

# 10-701: Introduction to Machine Learning

## Lecture 8 – Regularization

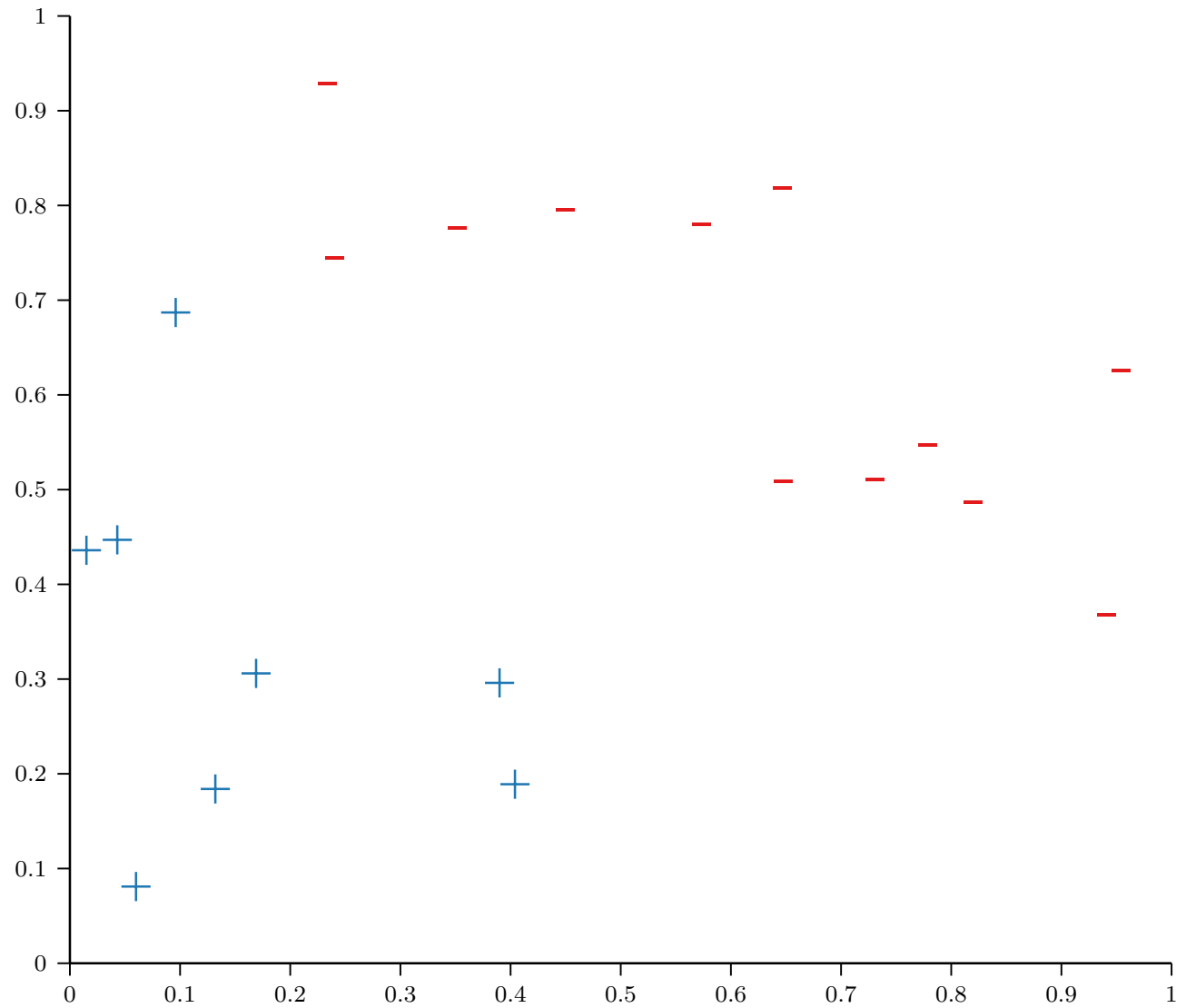
Henry Chai

2/12/24

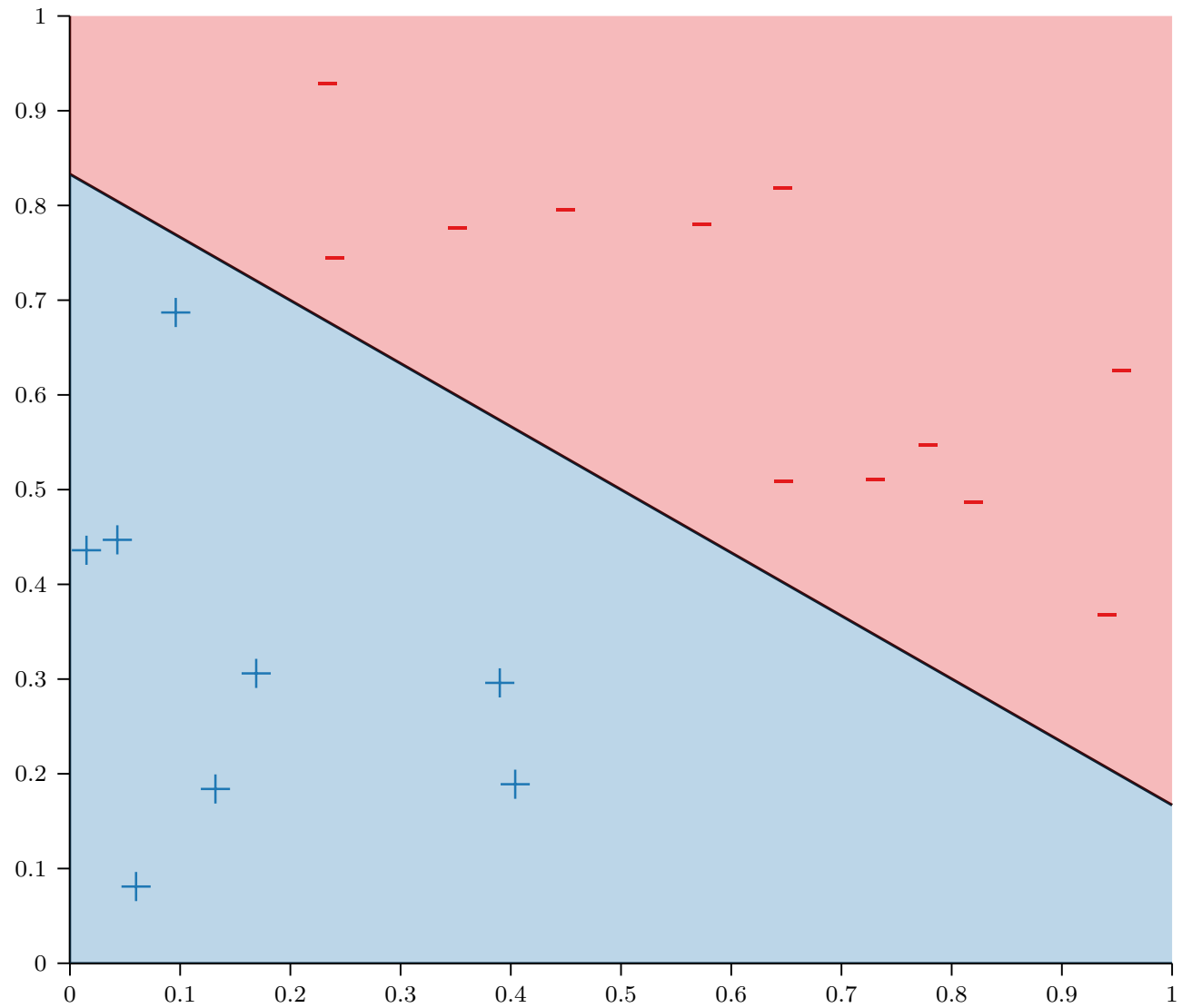
# Front Matter

- Announcements:
  - HW2 released 2/7, due **2/19** (previously 2/16) at 11:59 PM
  - HW3 released **2/19** (previously 2/16), due **2/28** (previously 2/26) at 11:59 PM
  - Lecture schedule has been updated, [see the course website](#) for full details
    - Lecture on 2/21 (Wednesday) and Recitation on 2/23 (Friday) have been swapped
- Recommended Readings:
  - Murphy, [Sections 7.5 & 14.4](#)

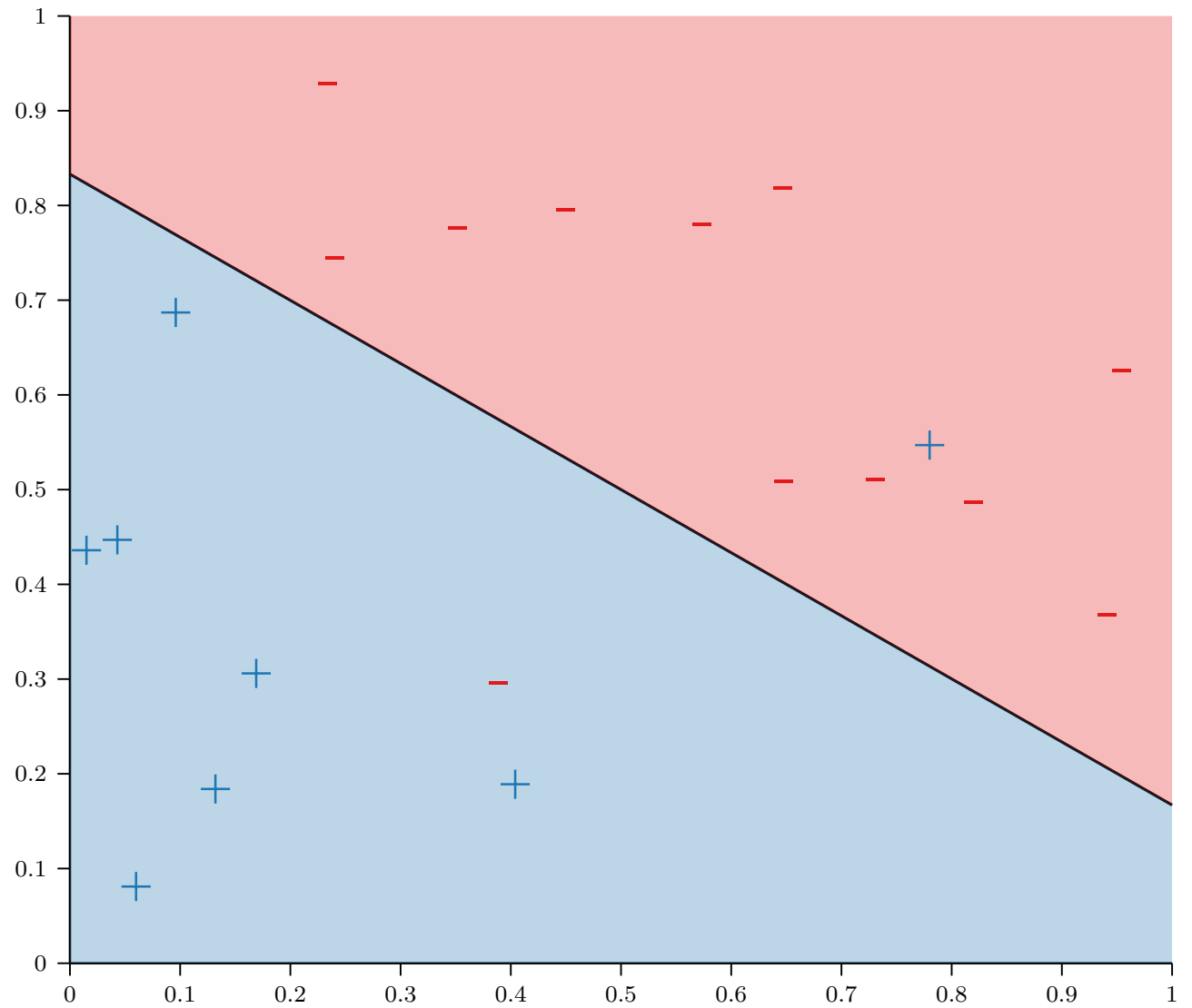
# Linear Models



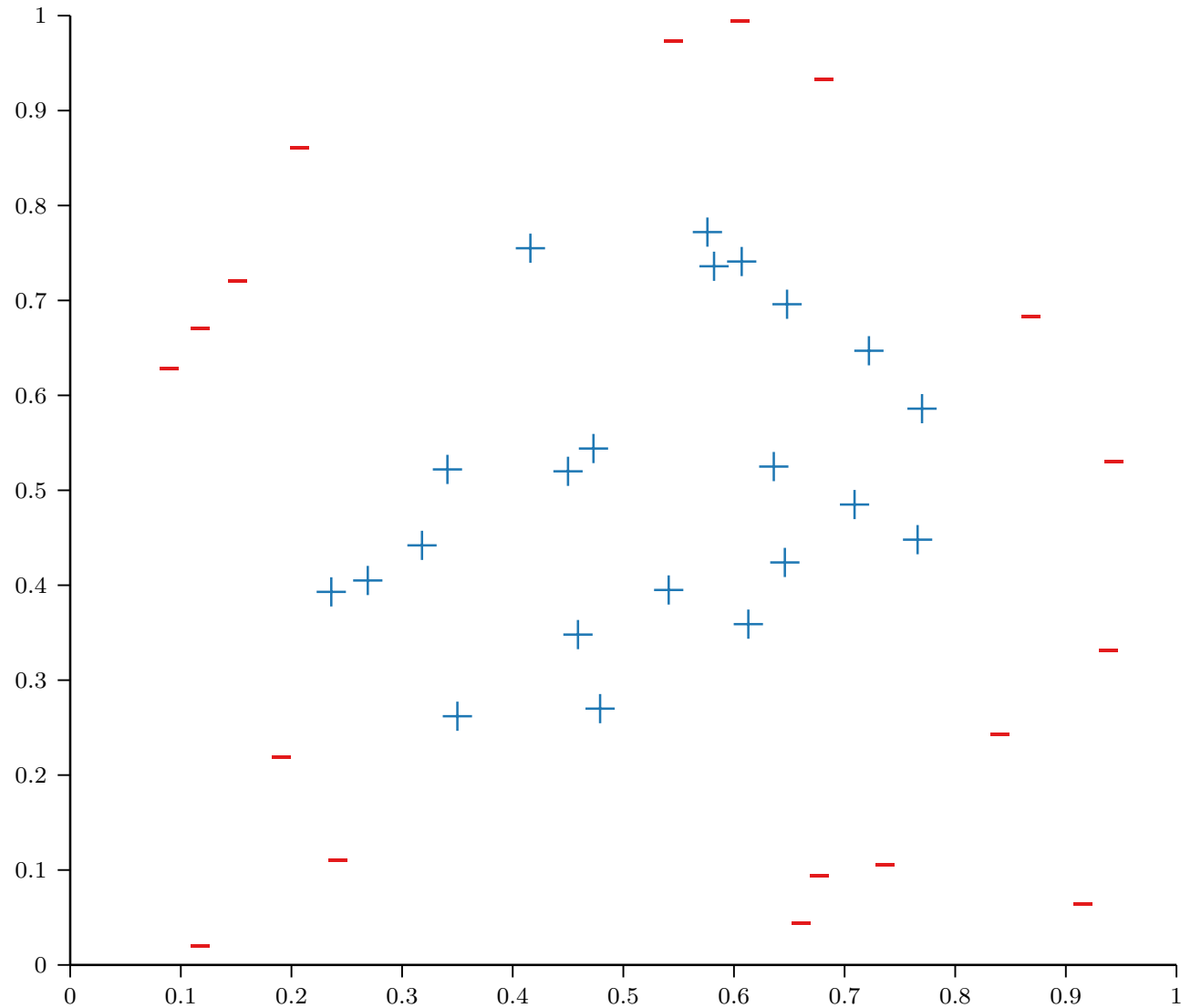
# Linear Models



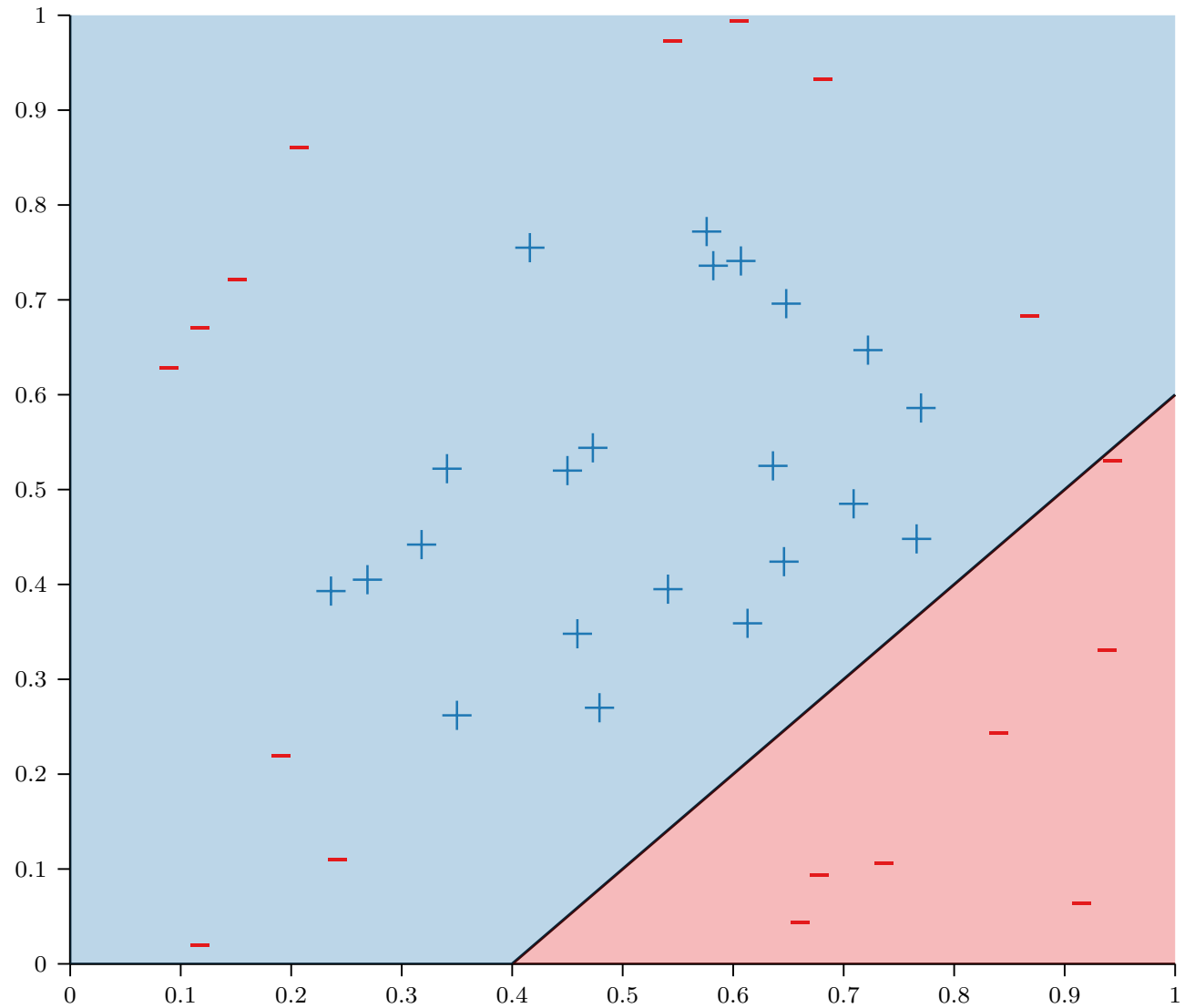
# Linear Models



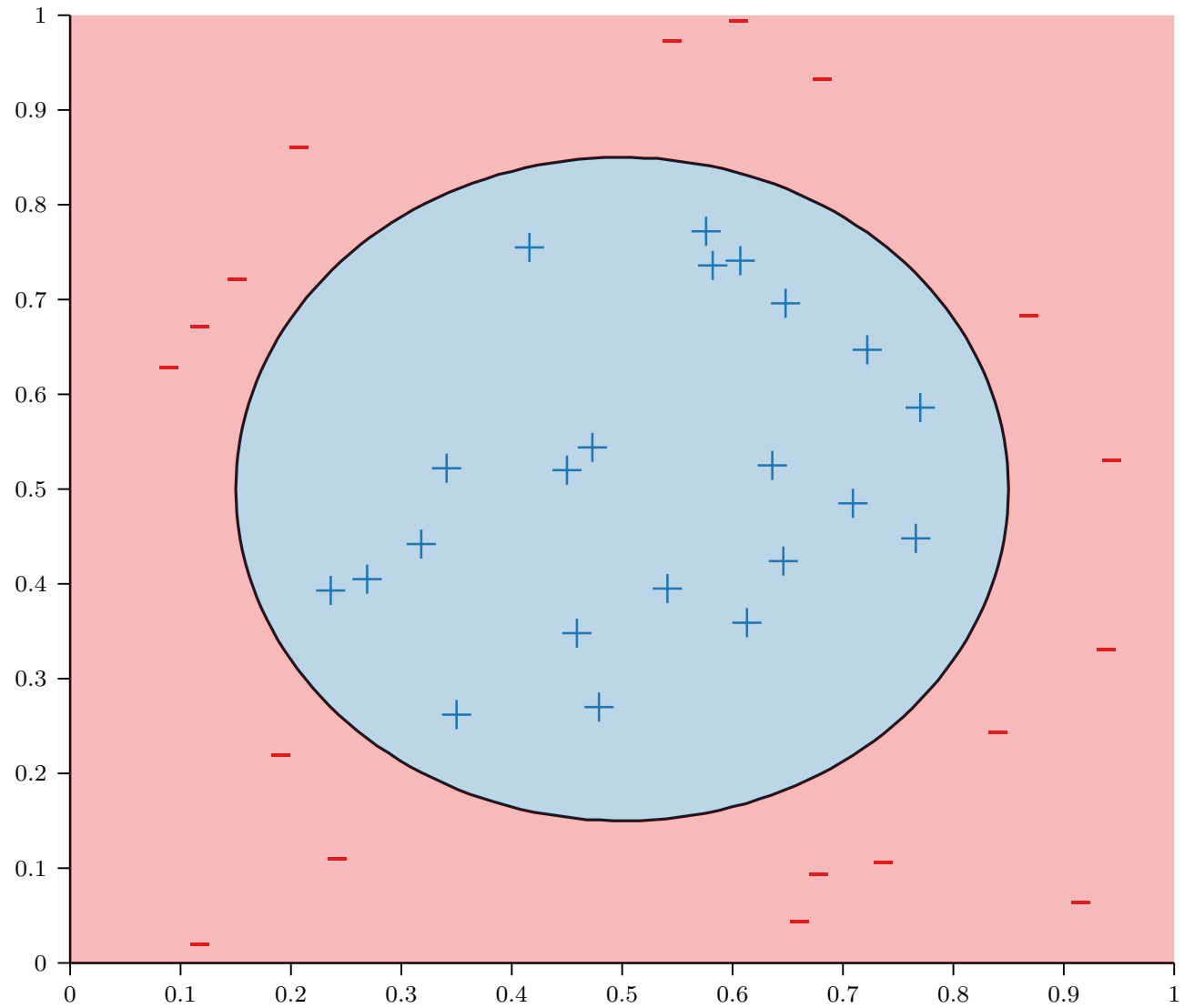
# Linear Models?



# Linear Models?



# Nonlinear Models



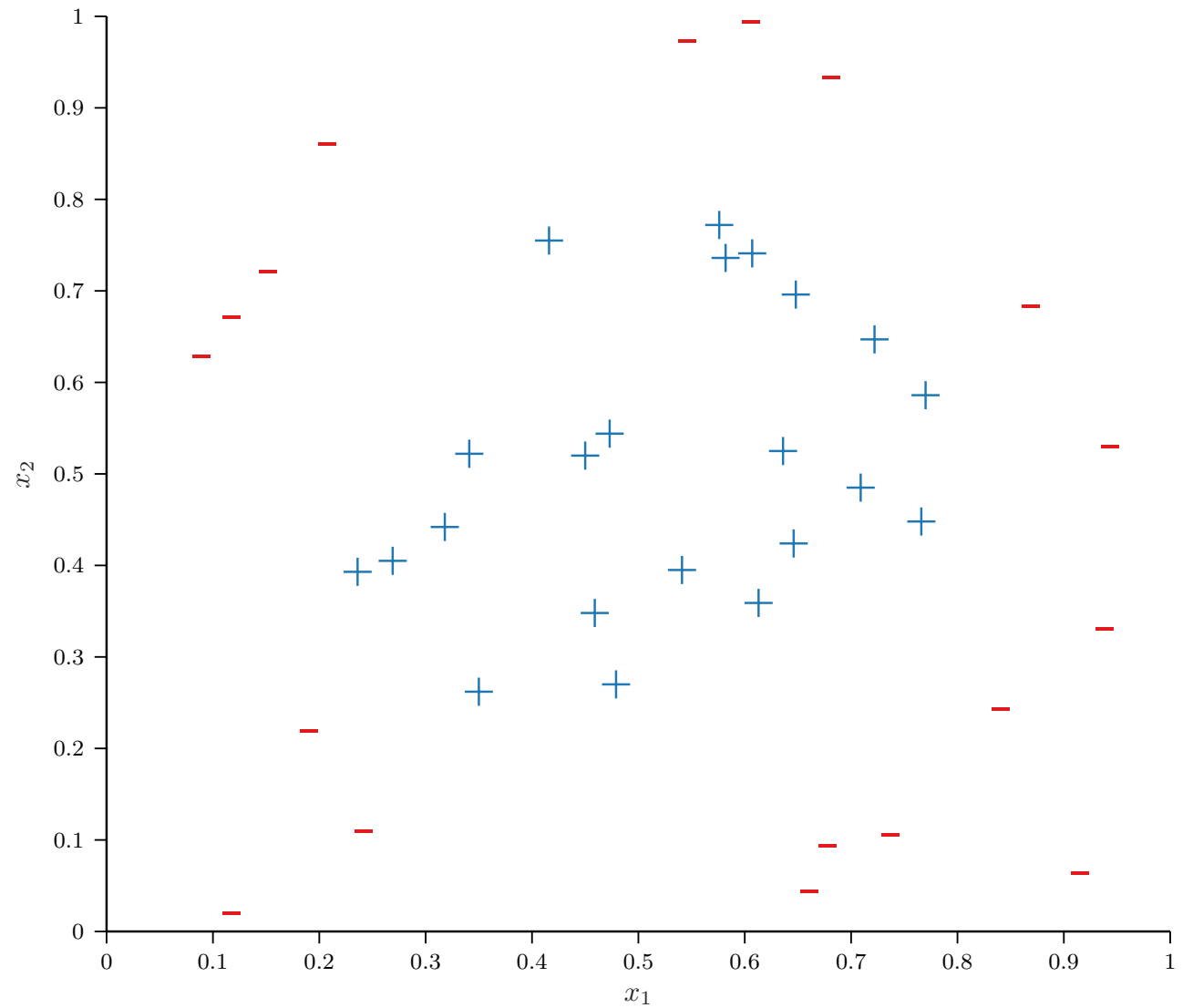


# Feature Transforms

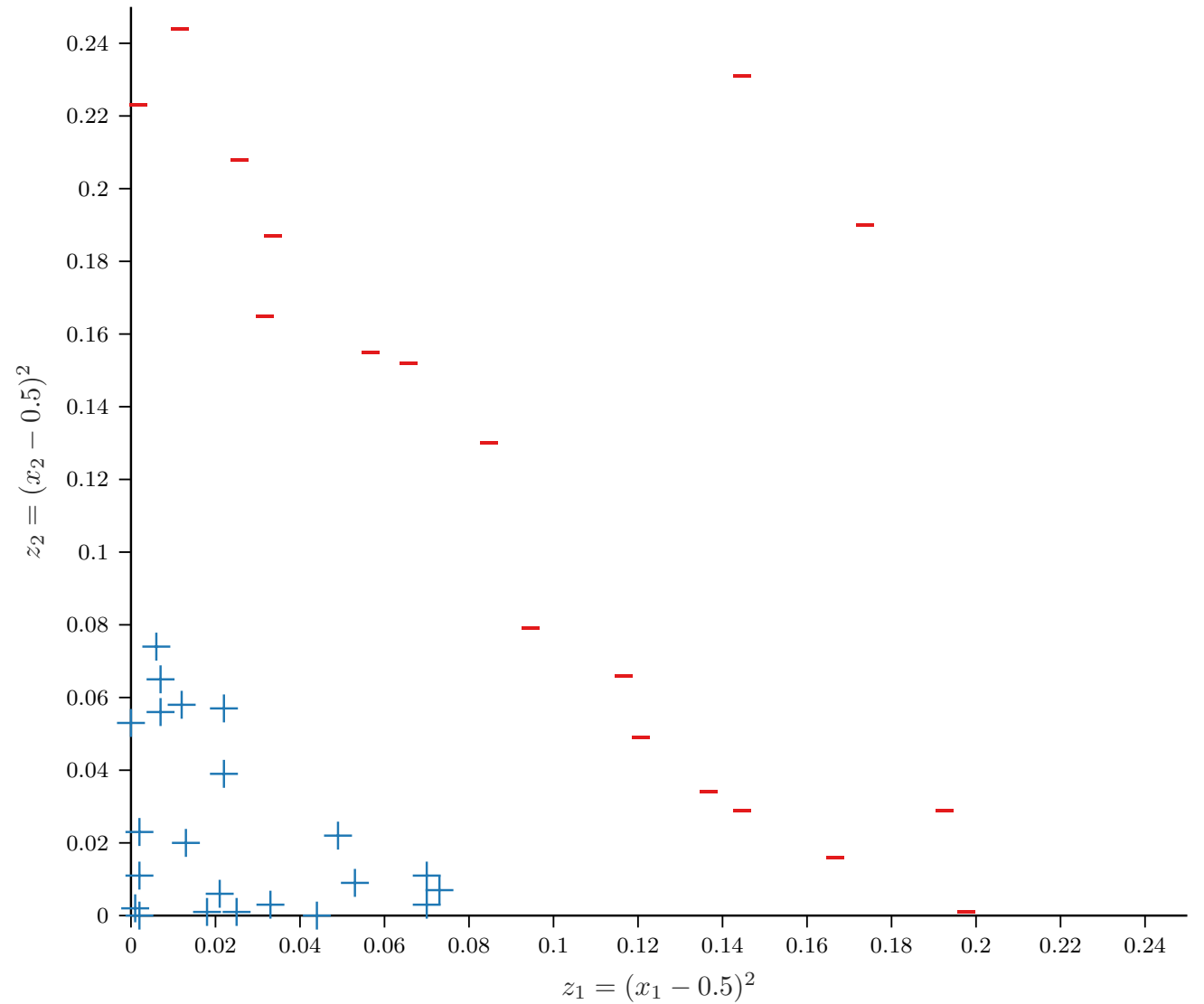
- Given  $D$ -dimensional inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , first compute some transformation of our input, e.g.,

$$\phi([x_1, x_2]) = [z_1 = (x_1 - 0.5)^2, z_2 = (x_2 - 0.5)^2]$$

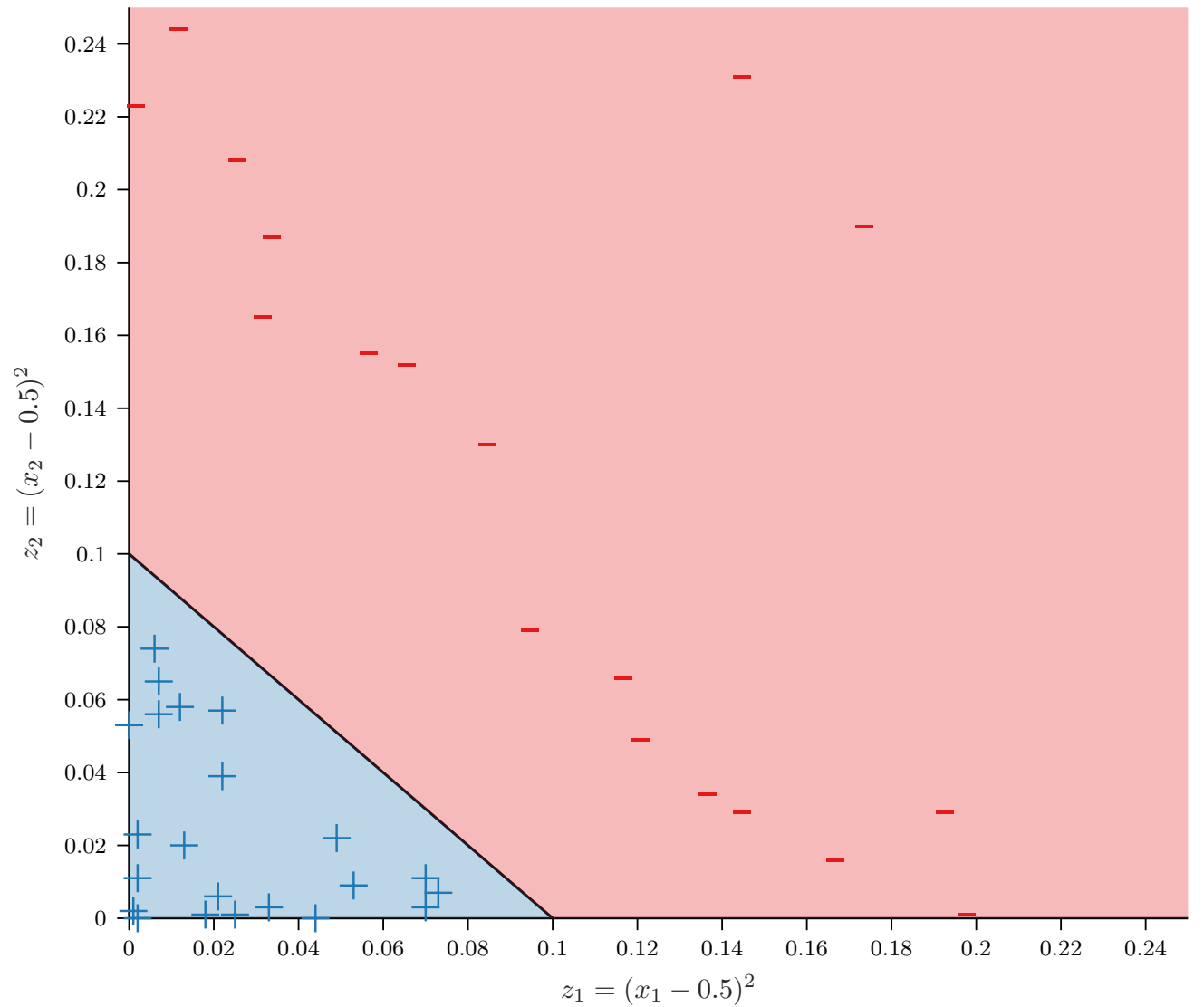
# Nonlinear Models



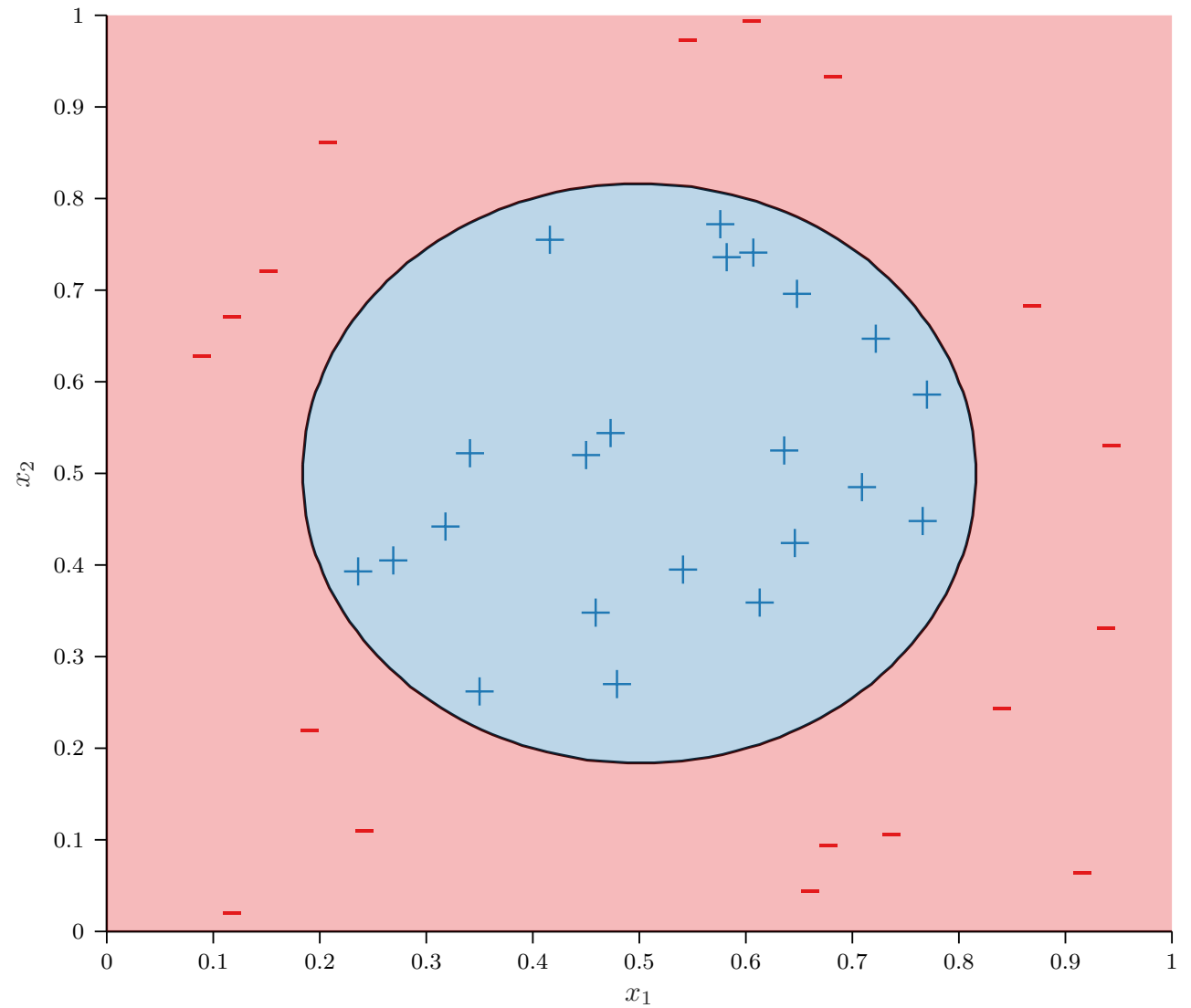
# Nonlinear Models



# Nonlinear Models



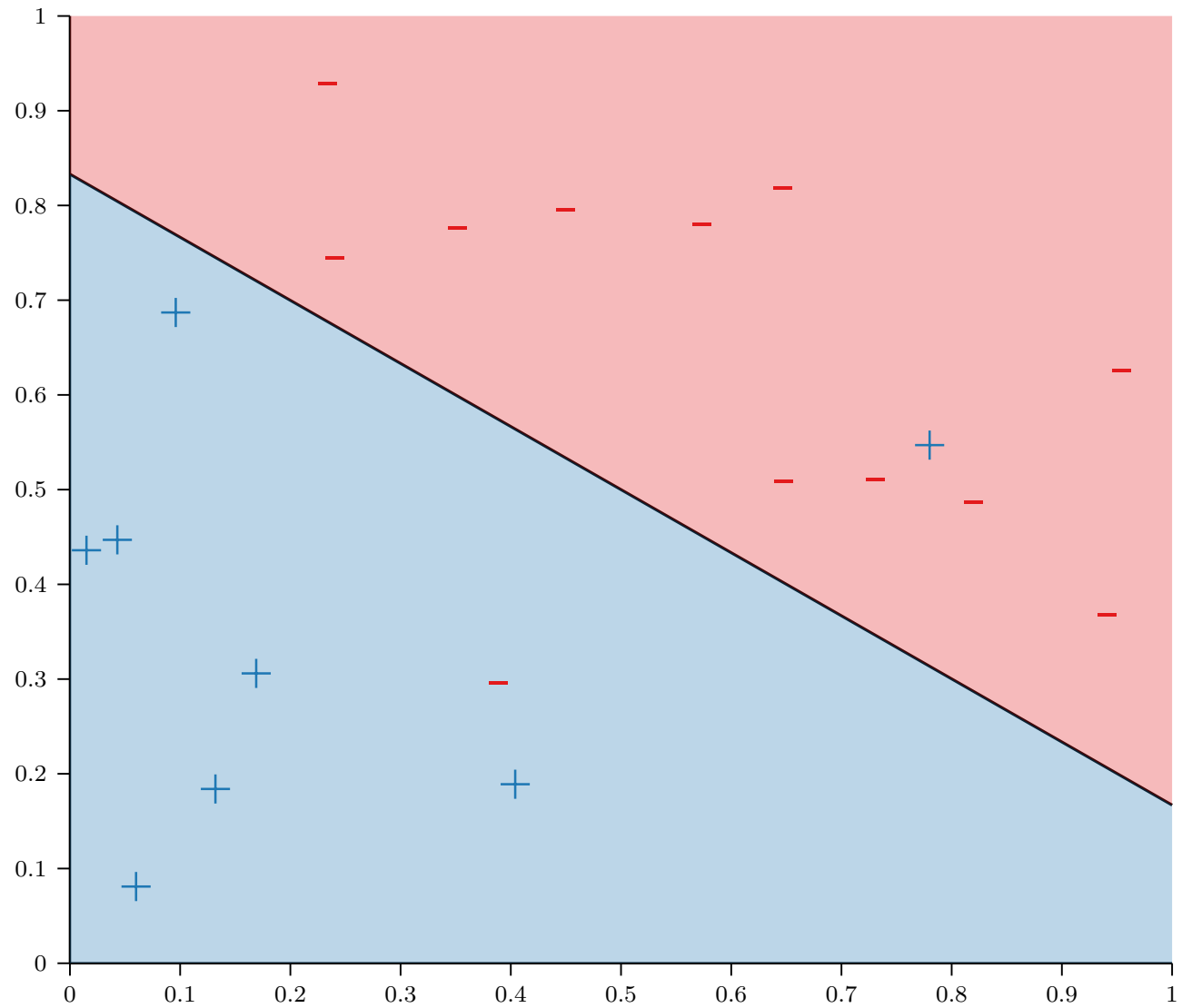
# Nonlinear Models



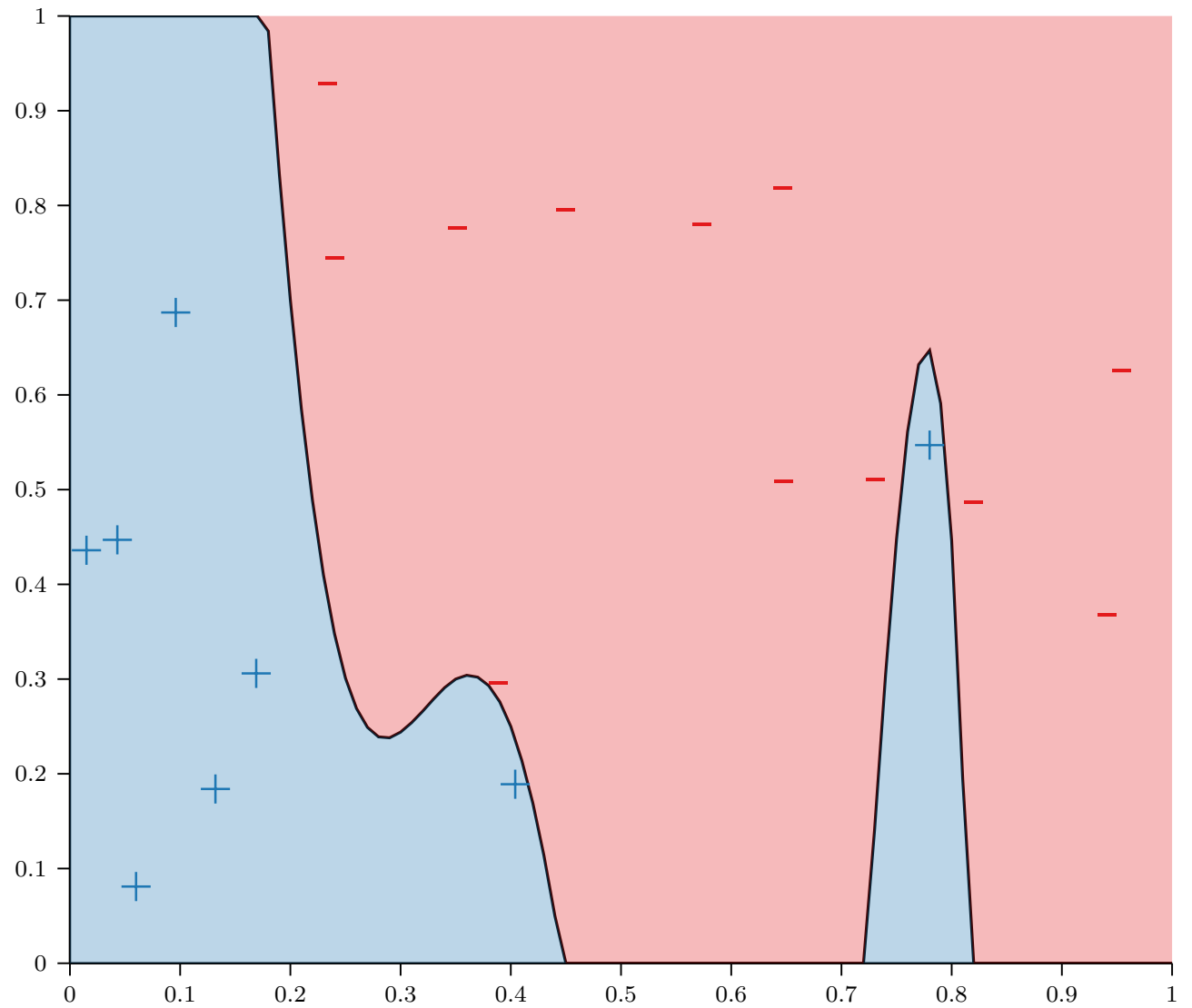
# General $Q^{th}$ -order Transforms

- $\phi_{2,2}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2]$
- $\phi_{2,3}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3]$
- $\phi_{2,4}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3, x_1^4, x_1^3x_2, x_1^2x_2^2, x_1x_2^3, x_2^4]$
- $\phi_{2,Q}$  maps a 2-D input to a  $O(Q^2)$ -D output
- Scales even worse for higher-dimensional inputs...

# Linear Models



# Nonlinear Models?





# Feature Transforms: Tradeoffs

	Low-Dimensional Input Space	High-Dimensional Input Space
Training Error	High	Low
Generalization	Good	Bad

# Feature Transforms: Experiment

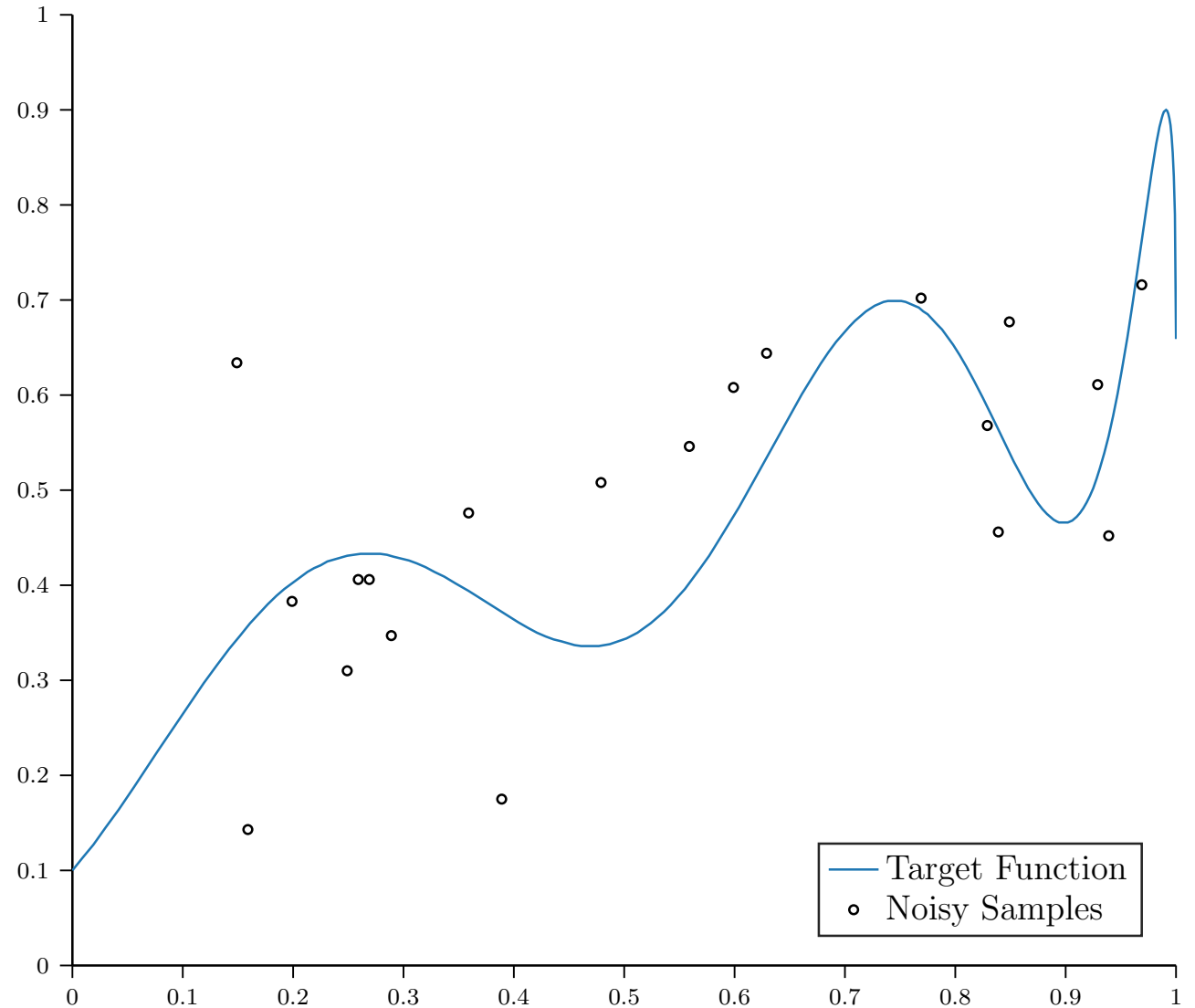
- $x \in \mathbb{R}, y \in \mathbb{R}$  and  $N = 20$
- Targets are generated by a 10<sup>th</sup>-order polynomial in  $x$  with additive Gaussian noise:

$$y = \sum_{d=0}^{10} a_d x^d + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2)$$

- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomials
  - $\phi_{1,2}(x) = [x, x^2]$
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomials
  - $\phi_{1,10}(x) = [x, x^2, x^3, x^4, x^5, x^6, x^7, x^8, x^9, x^{10}]$

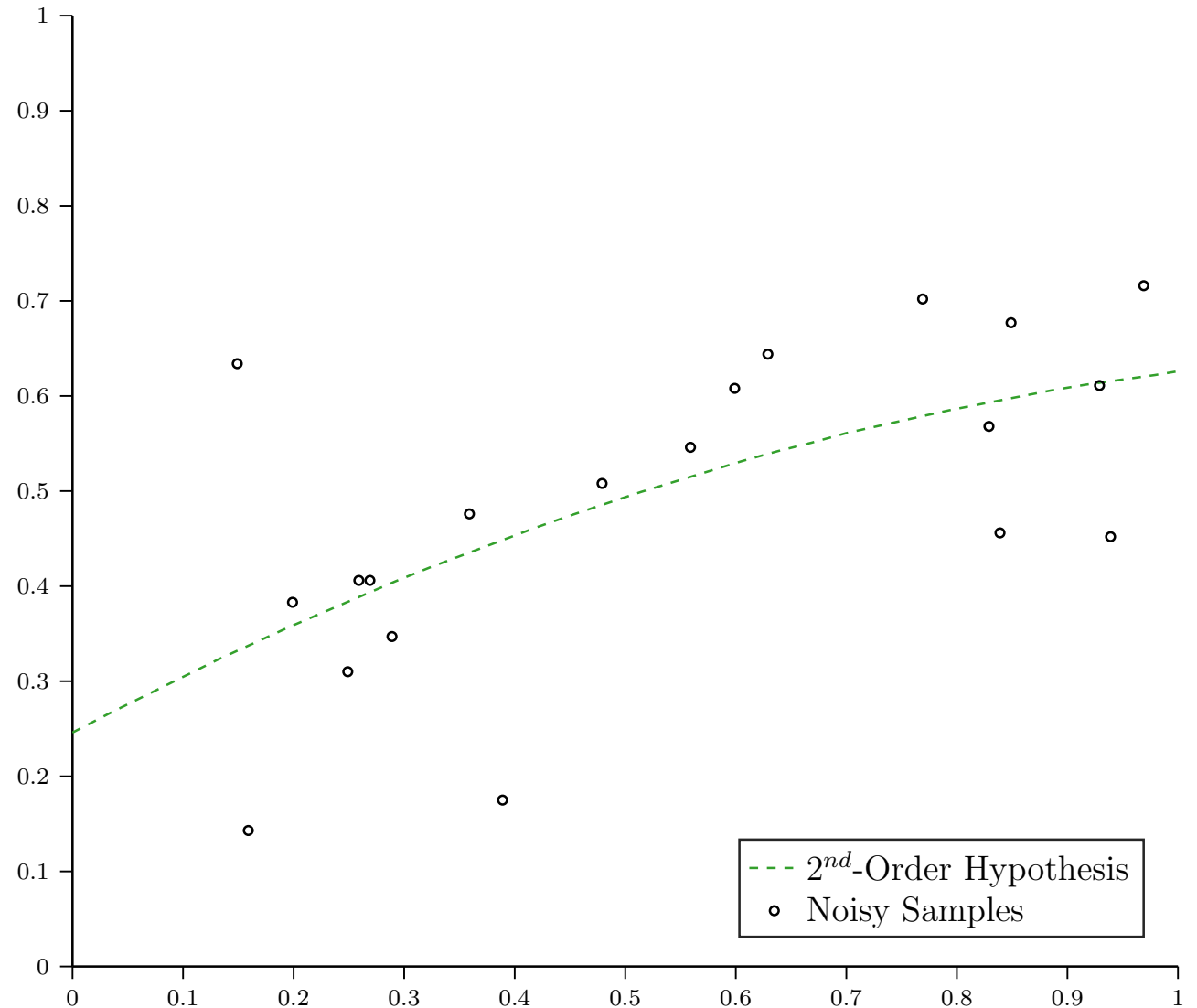
# Noisy Targets

- 10-dimensional target function with additive Gaussian noise
- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomial
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomial



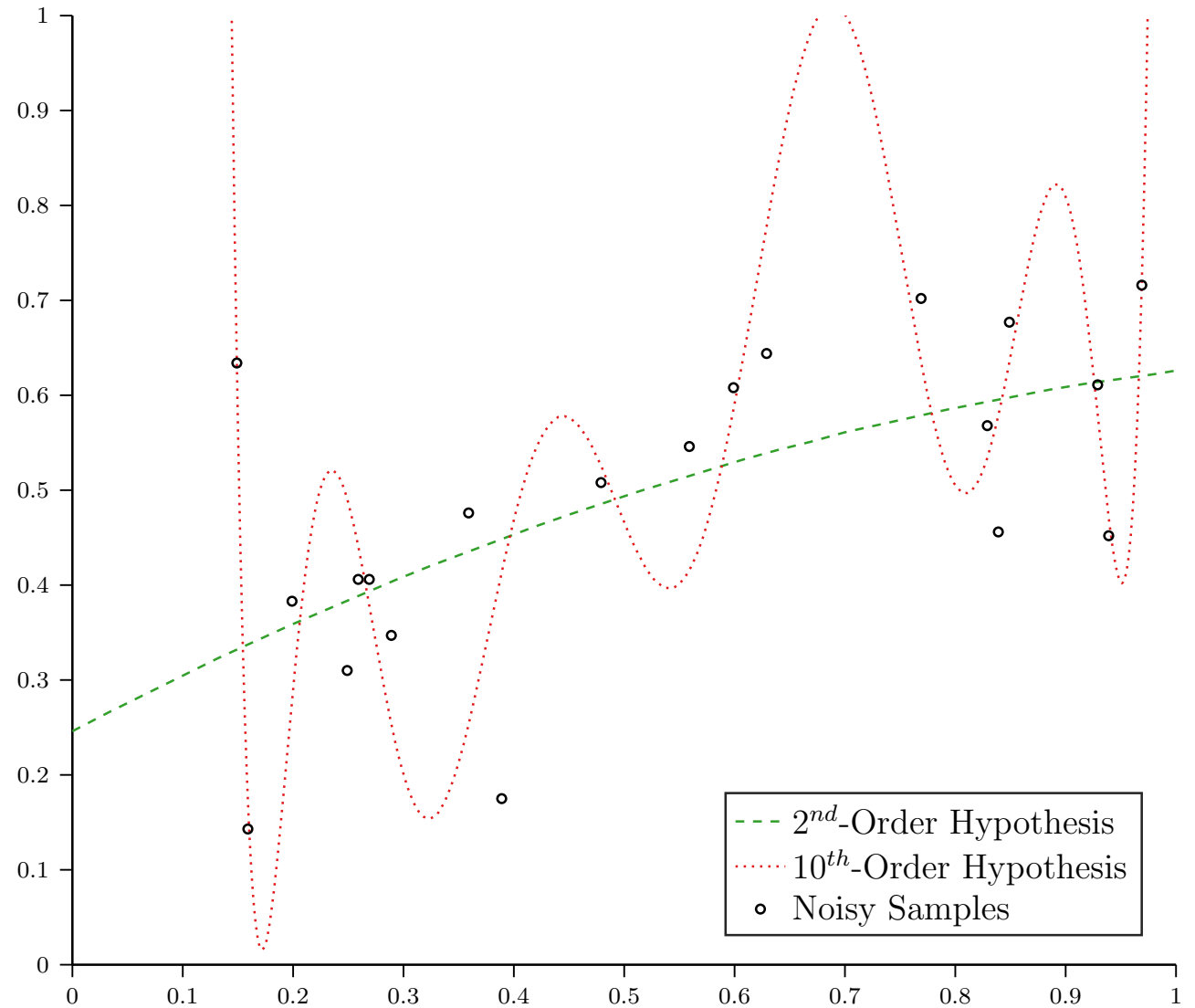
# Noisy Targets

- 10-dimensional target function with additive Gaussian noise
- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomial
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomial



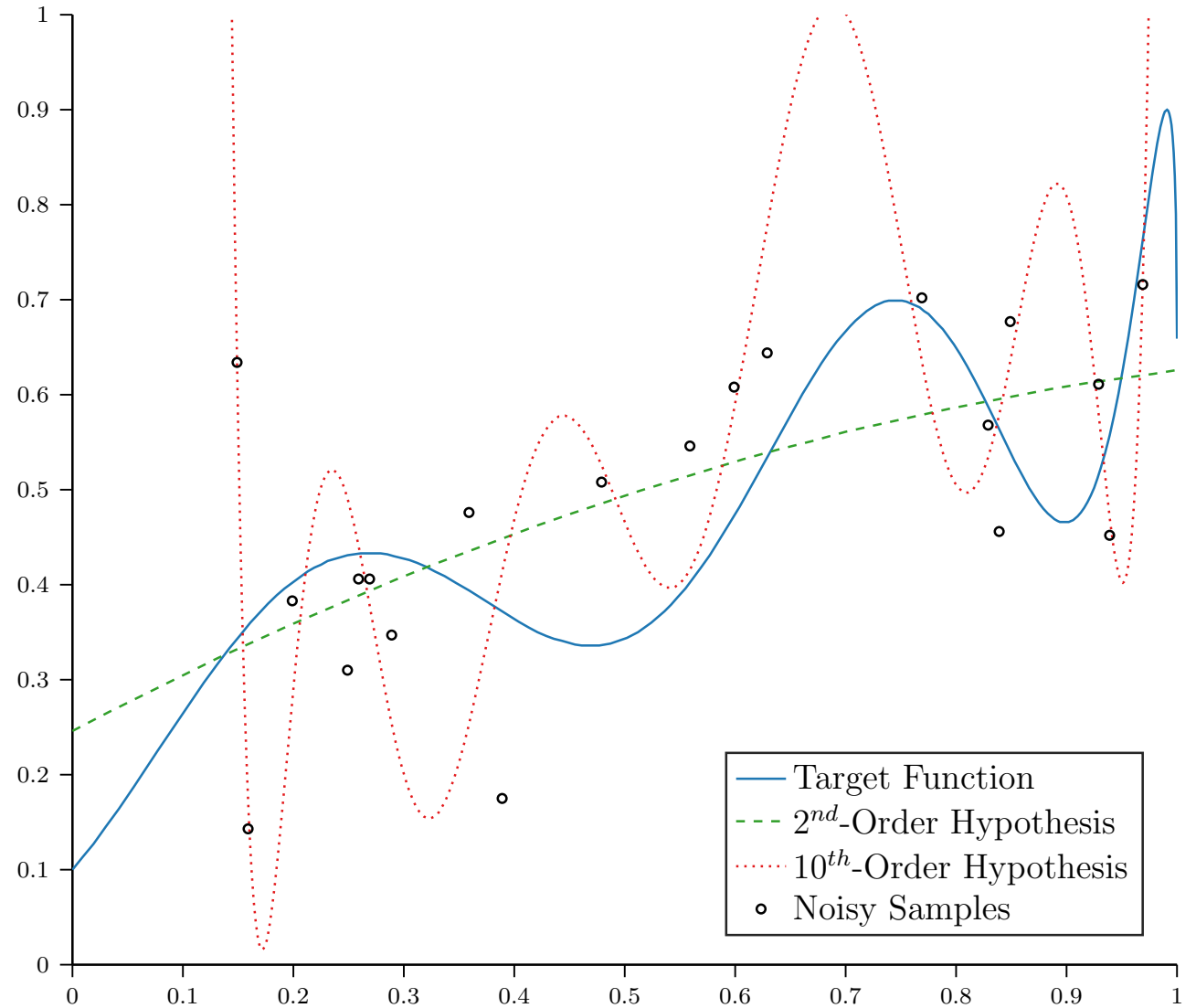
# Noisy Targets

- 10-dimensional target function with additive Gaussian noise
- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomial
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomial



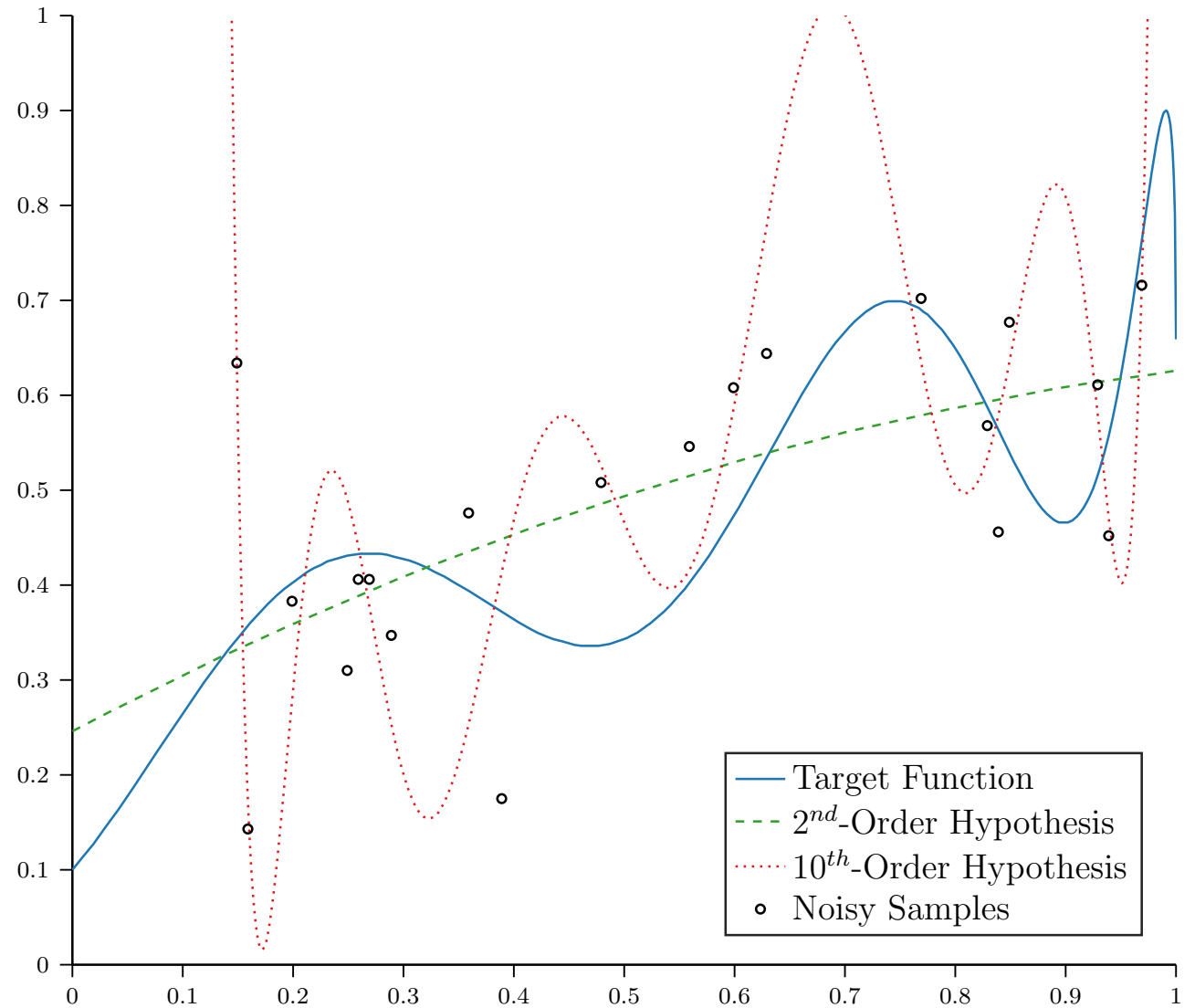
# Noisy Targets

- 10-dimensional target function with additive Gaussian noise
- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomial
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomial



# Noisy Targets

	$\mathcal{H}_2$	$\mathcal{H}_{10}$
Training Error	0.016	0.011
True Error	0.009	3797



# Feature Transforms: Experiment

- $x \in \mathbb{R}, y \in \mathbb{R}$  and  $N = 100$
- Targets are generated by a 10<sup>th</sup>-order polynomial in  $x$  with additive Gaussian noise:

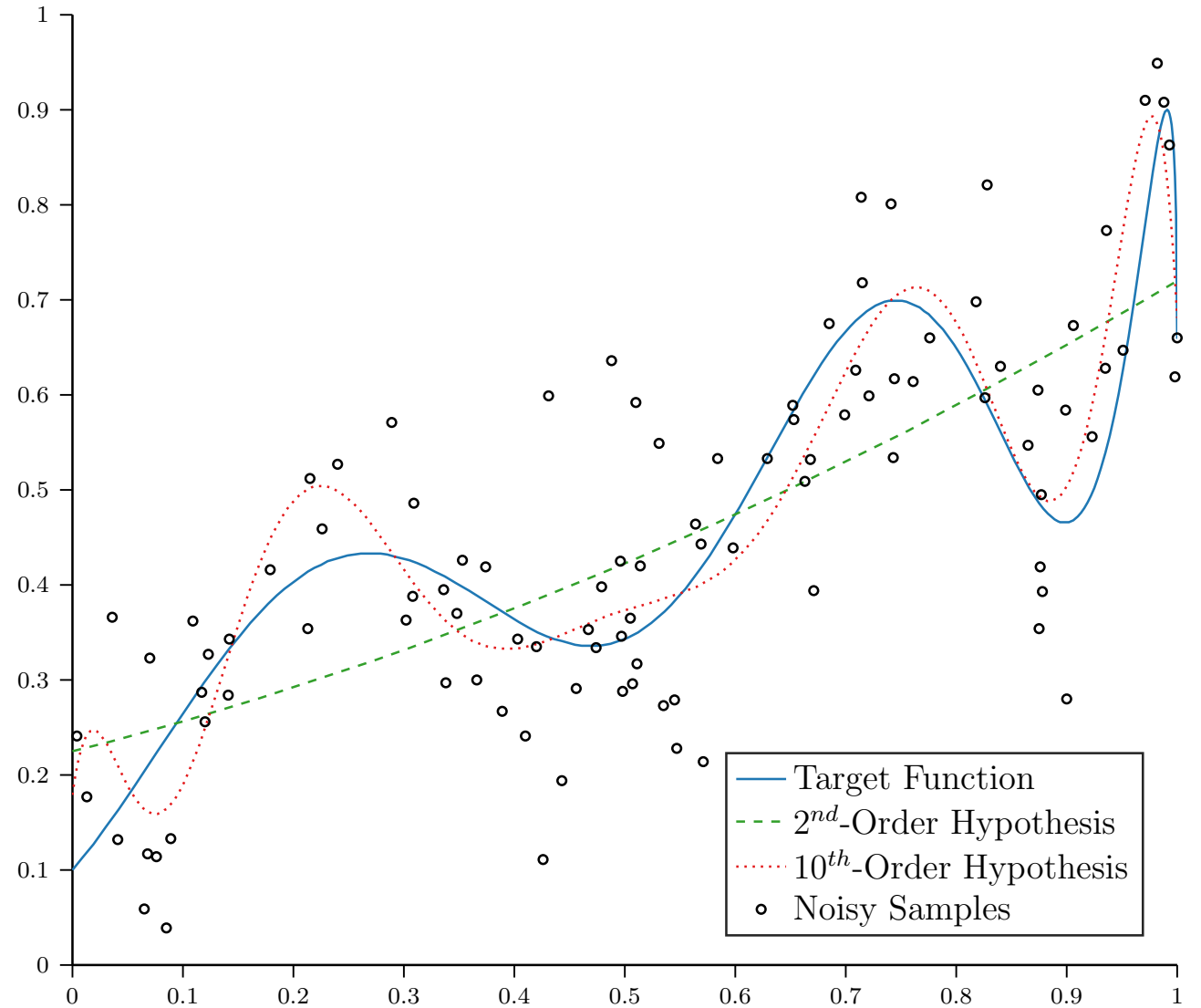
$$y = \sum_{d=0}^{10} a_d x^d + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2)$$

- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomials
  - $\phi_{1,2}(x) = [x, x^2]$
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomials
  - $\phi_{1,10}(x) = [x, x^2, x^3, x^4, x^5, x^6, x^7, x^8, x^9, x^{10}]$



# Noisy Targets

	$\mathcal{H}_2$	$\mathcal{H}_{10}$
Training Error	0.018	0.010
True Error	0.009	0.003



# Regularization

- Constrain models to prevent them from overfitting
- Learning algorithms are optimization problems and regularization imposes constraints on the optimization

# Hard Constraints

- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomials
  - $\phi_{1,10}(x) = [x, x^2, x^3, x^4, x^5, x^6, x^7, x^8, x^9, x^{10}]$

- Given  $X = \begin{bmatrix} 1 & \phi_{1,10}(x^{(1)}) \\ 1 & \phi_{1,10}(x^{(2)}) \\ \vdots & \vdots \\ 1 & \phi_{1,10}(x^{(N)}) \end{bmatrix}$  and  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$  find

$\boldsymbol{\omega} = [\omega_0, \omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6, \omega_7, \omega_8, \omega_9, \omega_{10}]$   
that minimizes

$$\frac{1}{N} (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$$

- Subject to

$$\omega_3 = \omega_4 = \omega_5 = \omega_6 = \omega_7 = \omega_8 = \omega_9 = \omega_{10} = 0$$

# Hard Constraints

- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomials
  - $\phi_{1,10}(x) = [x, x^2, x^3, x^4, x^5, x^6, x^7, x^8, x^9, x^{10}]$

- Given  $X = \begin{bmatrix} 1 & \phi_{1,10}(x^{(1)}) \\ 1 & \phi_{1,10}(x^{(2)}) \\ \vdots & \vdots \\ 1 & \phi_{1,10}(x^{(N)}) \end{bmatrix}$  and  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$  find

$\boldsymbol{\omega} = [\omega_0, \omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6, \omega_7, \omega_8, \omega_9, \omega_{10}]$   
that minimizes

$$\frac{1}{N} \sum_{n=1}^N \left( \left( \sum_{d=0}^{10} x_d^{(n)} \omega_d \right) - y^{(n)} \right)^2$$

- Subject to

$$\omega_3 = \omega_4 = \omega_5 = \omega_6 = \omega_7 = \omega_8 = \omega_9 = \omega_{10} = 0$$

# Hard Constraints

- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomials
  - $\phi_{1,10}(x) = [x, x^2, x^3, x^4, x^5, x^6, x^7, x^8, x^9, x^{10}]$

- Given  $X = \begin{bmatrix} 1 & \phi_{1,10}(x^{(1)}) \\ 1 & \phi_{1,10}(x^{(2)}) \\ \vdots & \vdots \\ 1 & \phi_{1,10}(x^{(N)}) \end{bmatrix}$  and  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$  find

$\boldsymbol{\omega} = [\omega_0, \omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6, \omega_7, \omega_8, \omega_9, \omega_{10}]$   
that minimizes

$$\frac{1}{N} \sum_{n=1}^N \left( \left( \sum_{d=0}^2 x_d^{(n)} \omega_d \right) - y^{(n)} \right)^2$$

- Subject to nothing!

# Hard Constraints

- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomials

- $\phi_{1,2}(x) = [x, x^2]$

- Given  $X = \begin{bmatrix} 1 & \phi_{1,2}(x^{(1)}) \\ 1 & \phi_{1,2}(x^{(2)}) \\ \vdots & \vdots \\ 1 & \phi_{1,2}(x^{(N)}) \end{bmatrix}$  and  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$  find

$$\boldsymbol{\omega} = [\omega_0, \omega_1, \omega_2]$$

that minimizes

$$\frac{1}{N} (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$$

- Subject to nothing!

# Soft Constraints

- More generally,  $\phi$  can be any nonlinear transformation, e.g., exp, log, sin, sqrt, etc...

- Given  $X = \begin{bmatrix} 1 & \phi_1(\mathbf{x}^{(1)}) & \cdots & \phi_m(\mathbf{x}^{(1)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}^{(N)}) & \cdots & \phi_m(\mathbf{x}^{(N)}) \end{bmatrix}$  and  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$ ,

find  $\boldsymbol{\omega}$  that minimizes

$$\frac{1}{N} (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$$

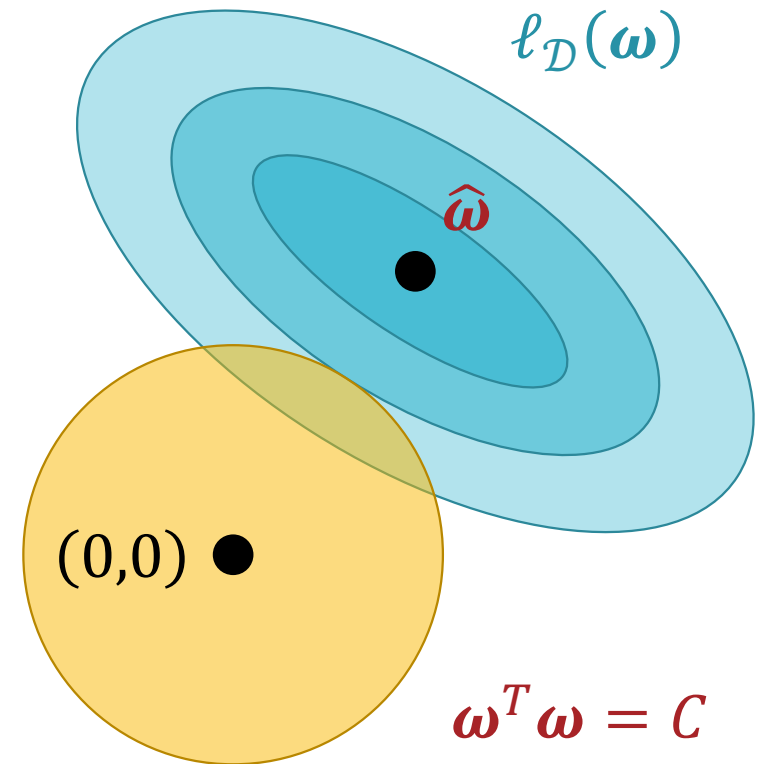
- Subject to:

$$\|\boldsymbol{\omega}\|_2^2 = \boldsymbol{\omega}^T \boldsymbol{\omega} = \sum_{d=0}^D \omega_d^2 \leq C$$

# Soft Constraints

minimize  $\ell_{\mathcal{D}}(\boldsymbol{\omega}) = (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$

subject to  $\boldsymbol{\omega}^T \boldsymbol{\omega} \leq C$

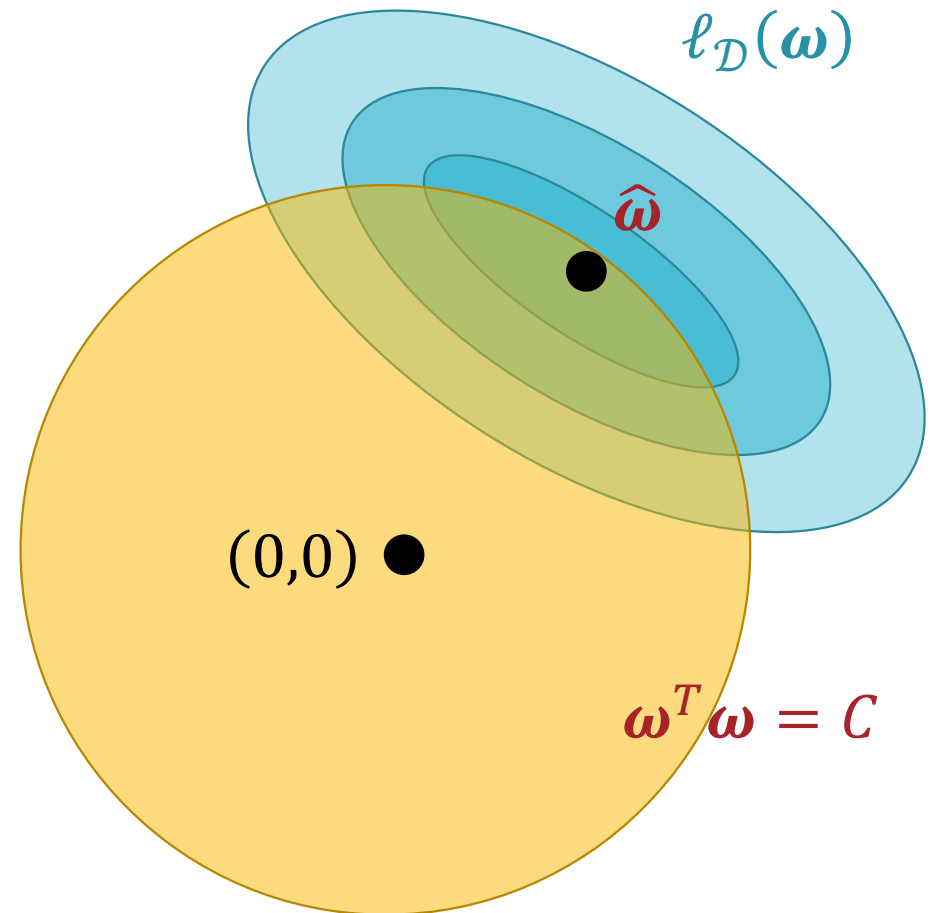




# Soft Constraints

minimize  $\ell_{\mathcal{D}}(\boldsymbol{\omega}) = (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$

subject to  $\boldsymbol{\omega}^T \boldsymbol{\omega} \leq C$



# Soft Constraints

minimize  $\ell_{\mathcal{D}}(\boldsymbol{\omega}) = (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$

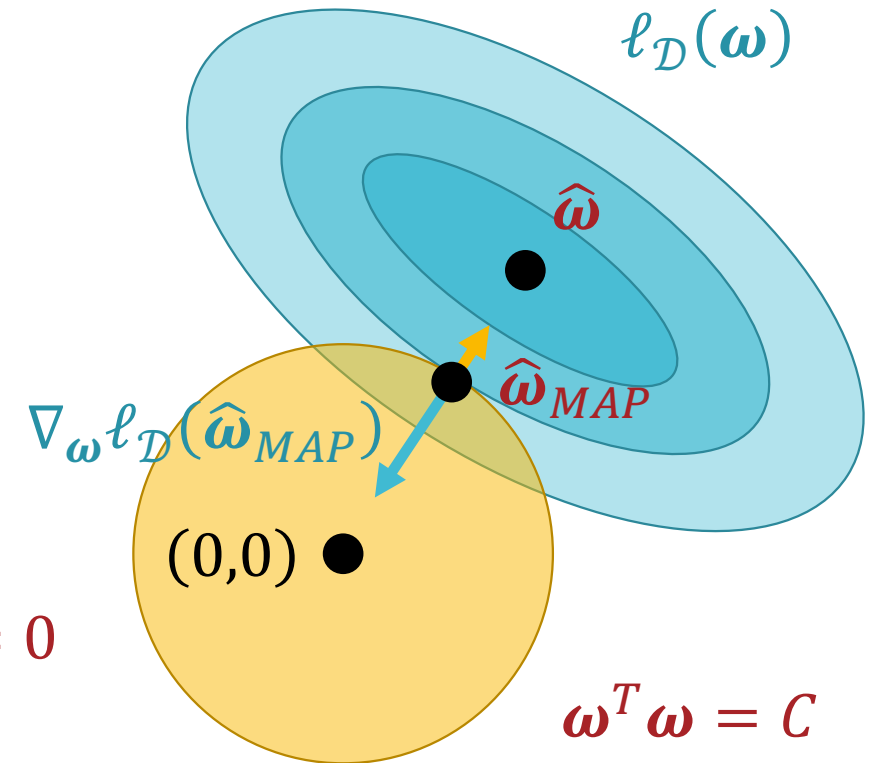
subject to  $\boldsymbol{\omega}^T \boldsymbol{\omega} \leq C$

$$\nabla_{\boldsymbol{\omega}} \ell_{\mathcal{D}}(\hat{\boldsymbol{\omega}}_{MAP}) \propto -2\hat{\boldsymbol{\omega}}_{MAP}$$

$$\nabla_{\boldsymbol{\omega}} \ell_{\mathcal{D}}(\hat{\boldsymbol{\omega}}_{MAP}) = -2\lambda_C \hat{\boldsymbol{\omega}}_{MAP}$$

$$\nabla_{\boldsymbol{\omega}} \ell_{\mathcal{D}}(\hat{\boldsymbol{\omega}}_{MAP}) + 2\lambda_C \hat{\boldsymbol{\omega}}_{MAP} = 0$$

$$\nabla_{\boldsymbol{\omega}} (\ell_{\mathcal{D}}(\hat{\boldsymbol{\omega}}_{MAP}) + \lambda_C (\hat{\boldsymbol{\omega}}_{MAP})^T \hat{\boldsymbol{\omega}}_{MAP}) = 0$$



Soft  
Constraints:  
Solving for  $\hat{\omega}_{MAP}$

$$\text{minimize } \ell_{\mathcal{D}}(\omega) = (X\omega - \mathbf{y})^T (X\omega - \mathbf{y})$$

$$\text{subject to } \omega^T \omega \leq C$$



$$\text{minimize } \ell_{\mathcal{D}}^{AUG}(\omega) = \ell_{\mathcal{D}}(\omega) + \lambda_C \omega^T \omega$$

# Ridge Regression

$$\text{minimize } \ell_D^{AUG}(\boldsymbol{\omega}) = \ell_D(\boldsymbol{\omega}) + \lambda_C \boldsymbol{\omega}^T \boldsymbol{\omega}$$

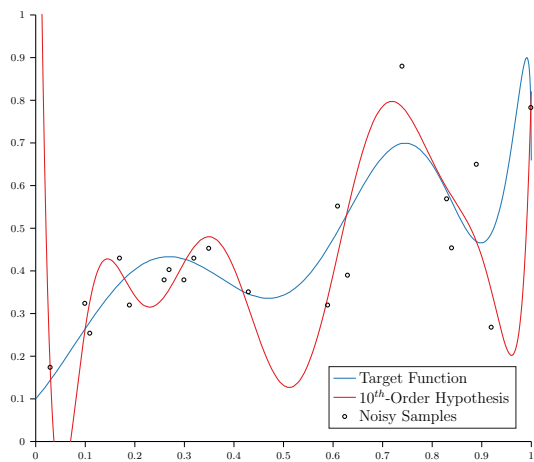
$$\nabla_{\boldsymbol{\omega}} \ell_D^{AUG}(\boldsymbol{\omega}) = 2(X^T X \boldsymbol{\omega} - X^T \mathbf{y} + \lambda_C \boldsymbol{\omega})$$

$$2(X^T X \hat{\boldsymbol{\omega}}_{MAP} - X^T \mathbf{y} + \lambda_C \hat{\boldsymbol{\omega}}_{MAP}) = 0$$

$$(X^T X + \lambda_C I_{D+1}) \hat{\boldsymbol{\omega}}_{MAP} = X^T \mathbf{y}$$

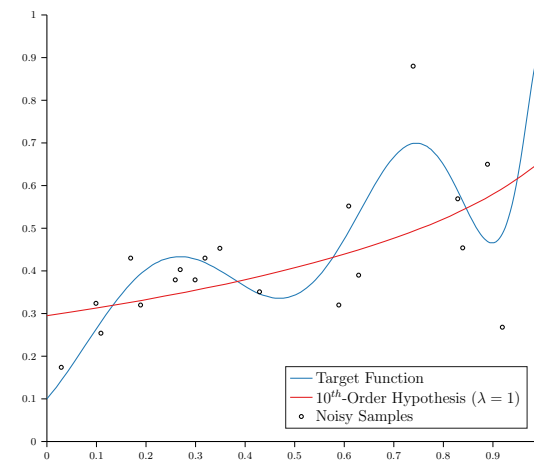
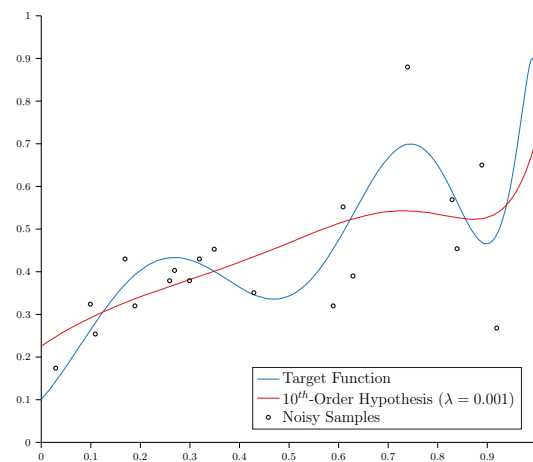
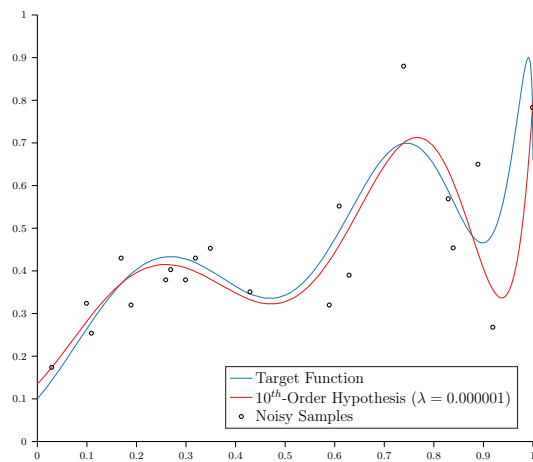
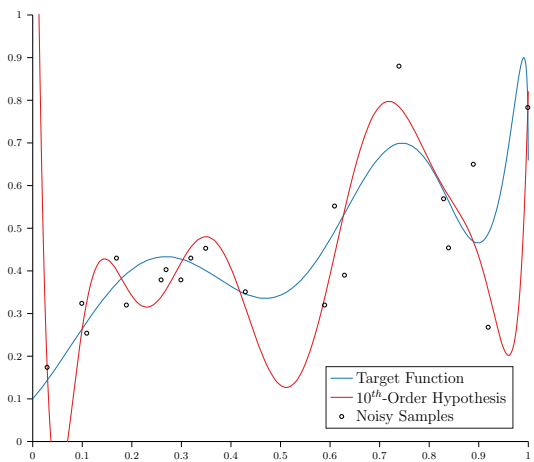
$$\hat{\boldsymbol{\omega}}_{MAP} = \underbrace{(X^T X + \lambda_C I_{D+1})^{-1}} X^T \mathbf{y}$$

Adding this positive ( $\lambda_C \geq 0$ ) diagonal matrix can help if  $X^T X$  is not invertible!



# Ridge Regression

- 10-dimensional target function with additive Gaussian noise
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomial



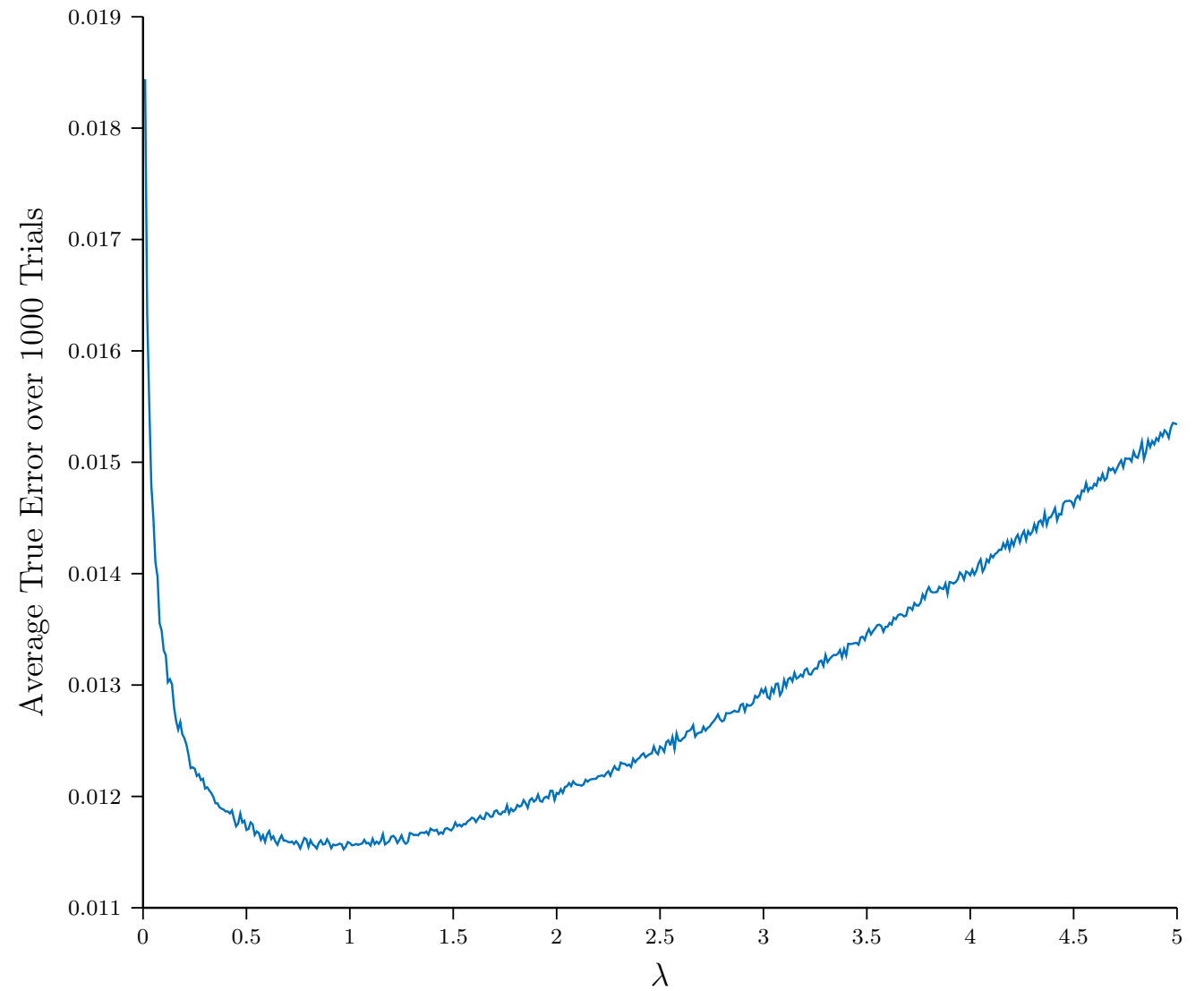
# Ridge Regression

$$\lambda_c = 0$$

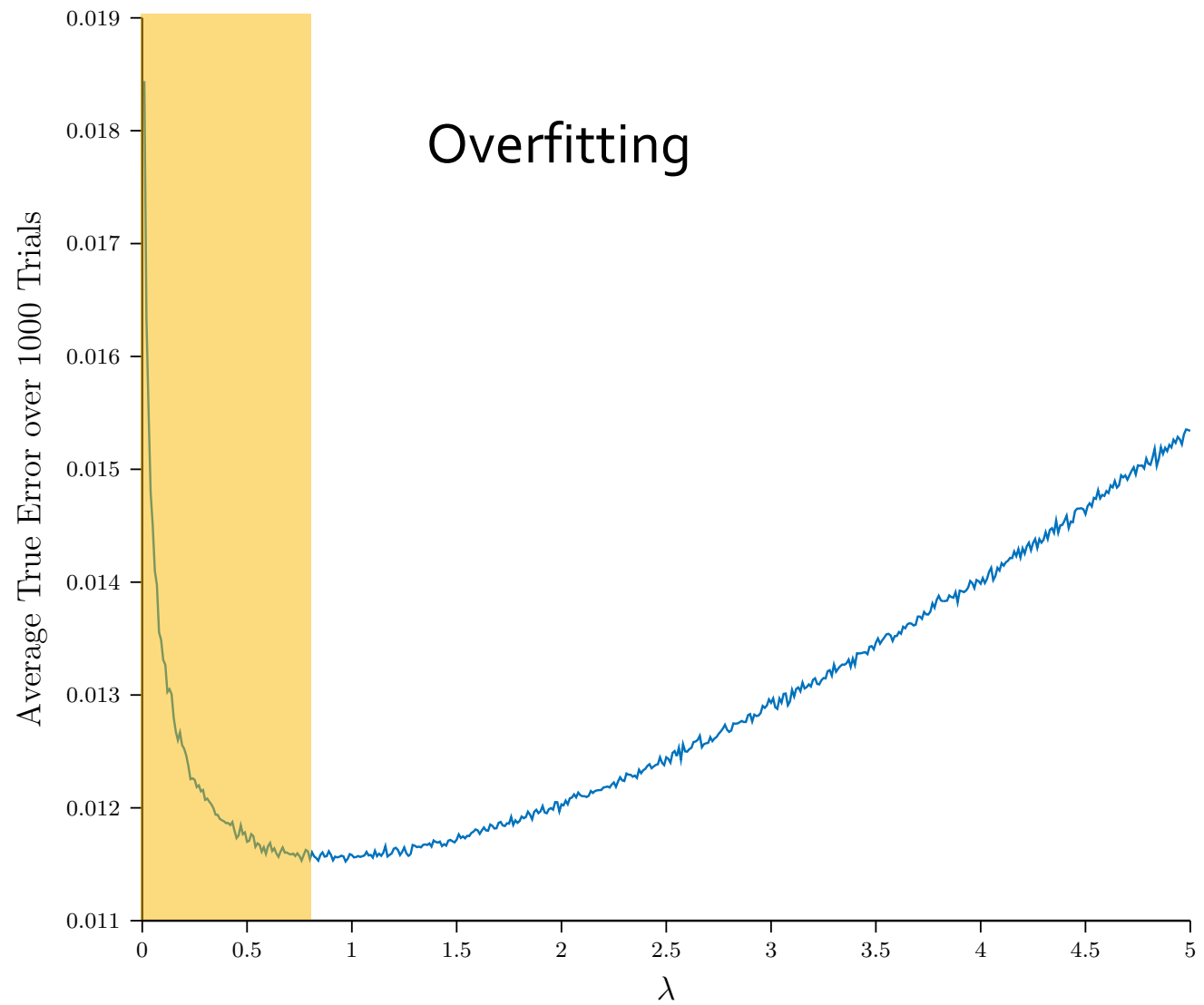
True Error  
0.059

Overfit

# Setting $\lambda$

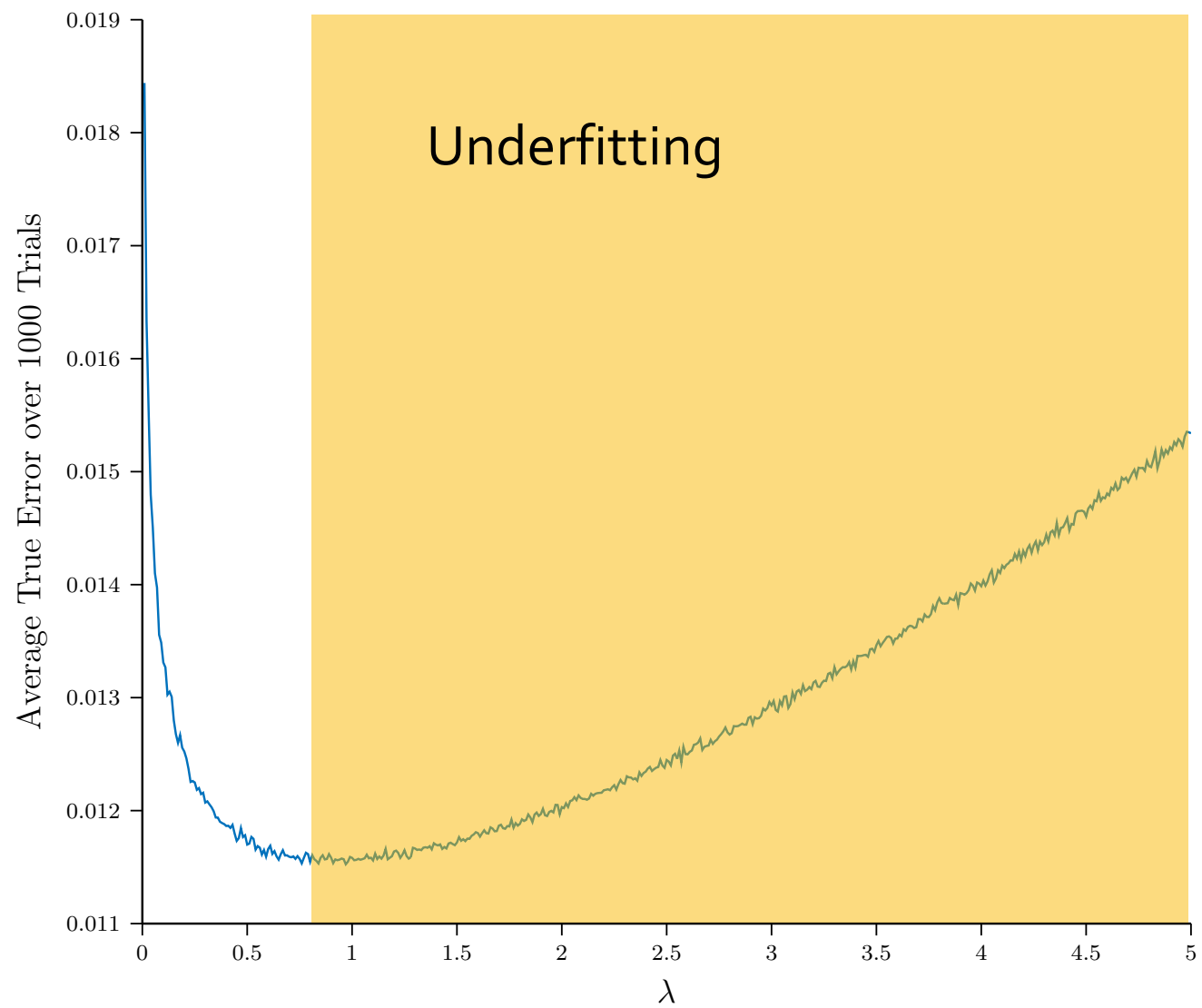


# Setting $\lambda$

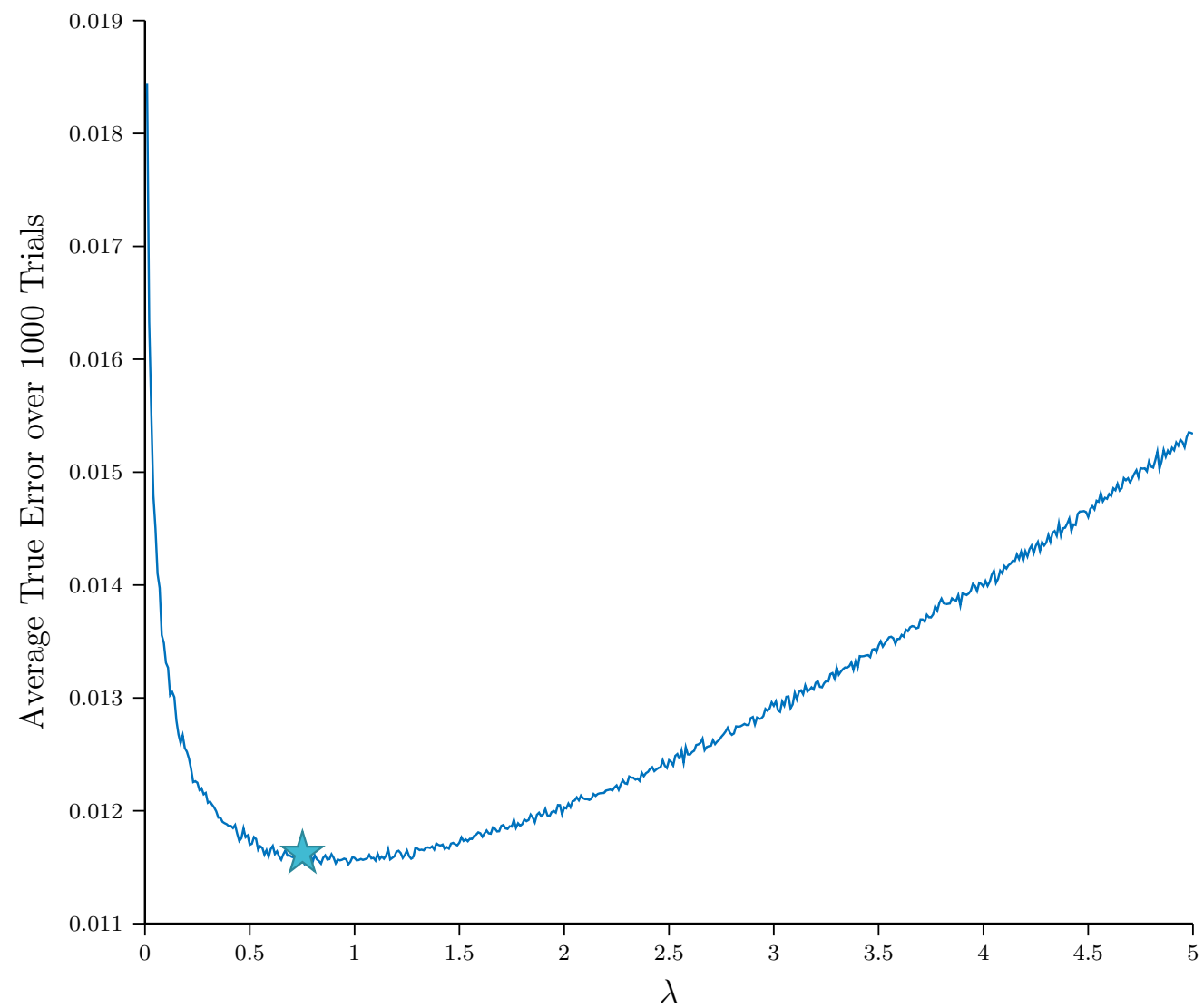




# Setting $\lambda$

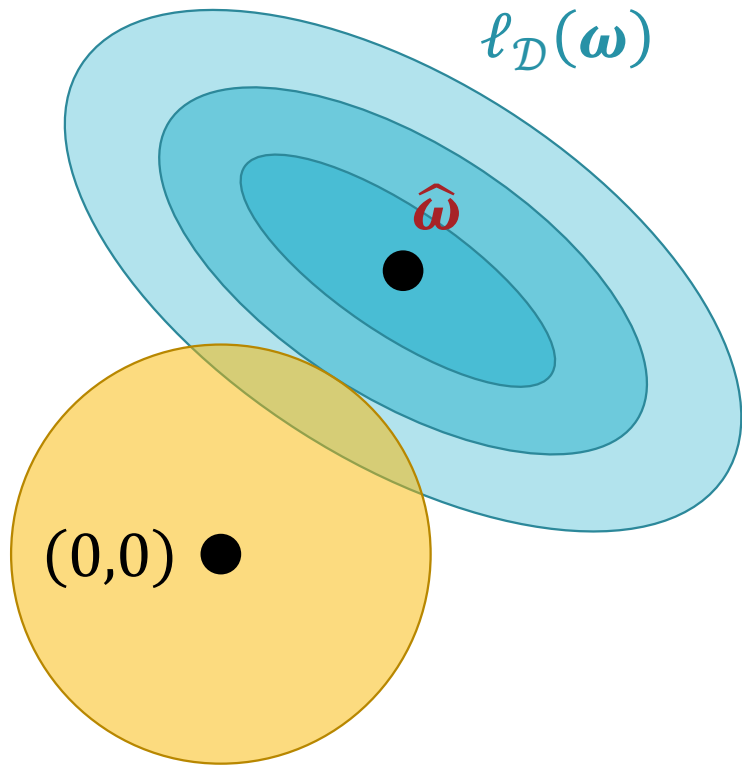


# Setting $\lambda$

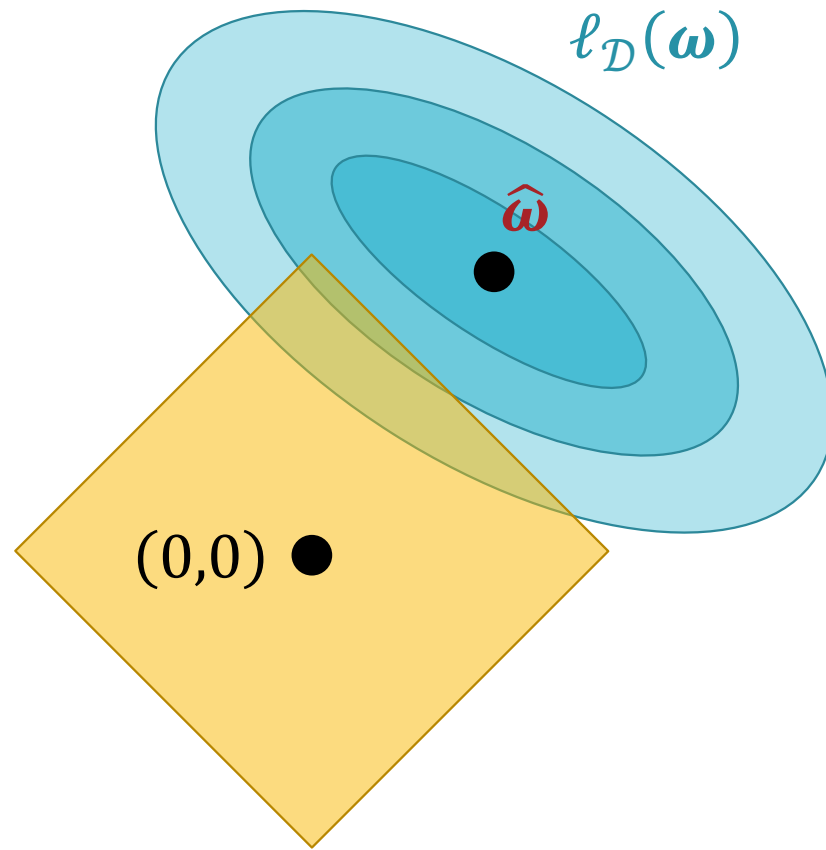


## Other Regularizers

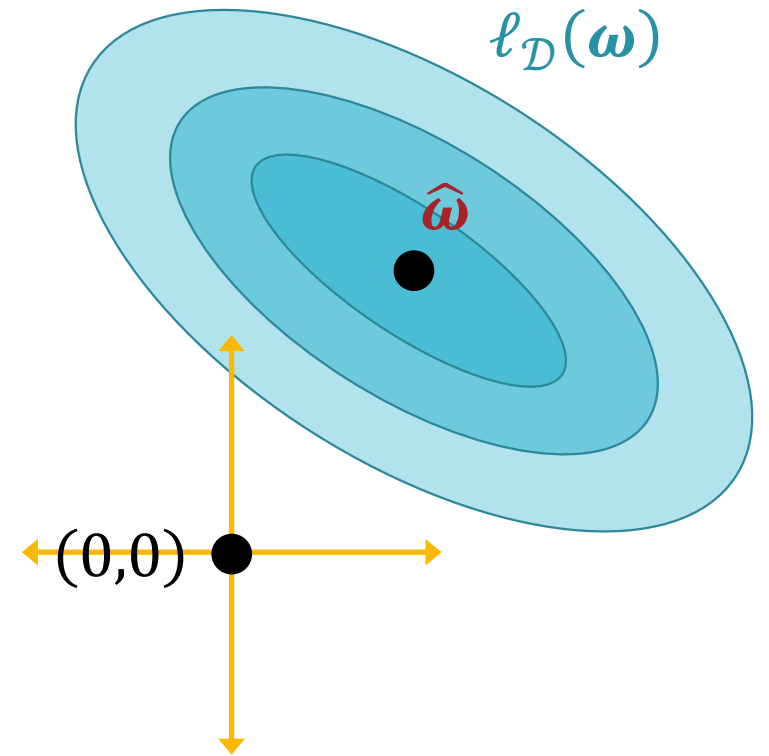
$\ell_{\mathcal{D}}(\boldsymbol{\omega}) + \lambda r(\boldsymbol{\omega})$		
Ridge or $L2$	$r(\boldsymbol{\omega}) = \ \boldsymbol{\omega}\ _2^2 = \sum_{d=0}^D \omega_d^2$	Encourages small weights
Lasso or $L1$	$r(\boldsymbol{\omega}) = \ \boldsymbol{\omega}\ _1 = \sum_{d=0}^D  \omega_d $	Encourages sparsity
$L0$	$r(\boldsymbol{\omega}) = \ \boldsymbol{\omega}\ _0 = \sum_{d=0}^D \mathbb{1}(\omega_d \neq 0)$	Encourages sparsity (intractable)



Ridge or  $L_2$



Lasso or  $L_1$



$L_0$

## Other Regularizers

# M(C)LE for Linear Regression

- If we assume a linear model with additive Gaussian noise

$$y = \boldsymbol{\omega}^T \boldsymbol{x} + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2) \rightarrow y \sim N(\boldsymbol{\omega}^T \boldsymbol{x}, \sigma^2)$$

- Then given  $X = \begin{bmatrix} 1 & \boldsymbol{x}^{(1)} \\ 1 & \boldsymbol{x}^{(2)} \\ \vdots & \vdots \\ 1 & \boldsymbol{x}^{(N)} \end{bmatrix}$  and  $\boldsymbol{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$  the MLE of  $\boldsymbol{\omega}$  is

$$\hat{\boldsymbol{\omega}} = \underset{\boldsymbol{\omega}}{\operatorname{argmax}} \log P(\boldsymbol{y}|X, \boldsymbol{\omega})$$

$$\vdots$$

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

# MAP for Linear Regression

- If we assume a linear model with additive Gaussian noise

$$y = \boldsymbol{\omega}^T \mathbf{x} + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2) \rightarrow y \sim N(\boldsymbol{\omega}^T \mathbf{x}, \sigma^2)$$

and independent Gaussian priors on all the weights...

$$\omega_d \sim N\left(0, \frac{\sigma^2}{\lambda}\right)$$

- ... then, the MAP of  $\boldsymbol{\omega}$  is the ridge regression solution!

$$\hat{\boldsymbol{\omega}} = \underset{\boldsymbol{\omega}}{\operatorname{argmax}} \log P(\boldsymbol{\omega}|X, \mathbf{y}) = \underset{\boldsymbol{\omega}}{\operatorname{argmax}} \log P(\mathbf{y}|X, \boldsymbol{\omega})P(\boldsymbol{\omega})$$

⋮

$$= (X^T X + \lambda_C I_{D+1})^{-1} X^T \mathbf{y}$$

# MAP for Linear Regression

- If we assume a linear model with additive Gaussian noise  $y = \boldsymbol{\omega}^T \mathbf{x} + \epsilon$  where  $\epsilon \sim N(0, \sigma^2) \rightarrow y \sim N(\boldsymbol{\omega}^T \mathbf{x}, \sigma^2)$  and independent *Laplace* priors on all the weights...

$$\omega_d \sim \text{Laplace} \left( 0, \frac{2\sigma^2}{\lambda} \right)$$

- ... then, the MAP of  $\boldsymbol{\omega}$  is the Lasso regression solution!
- No closed form solution exists but we can solve via (sub-)gradient descent

# Nonlinear Transforms

- Decide on some transformation  $\Phi: \mathcal{X} \rightarrow \mathcal{Z}$
- Given  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ , learn a hypothesis,  $\tilde{h}(\mathbf{z})$ ,  
using  $\tilde{\mathcal{D}} = \{(\mathbf{z}^{(n)} = \Phi(\mathbf{x}^{(n)}), y^{(n)})\}_{n=1}^N$
- Return the corresponding predictor in the original space:  
 $h(\mathbf{x}) = \tilde{h}(\Phi(\mathbf{x}))$



# Efficiency

- Depending on the transformation  $\Phi$  and the dimensionality of the original input space  $\mathcal{X}$ , computing  $\Phi(\mathbf{x})$  can be prohibitively computationally expensive
  - Computing  $\Phi_2(\mathbf{x}) = [x_1, x_2, \dots, x_D, x_1^2, x_1x_2, \dots, x_D^2]$  for  $\mathbf{x} \in \mathbb{R}^D$  requires  $D + \binom{D}{2} + D = \frac{D^2+3D}{2} = O(D^2)$  time
  - Computing  $\Phi_{10}(\mathbf{x})$  requires  $O(D^{10})$  time
- Tradeoff:
  - High-dimensional transformations can result in good hypotheses (as long as they don't overfit) but...
  - High-dimensional transformations are expensive

# Inner Product Methods

- Insight: the predictions of many machine learning models can be expressed as a function of inner products between feature

vectors i.e., given  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$

$$h_{\mathcal{D}}(\mathbf{x}) = g\left(\{\mathbf{x}^T \mathbf{x}^{(n)}\}_{n=1}^N, \{\mathbf{x}^{(n)T} \mathbf{x}^{(m)}\}_{n,m=1}^N\right)$$

- Crucially, feature vectors **only** appear in inner products with other feature vectors
- Applying a feature transformation  $\Phi$  gives

$$h_{\mathcal{D}}(\Phi(\mathbf{x})) = g\left(\{\Phi(\mathbf{x})^T \Phi(\mathbf{x}^{(n)})\}_{n=1}^N, \{\Phi(\mathbf{x}^{(n)})^T \Phi(\mathbf{x}^{(m)})\}_{n,m=1}^N\right)$$

# The Kernel Trick

- Idea: for inner product methods, instead of computing  $\Phi(\mathbf{x})$ , use some function  $K_\Phi$  s.t.  $K_\Phi(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$ 
  - $K_\Phi(\mathbf{x}, \mathbf{x}')$  should be cheaper to compute than  $\Phi(\mathbf{x})$

- Example:  $\Phi'_2(\mathbf{x}) = [x_1, \dots, x_D, x_1^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{D-1}x_D, x_D^2]$

$$\begin{aligned}\Phi'_2(\mathbf{x})^T \Phi'_2(\mathbf{x}') &= \sum_{i=1}^D x_i x'_i + \sum_{i=1}^D x_i^2 x'^2_i + \sum_{i=1}^D \sum_{j \neq i} 2x_i x'_i x_j x'_j \\ &= \sum_{i=1}^D x_i x'_i + \left( \sum_{i=1}^D x_i x'_i \right)^2 = \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2\end{aligned}$$

$$K_{\Phi'_2}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$$

- Computing  $\Phi'_2(\mathbf{x})^T \Phi'_2(\mathbf{x}')$  requires  $O(D^2)$  time whereas computing  $K_{\Phi'_2}(\mathbf{x}, \mathbf{x}')$  only takes  $O(D)$ !

# Common Kernels

- $K_{\Phi'_2}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$

- Implied feature transformation:

$$\Phi'_2(\mathbf{x}) = [x_1, \dots, x_D, x_1^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{D-1}x_D, x_D^2]$$

- Implied dimensionality:  $\frac{D^2+3D}{2}$

- $K_{\Phi_2^{(\gamma)}}(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^2 - 1$

- Implied feature transformation:

$$\Phi_2^{(\gamma)}(\mathbf{x}) = [\sqrt{2\gamma}x_1, \dots, \sqrt{2\gamma}x_D, \gamma x_1^2, \gamma x_1x_2, \dots, \gamma x_D^2]$$

- $\gamma$  affects the geometry of the transform

- Implied dimensionality:  $\frac{D^2+3D}{2}$

# Common Kernels

- Polynomial Kernel:  $K_{\Phi_Q^{(\gamma)}}(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^Q - 1$ 
  - Implied dimensionality:  $O(D^Q)$
  - $\gamma$  affects the geometry of the transform

- Gaussian-RBF Kernel:  $K_{\Phi_r}(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2r}}$

- Implied feature transformation:  $\Phi_r(\mathbf{x}) = \left[ e^{-\frac{x_1^2}{2r}}, \dots, e^{-\frac{x_D^2}{2r}}, \right.$

$$\left. e^{-\frac{x_1^2}{2r}} \sqrt{\frac{(x_1)^2}{1!r^1}}, \dots, e^{-\frac{x_D^2}{2r}} \sqrt{\frac{(x_D)^2}{1!r^1}}, e^{-\frac{x_1^2}{2r}} \sqrt{\frac{(x_1^2)^2}{2!r^2}}, \dots, e^{-\frac{x_D^2}{2r}} \sqrt{\frac{(x_D^2)^2}{2!r^2}}, \dots \right]$$

# Common Kernels

- Polynomial Kernel:  $K_{\Phi_Q^{(\gamma)}}(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^Q - 1$ 
  - Implied dimensionality:  $O(D^Q)$
  - $\gamma$  affects the geometry of the transform
- Gaussian-RBF Kernel:  $K_{\Phi_r}(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2r}}$ 
  - Implied feature transformation:  $\Phi_r(\mathbf{x}) = \left[ \left[ e^{-\frac{x_1^2}{2r}} \sqrt{\frac{(x_1^d)^2}{d!r^d}}, \dots, e^{-\frac{x_D^2}{2r}} \sqrt{\frac{(x_1^d)^2}{d!r^d}} \right] : d \in \mathbb{N} \right]$
  - Implied dimensionality:  $\infty!$

# Kernels Everywhere!

- Any method that only depends on the Euclidean distance between data points is an inner product method:

$$\|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{(\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}')} = \sqrt{\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{x}' + \mathbf{x}'^T \mathbf{x}'}$$

- We can kernelize  $k$ NN!

# Kernels Everywhere!

- We can also kernelize linear/ridge regression!

$$\begin{aligned}\hat{\boldsymbol{\omega}}_{MAP} &= (X^T X + \lambda_C I_{D+1})^{-1} X^T \mathbf{y} \\ &= X^T (X X^T + \lambda_C I_N)^{-1} \mathbf{y}\end{aligned}$$

$$X X^T = \begin{bmatrix} 1 & \mathbf{x}^{(1)T} \\ 1 & \mathbf{x}^{(2)T} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(N)T} \end{bmatrix} \begin{bmatrix} 1 & 1 & \cdots & 1 \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \cdots & \mathbf{x}^{(N)} \end{bmatrix}$$

Let  $\boldsymbol{\alpha} = (X X^T + \lambda_C I_{D+1})^{-1} \mathbf{y} \rightarrow \hat{\boldsymbol{\omega}}_{MAP} = X^T \boldsymbol{\alpha}$

$$\boldsymbol{\alpha} \in \mathbb{R}^N \rightarrow \hat{\boldsymbol{\omega}}_{MAP} = \sum_{i=1}^N \alpha_i \begin{bmatrix} 1 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

$$\rightarrow h(\mathbf{x}) = \hat{\boldsymbol{\omega}}_{MAP}^T \mathbf{x} = \sum_{i=1}^N \alpha_i \begin{bmatrix} 1 & \mathbf{x}^{(i)T} \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$



# Valid Kernels

- Any function  $K$  is a valid kernel if and only if:
  - $\exists$  a transformation  $\Phi$  s.t.

$$K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') \quad \forall \mathbf{x}, \mathbf{x}'$$



- the Gram matrix

$$K = \begin{bmatrix} K(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & \dots & K(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ K(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & K(\mathbf{x}^{(2)}, \mathbf{x}^{(2)}) & \dots & K(\mathbf{x}^{(2)}, \mathbf{x}^{(N)}) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & K(\mathbf{x}^{(N)}, \mathbf{x}^{(2)}) & \dots & K(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix}$$

is symmetric and positive semi-definite  $\forall$  sets

$$\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$$

# Key Takeaways

- Polynomial/non-linear feature transformations allow for learning non-linear functions/decision boundaries
  - Can lead to overfitting...
    - Address with regularization!
    - Analogous to constrained optimization, solve via method of Lagrange multipliers
    - Regularization level is a hyperparameter
  - Can be computationally expensive...
    - Address with kernels!
    - Alternative to explicitly computing feature transformations for inner product methods