

# 10-315: Introduction to Machine Learning

## Lecture 7 – Regularization

Henry Chai

2/9/22

# Linear Regression

- Given *design matrix*

$$\mathbf{A} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} = \begin{bmatrix} 1 & X_1^{(1)} & \cdots & X_1^{(p)} \\ 1 & X_2^{(1)} & \cdots & X_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_n^{(1)} & \cdots & X_n^{(p)} \end{bmatrix} \in \mathbb{R}^{n \times p+1}$$

- and *target vector*

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbb{R}^n$$

- the goal of linear regression is to find

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} J(\beta) = \underset{\beta}{\operatorname{argmin}} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

Poll Review: Is  $J(\beta)$  convex in  $\beta$ ?

$$J(\beta) = \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

- A) Convex, quadratic in  $\beta$
- B) Non-convex,  $\mathbf{A}$  may not be positive semi-definite
- C) Depends on conditioning (ratio of max:min eigenvalues) of  $\mathbf{A}^T \mathbf{A}$
- ☒ D) Convex,  $\mathbf{A}^T \mathbf{A}$  is positive semi-definite

## Minimizing the Mean Squared Error

$$J(\beta) = \frac{1}{n} (A\beta - Y)^T (A\beta - Y)$$

$$J(\beta) = \frac{1}{n} (\beta^T A^T A \beta - 2\beta^T A^T Y + Y^T Y)$$

$$\nabla_{\beta} J(\beta) = \frac{1}{n} (2A^T A \beta - 2A^T Y)$$

$$\frac{1}{n} (2A^T A \hat{\beta} - 2A^T Y) = 0$$

$$A^T A \hat{\beta} - A^T Y = 0$$

$$\hat{\beta} = (A^T A)^{-1} A^T Y$$

ordinary least squares or OLS

# Minimizing the Mean Squared Error

$$\hat{\beta} = (A^T A)^{-1} A^T Y$$

---

1. Is  $A^T A$  invertible?

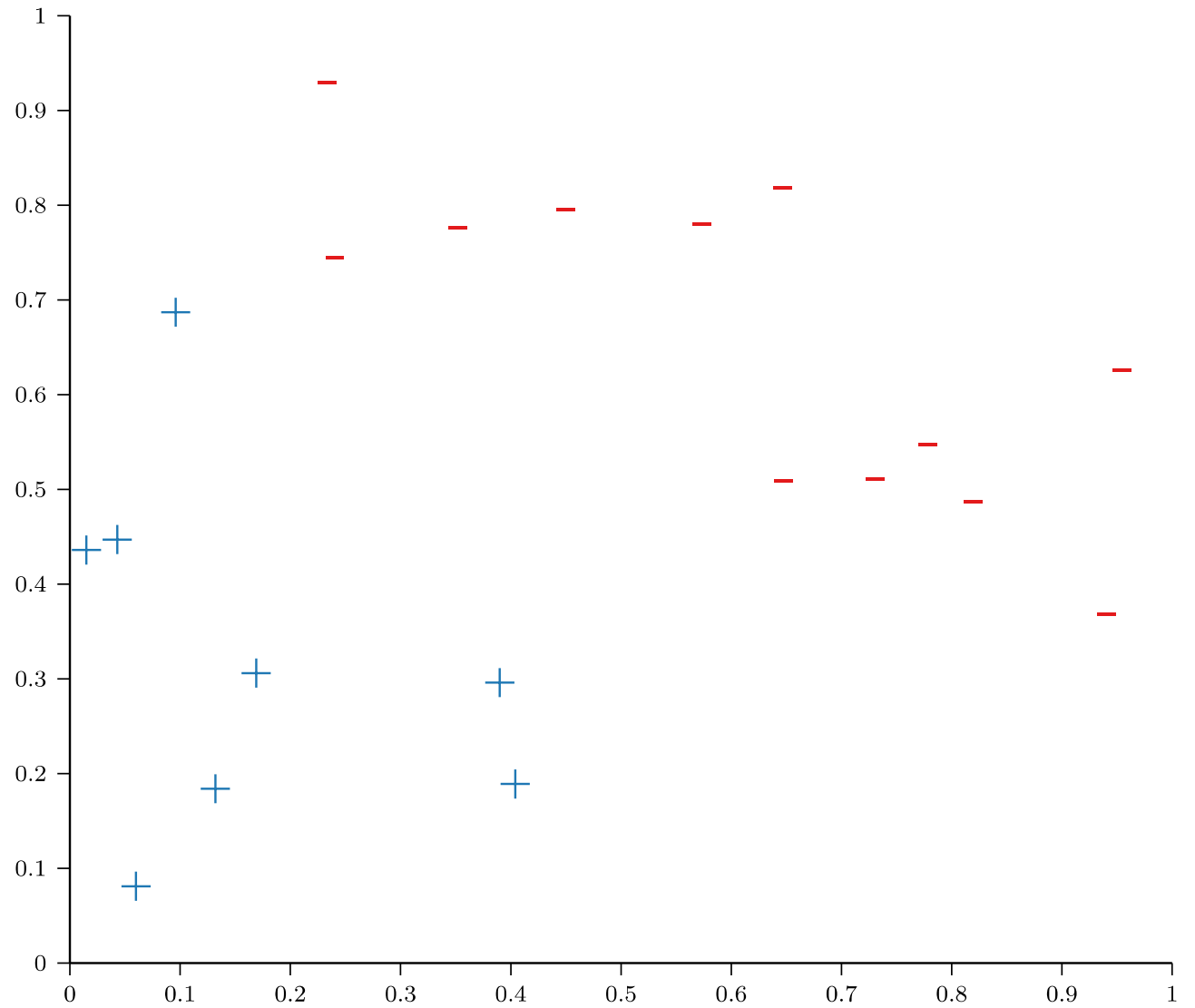
In practice, almost always  
if  $n \gg p$

2. If so, how computationally expensive is inverting  $A^T A$ ?

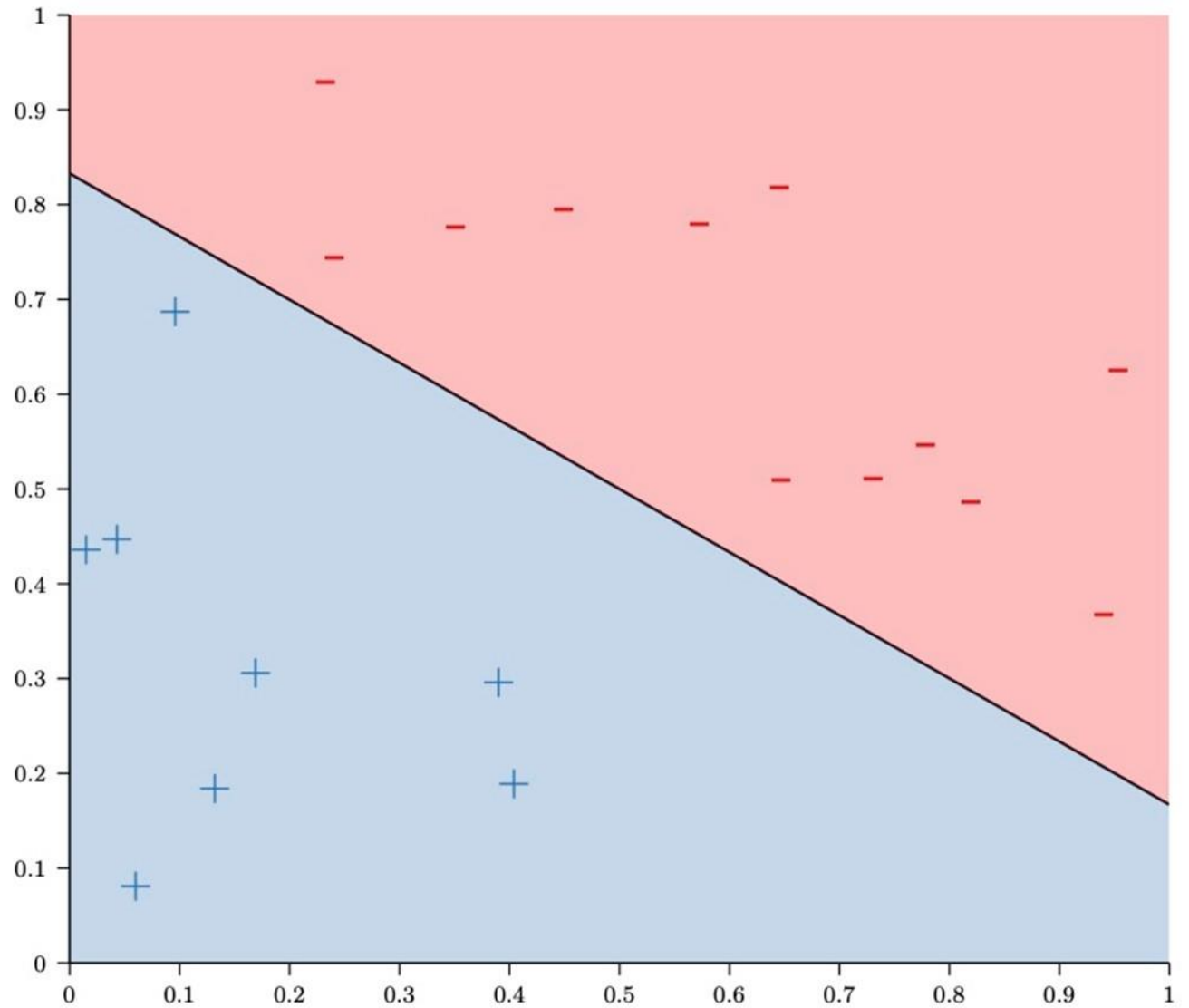
$$A^T A \in \mathbb{R}^{(p+1) \times (p+1)} \quad O(p^3)$$

$$\beta^{t+1} \leftarrow \beta^t - \eta \nabla_{\beta} J(\beta) = \beta^t - \frac{2\eta}{n} (A^T A \beta - A^T Y)$$

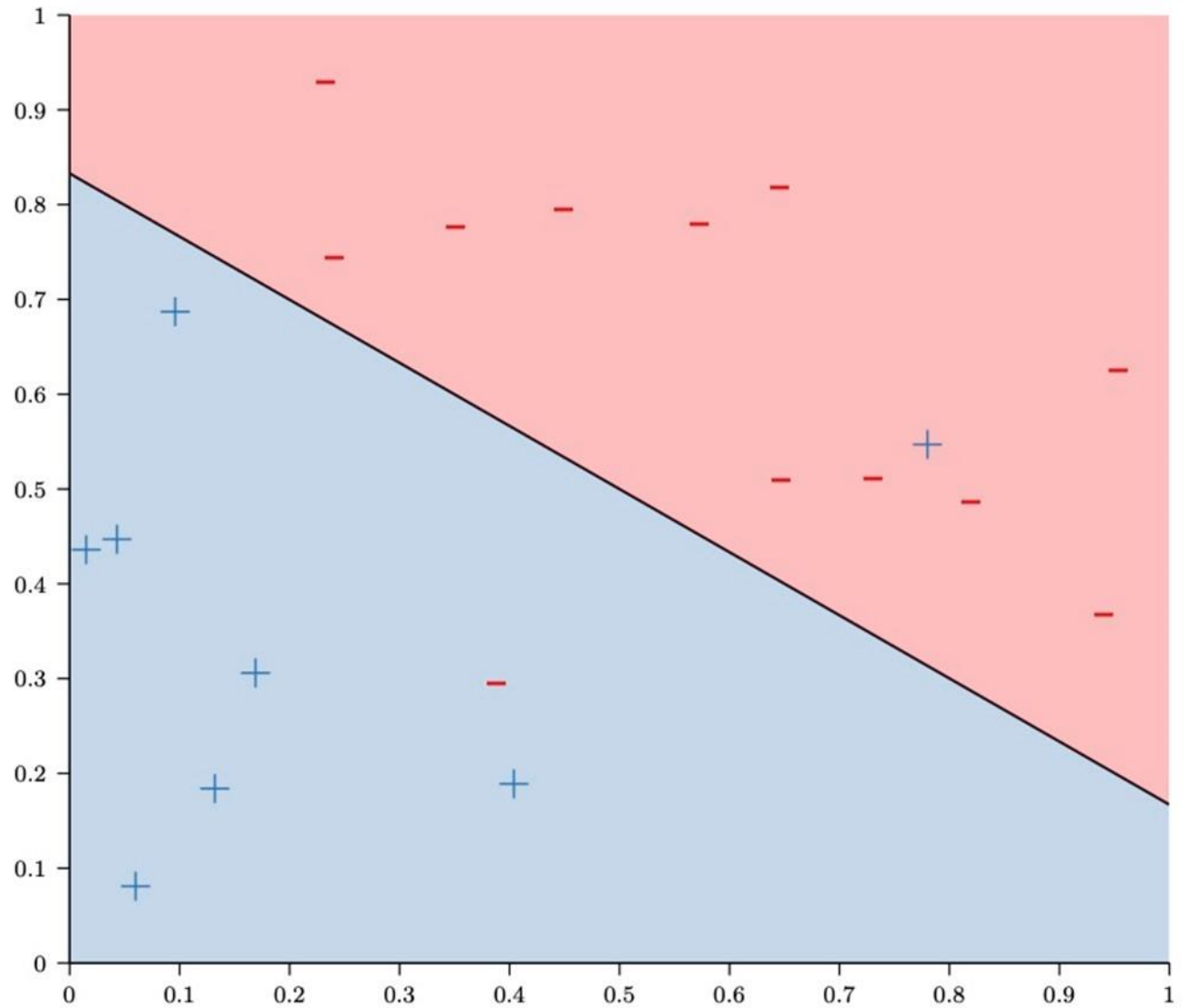
# Linear Models



# Linear Models

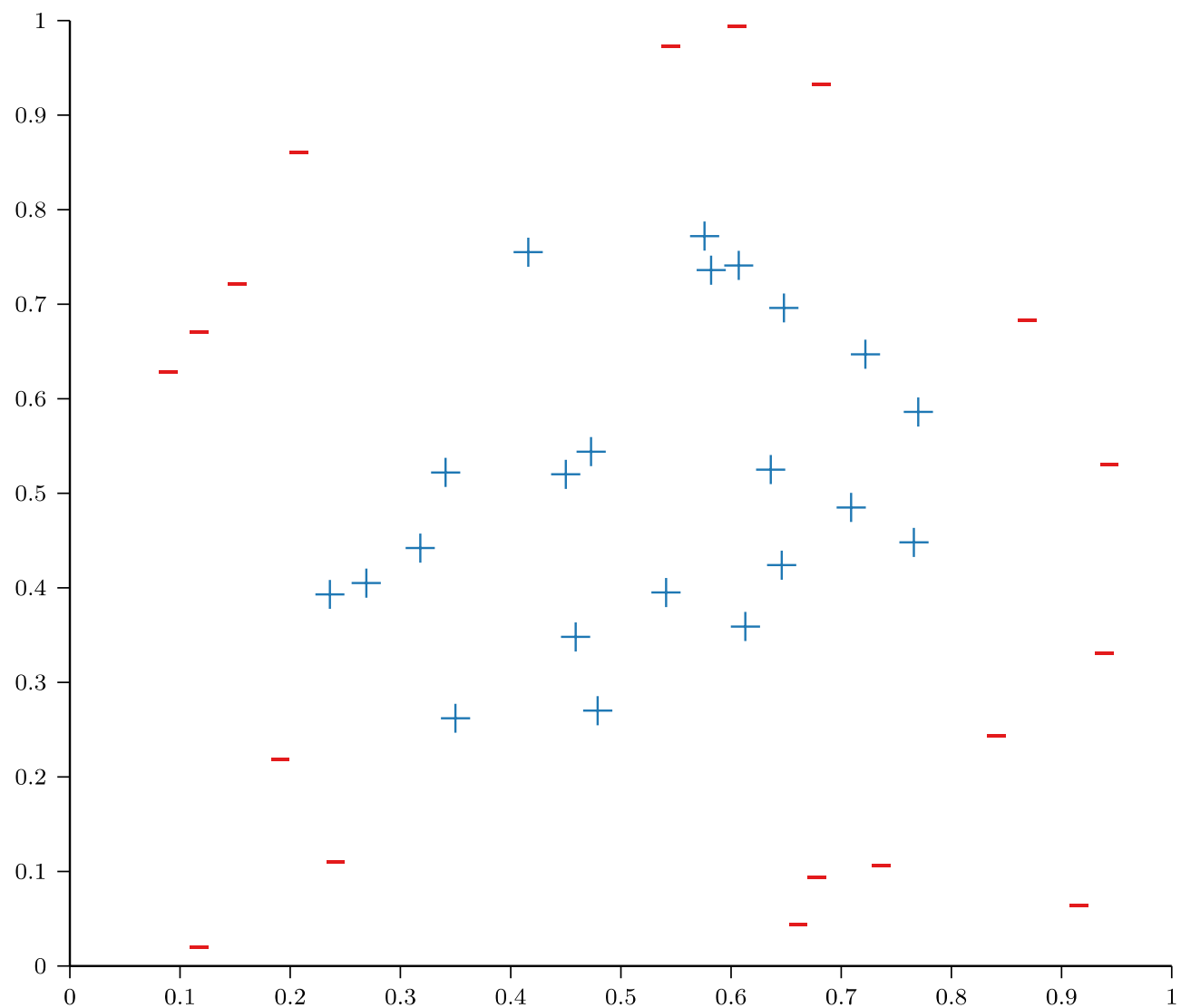


# Linear Models

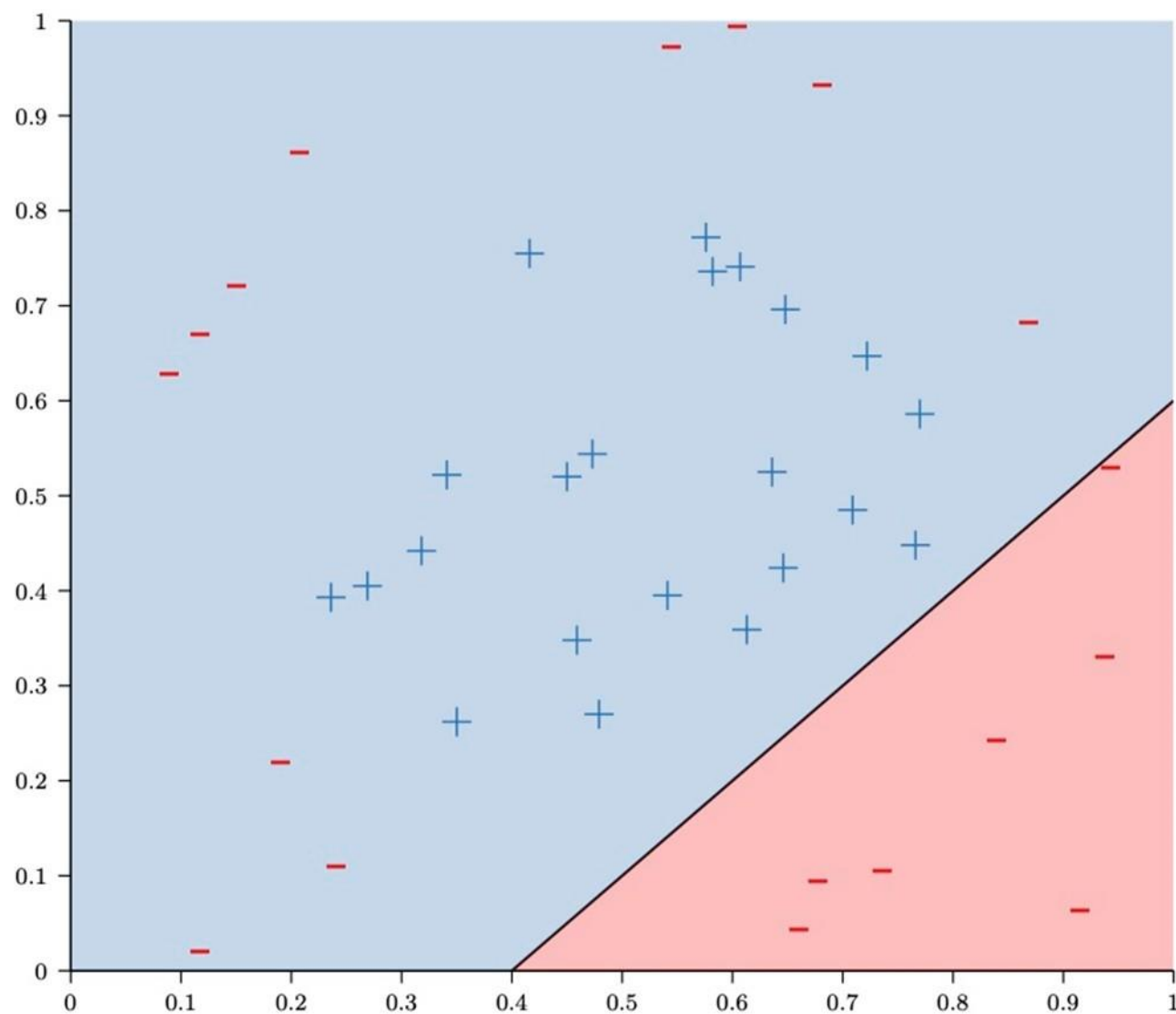




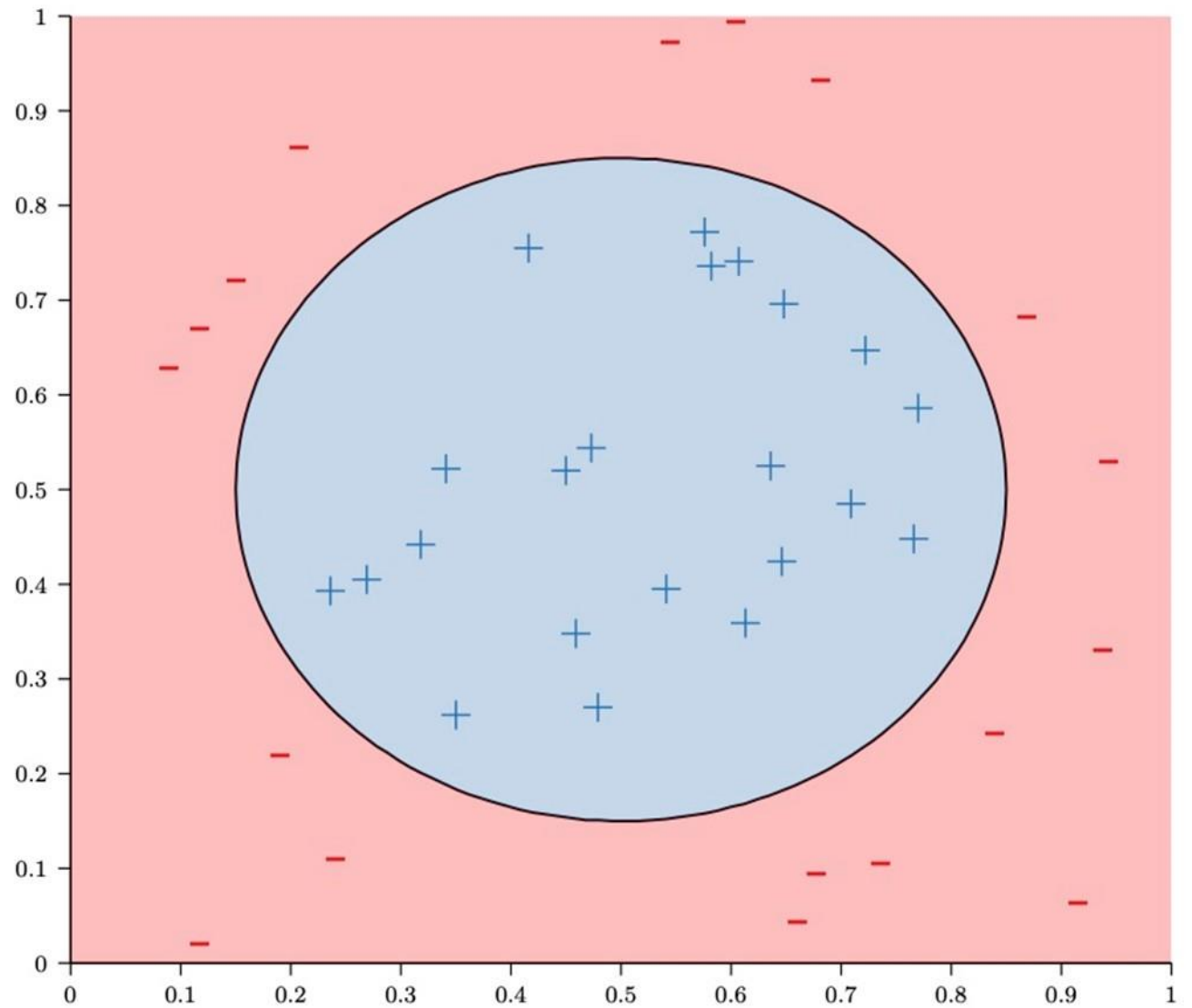
# Linear Models?



# Linear Models?



# Nonlinear Models

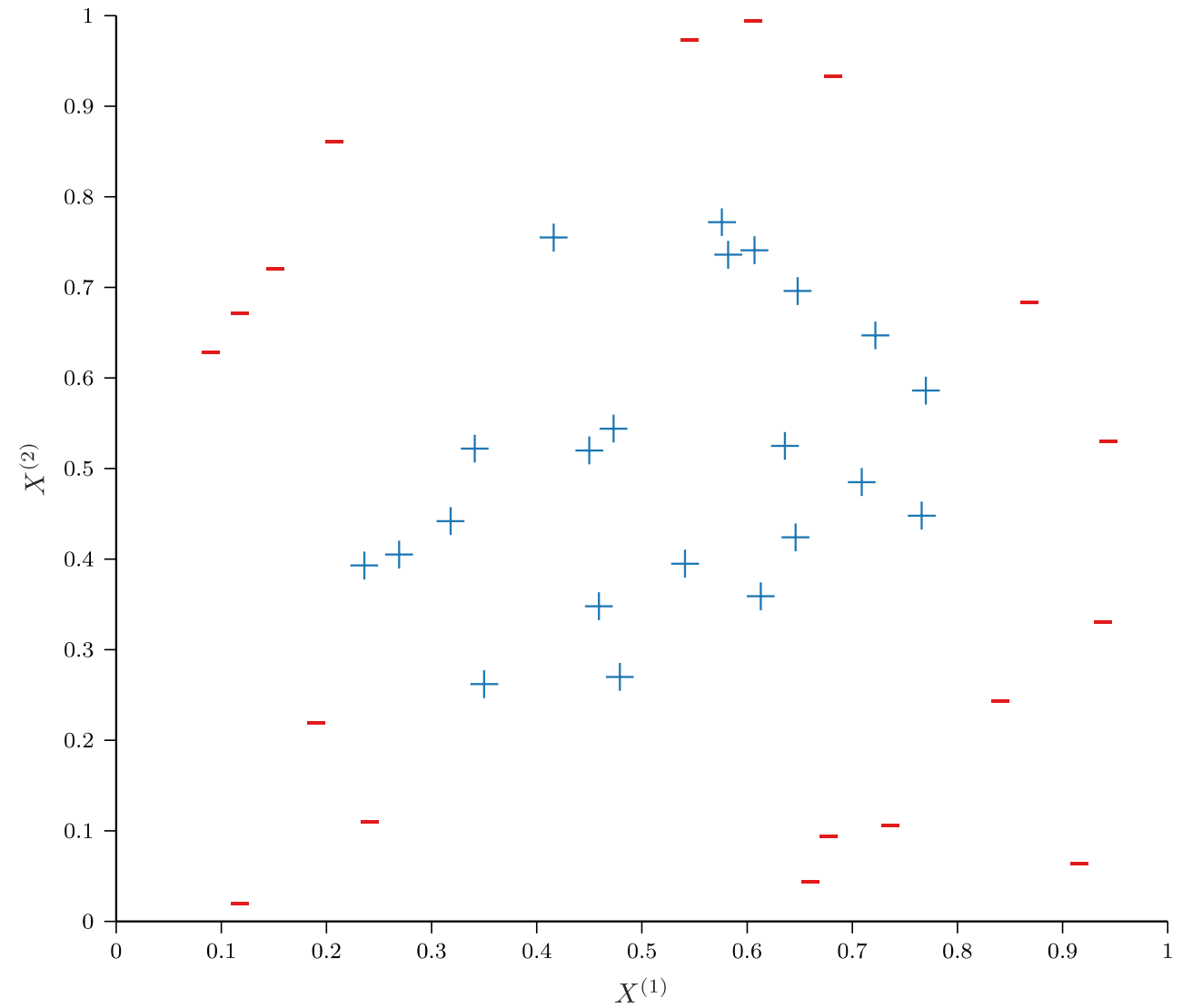


# Feature Transforms

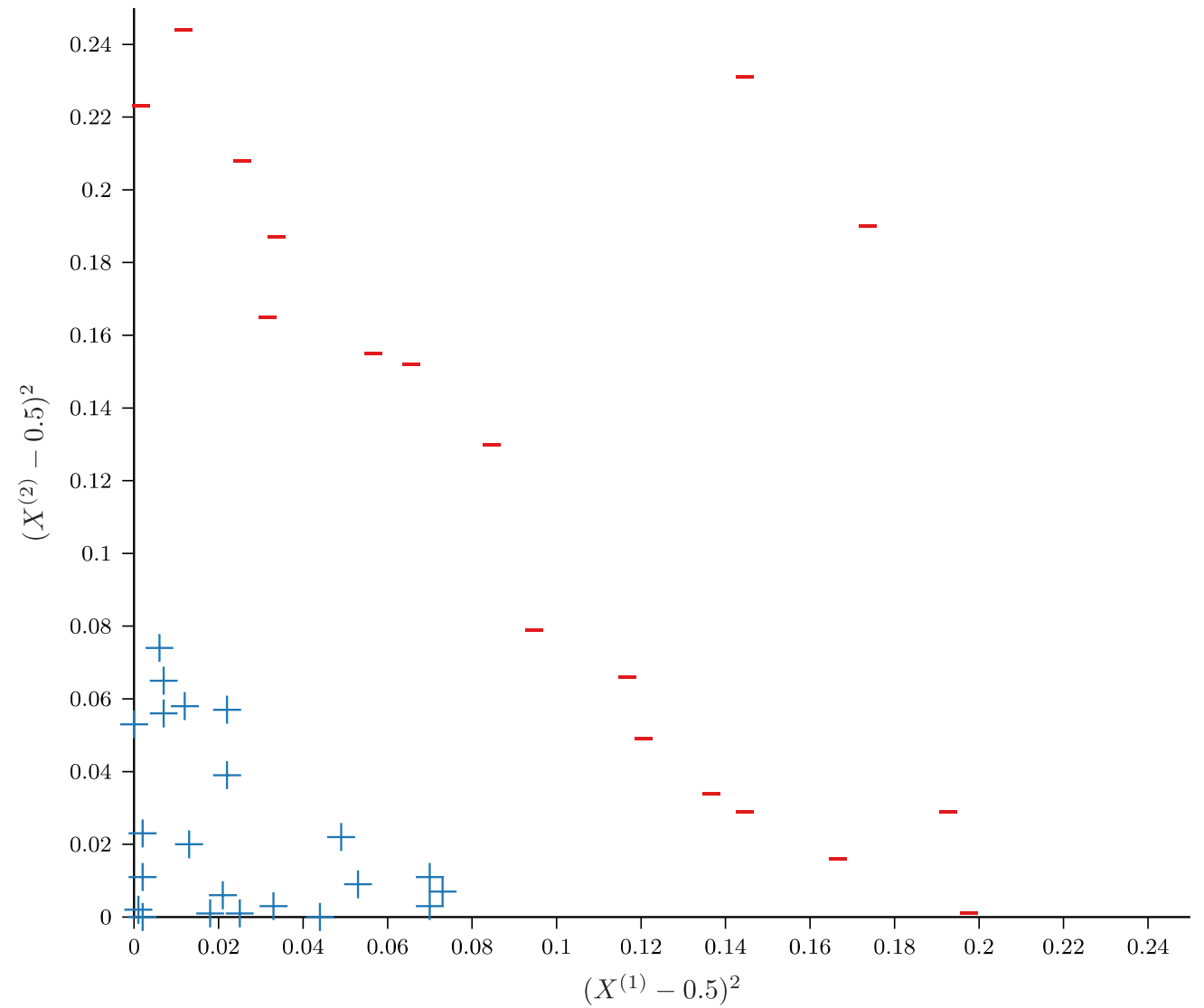
- Given  $p$ -dimensional inputs  $X = [X^{(1)}, \dots, X^{(p)}]$ , first compute some transformation of our input, e.g.,

$$\phi([X^{(1)}, X^{(2)}]) = [(X^{(1)} - 0.5)^2, (X^{(2)} - 0.5)^2]$$

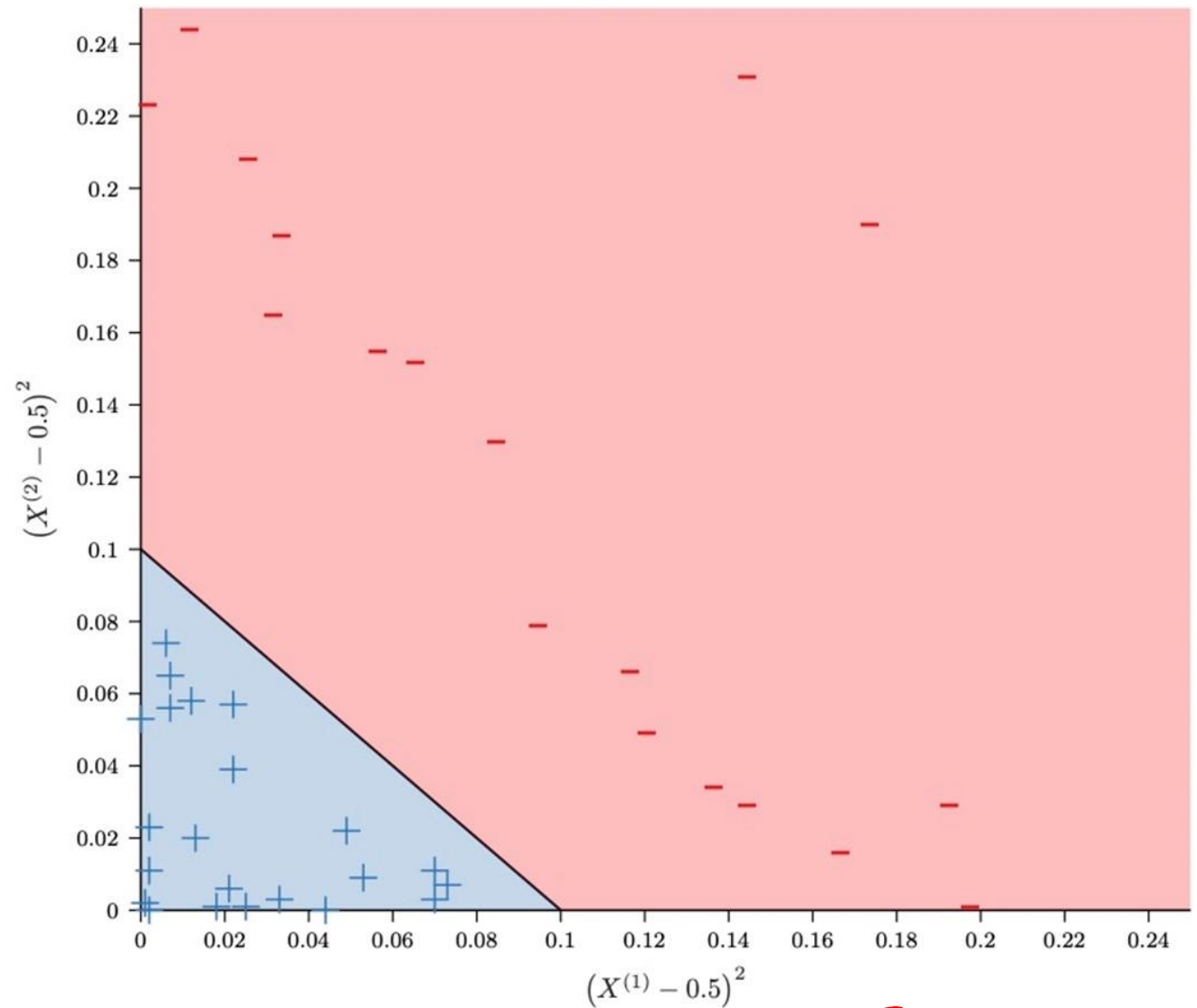
# Nonlinear Models



# Nonlinear Models

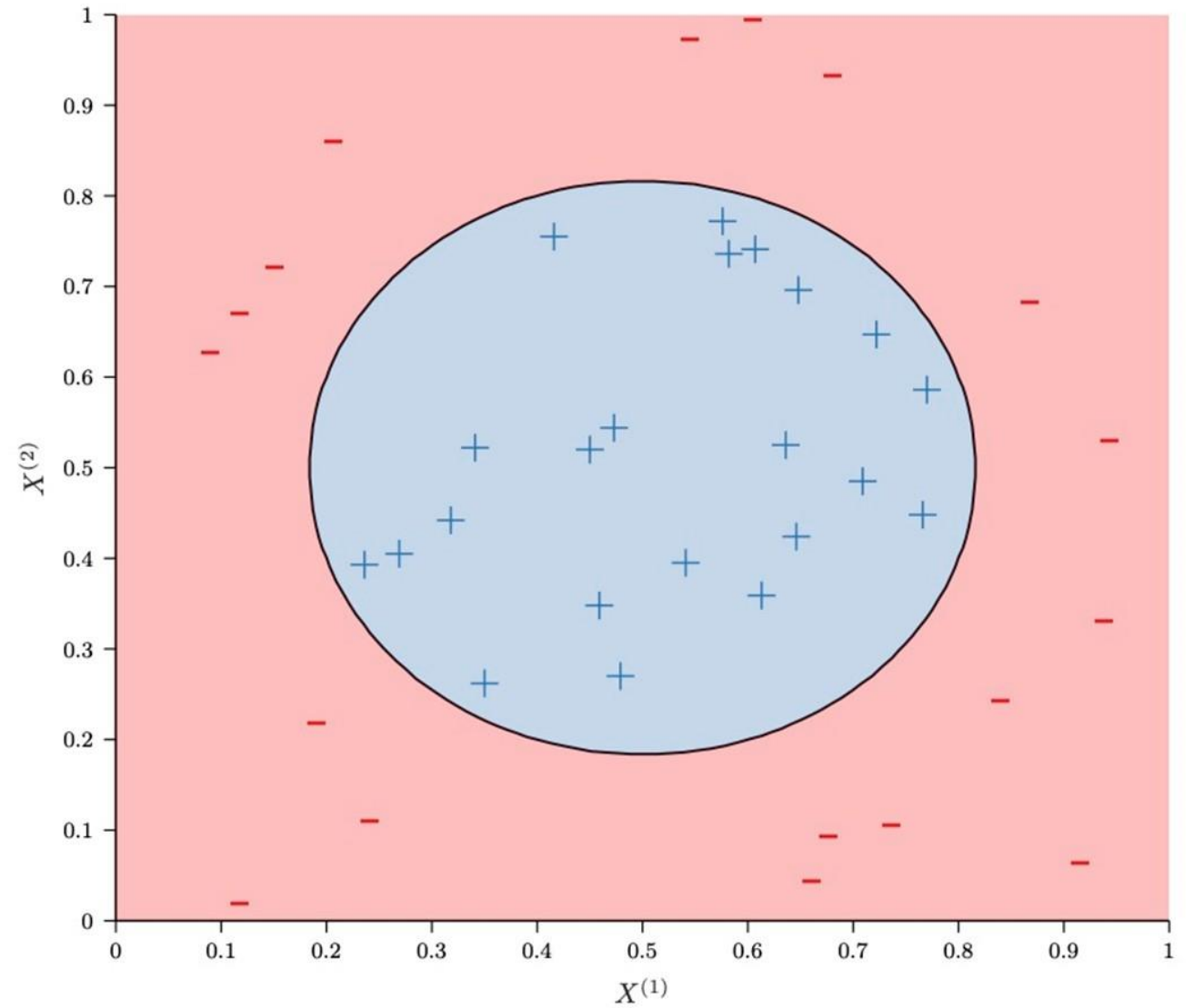


# Nonlinear Models



$$(X^{(1)} - 0.5)^2 + (X^{(2)} - 0.5)^2 = 0.1$$

# Nonlinear Models



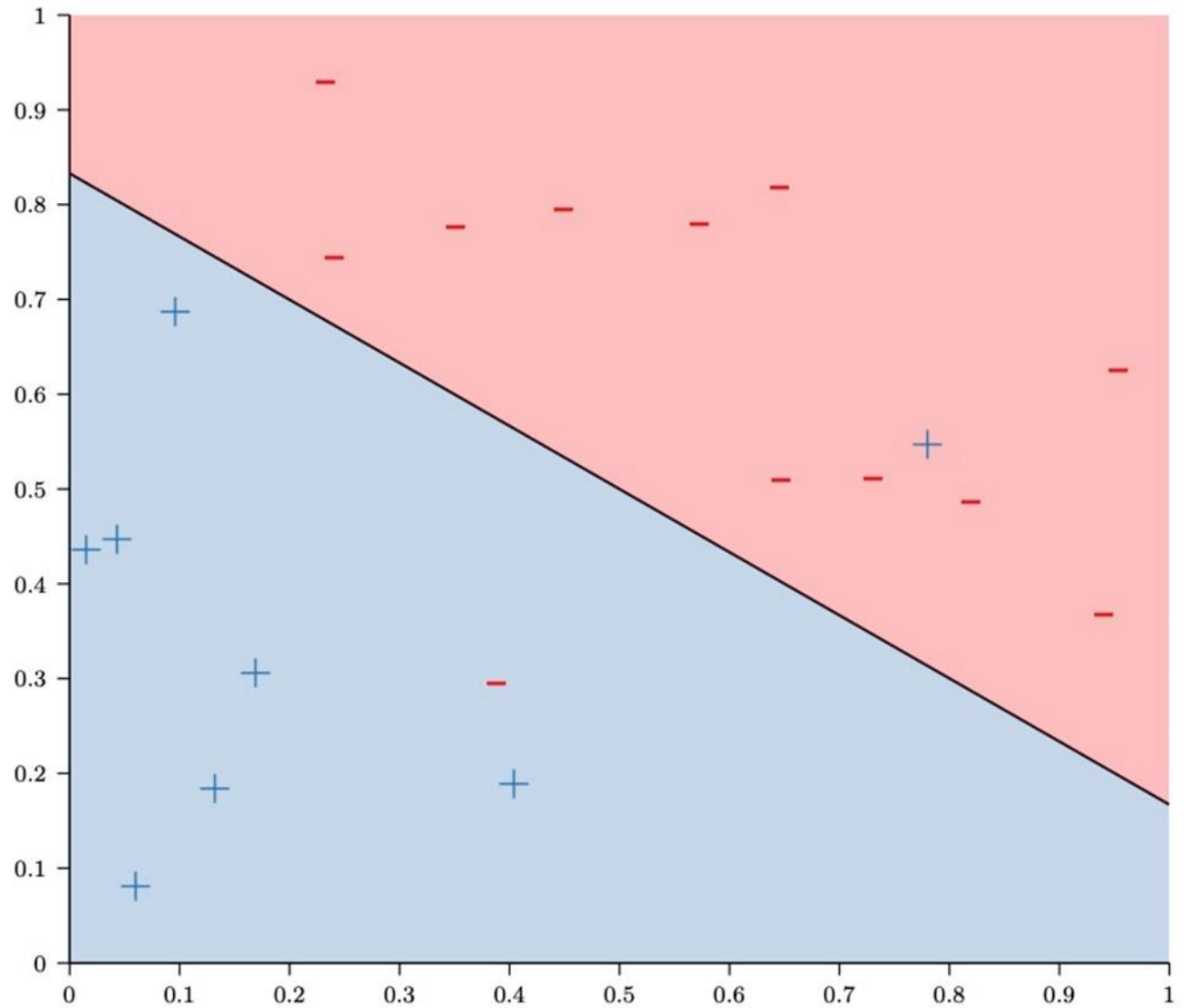


# General $k^{th}$ -order Transforms

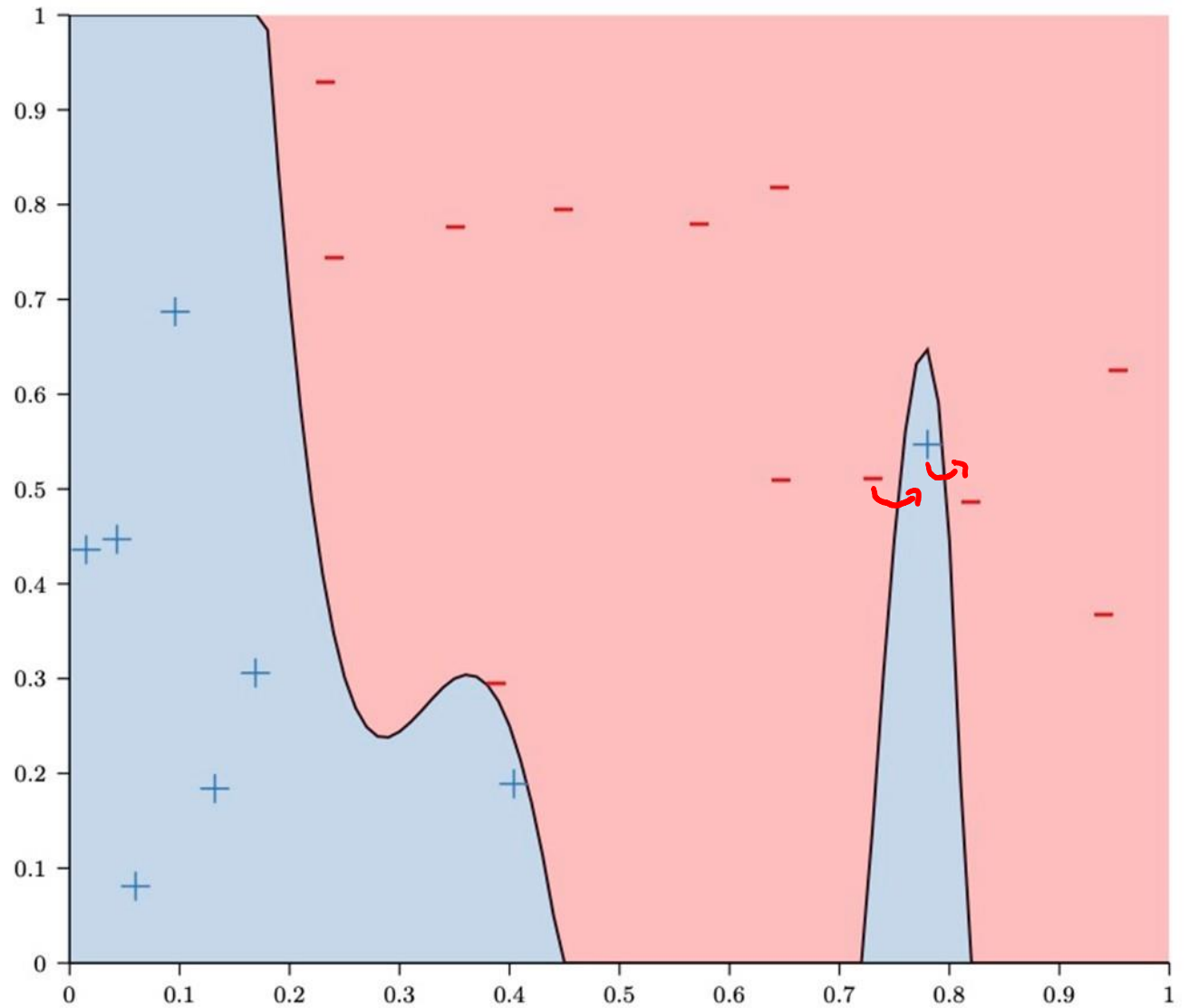
- $\phi_{2,2}([X^{(1)}, X^{(2)}]) = [X^{(1)}, X^{(2)}, X^{(1)2}, X^{(1)}X^{(2)}, X^{(2)2}]$
- $\phi_{2,3}([X^{(1)}, X^{(2)}]) =$   
 $[ \phi_{2,2}([X^{(1)}, X^{(2)}]), X^{(1)3}, X^{(1)2}X^{(2)}, X^{(1)}X^{(2)2}, X^{(2)3} ]$
- $\phi_{2,4}([X^{(1)}, X^{(2)}]) =$   
 $[ \phi_{2,3}([X^{(1)}, X^{(2)}]), X^{(1)4}, X^{(1)3}X^{(2)}, X^{(1)2}X^{(2)2}, X^{(1)}X^{(2)3}, X^{(2)4} ]$
- $\phi_{2,Q}$  maps a 2-dimensional input to a  $\frac{Q(Q+3)}{2}$ -dimensional output
- Scales even worse for higher-dimensional inputs...

Solution: Kernel methods

# Linear Models



# Nonlinear Models?



# Feature Transforms: Tradeoffs

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

Overfitting



# 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:

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

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

Poll 1: Which model do you think will have lower *training* error?

- A.  $\mathcal{F}_2$
- B.  $\mathcal{F}_{10}$

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

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

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

Poll 2: Which model do you think will have lower *true* error?

- A.  $\mathcal{F}_2$
- B.  $\mathcal{F}_{10}$

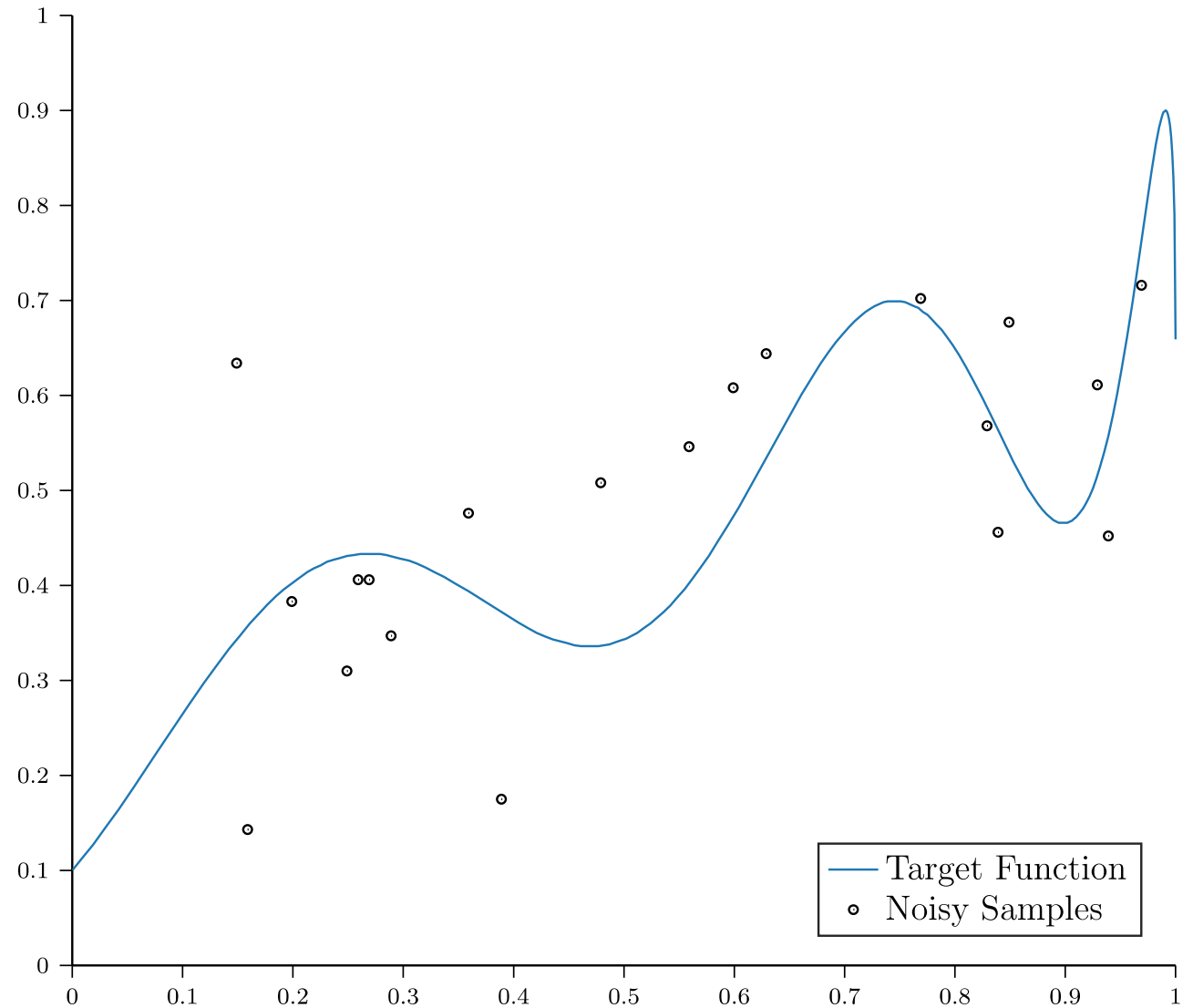
- $X \in \mathbb{R}, Y \in \mathbb{R}$  and  $n = 20$
- Targets are generated by a 10<sup>th</sup>-order polynomial in  $X$  with additive Gaussian noise:

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

- $\mathcal{F}_2 = 2^{\text{nd}}$ -order polynomials
  - $\phi_{1,2}(X) = [X, X^2]$
- $\mathcal{F}_{10} = 10^{\text{th}}$ -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

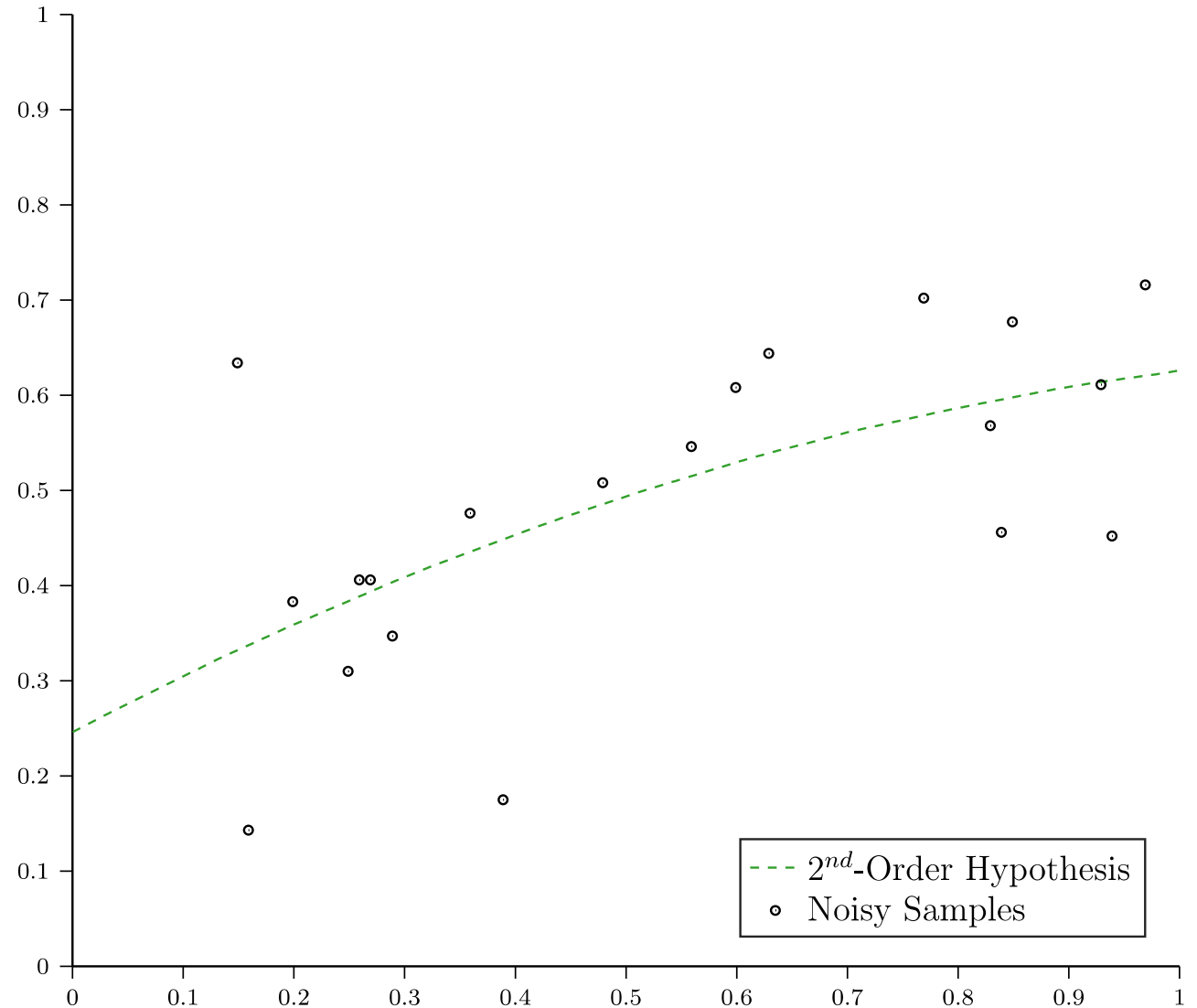
- 10-dimensional target function with additive Gaussian noise
- $\mathcal{F}_2 = 2^{\text{nd}}$ -order polynomial
- $\mathcal{F}_{10} = 10^{\text{th}}$ -order polynomial





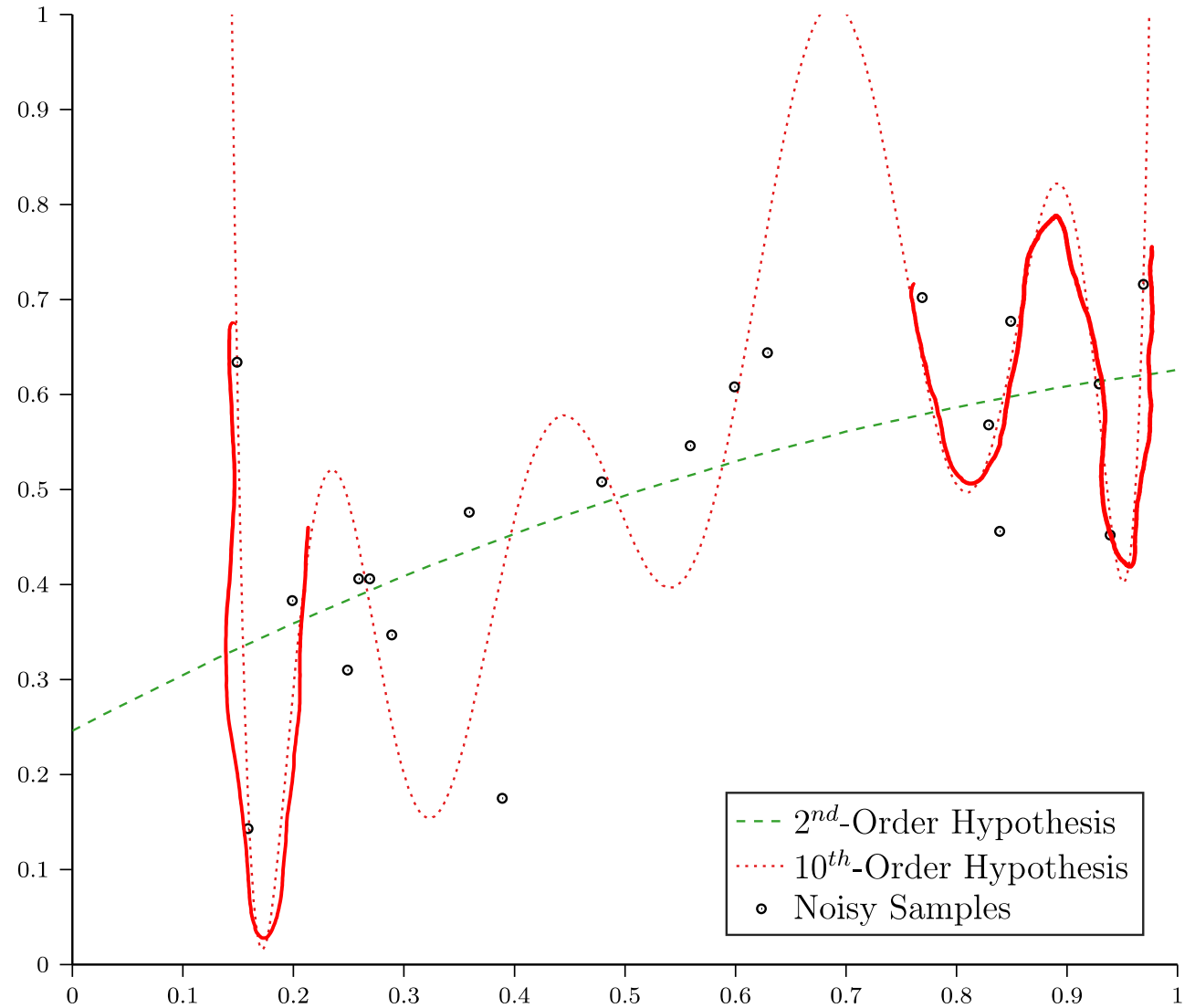
# Noisy Targets

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



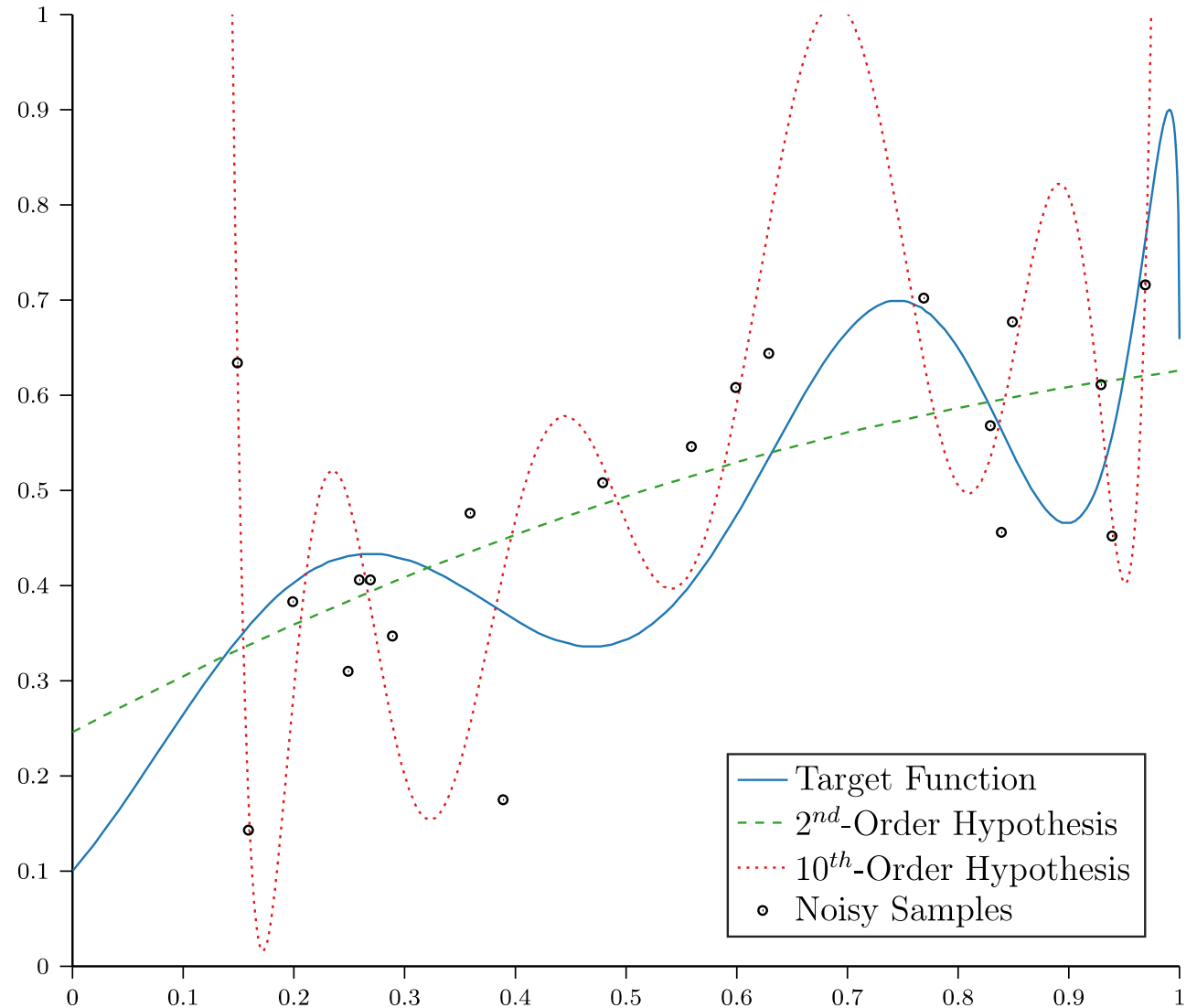
# Noisy Targets

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



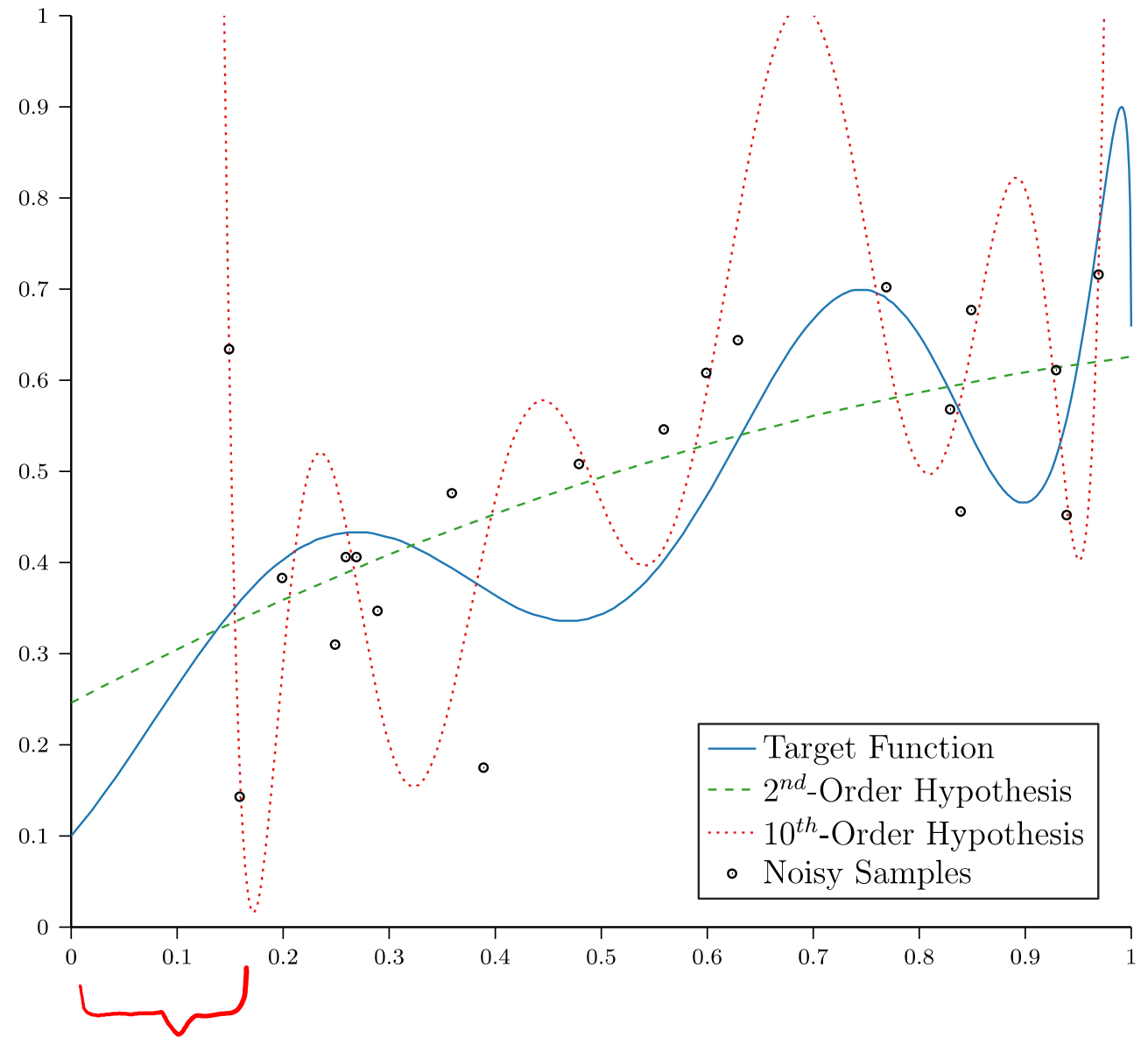
# Noisy Targets

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



# Noisy Targets

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



# Regularization

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

# Hard Constraints

- $\mathcal{F}_{10} = 10^{\text{th}}$ -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  $\mathbf{A} = \begin{bmatrix} 1 & \phi_{1,2}(X_1) \\ 1 & \phi_{1,2}(X_2) \\ \vdots & \vdots \\ 1 & \phi_{1,2}(X_n) \end{bmatrix}$  and  $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$  find  $\beta = [\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7, \beta_8, \beta_9, \beta_{10}]$  that minimizes
$$\frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$
- Subject to
$$\beta_3 = \beta_4 = \beta_5 = \beta_6 = \beta_7 = \beta_8 = \beta_9 = \beta_{10} = 0$$

# Hard Constraints

- $\mathcal{F}_{10} = 10^{\text{th}}$ -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  $\mathbf{A} = \begin{bmatrix} 1 & \phi_{1,2}(X_1) \\ 1 & \phi_{1,2}(X_2) \\ \vdots & \vdots \\ 1 & \phi_{1,2}(X_n) \end{bmatrix}$  and  $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$  find  $\beta = [\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7, \beta_8, \beta_9, \beta_{10}]$  that minimizes 
$$\frac{1}{n} \sum_{i=1}^n \left( \left( \sum_{d=0}^{10} X_i^{(d)} \beta_d \right) - Y_i \right)^2$$
- Subject to 
$$\beta_3 = \beta_4 = \beta_5 = \beta_6 = \beta_7 = \beta_8 = \beta_9 = \beta_{10} = 0$$

# Hard Constraints

- $\mathcal{F}_{10} = 10^{\text{th}}$ -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  $\mathbf{A} = \begin{bmatrix} 1 & \phi_{1,2}(X_1) \\ 1 & \phi_{1,2}(X_2) \\ \vdots & \vdots \\ 1 & \phi_{1,2}(X_n) \end{bmatrix}$  and  $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$  find  $\beta = [\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7, \beta_8, \beta_9, \beta_{10}]$  that minimizes 
$$\frac{1}{n} \sum_{i=1}^n \left( \left( \sum_{d=0}^2 X_i^{(d)} \beta_d \right) - Y_i \right)^2$$
- Subject to nothing!



# Hard Constraints

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

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

- Given  $\mathbf{A} = \begin{bmatrix} 1 & \phi_{1,2}(X_1) \\ 1 & \phi_{1,2}(X_2) \\ \vdots & \vdots \\ 1 & \phi_{1,2}(X_n) \end{bmatrix}$  and  $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$  find

$\beta = [\beta_0, \beta_1, \beta_2]$  that minimizes

$$\frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

- Subject to nothing!

# Soft Constraints

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

- Given  $\mathbf{A} = \begin{bmatrix} 1 & \phi_1(X_1) & \cdots & \phi_m(X_1) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(X_n) & \cdots & \phi_m(X_n) \end{bmatrix}$  and  $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$ ,

find  $\beta$  that minimizes

$$\frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

- Subject to:

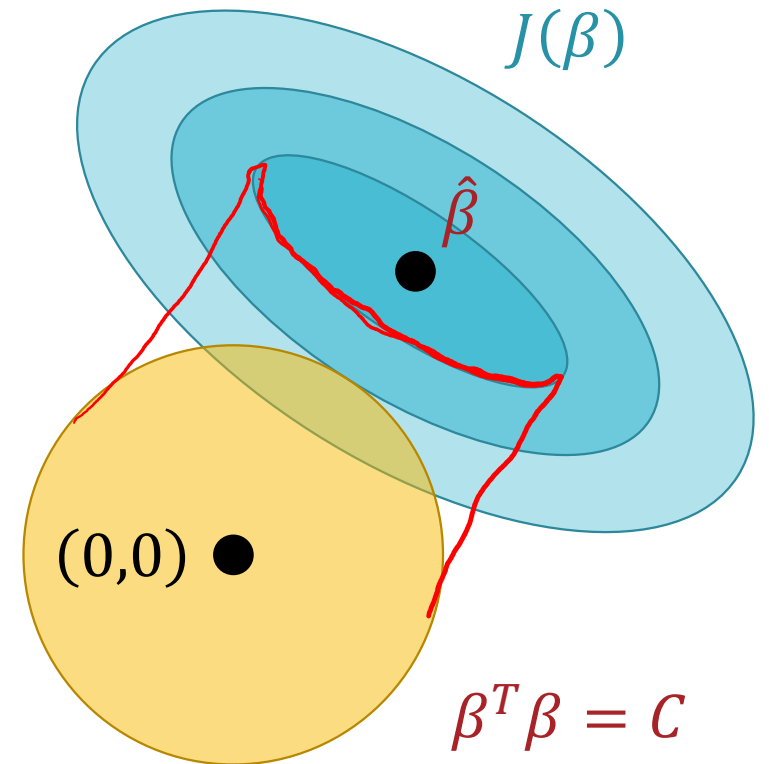
$$\underline{\|\beta\|_2^2} = \beta^T \beta = \sum_{d=0}^m \beta_d^2 \leq C$$

# Soft Constraints

$$\text{minimize } J(\beta) = \frac{1}{n} (A\beta - Y)^T (A\beta - Y)$$

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

$$\beta_1^2 + \beta_2^2 \leq C$$

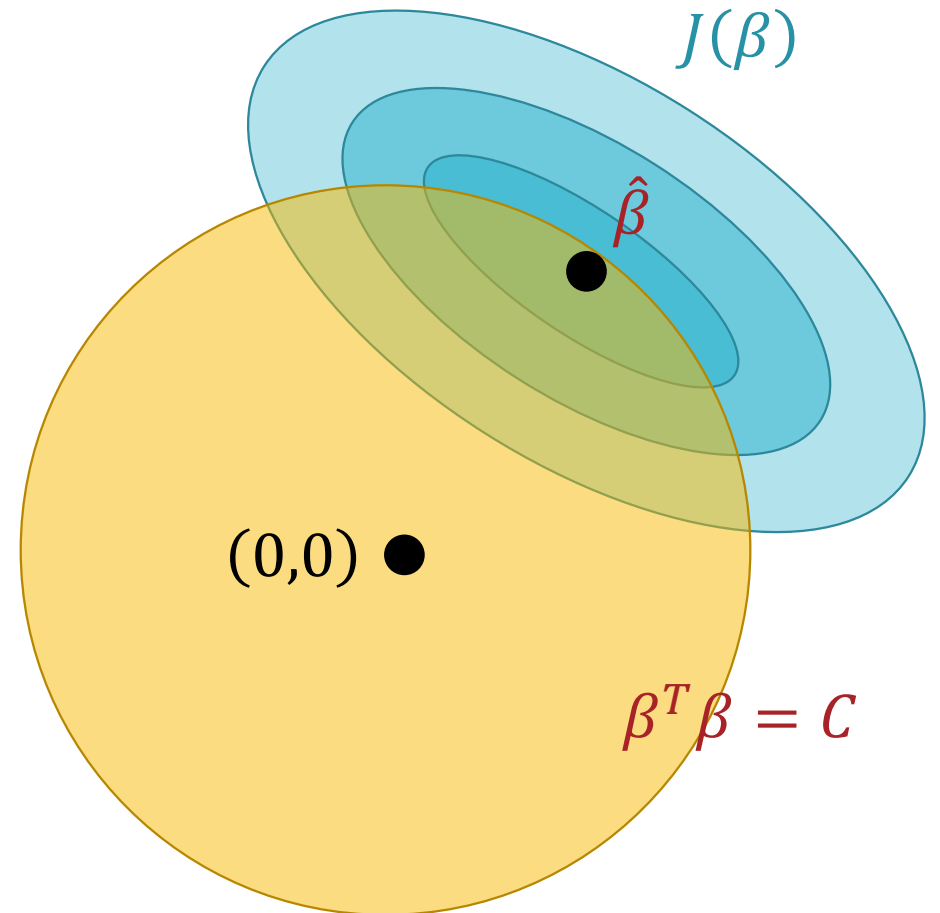


# Soft Constraints

$$\text{minimize } J(\beta) = \frac{1}{n} (A\beta - Y)^T (A\beta - Y)$$

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

$\hat{\beta}$  is optimal



# Soft Constraints

$$\text{minimize } J(\beta) = \frac{1}{n} (A\beta - Y)^T (A\beta - Y)$$

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

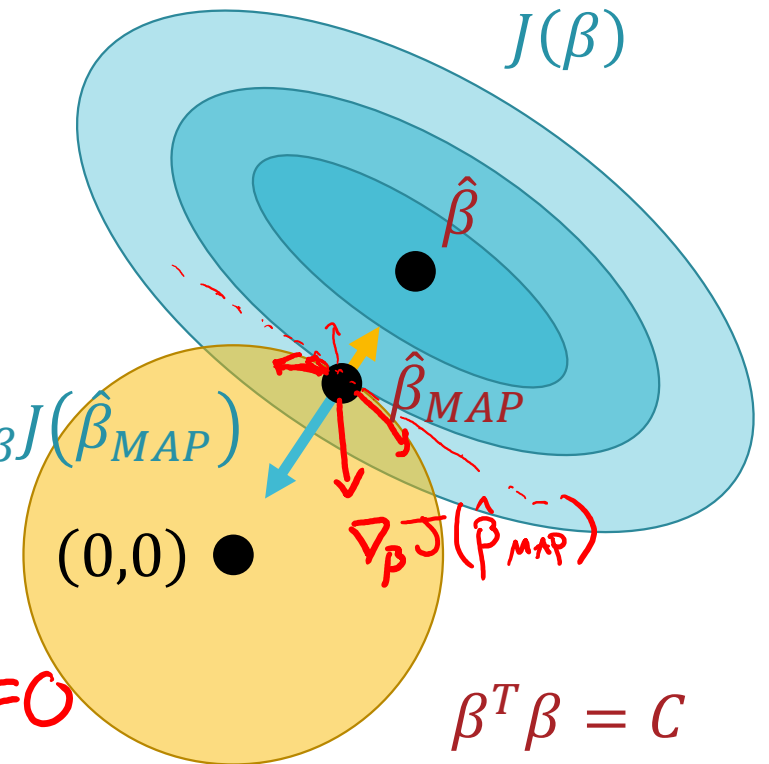
$$\nabla_{\beta} J(\hat{\beta}_{MAP}) \propto -\frac{2}{n} \hat{\beta}_{MAP}$$

$$\nabla_{\beta} J(\hat{\beta}_{MAP}) = -\frac{2}{n} \lambda_c \hat{\beta}_{MAP} \quad \nabla_{\beta} J(\hat{\beta}_{MAP})$$

$$\text{where } \lambda_c \geq 0$$

$$\nabla_{\beta} J(\hat{\beta}_{MAP}) + \frac{2}{n} \lambda_c \hat{\beta}_{MAP} = 0$$

$$\nabla_{\beta} (J(\hat{\beta}_{MAP}) + \frac{1}{n} \lambda_c \hat{\beta}_{MAP}^T \hat{\beta}_{MAP}) = 0$$



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

$$\text{minimize } J(\beta) = \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

$$\text{subject to } \beta^T \beta \leq c$$



$$\text{minimize } J_{AUG}(\beta) = \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) + \frac{\lambda_c}{n} \beta^T \beta$$

## Ridge Regression

$$\text{minimize } J_{AUG}(\beta) = \frac{1}{n} (A\beta - Y)^T (A\beta - Y) + \frac{\lambda_c}{n} \beta^T \beta$$

$$\nabla_{\beta} (J_{AUG}(\beta)) = \frac{1}{n} (2A^T A \beta - 2A^T Y + 2\lambda_c \beta)$$

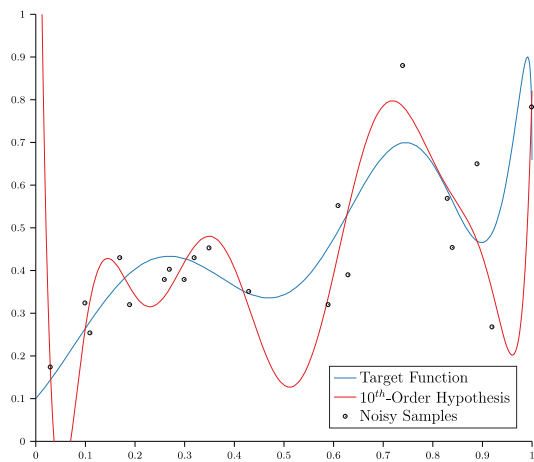
$$\Rightarrow \frac{1}{n} (2A^T A \hat{\beta}_{MAP} - 2A^T Y + 2\lambda_c \hat{\beta}_{MAP}) = 0$$

$$A^T A \hat{\beta}_{MAP} + \lambda_c \hat{\beta}_{MAP} = A^T Y$$

$$(A^T A + \lambda_c I_{p+1}) \hat{\beta}_{MAP} = A^T Y$$

$$\hat{\beta}_{MAP} = \underbrace{(A^T A + \lambda_c I_{p+1})^{-1}} A^T Y$$

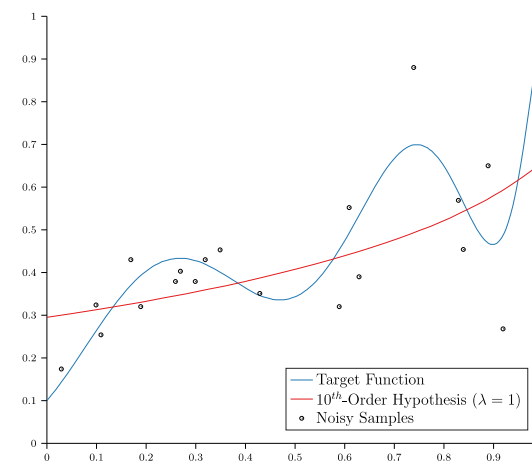
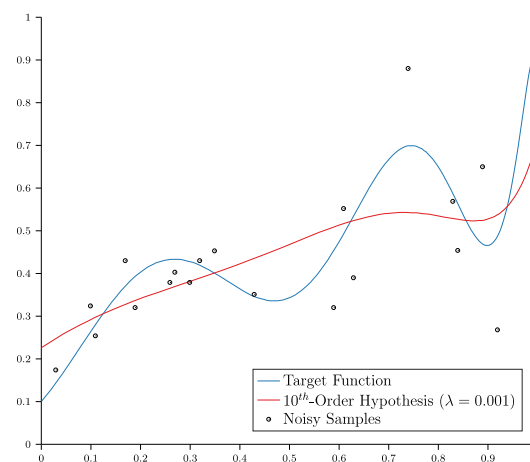
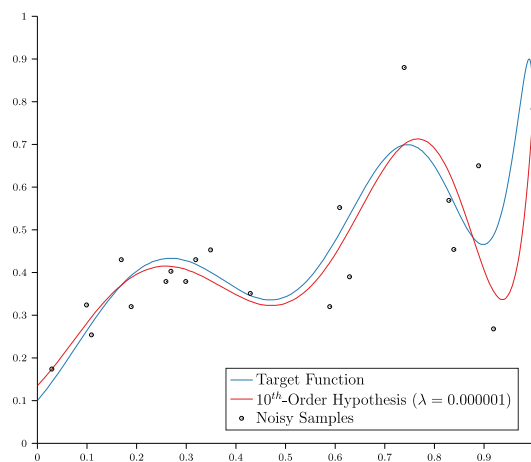
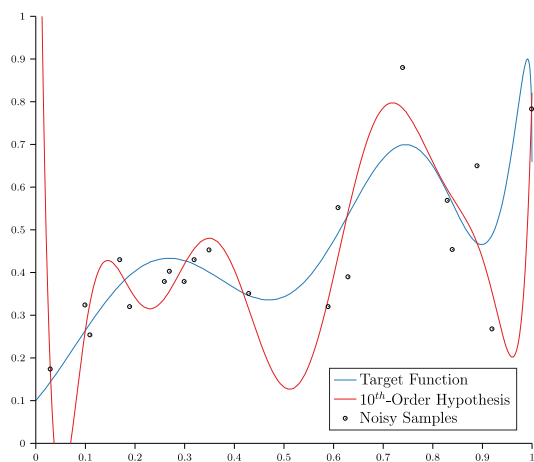
when  $\lambda_c \geq 0$ , this helps make  $A^T A$  invertible



# Ridge Regression

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





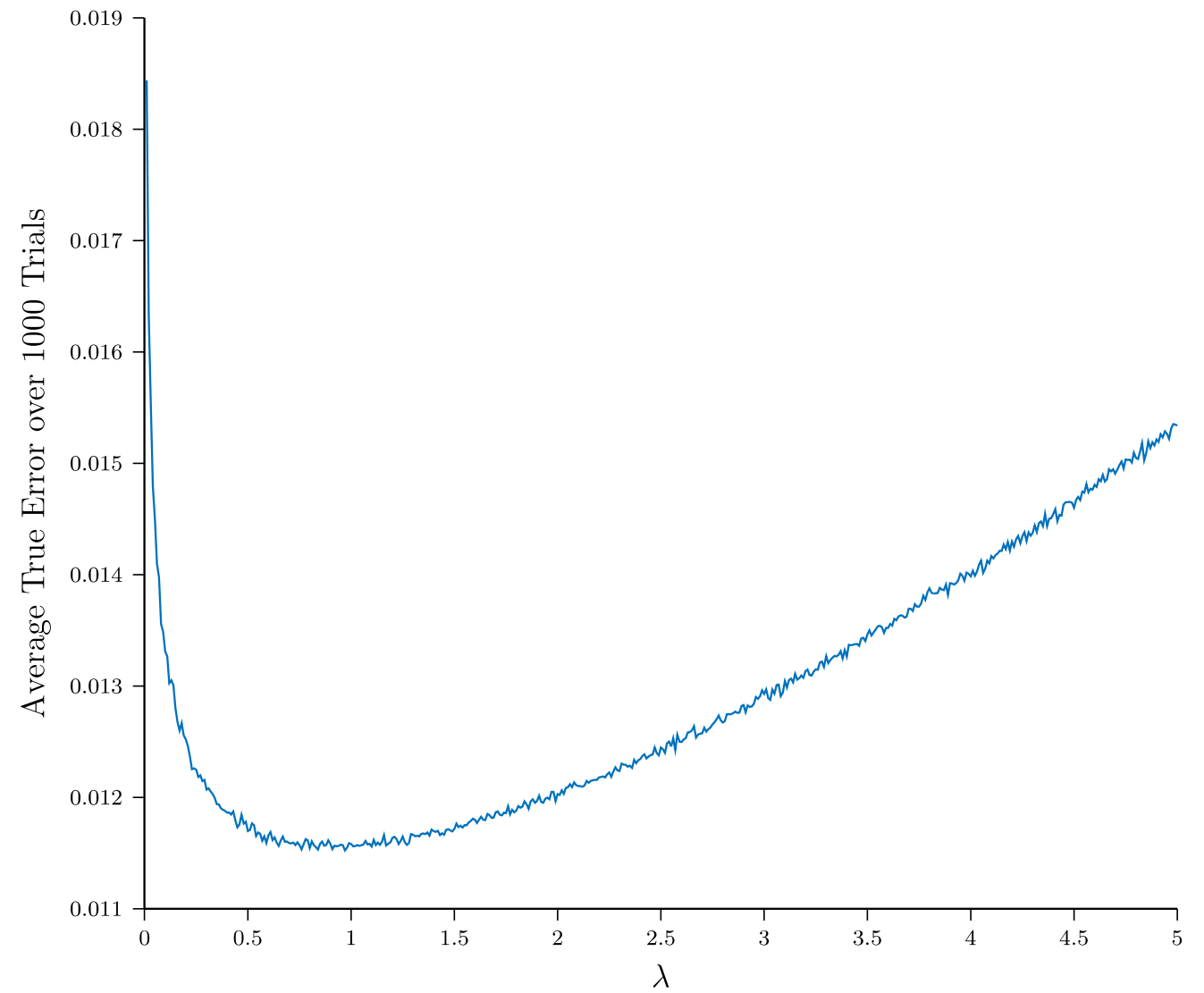
# Ridge Regression

$$\lambda_c = 0$$

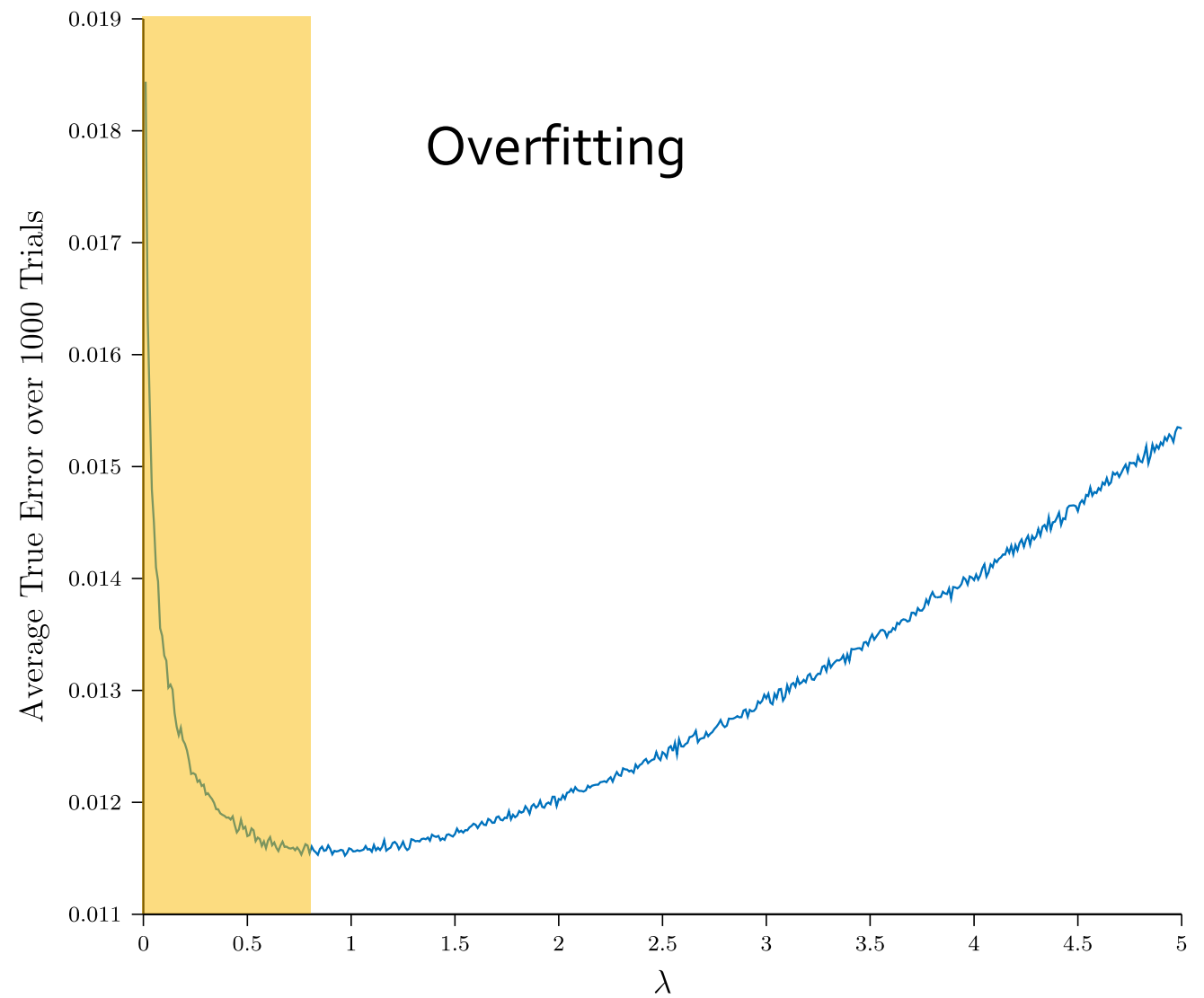
True  
Error 0.059

Overfit

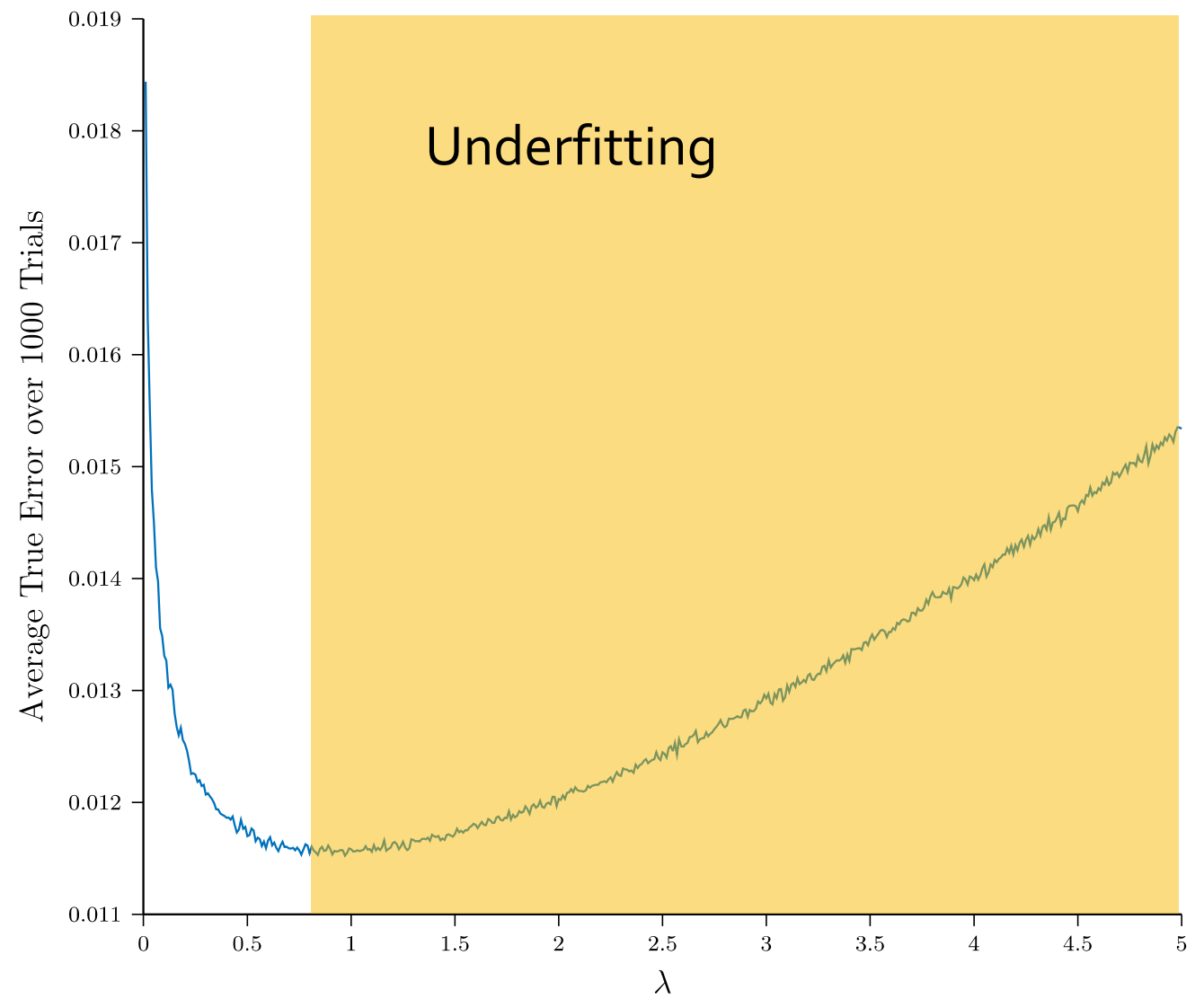
# Setting $\lambda$



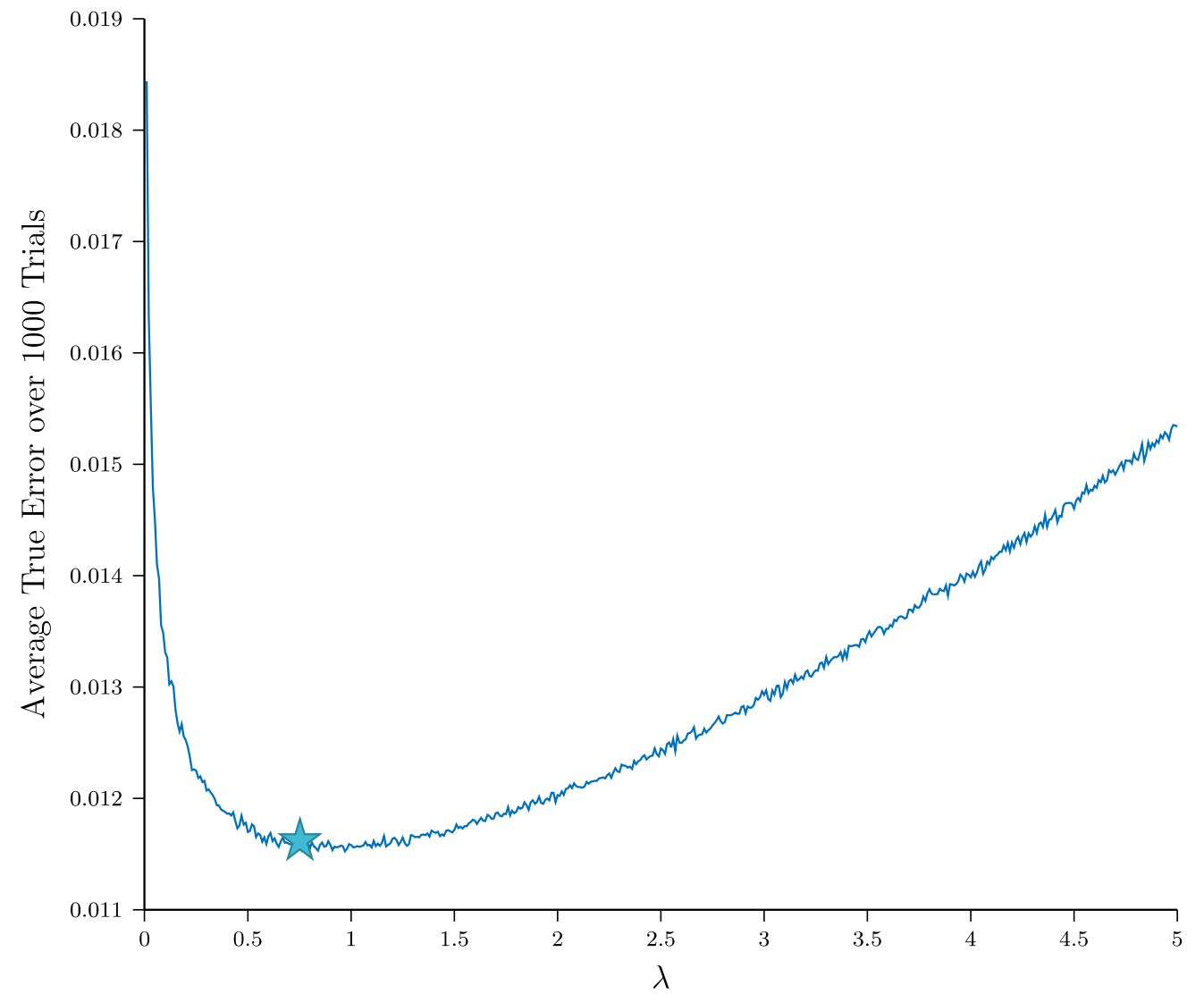
# Setting $\lambda$



# Setting $\lambda$

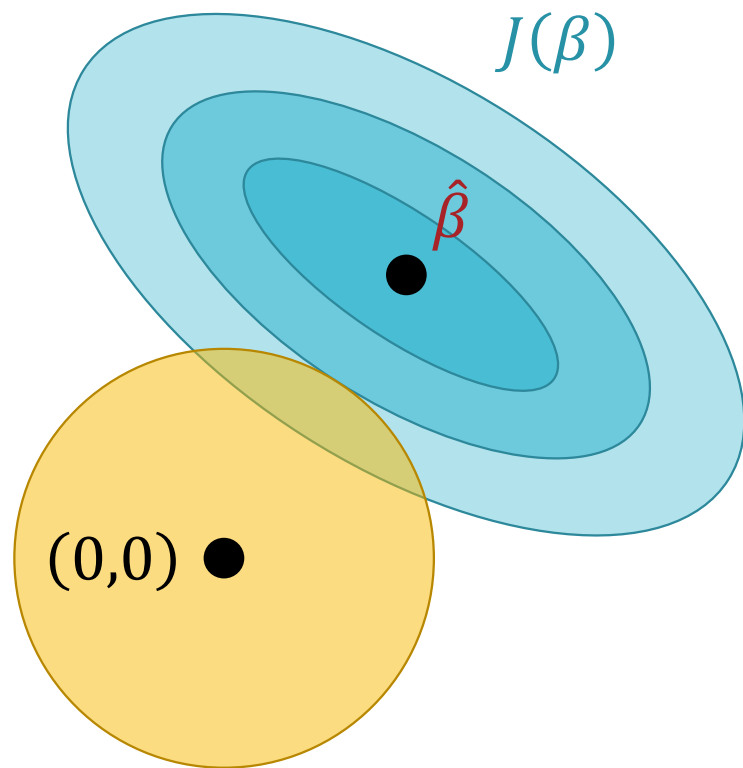


# Setting $\lambda$

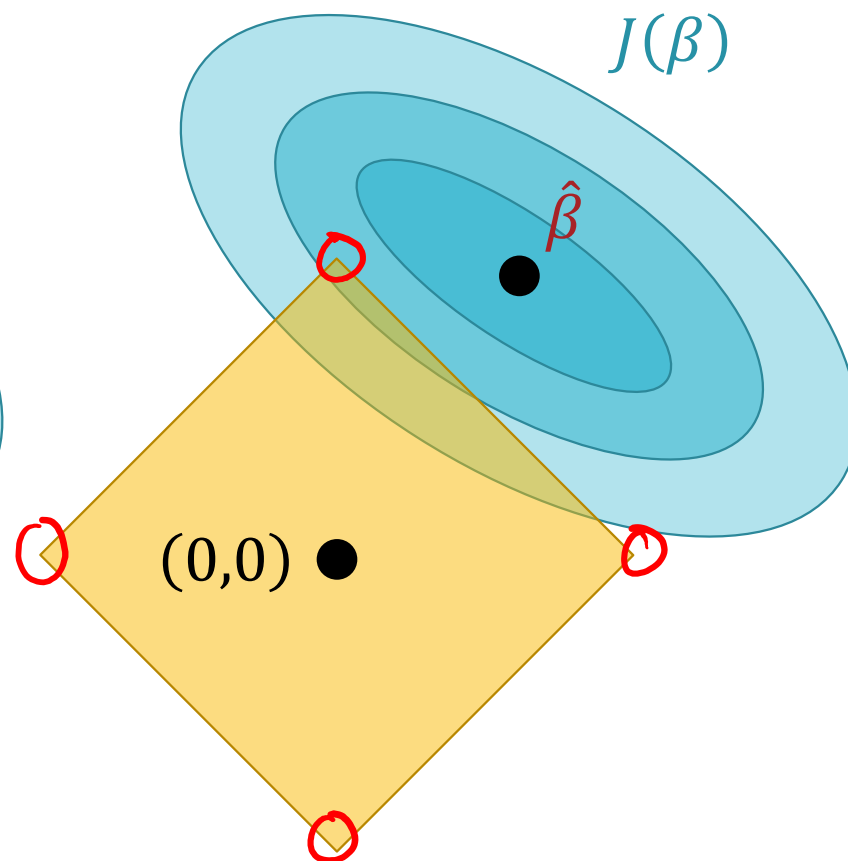


## Other Regularizers

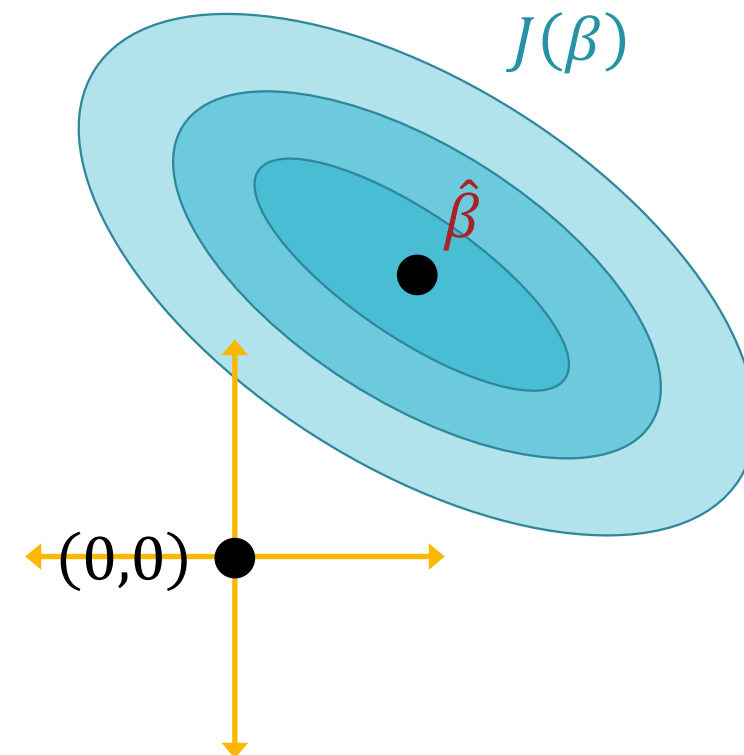
$J(\beta) + \lambda \text{pen}(\beta)$		
Ridge or $L2$	$\text{pen}(\beta) = \ \beta\ _2^2 = \sum_{d=0}^p \beta_d^2$	Encourages small weights
Lasso or $L1$	$\text{pen}(\beta) = \ \beta\ _1 = \sum_{d=0}^p  \beta_d $	Encourages sparsity
$L0$ <i>Do not use</i>	$\text{pen}(\beta) = \ \beta\ _0 = \sum_{d=0}^p \mathbb{1}(\beta_d \neq 0)$	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

$$Y = X\beta + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2) \rightarrow Y \sim N(X\beta, \sigma^2)$$

- Then given  $\mathbf{A} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$  and  $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$  the MLE of  $\beta$  is

$$\hat{\beta} = \operatorname{argmax}_{\beta} \log P(\mathbf{Y}|\mathbf{A}, \beta)$$

$$= \operatorname{argmax}_{\beta} \log \exp \left( -\frac{1}{2\sigma^2} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) \right)$$

$$= \operatorname{argmin}_{\beta} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$N(\mu, \sigma^2) = \frac{1}{\sqrt{2\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$$



# MAP for Linear Regression

- If we assume a linear model with additive Gaussian noise

$$Y = X\beta + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2) \rightarrow Y \sim N(X\beta, \sigma^2)$$

and a Gaussian prior on the weights...

$$\underline{\beta \sim N\left(0, \frac{\sigma^2}{\lambda}\right)} \rightarrow p(\beta) \propto \exp\left(-\frac{1}{2\sigma^2}(\lambda\beta^T\beta)\right)$$

- ... then, the MAP of  $\beta$  is the ridge regression solution!

$$\begin{aligned}\hat{\beta}_{MAP} &= \underset{\beta}{\operatorname{argmin}} (A\beta - Y)^T (A\beta - Y) + \lambda\beta^T\beta \\ &= (A^T A + \lambda I_{p+1})^{-1} A^T Y\end{aligned}$$

# MAP for Linear Regression

- If we assume a linear model with additive Gaussian noise

$$Y = X\beta + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2) \rightarrow Y \sim N(X\beta, \sigma^2)$$

and a Laplace prior on the weights...

$$\beta \sim \text{Laplace}\left(0, \frac{2\sigma^2}{\lambda}\right) \rightarrow p(\beta) \propto \exp\left(-\frac{1}{2\sigma^2}(\lambda\|\beta\|_1)\right)$$

- ... then, the MAP of  $\beta$  is the lasso regression solution!

$$\hat{\beta}_{MAP} = \underset{\beta}{\operatorname{argmin}} (A\beta - Y)^T (A\beta - Y) + \lambda\|\beta\|_1$$

- No closed form solution but can solve via sub-gradient descent