Polynomial Feature Transform

E.g. Polynomial Transform for 2-D $\mathcal{X}$-space

Degree 1: $(1, x_1, x_2) \overset{\boldsymbol{\Phi_1}}{\to} (1, x_1, x_2)$ $\Rightarrow \tilde{d}_{VC} = 3$

Degree 2: $\overset{\boldsymbol{\Phi_2}}{\to} (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$ $\Rightarrow \tilde{d}_{VC} = 6$

Degree 3: $\overset{\boldsymbol{\Phi_3}}{\to} (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2, x_1^3, x_1^2 x_2, x_1 x_2^2, x_2^3)$ $\Rightarrow \tilde{d}_{VC} = 10$

Degree 4: $\overset{\boldsymbol{\Phi_4}}{\to} (1, x_1, x_2, x_1^2, \dots, x_2^2, x_1^3, \dots, x_2^3, x_1^4, \dots, x_2^4)$ $\Rightarrow \tilde{d}_{VC} = 15$

In general:

For a Q-th order polynomial feature transform, with original $x$ of $d$ dimensions,

$$\tilde{d} = \binom{Q + d}{Q} - 1 \Rightarrow \tilde{d}_{VC} = \binom{Q + d}{Q}$$

# Pick Feature Transform BEFORE Seeing Data

Construct Features Carefully, Before Seeing the Data

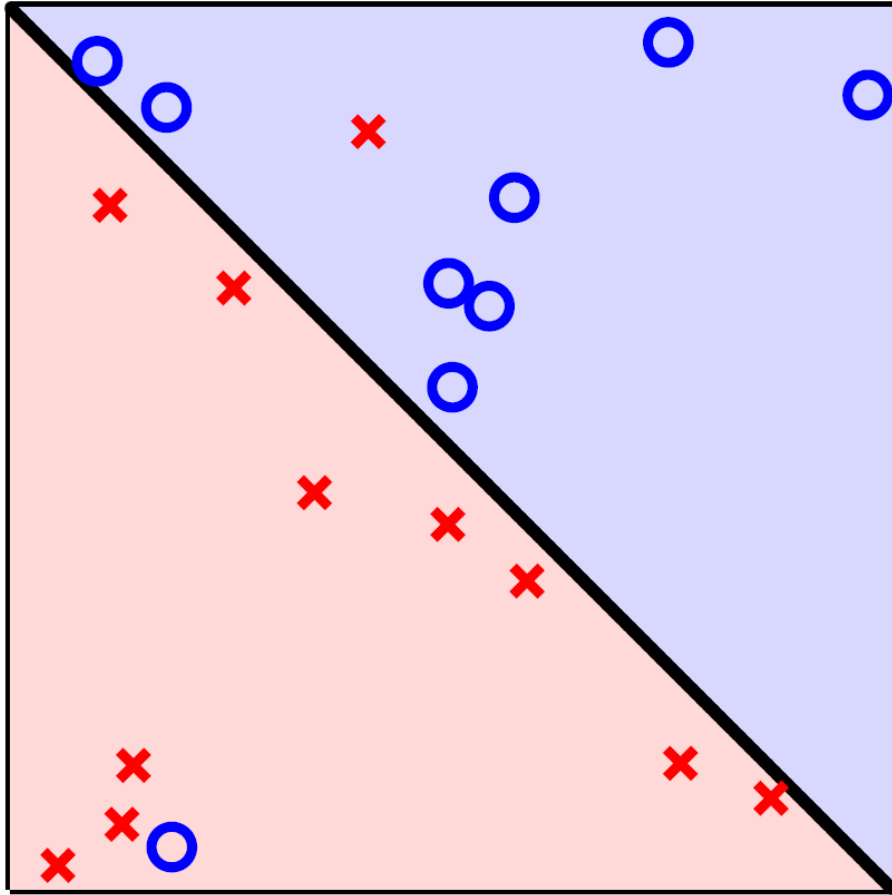- Do NOT Pick $\Phi$ after looking at the data

- Pick $\Phi$ BEFORE seeing the data

If you suspect linear model is insufficient to approximate the target,

Pick one of the standard feature transforms

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \xrightarrow{\Phi_2} \Phi_2(x) = \begin{bmatrix} 1 \\ \Phi_1(x) \\ \Phi_2(x) \\ \Phi_3(x) \\ \Phi_4(x) \\ \Phi_5(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix}$$

# Use Feature Transforms Responsibly

- Approximation – Generalization Tradeoff



$Q = 1$
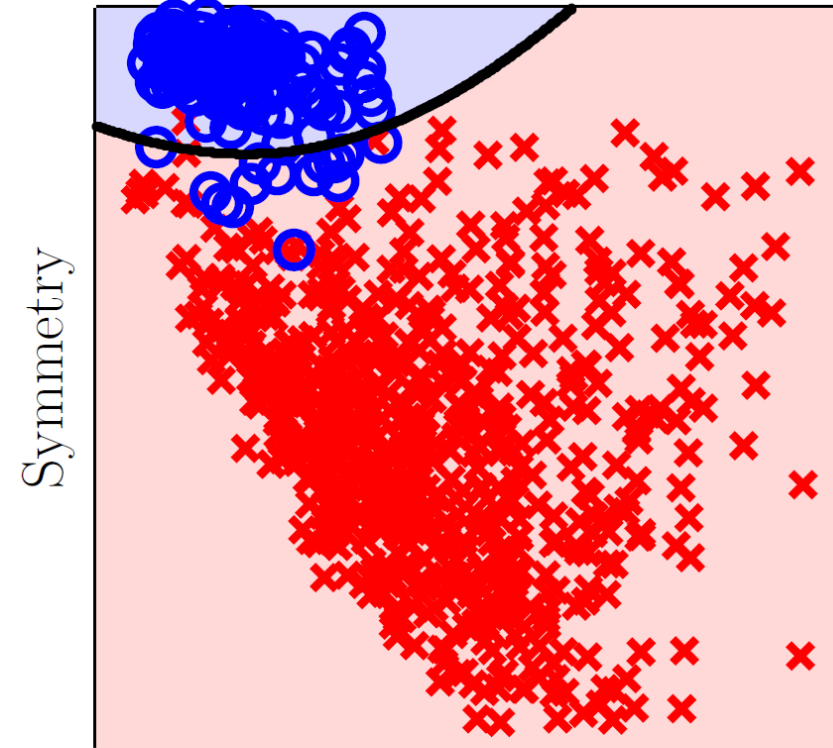
$Q = 4$

# Digits Data – "1" vs. "All"
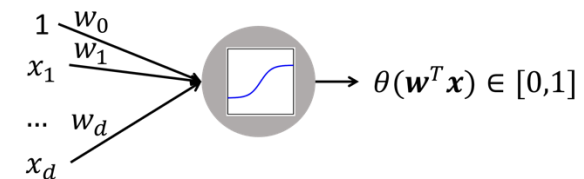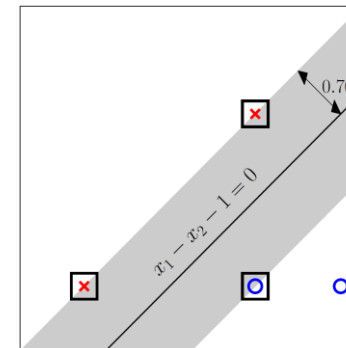


**Linear model**
$$E_{in} \approx 2.13\%$$
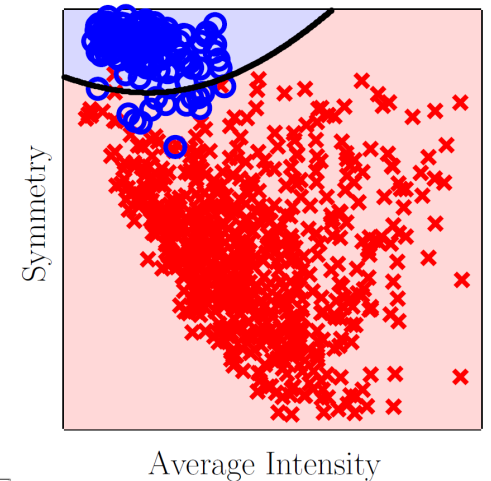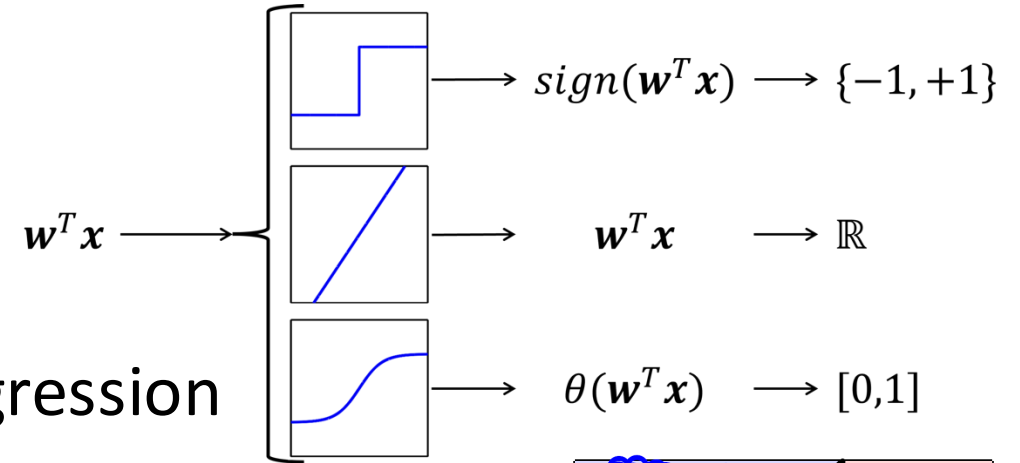$$E_{out} \approx 2.38\%$$

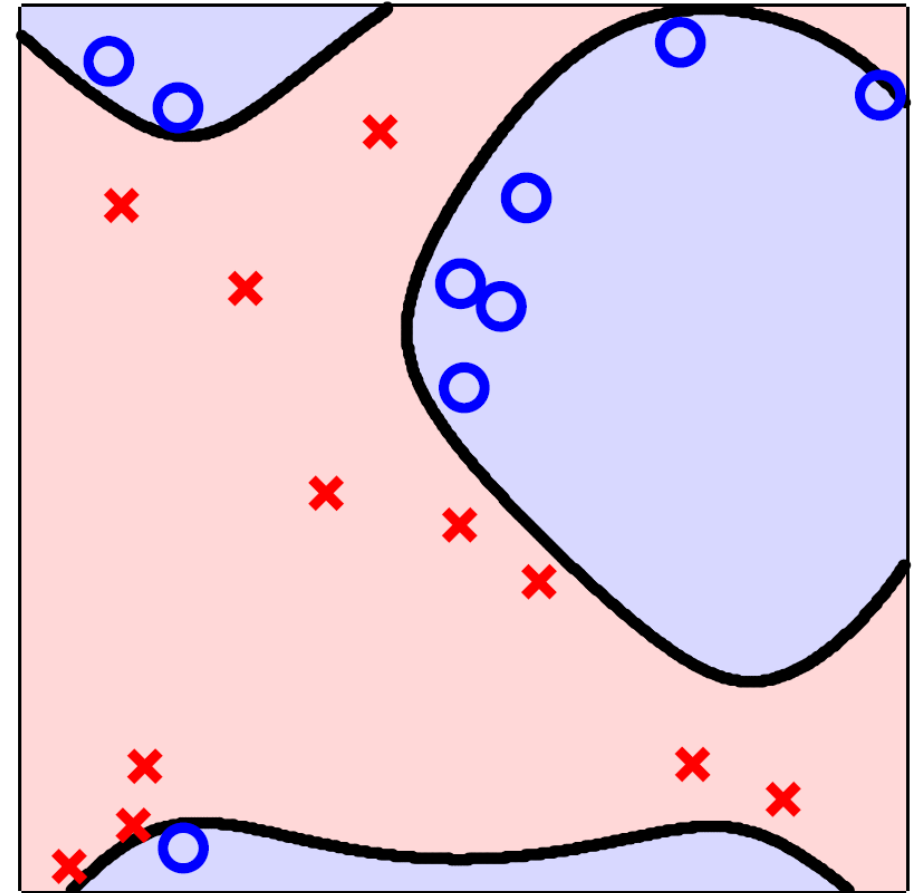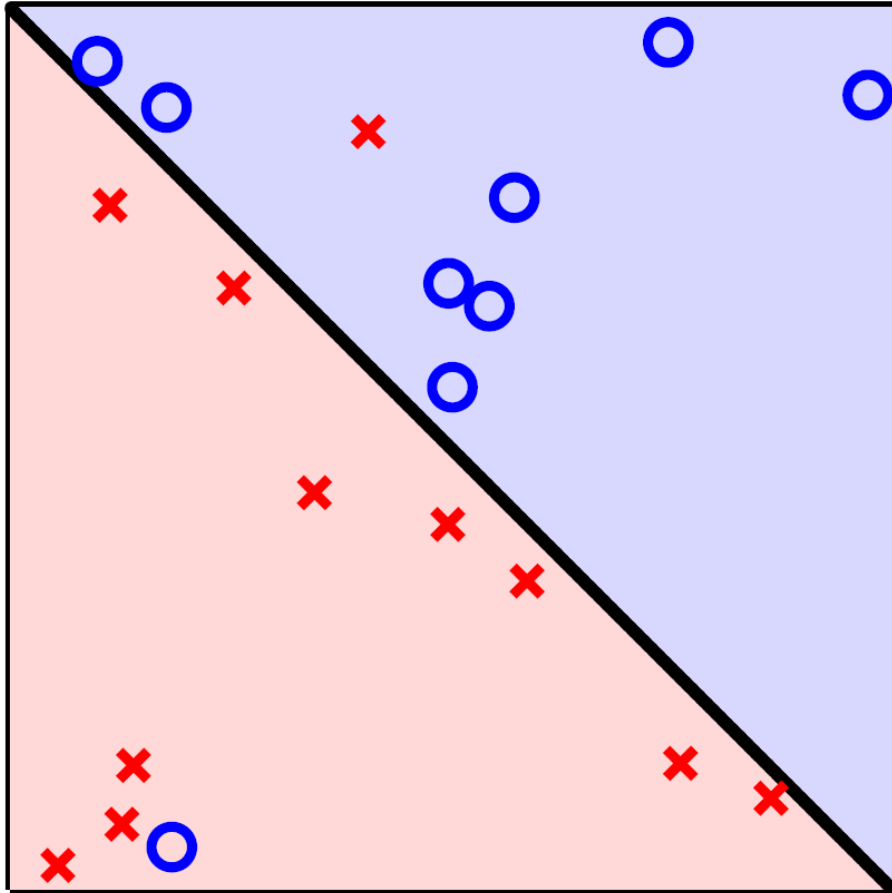**3$^{rd}$ Order Polynomial model**
$$E_{in} \approx 1.75\%$$
$$E_{out} \approx 1.87\%$$

# Use the Linear Model!

- Try it first – simple, robust, works
- Works for different kinds of problems:
  - Classification, Regression, Logistic Regression
- Works for nonlinear targets:
  - All the machinery and algorithms
    - Can tolerate errors
    - Easily extended to use nonlinear transforms
- Pick feature transform before seeing the data
- Linear Model is Fundamental!
  - Building block for more complex models

$$w^T x \longrightarrow$$

$$sign(w^T x) \longrightarrow \{-1, +1\}$$

$$w^T x \longrightarrow \mathbb{R}$$

$$\theta(w^T x) \longrightarrow [0,1]$$

Symmetry

Average Intensity

$$x_1 - x_2 - 1 = 0$$

0.707

$$\theta(w^T x) \in [0,1]$$

$1$ $w_0$
$x_1$ $w_1$
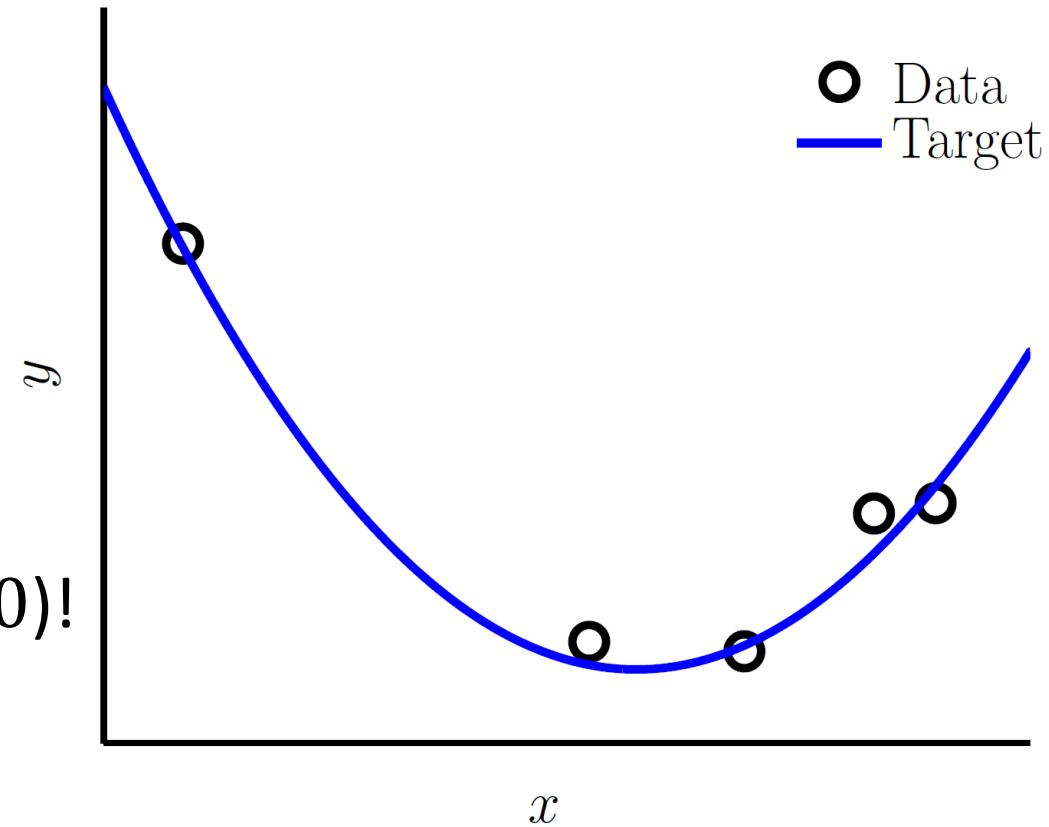$\ldots$ $w_d$
$x_d$

# Overfitting

# Another Simple Learning Problem

- Target $f$ is *quadratic*
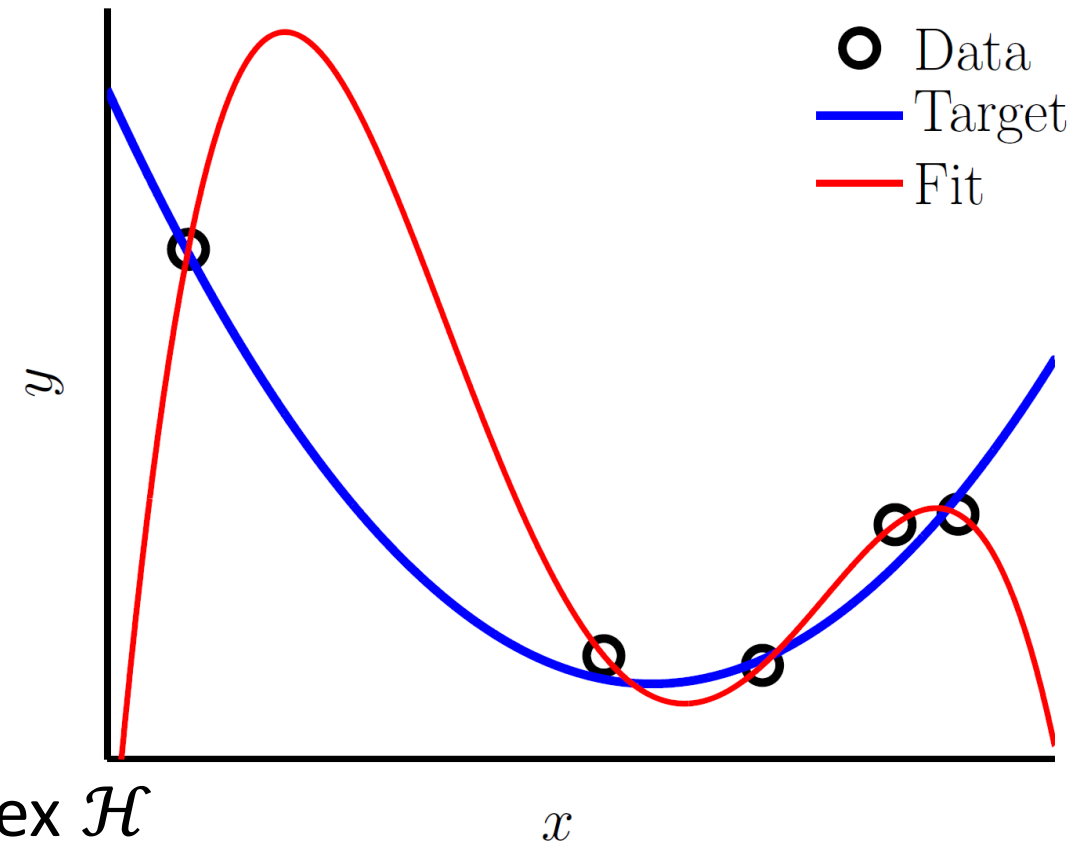- 5 data points $x_1, x_2, \ldots, x_5$
- A little noise (small measurement error)

- 5 data points?

$4^{\text{th}}$ order polynomial can fit it exactly ($E_{in} = 0$)!

# An Illustration of Overfitting

- Target $f$ is *quadratic*

- 5 data points $x_1, x_2, \ldots, x_5$

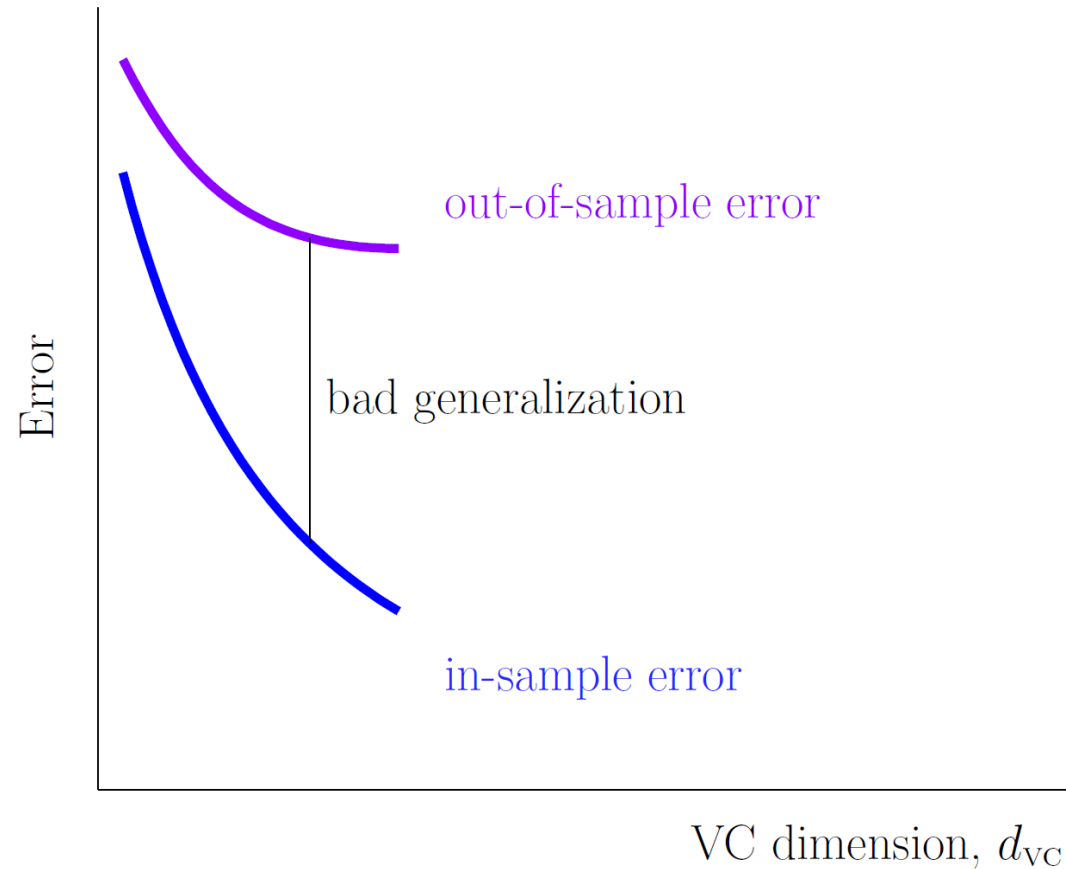- A little noise (small measurement error)

- Classic Overfitting

Learning a simple target with excessively complex $\mathcal{H}$

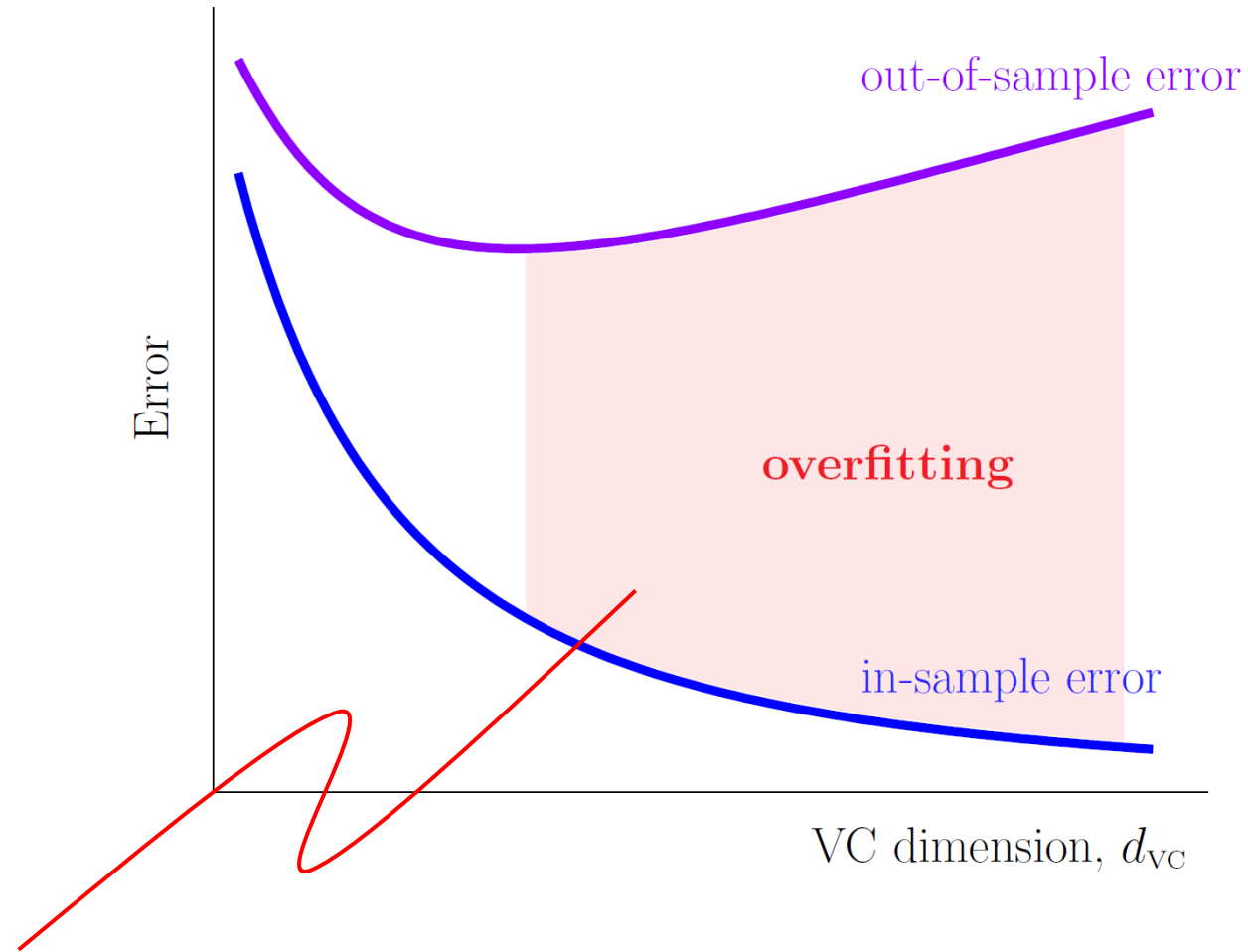Small $E_{in} \approx 0$, but High $E_{out} \gg 0$

Because we listened to the **noise**. (We fit the data more than was warranted)

# Overfitting is Not Just Bad Generalization



- VC analysis covers bad generalization
- Fitting the data, lowering $E_{in}$, is fine so long as $E_{out}$ is also lowered
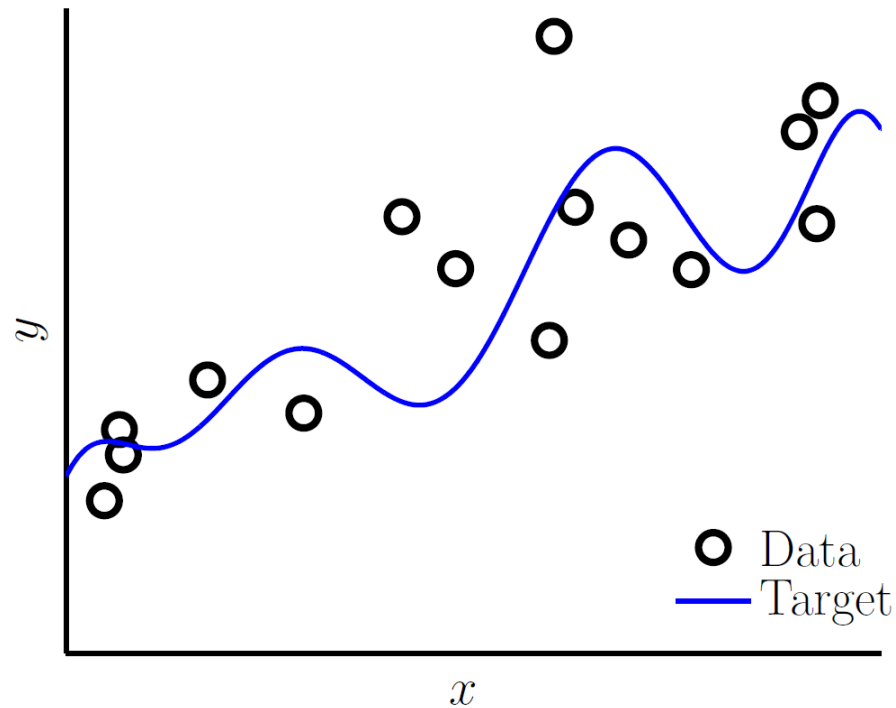
# Overfitting: When Lower $E_{in}$ results in *Higher $E_{out}$*



out-of-sample error

overfitting

in-sample error

Error

VC dimension, $d_{\mathrm{VC}}$

Lowering $E_{in}$ is no longer a good guide to lower $E_{out}$

# The Curious Case of the Simple Model
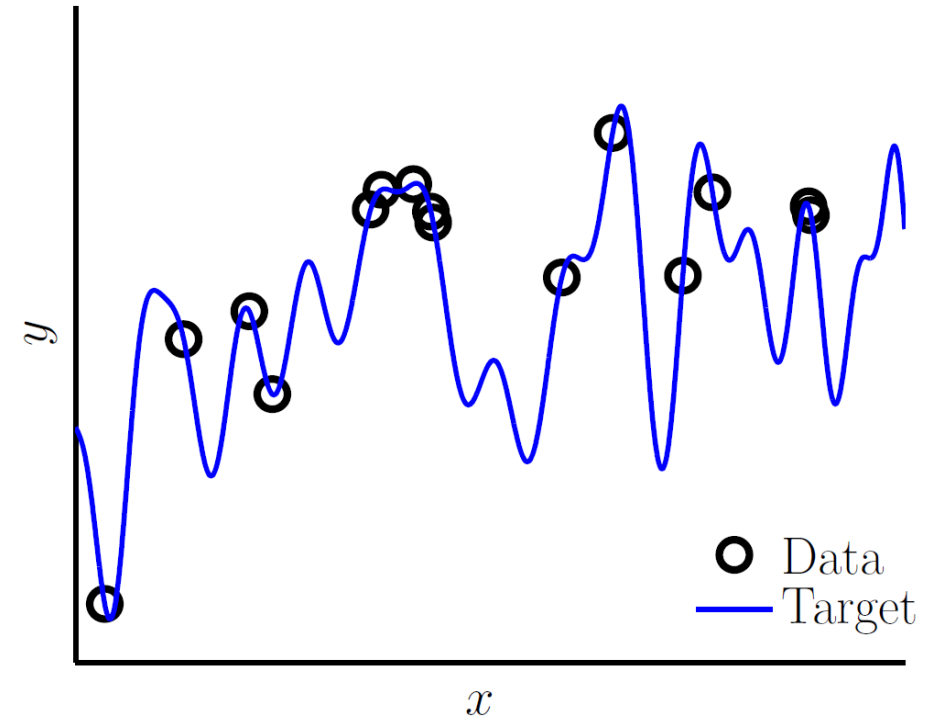## 2$^{nd}$ order vs. 10$^{th}$ Order Polynomial Fit
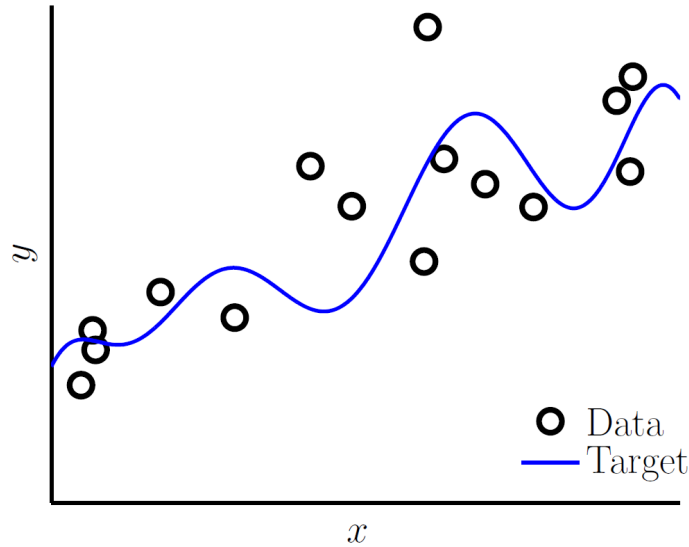## Two Learning Problems



Simple Noisy Target

Complex Noiseless Target

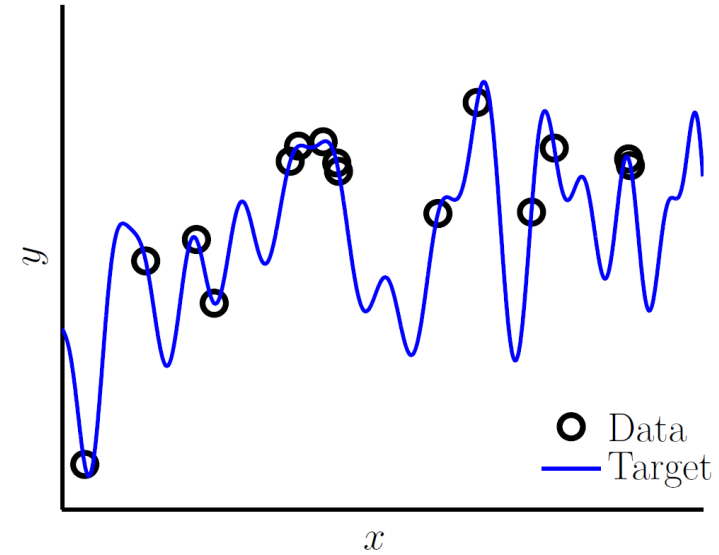10th order $f$ with noise.

50th order $f$ with no noise.

# Two Contenders: 2nd Order vs 10th Order Polynomial



10th order $f$ with noise.



50th order $f$ with no noise.
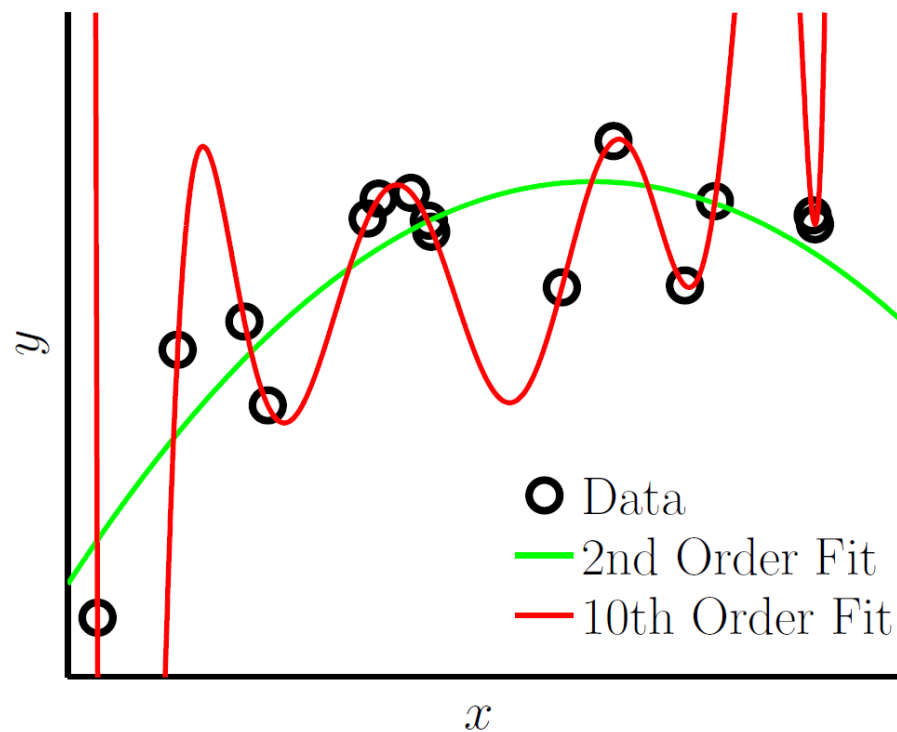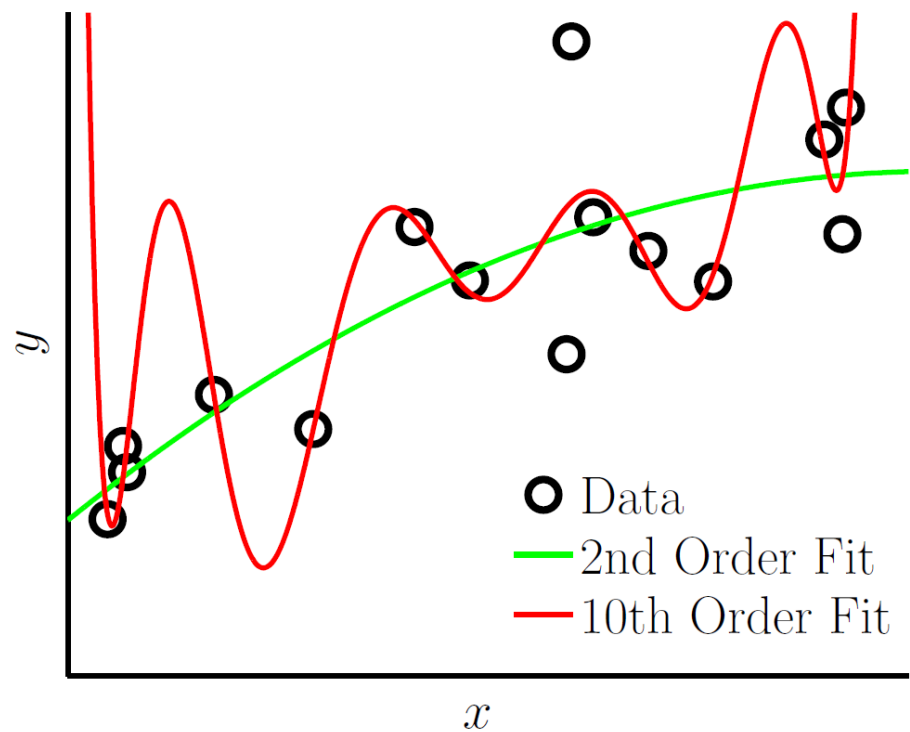
$$\mathcal{H}_2: h(\boldsymbol{x}) = w^T \Phi_2(\boldsymbol{x})$$

Learn linear model with 2nd order polynomial transform

$$\mathcal{H}_{10}: h(\boldsymbol{x}) = w^T \Phi_{10}(\boldsymbol{x})$$

Learn linear model with 10th order polynomial transform

# 2$^{nd}$ Order vs 10$^{th}$ Order Polynomial



**simple noisy target**

| | | 2$^{nd}$ order | 10$^{th}$ order |
|---|---|---|---|
| $E_{in}$ | \| | 0.050 | 0.034 |
| $E_{out}$ | \| | **0.127** | **9.00** |

**complex noiseless target** Small $N$

| | | 2$^{nd}$ order | 10$^{th}$ order |
|---|---|---|---|
| $E_{in}$ | \| | 0.029 | $10^{-5}$ |
| $E_{out}$ | \| | **0.120** | **7680** |

# When is $\mathcal{H}_2$ better than $\mathcal{H}_{10}$?

Learning curves for $\mathcal{H}_2$

Learning curves for $\mathcal{H}_{10}$



$\mathcal{H}$ should match quantity ($N$) and quality of data (noise), not the target ($f$)

Overfitting: When $E_{out}(\mathcal{H}_{10}) > E_{out}(\mathcal{H}_2)$

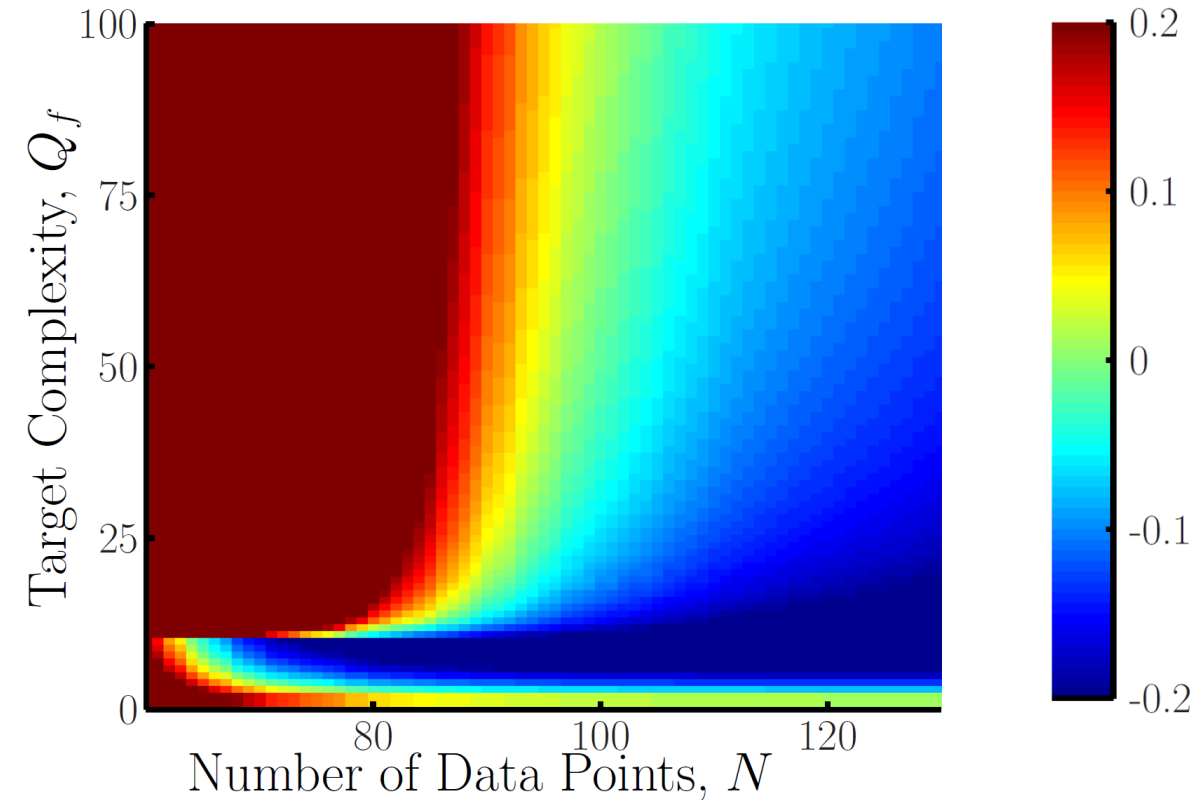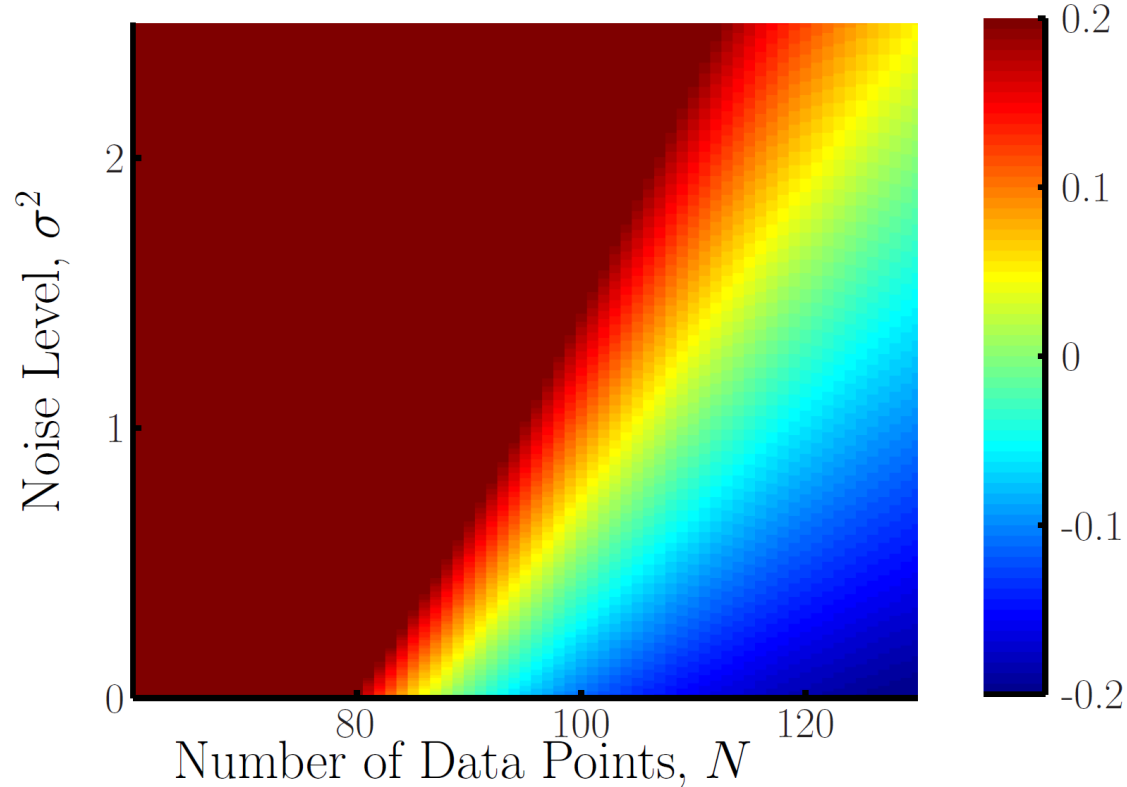# Overfit Measure: $E_{out}(\mathcal{H}_{10}) - E_{out}(\mathcal{H}_2)$

Learning Scenario Parameters:

- Noise level $\sigma^2$ (variance)

- Target complexity $Q_f$

- Number of Data Points $N$

# Overfit Measure: $E_{out}(\mathcal{H}_{10}) - E_{out}(\mathcal{H}_2)$



$E_{out}(\mathcal{H}_{10}) > E_{out}(\mathcal{H}_2)$

$E_{out}(\mathcal{H}_2) > E_{out}(\mathcal{H}_{10})$

| | | |
|---|---|---|
| Noise level $\sigma^2$ (variance) | ↑ ↑ | Overfitting |
| Target complexity $Q_f$ | ↑ ↑ | Overfitting |
| Number of Data Points $N$ | ↑ ↓ | Overfitting |

27

# Noise: The Part of $y$ we <span style="color:red">cannot model</span>

- Towards a unified view of overfitting


- Noise:
  - Stochastic noise: measurement error
  - Deterministic noise (?)

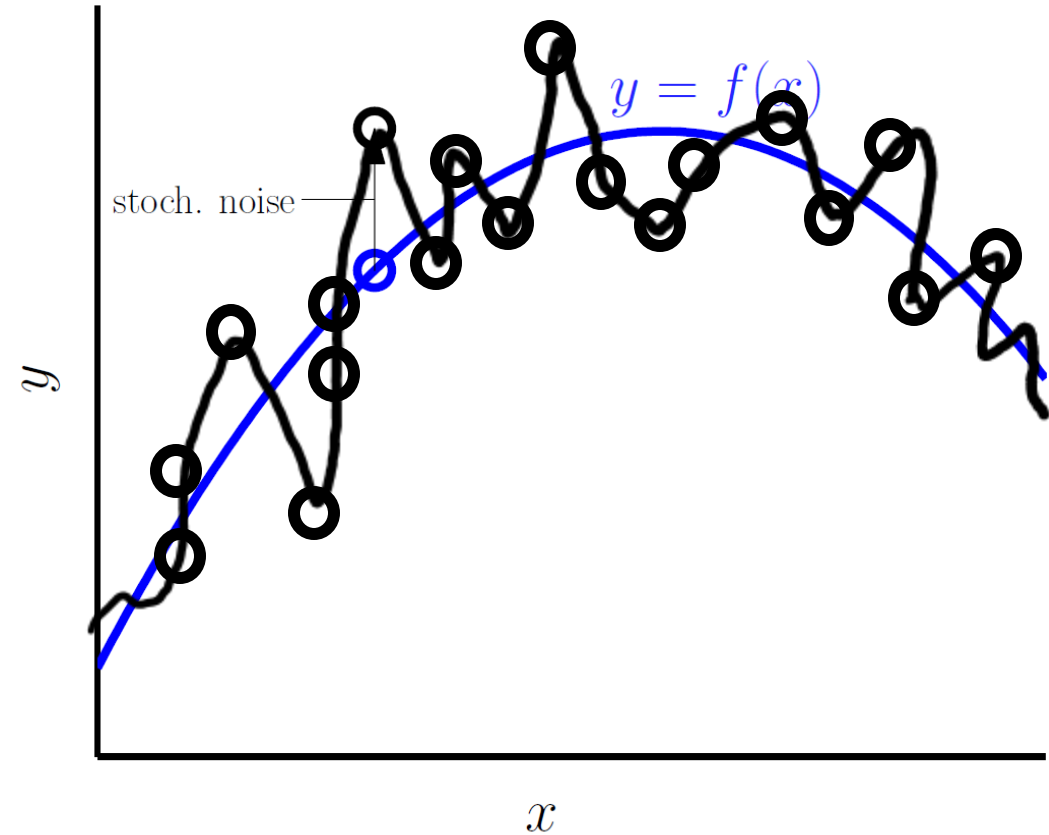# Stochastic Noise: A Part of $y$ we cannot model

Want to learn from **o**: $y = f(x)$

Unfortunately we only observe **o**:
$y = f(x) +$ stochastic noise

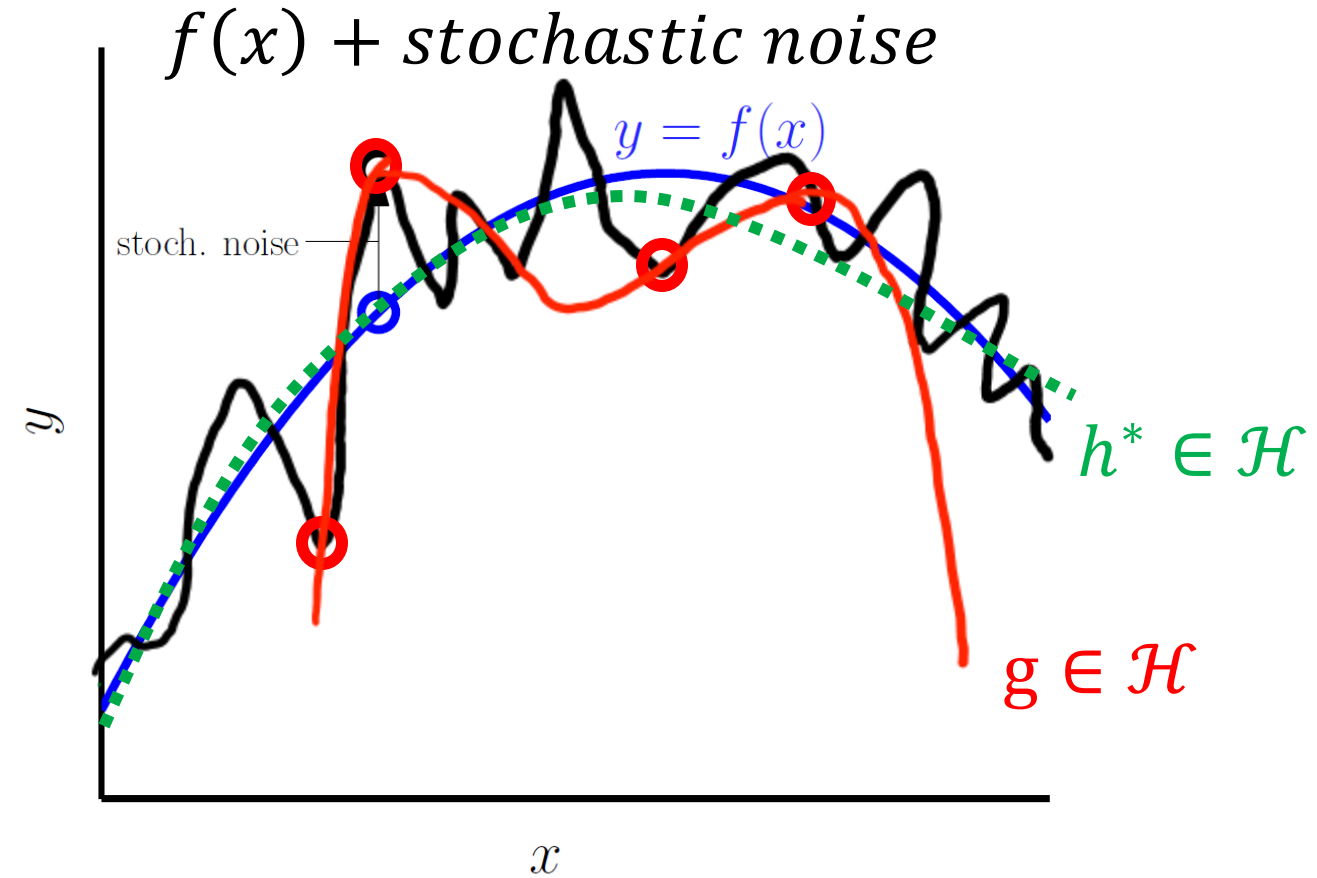We have no way to model stochastic noise



stoch. noise

$y = f(x)$

$y$

$x$

# Stochastic Noise: A Part of $y$ we cannot model

Want to learn from **o**: $y = f(x)$

Unfortunately we only observe **o**:
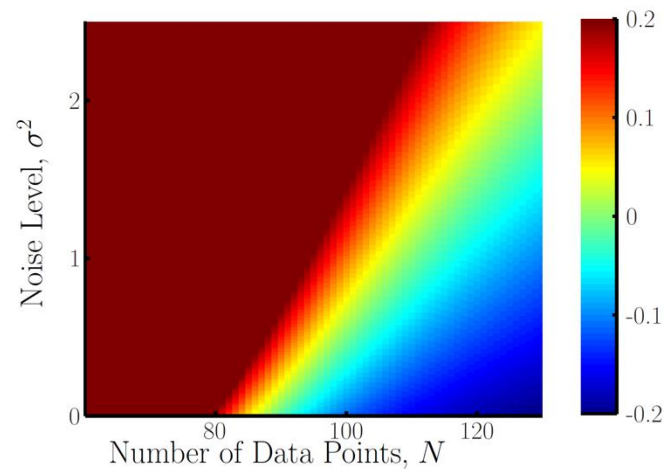$y = f(x) +$ stochastic noise

We have no way to model stochastic noise



$f(x) + stochastic\ noise$

$y = f(x)$

stoch. noise

$h^* \in \mathcal{H}$

$g \in \mathcal{H}$

Noise level $\sigma^2$ (variance) ↑↑  Overfitting

# Stochastic and Deterministic Noise

Stochastic noise:
measurement errors we cannot model

Deterministic noise:
the other part of $y$ we cannot model
the best we can possibly do with $\mathcal{H}$

# Deterministic Noise

# Deterministic Noise: Model Error

Want: Pick $h^*$

the best approximation of $f$ in $\mathcal{H}$

Want: To learn from **o**

Unfortunately, we observe **o**

$$y_n = f(x_n)$$
$$y_n = h^*(x_n) + \text{deterministic noise}$$

Target complexity $Q_f$ ↑ ↑ Overfitting

# Stochastic and Deterministic Noise Hurt Learning

**Stochastic Noise**

random measurement errors

- Measuring $y_n$ again stochastic noise changes

- Independent of $\mathcal{H}$

**Deterministic Noise**

$\mathcal{H}$ cannot model $f$

- Measuring $y_n$ again deterministic noise constant

- Depends on $\mathcal{H}$



$y = f(x)$

$y = f(x)$ + stoch. noise



$h^* \in \mathcal{H}$

$y = h^*(x)$ + det. noise

# Bias Variance Decomposition with Stochastic Noise

$$\mathbb{E}[E_{out}(\boldsymbol{x})] = bias(\boldsymbol{x}) \;\; + \;\;\; var(\boldsymbol{x})$$
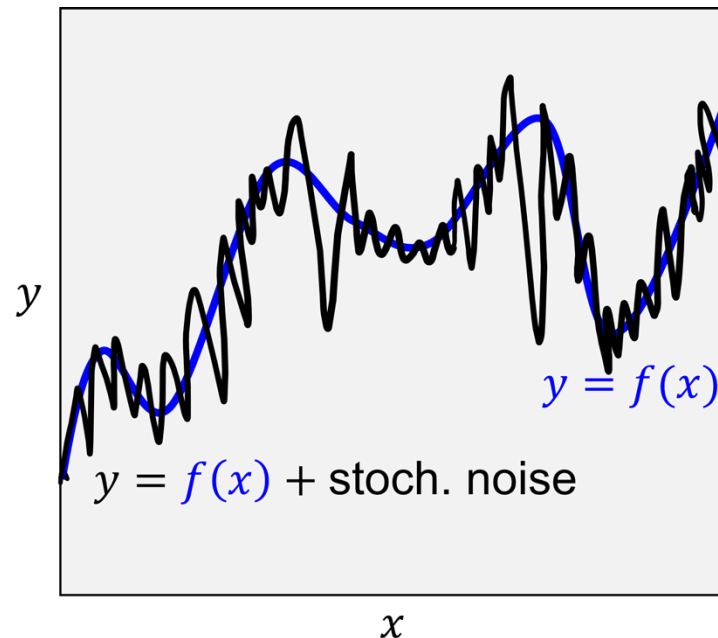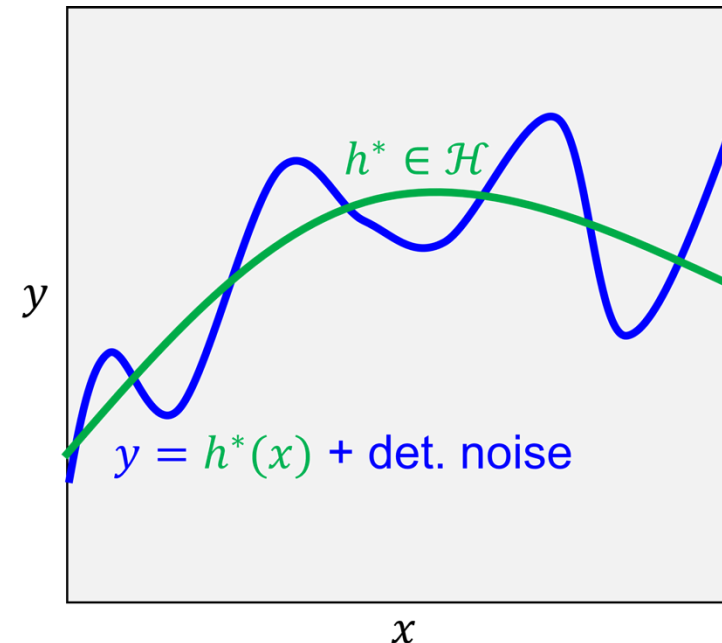
$$bias(\boldsymbol{x}) = \left(\bar{g}(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2$$

$$var(\boldsymbol{x}) = var(g(\boldsymbol{x}))$$

$$E_{out}(\boldsymbol{x}) = \left(g(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2$$

With noise:

$$
\begin{aligned}
E_{out}(\boldsymbol{x}) \;\; &= (g(\boldsymbol{x}) - y)^2, \text{ where } y = f(\boldsymbol{x}) + \epsilon \\
&= (g(\boldsymbol{x}) - f(\boldsymbol{x}) - \epsilon)^2 \\
&= \left(g(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2 - 2\epsilon\left(g(\boldsymbol{x}) - f(\boldsymbol{x})\right) + \epsilon^2
\end{aligned}
$$

# Bias Variance Decomposition with Stochastic Noise

$$E_{out}(\boldsymbol{x}) \quad = \quad \big(g(\boldsymbol{x}) - f(\boldsymbol{x})\big)^2 - 2\epsilon\big(g(\boldsymbol{x}) - f(\boldsymbol{x})\big) + \epsilon^2$$

$$\mathbb{E}_{\boldsymbol{x}}[E_{out}] \ = \mathbb{E}_{\boldsymbol{x}}\left[\big(g(\boldsymbol{x}) - f(\boldsymbol{x})\big)^2 - 2\epsilon\big(g(\boldsymbol{x}) - f(\boldsymbol{x})\big) + \epsilon^2\right]$$

$$\mathbb{E}[\epsilon] = 0 \Rightarrow$$

$$\mathbb{E}_{\boldsymbol{x}}[E_{out}] \ = \mathbb{E}_{\boldsymbol{x}}\left[\big(g(\boldsymbol{x}) - f(\boldsymbol{x})\big)^2\right] + \mathbb{E}_{\boldsymbol{x}}[\epsilon^2]$$

$$\mathbb{E}_{\boldsymbol{x}}[E_{out}] \ = \mathbb{E}_{\boldsymbol{x}}[bias(\boldsymbol{x}) \ + var(\boldsymbol{x}) \ + \epsilon^2]$$

$$E_{out} \quad = \quad bias \quad + \quad \sigma^2 \quad + \quad var \qquad \text{depends on } N$$

det. noise      stoch. noise    indirect impact of noise