

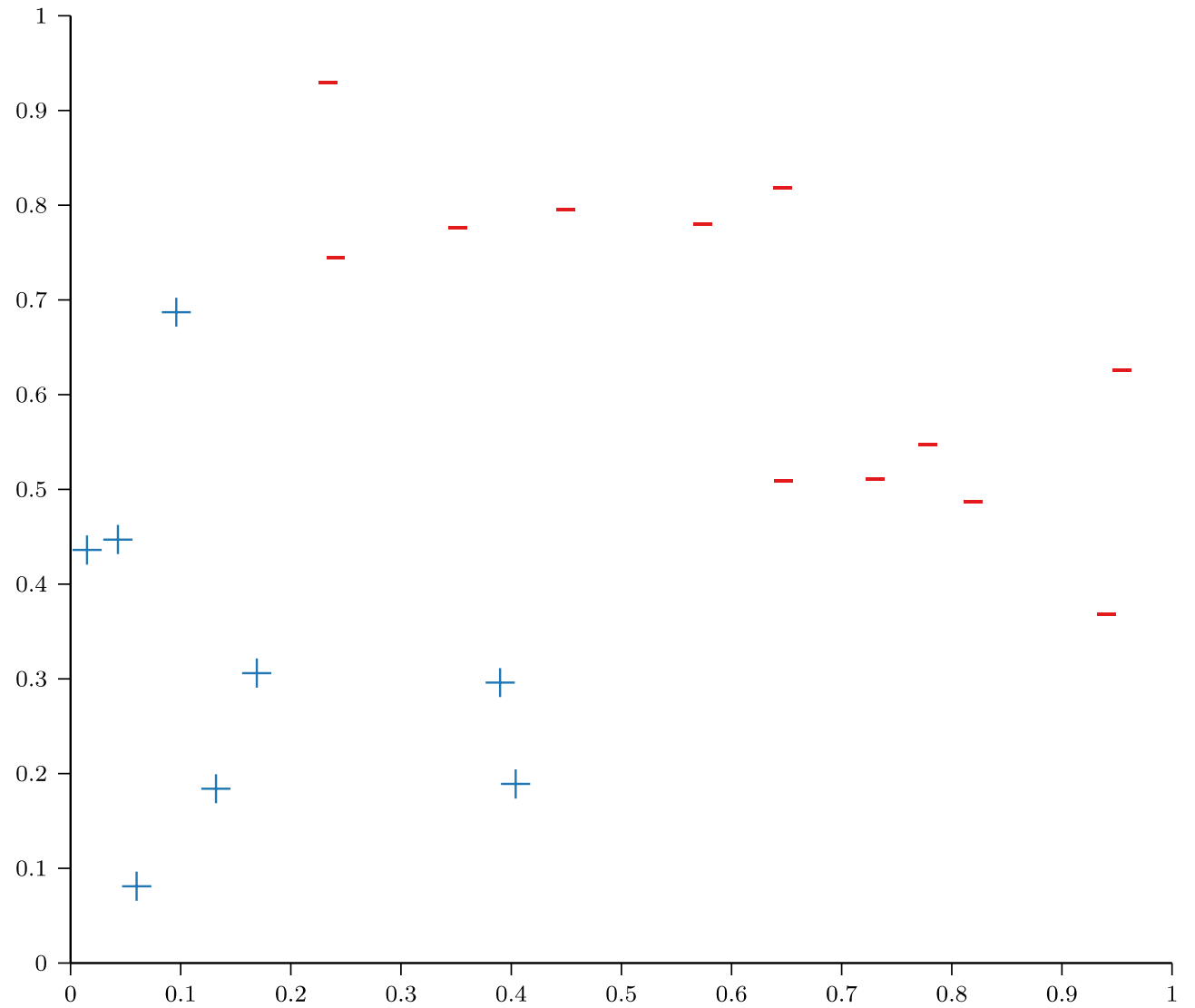
10-701: Introduction to Machine Learning

Lecture 8 – Regularization

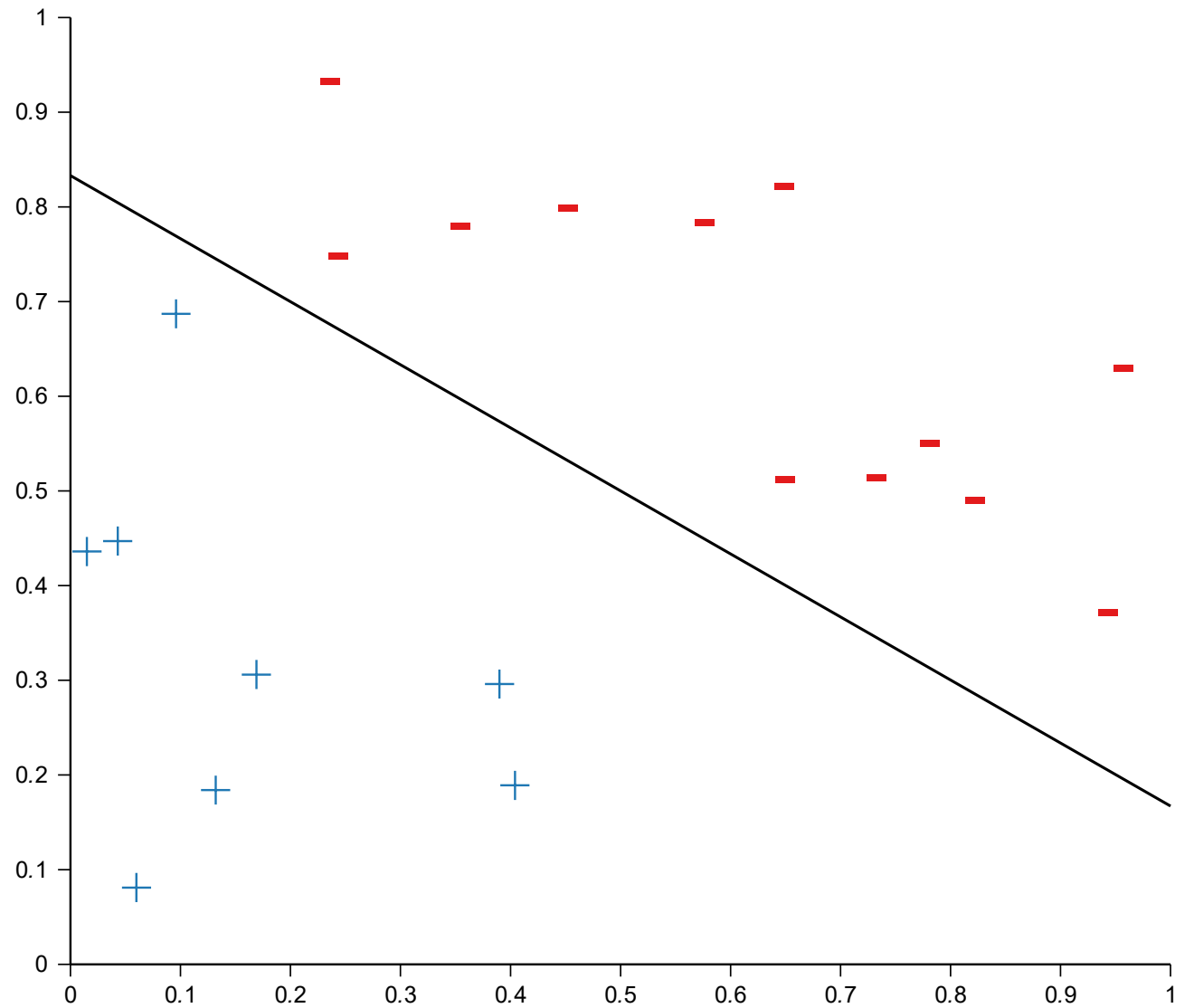
Hoda Heidari

* Slides adopted from F24 offering of 10701 by Henry Chai.

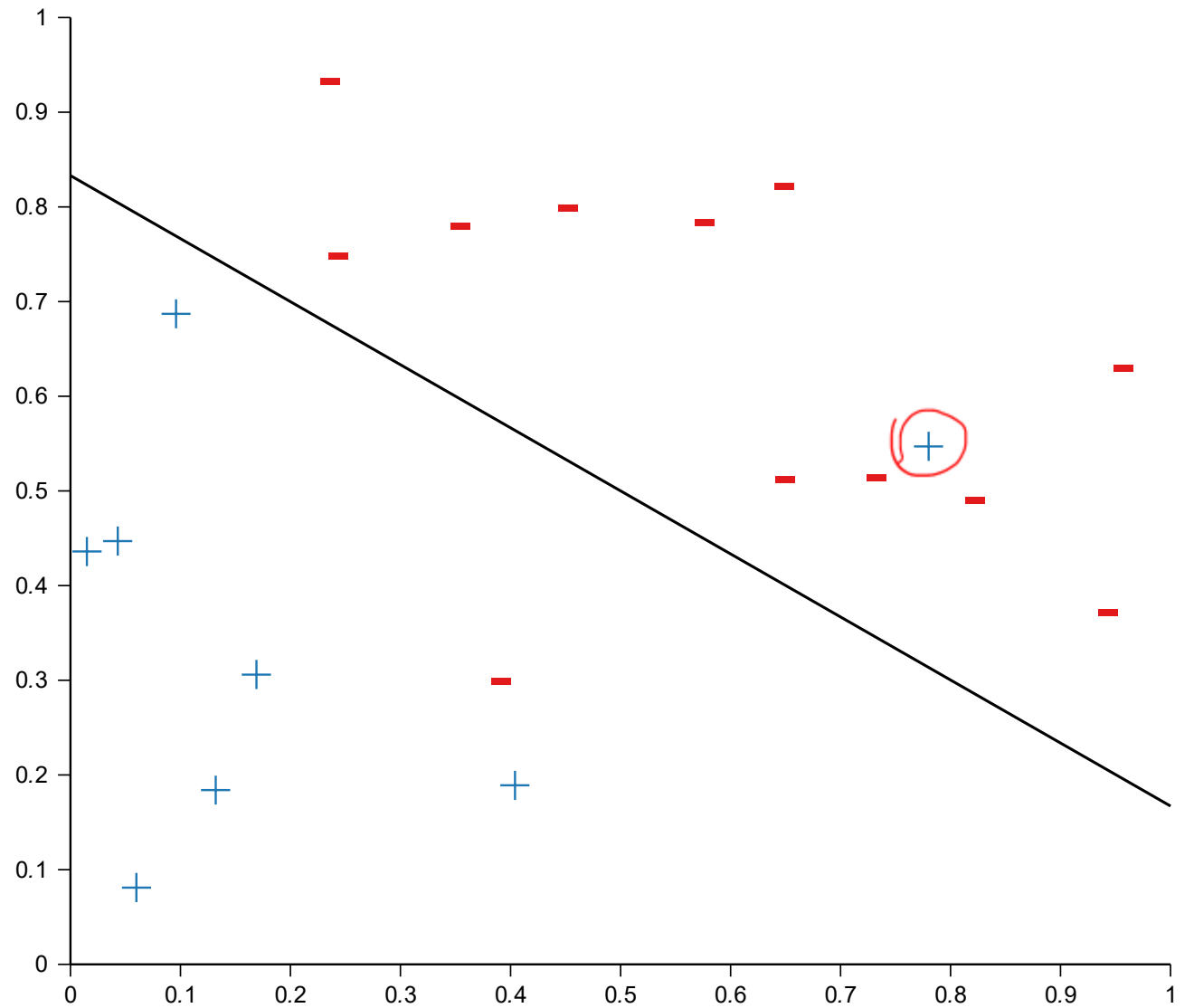
Linear Models



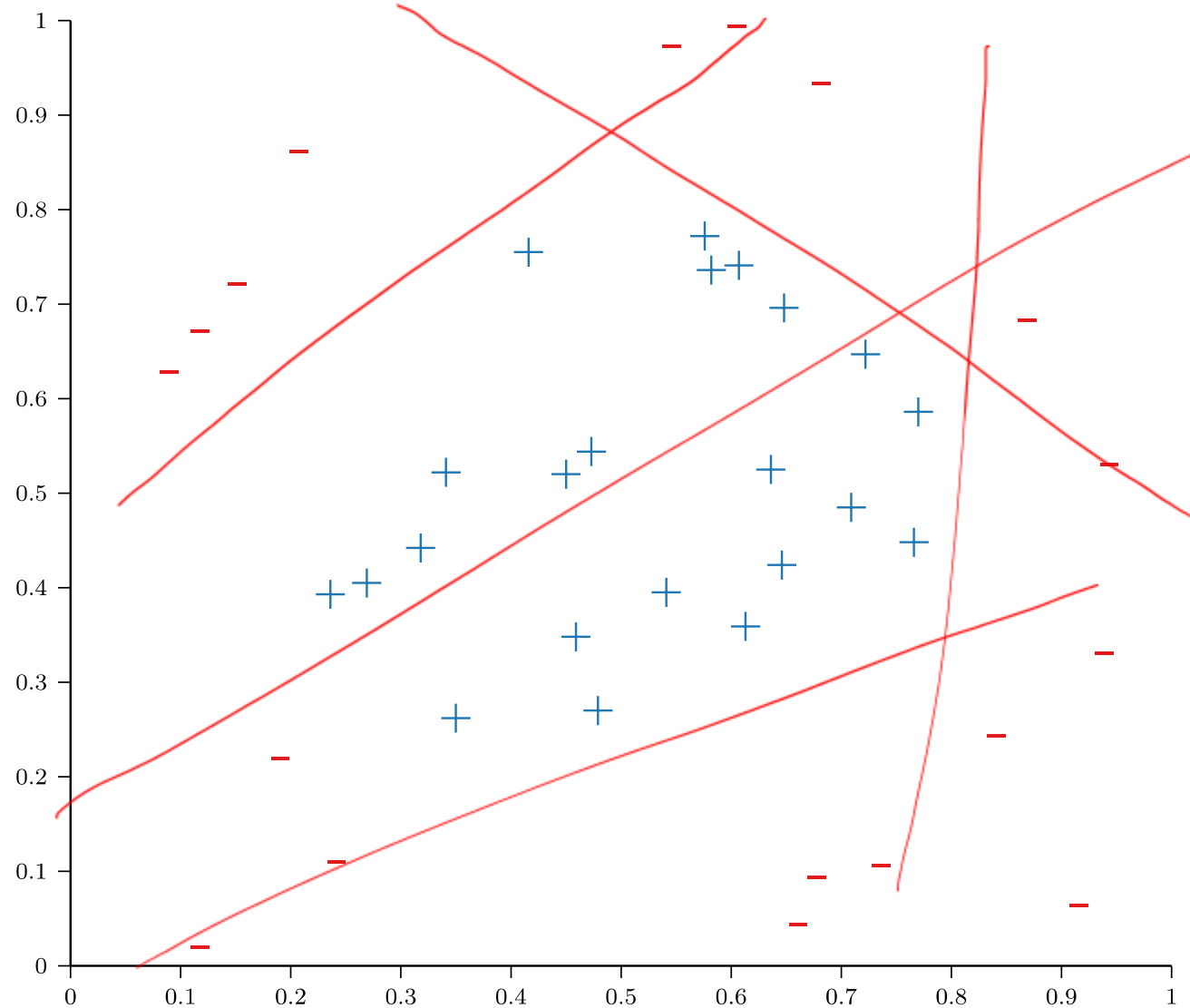
Linear Models



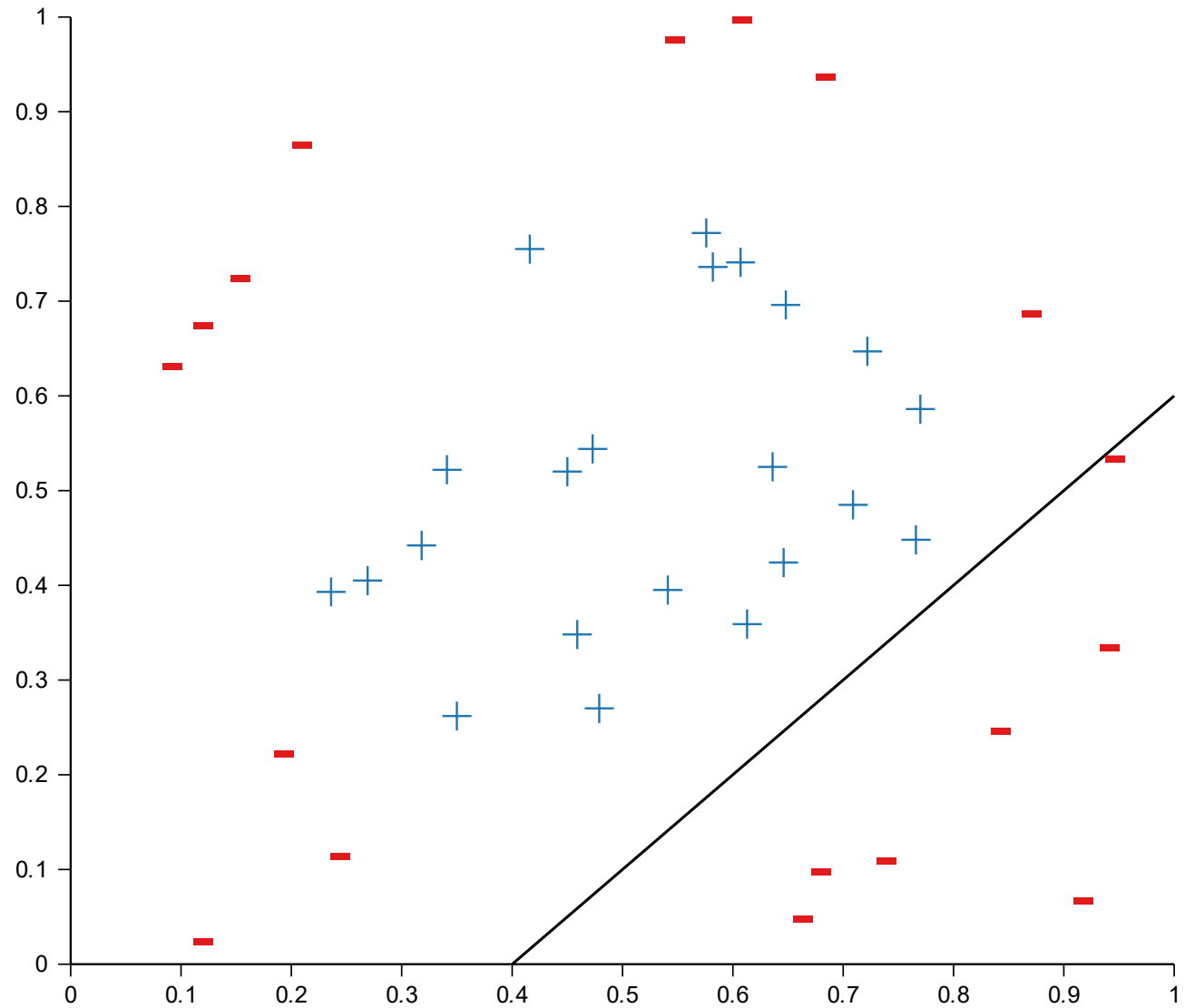
Linear Models



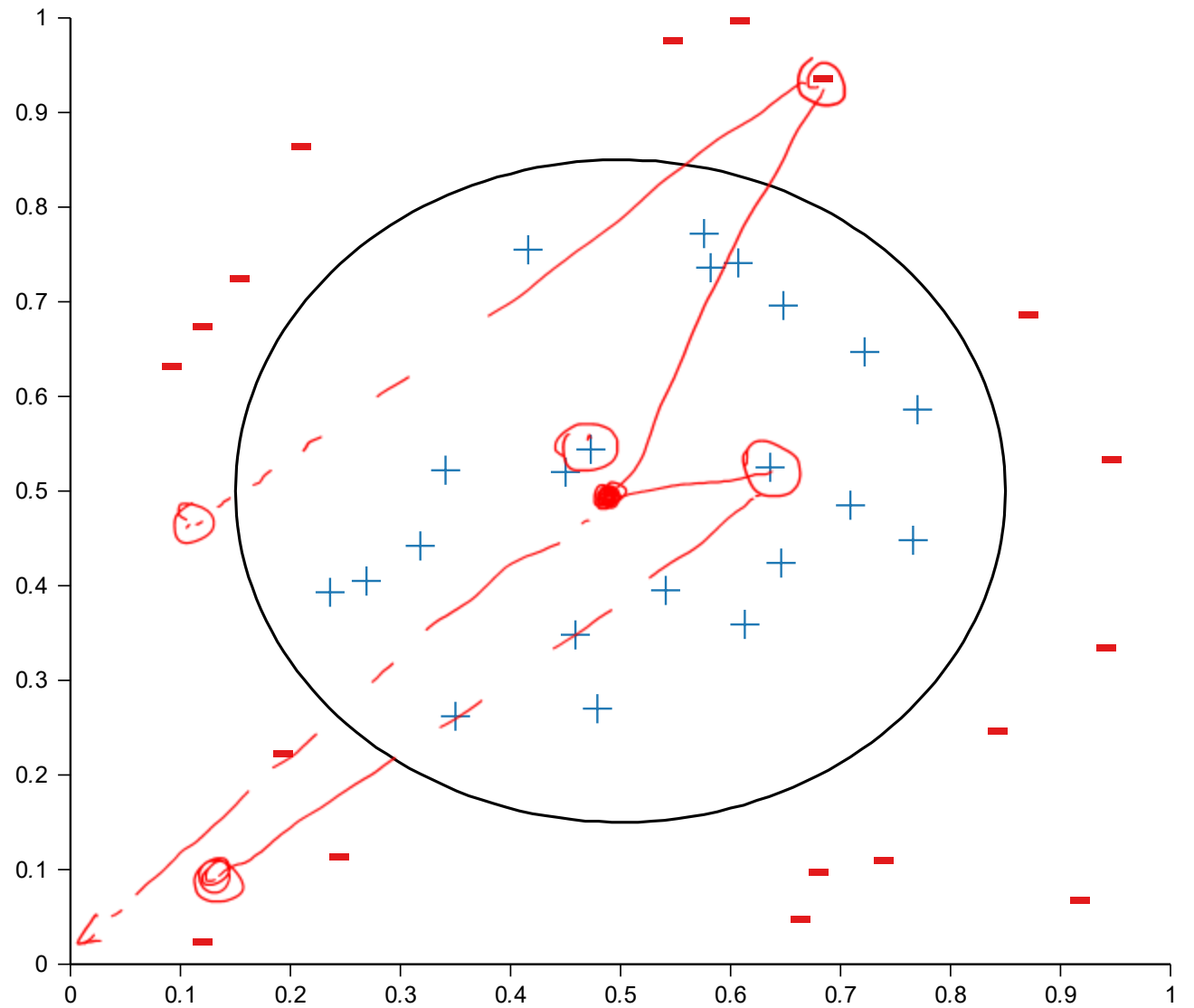
Linear Models?



Linear Models?



Nonlinear Models



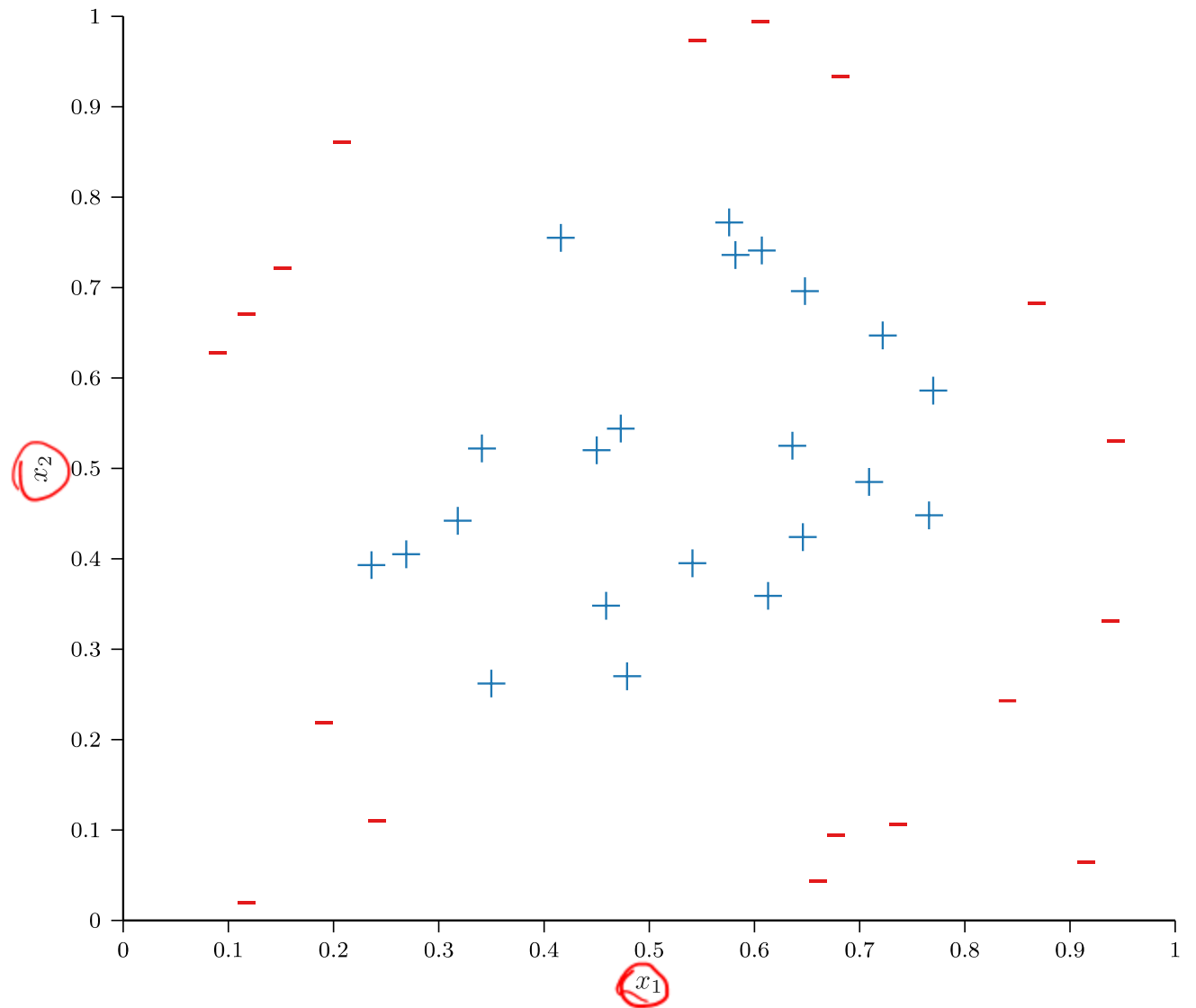
Feature Transforms

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

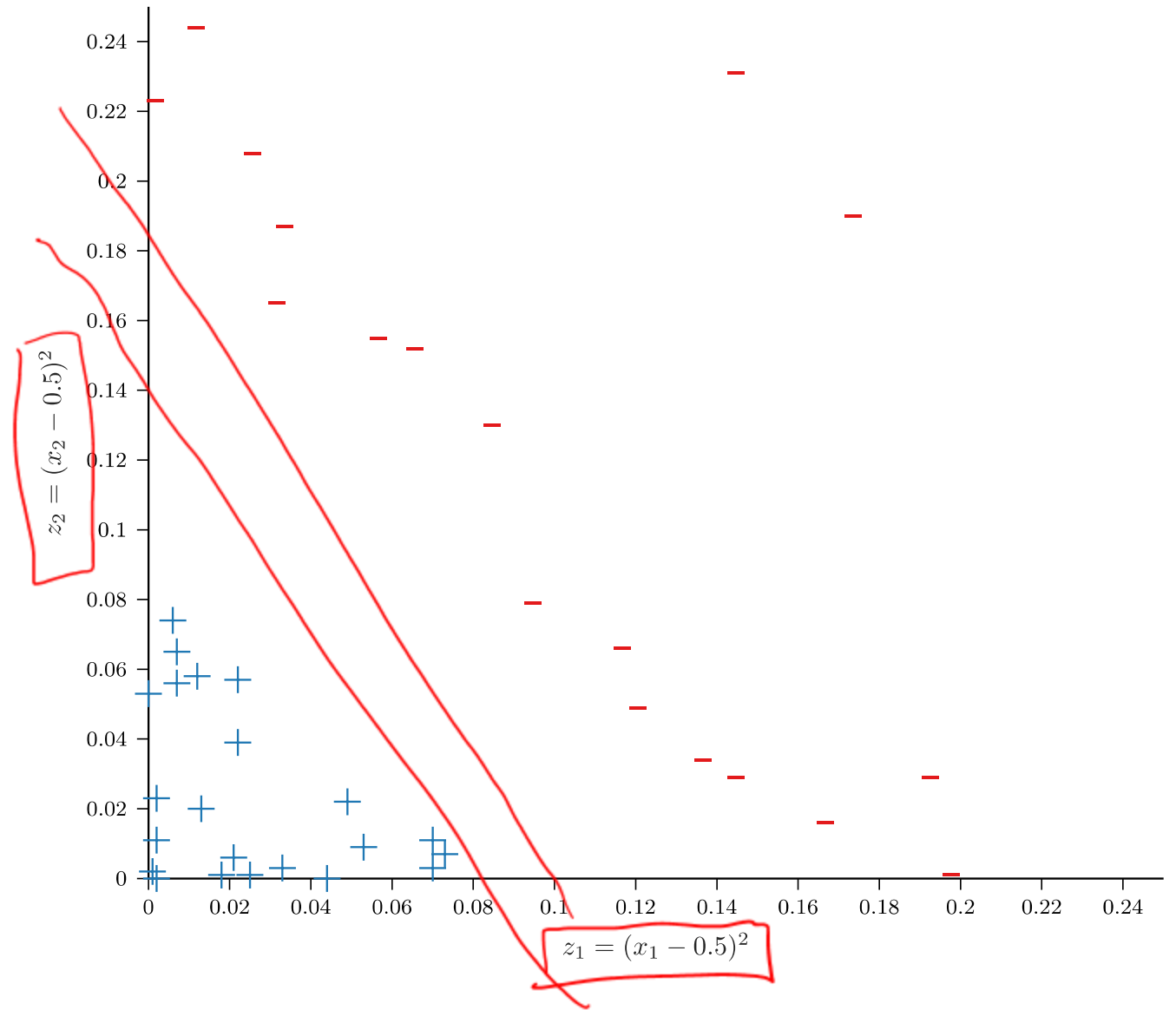
$$\phi([x_1, x_2]) = [z_1 = (x_1 - 0.5)^2, z_2 = (x_2 - 0.5)^2]$$

$(0.5, 0.5)$ is the center of ~~the~~ ^{red} circle.

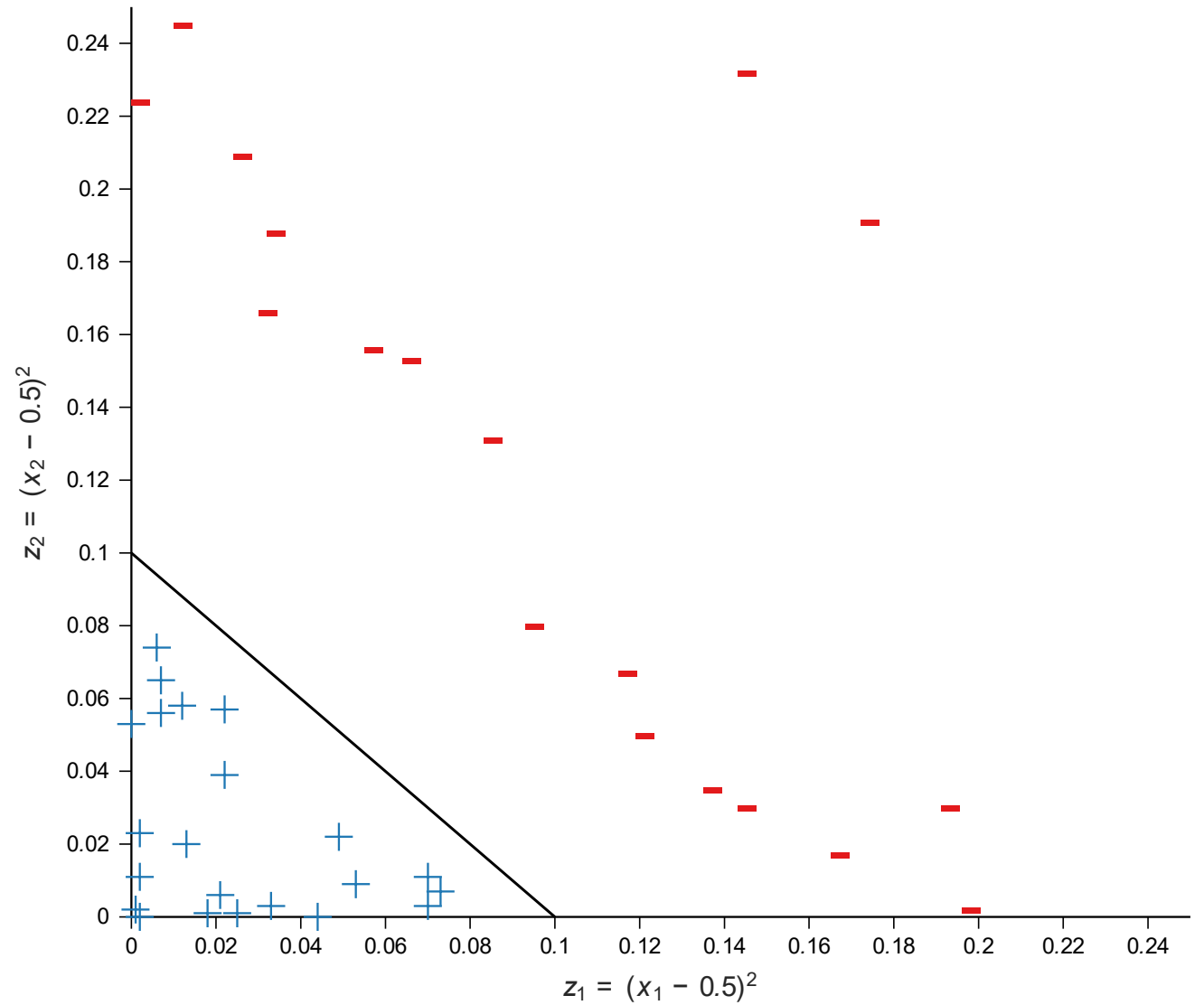
Nonlinear Models



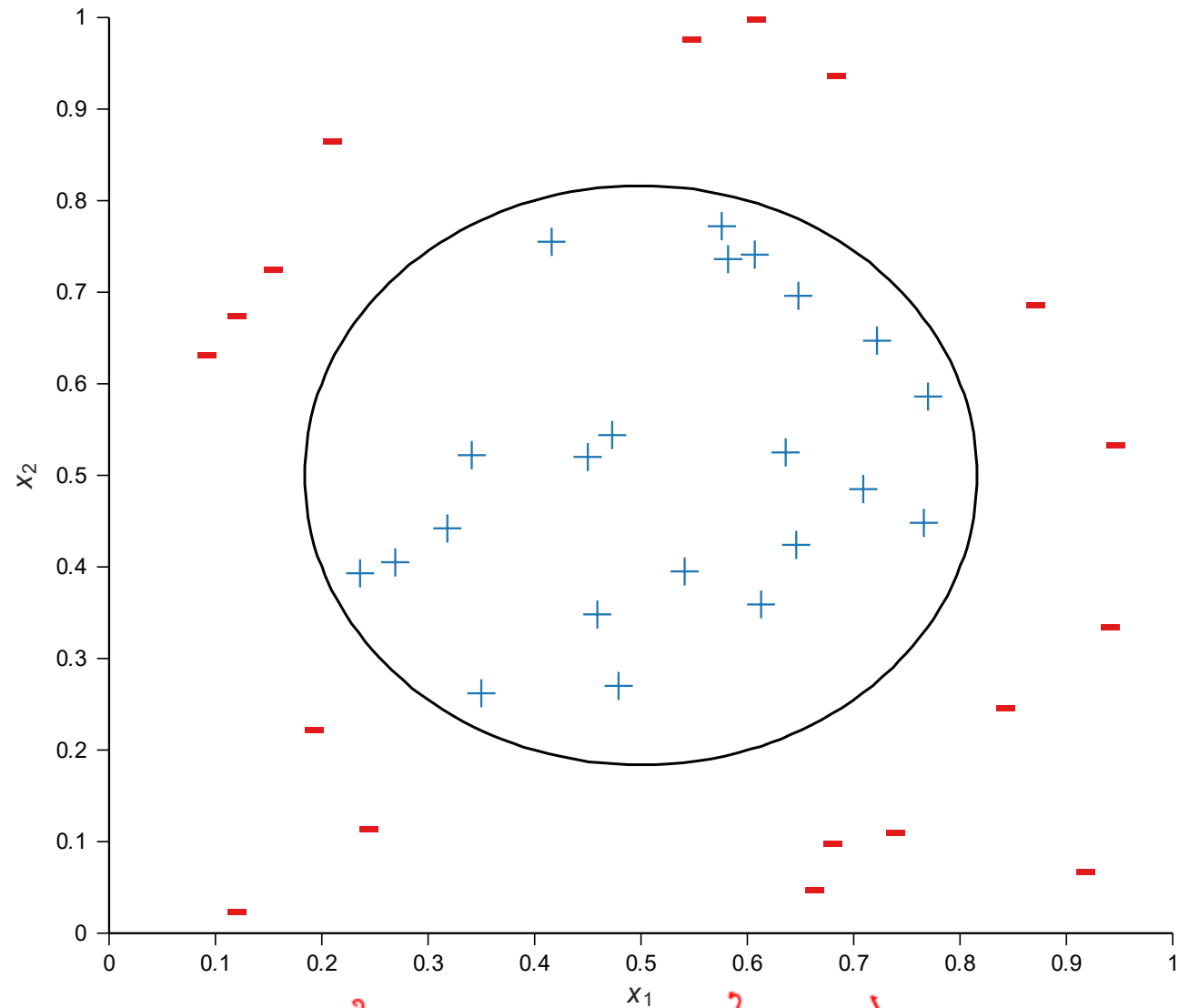
Nonlinear Models



Nonlinear Models



Nonlinear Models



$$\underbrace{w_1(x_1^2 + \dots)}_{z_1} + \underbrace{w_2(x_2^2 + \dots)}_{z_2} + w_0$$

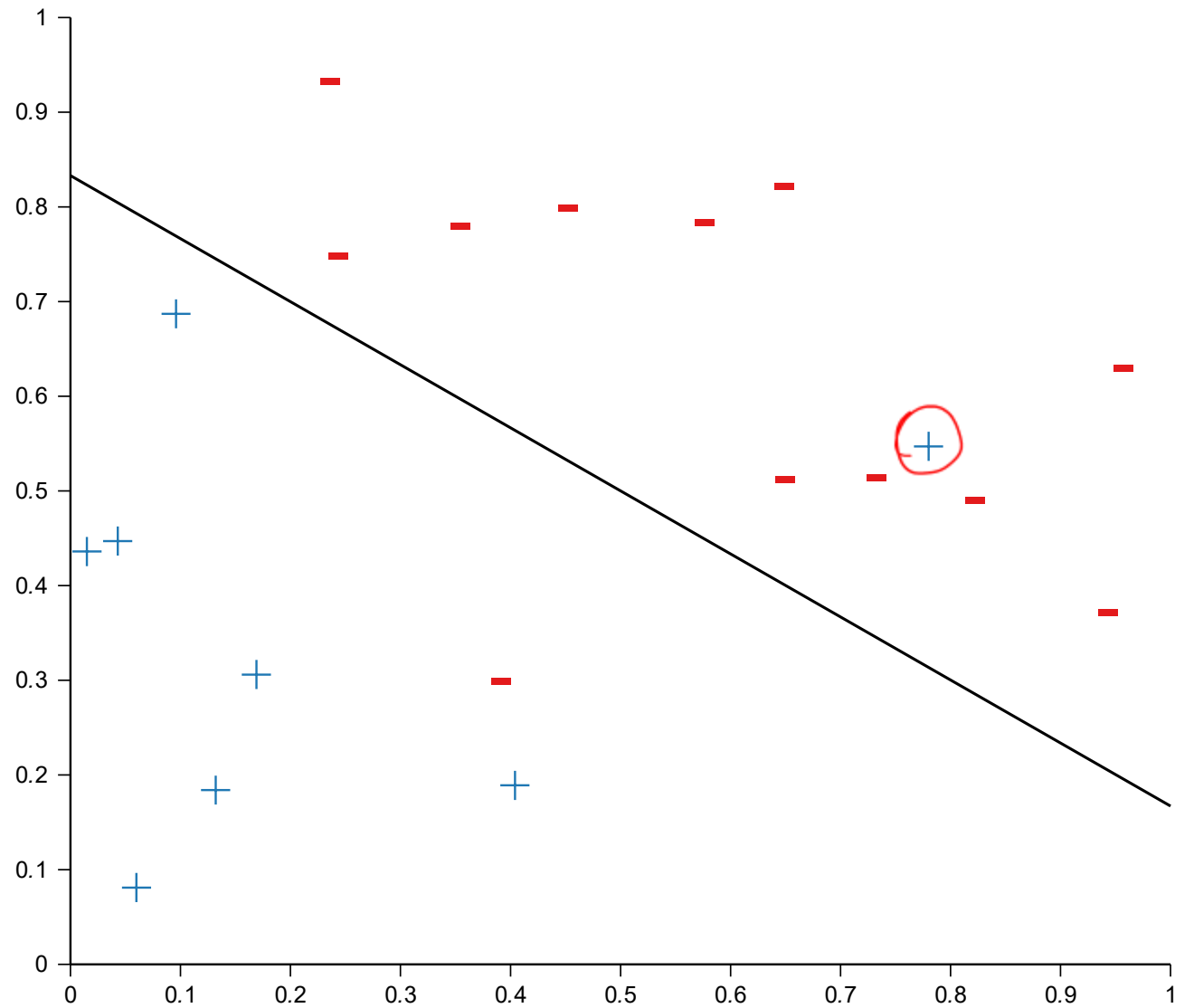
General Q^{th} -order Transforms

of features or D highest degree polynomial term allowed

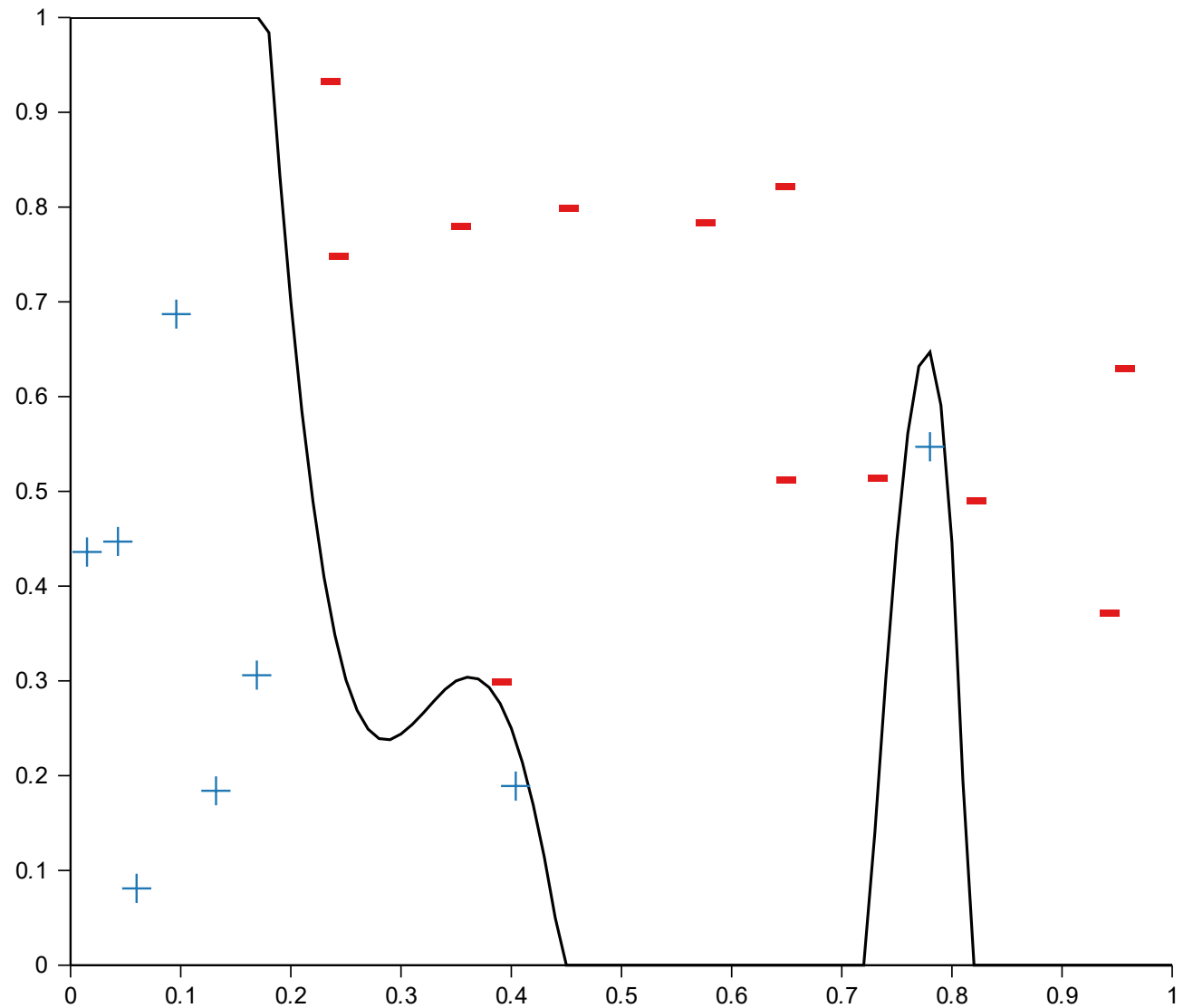
- $\phi_{2,2}([x_1, x_2]) = [\underbrace{x_1}, \underbrace{x_2}, \underbrace{x_1^2}, \underbrace{x_1 x_2}, \underbrace{x_2^2}]$
- $\phi_{2,3}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1 x_2, x_2^2, \underbrace{x_1^3}, \underbrace{x_1^2 x_2}, \underbrace{x_1 x_2^2}, \underbrace{x_2^3}]$
- $\phi_{2,4}([x_1, x_2])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

If feature space is D -dimensional, $\phi_{D,Q}$ is defined similarly
 how many features are in $\phi_{D,Q}$? $O(Q^D)$ $\binom{Q+D}{D}$

Linear Models



Nonlinear Models?



Feature Transforms: Tradeoffs

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

Test Error

Feature Transforms: Experiment

1-dim feature space
↘

- $x \in \mathbb{R}, y \in \mathbb{R}$ and $N = 20$ $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_{20}, y_{20})\}$
- Targets are generated by a 10th-order polynomial in x with additive Gaussian noise:

$$y_i = \underbrace{\sum_{d=0}^{10} a_d x_i^d}_{\text{polynomial}} + \underbrace{\epsilon}_{\text{noise}} \text{ where } \epsilon \sim \underbrace{N(0, \sigma^2)}_{\text{Gaussian noise}}$$

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

In-Class Poll:

- $x \in \mathbb{R}, y \in \mathbb{R}$ and $N = 20$

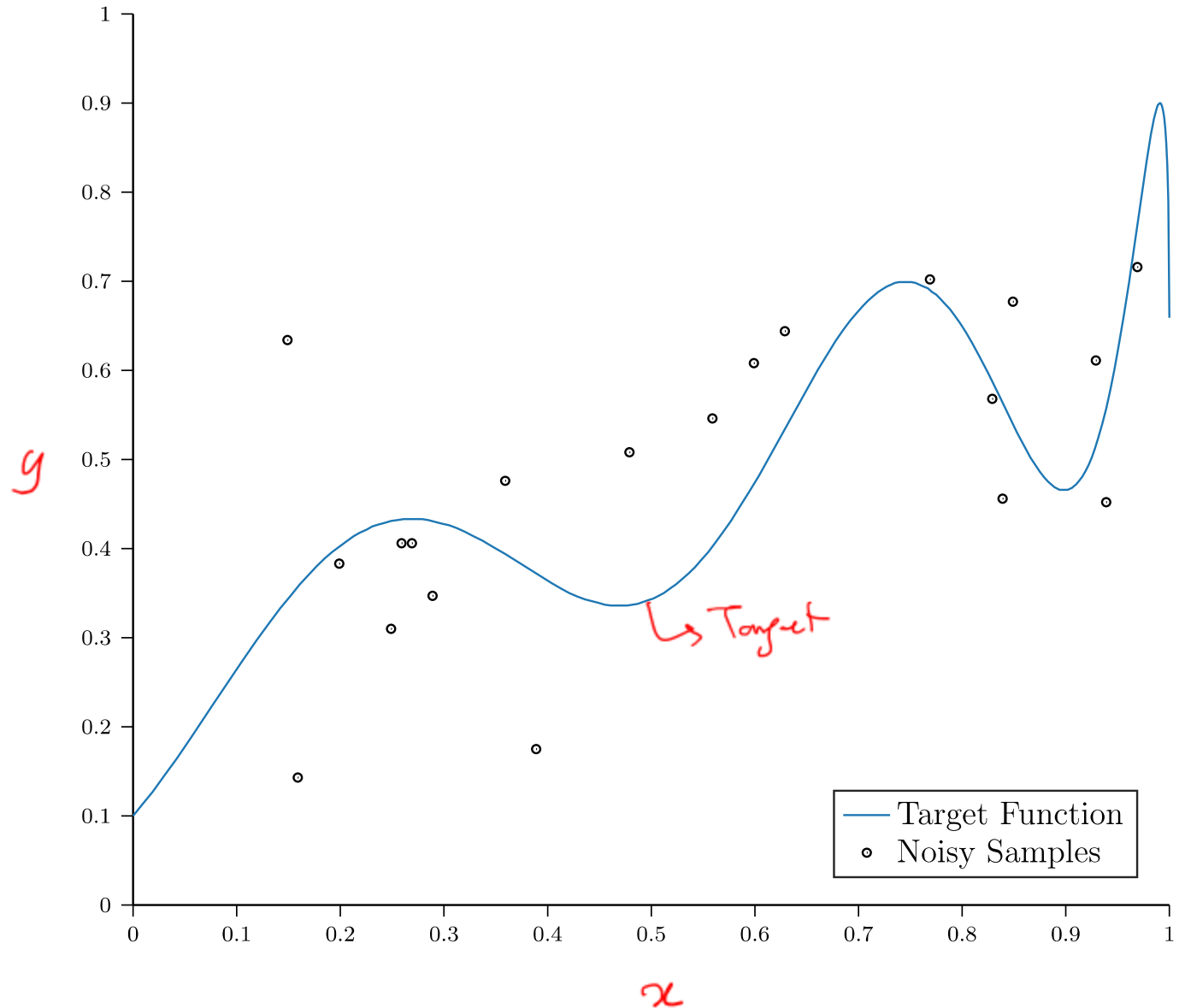
$$y = \sum_{d=0}^{10} a_d x^d + \epsilon \text{ where } \epsilon \sim N(0, \sigma^2)$$

- $\mathcal{H}_2 = 2^{\text{nd}}$ -order polynomials
- $\mathcal{H}_{10} = 10^{\text{th}}$ -order polynomials

Which hypothesis space (\mathcal{H}_2 or \mathcal{H}_{10}) is a better choice to fit to this data set?

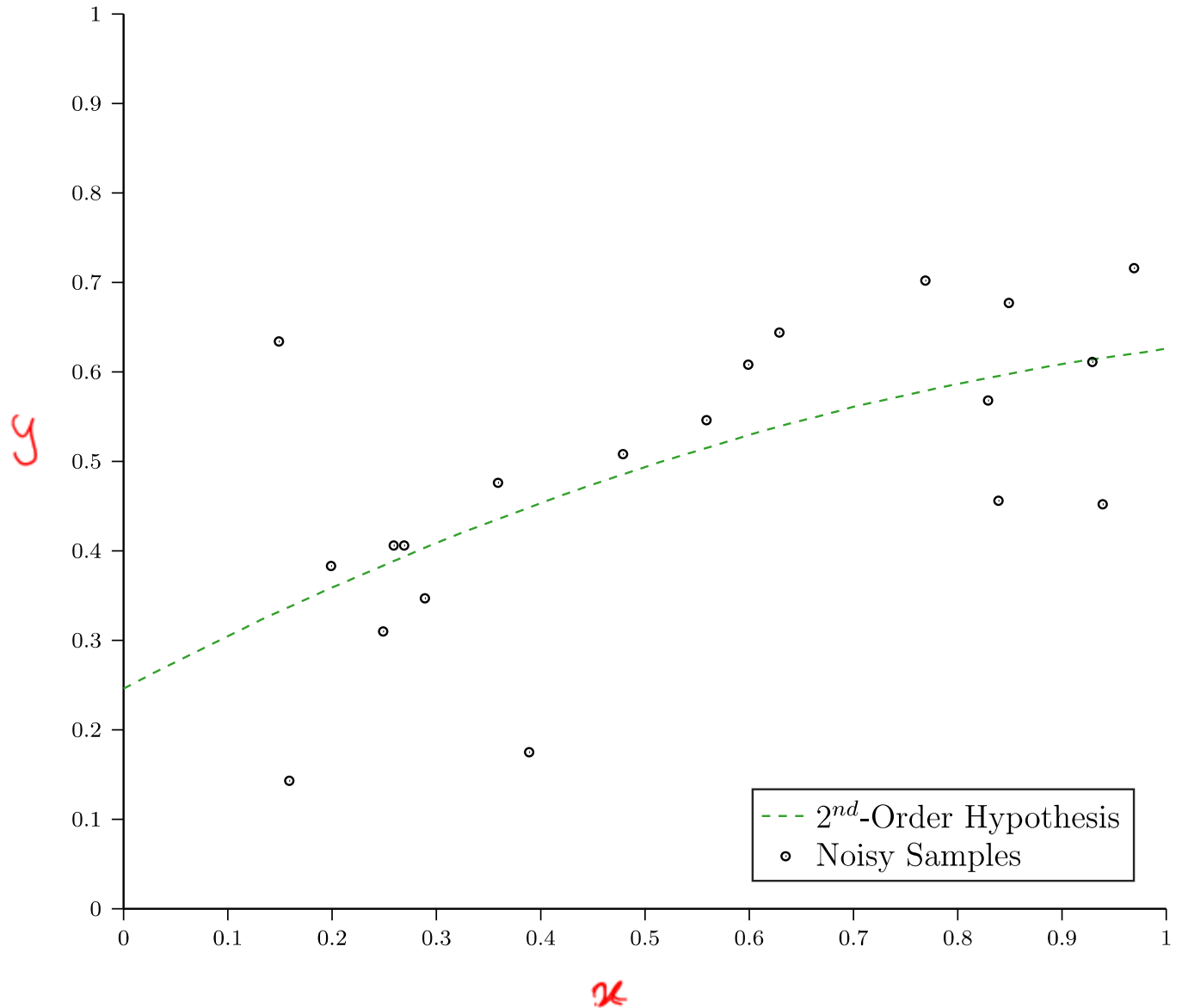
Noisy Targets

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



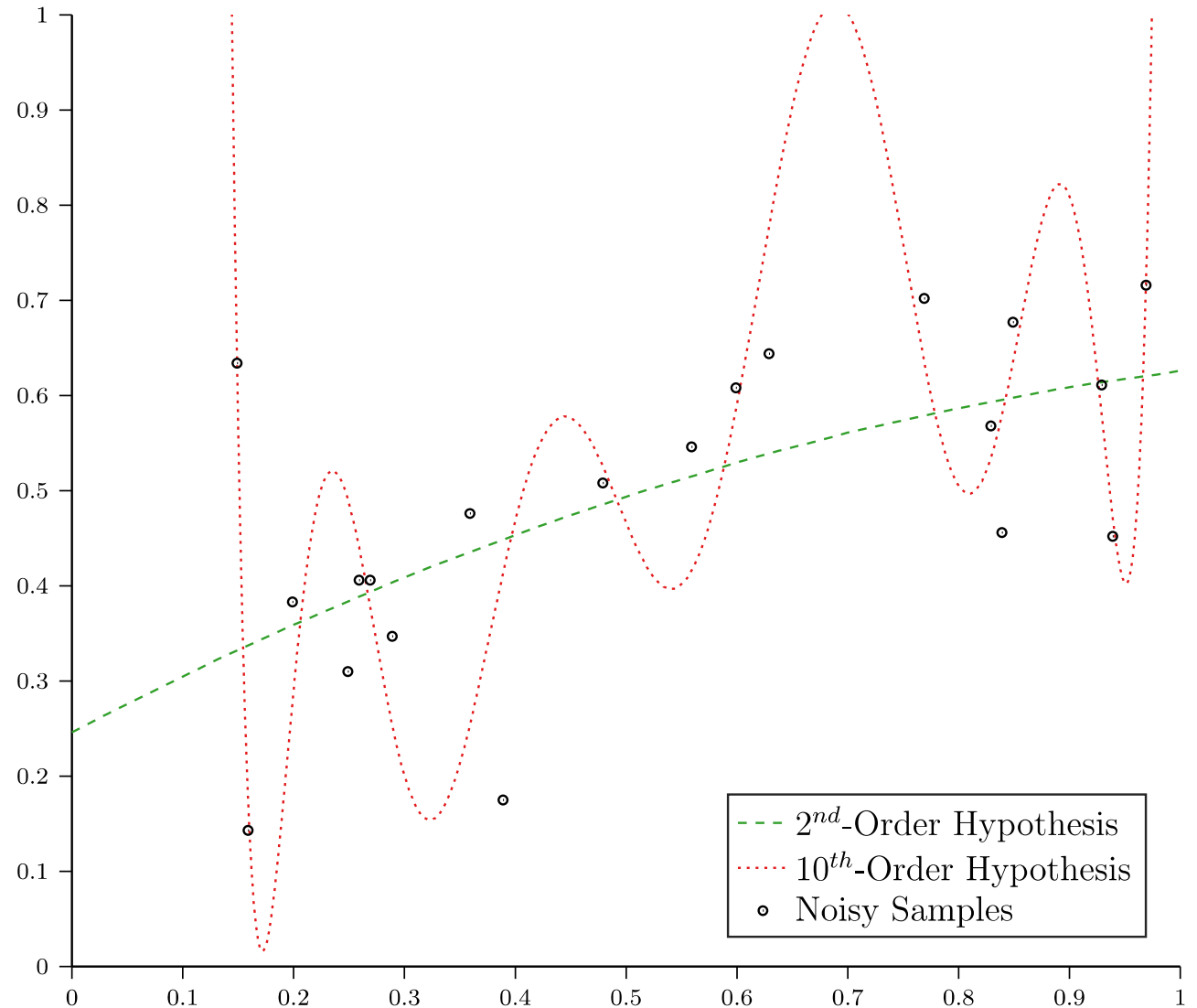
Noisy Targets

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



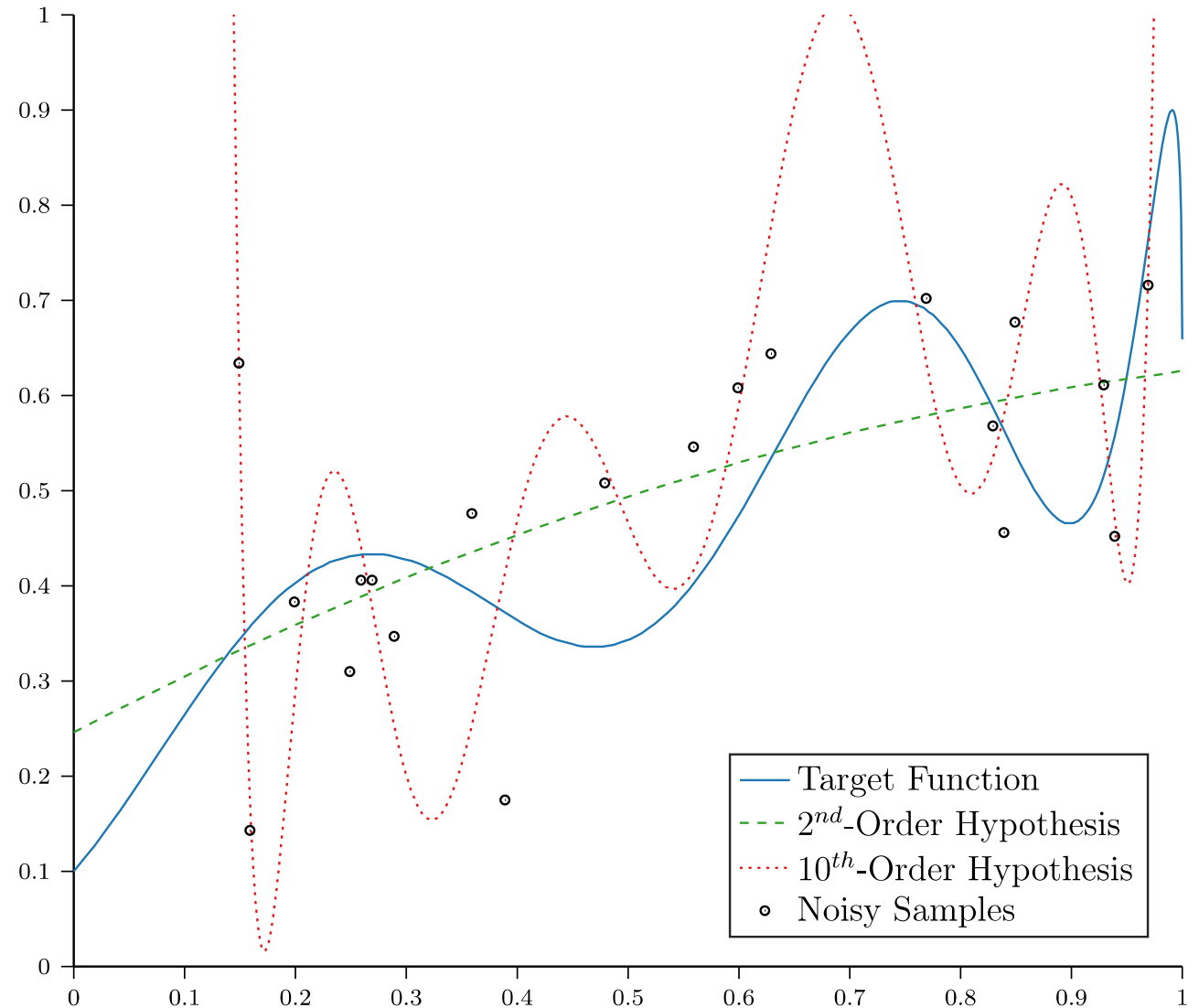
Noisy Targets

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



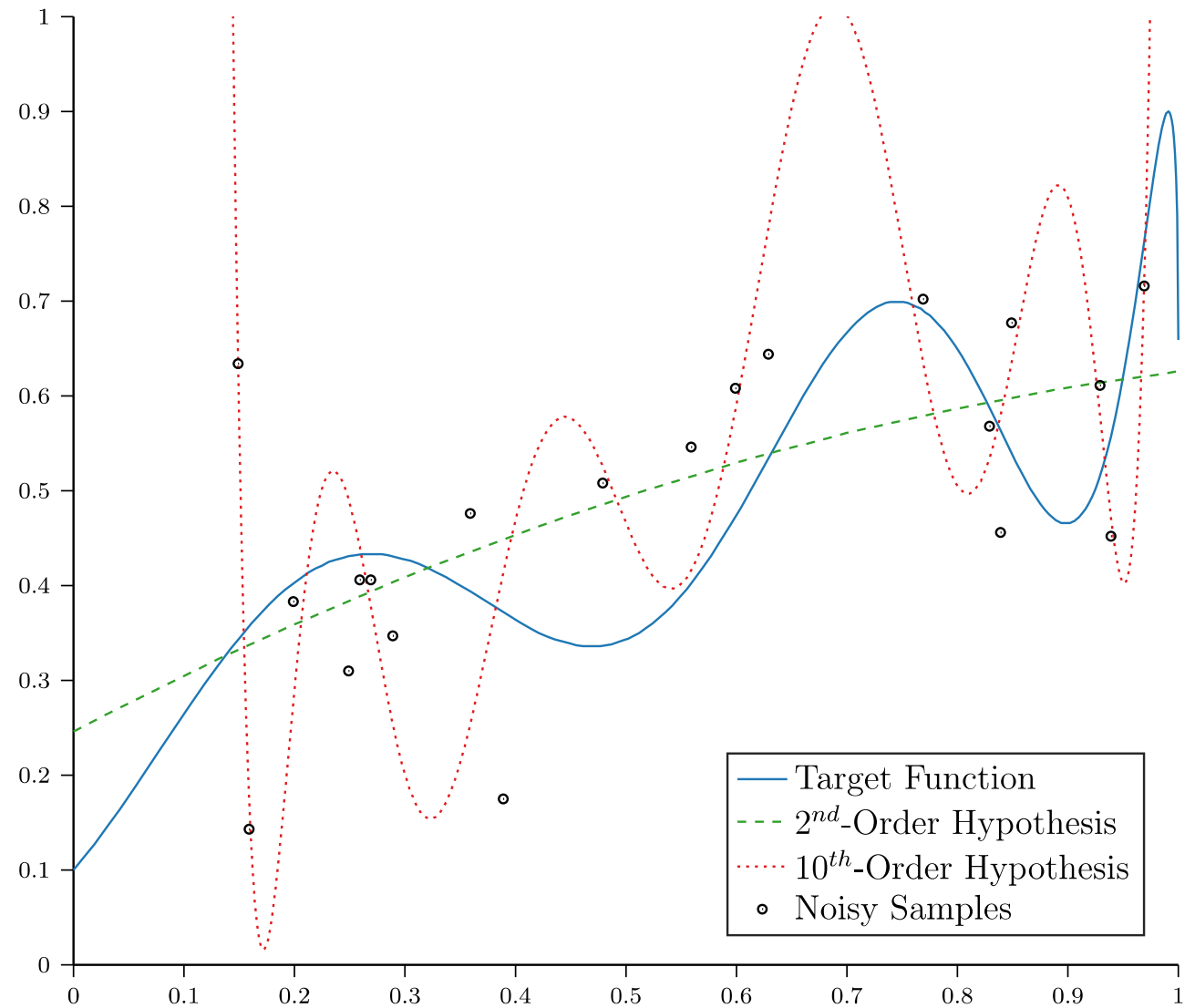
Noisy Targets

- 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

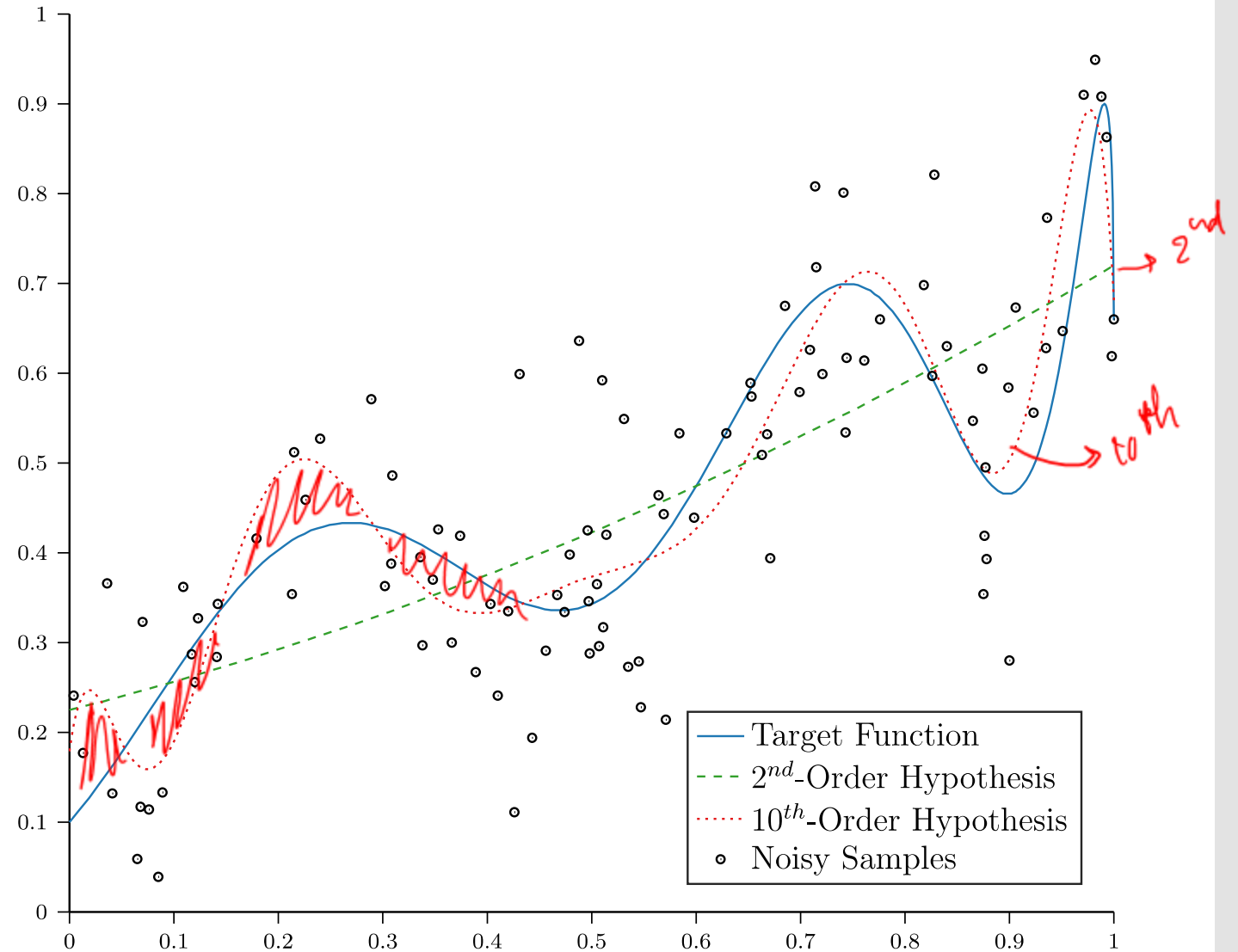
	\mathcal{H}_2	\mathcal{H}_{10}
Training Error	0.016	0.011
True Error	0.009	3797



More Data

	\mathcal{H}_2	\mathcal{H}_{10}
Training Error	<u>0.018</u>	<u>0.010</u>
True Error	<u>0.009</u>	<u>0.003</u>

test error



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 $\mathcal{D} = \{(\underline{x_1}, y_1), \dots, (\underline{x_{20}}, y_{20})\}$
- $\phi_{1,10}(x) = [x, x^2, x^3, x^4, x^5, x^6, x^7, x^8, x^9, x^{10}]$

- 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

$$\min_{\boldsymbol{\omega}} \ell_{\mathcal{D}}(\boldsymbol{\omega}) = \frac{1}{N} (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})$$

- Subject to

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

$$\frac{1}{N} \sum_{n=1}^N \left(\left(\sum_{d=0}^9 x_d^{(n)} \omega_d \right) - y^{(n)} \right)^2$$

- 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

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

- Subject to nothing!

Soft Constraints

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

transformed design matrix

←

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

find $\boldsymbol{\omega}$ that minimizes

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

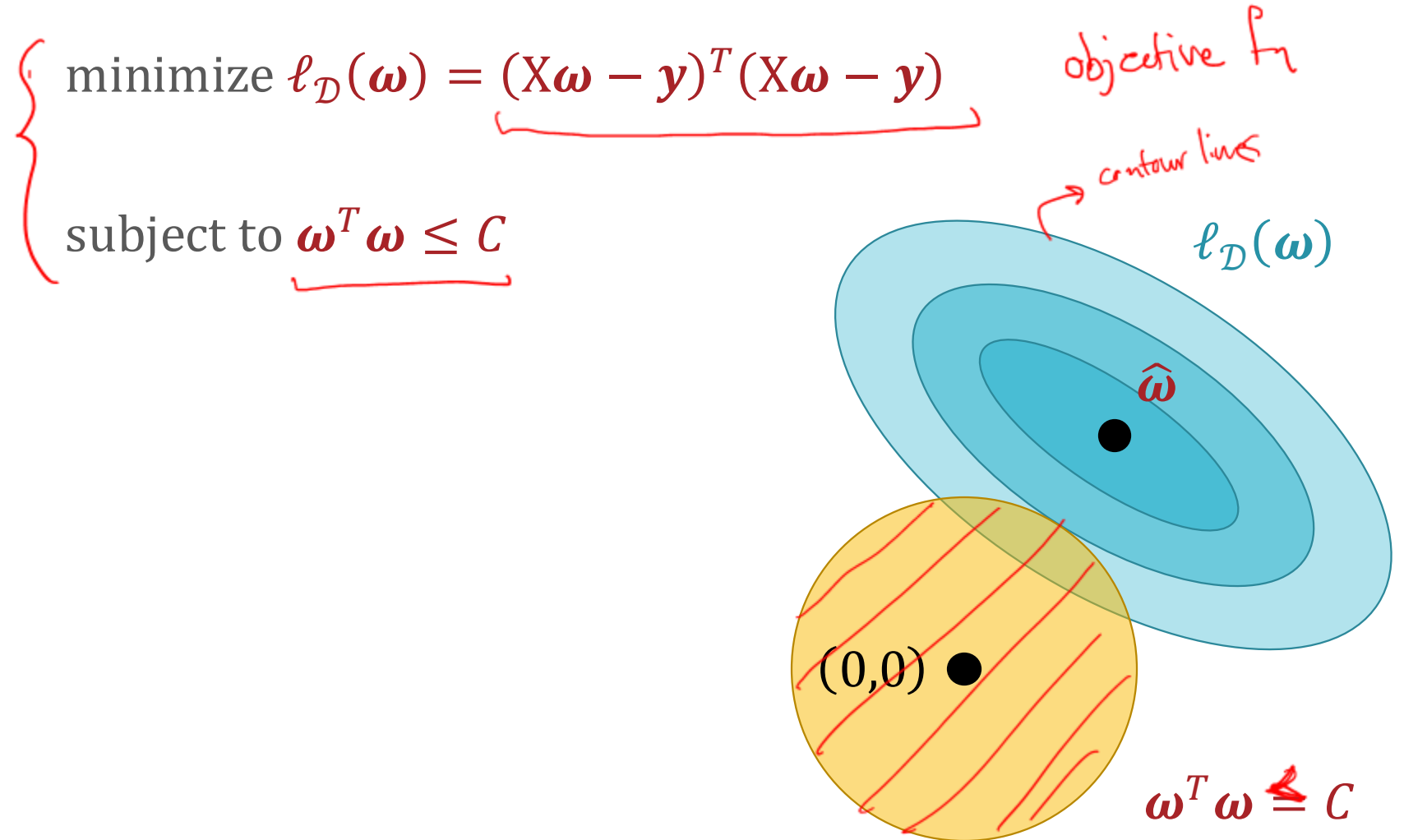
- Subject to: Some constraint on $\boldsymbol{\omega}$, e.g.,

Ridge ~~regression~~ regularizer

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

Soft Constraints

constrained optimization

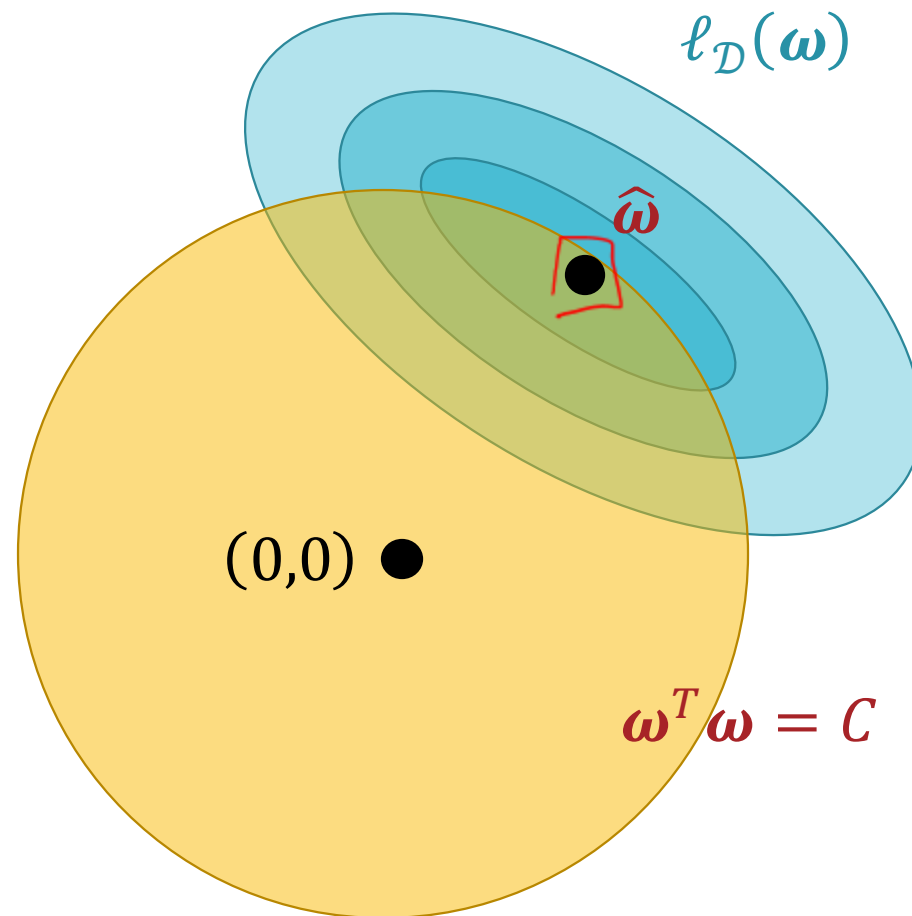


Soft Constraints

$$\ell_D(\omega) + 0.00001 \lambda^T \omega$$

$$\text{minimize } \ell_D(\omega) = (X\omega - y)^T (X\omega - y)$$

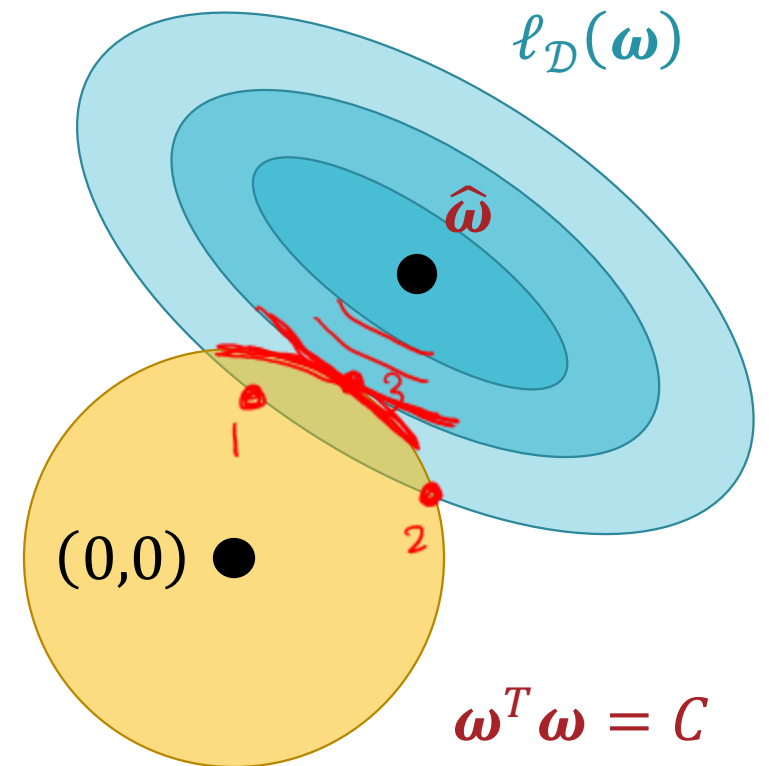
$$\text{subject to } \omega^T \omega \leq C$$



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

$$\left\{ \begin{array}{l} \text{minimize } \ell_D(\omega) = (X\omega - y)^T(X\omega - y) \\ \text{subject to } \omega^T \omega \leq C \end{array} \right.$$

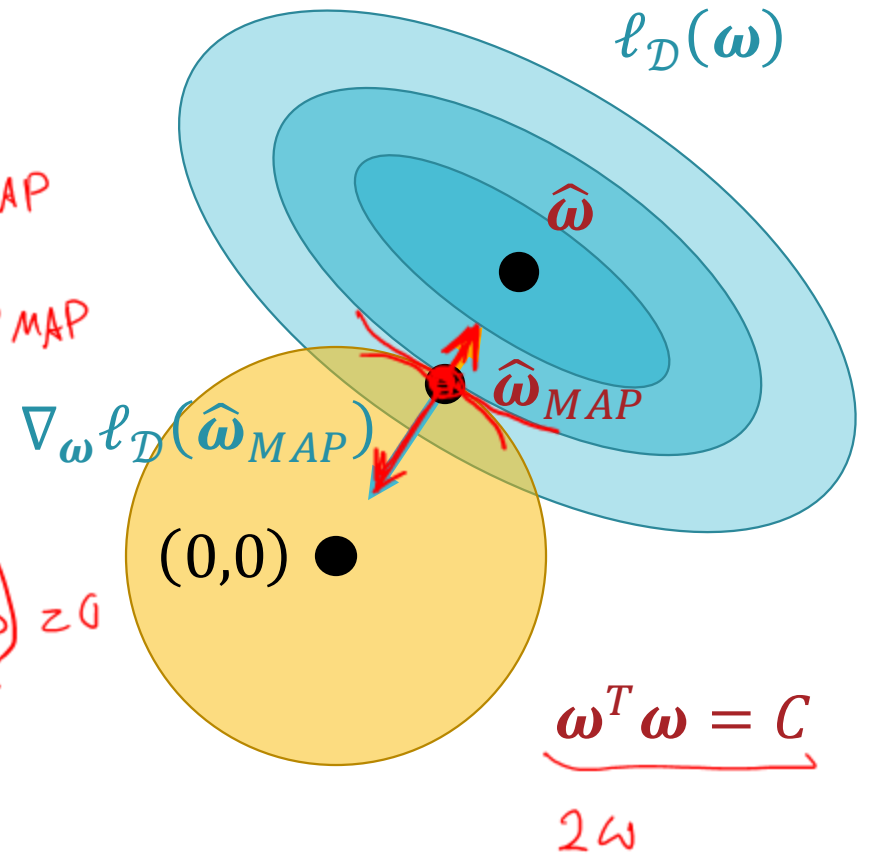
$$\nabla_{\omega} \ell_D(\hat{\omega}_{MAP}) \propto -2\hat{\omega}_{MAP}$$

$$\Leftrightarrow \nabla_{\omega} \ell_D(\hat{\omega}_{MAP}) = -2\lambda_C \hat{\omega}_{MAP}$$

$$\Leftrightarrow \nabla_{\omega} \ell_D(\hat{\omega}_{MAP}) + 2\lambda_C \hat{\omega}_{MAP} = 0$$

$$\Leftrightarrow \nabla_{\omega} \left(\ell_D(\hat{\omega}_{MAP}) + \lambda_C \hat{\omega}_{MAP}^T \hat{\omega}_{MAP} \right) = 0$$

new objective function
for unconstrained problem.



Soft Constraints

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

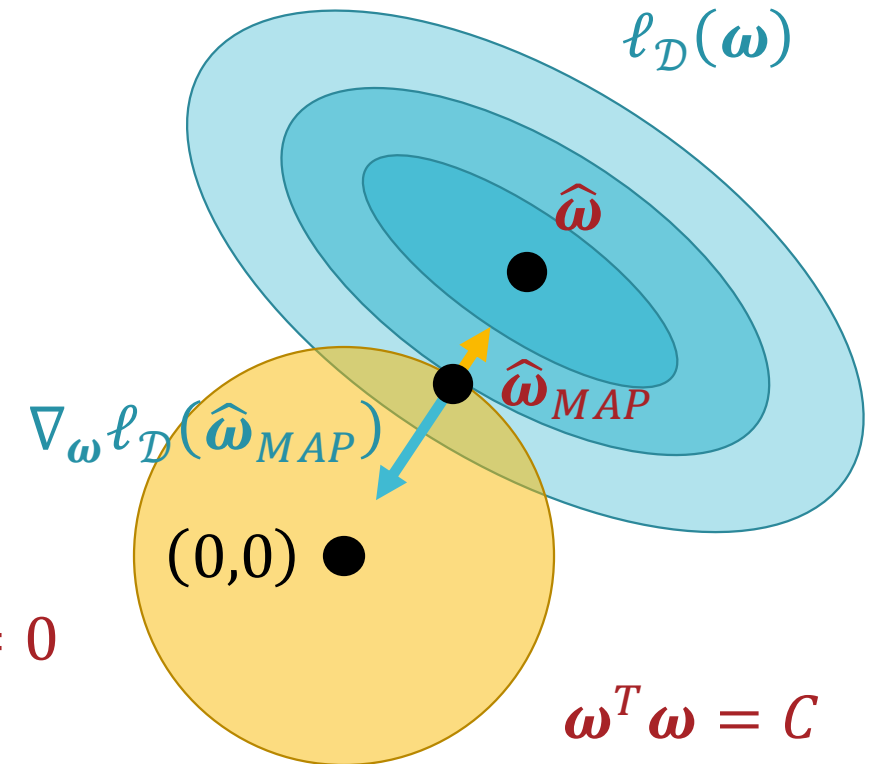
$$\text{subject to } \boldsymbol{\omega}^T \boldsymbol{\omega} \leq C$$

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

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

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

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



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

$$\left\{ \begin{array}{l} \text{minimize } \ell_D(\omega) = (X\omega - y)^T(X\omega - y) \\ \text{subject to } \omega^T \omega \leq C \end{array} \right. \quad (\star)$$



$$\text{minimize } \ell_D^{AUG}(\omega) = \ell_D(\omega) + \lambda_c \omega^T \omega \quad (\star) \quad \leftarrow$$

$$\nabla_{\omega} \ell_D^{AUG}(\omega) = \left(2X^T X \omega \quad -2X^T y + 2\lambda_c \omega \right)$$

$$\Rightarrow 2X^T X \hat{\omega}_{MAP} - 2X^T y + 2\lambda_c \hat{\omega}_{MAP} = 0$$

$$\Rightarrow \underbrace{X^T X \hat{\omega}_{MAP}} + \underbrace{\lambda_c \hat{\omega}_{MAP}} = X^T y \quad \rightarrow \equiv \lambda_c I_{D+1} \hat{\omega}_{MAP}$$

$$\Rightarrow \underbrace{(X^T X + \lambda_c I_{D+1})}_{\wedge} \hat{\omega}_{MAP} = X^T y$$

$$\Rightarrow \hat{\omega}_{MAP} = (X^T X + \lambda_c I_{D+1})^{-1} X^T y$$

Ridge Regression

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

Ridge Regression

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

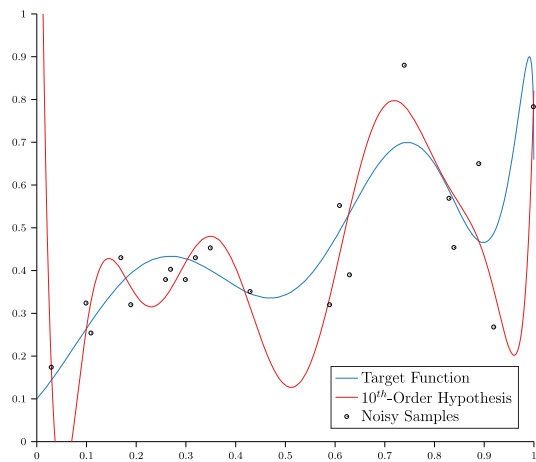
$$\nabla_{\boldsymbol{\omega}} \ell_D^{AUG}(\boldsymbol{\omega}) = 2(X^T X \boldsymbol{\omega} - X^T \mathbf{y} + \lambda_C \boldsymbol{\omega})$$

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

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

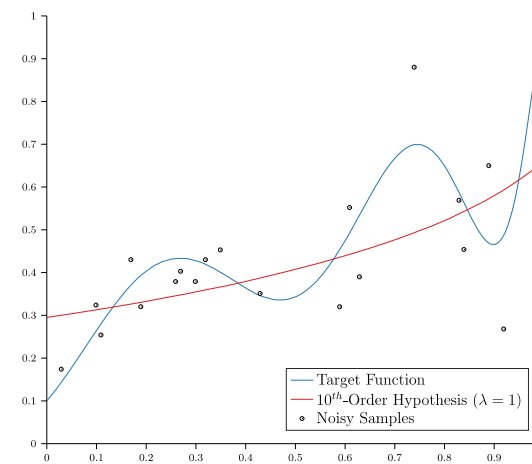
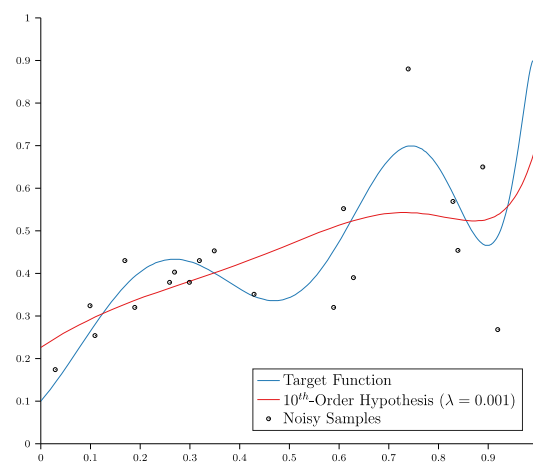
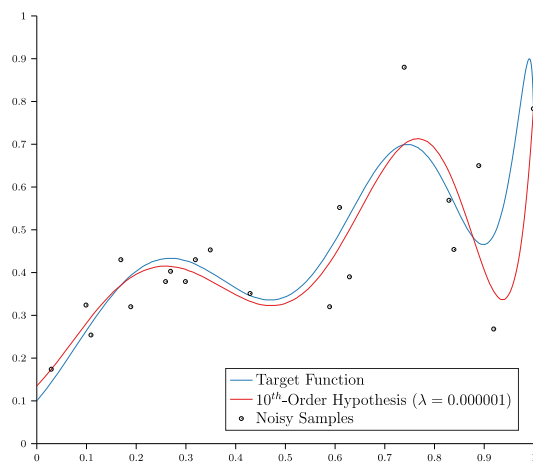
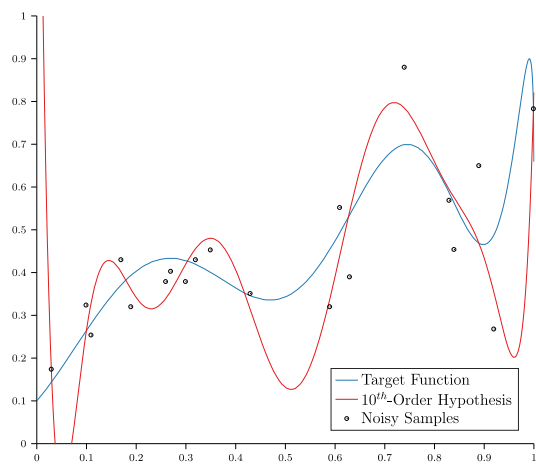
$$\hat{\boldsymbol{\omega}}_{MAP} = \underbrace{(X^T X + \lambda_C I_{D+1})^{-1}} X^T \mathbf{y}$$

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



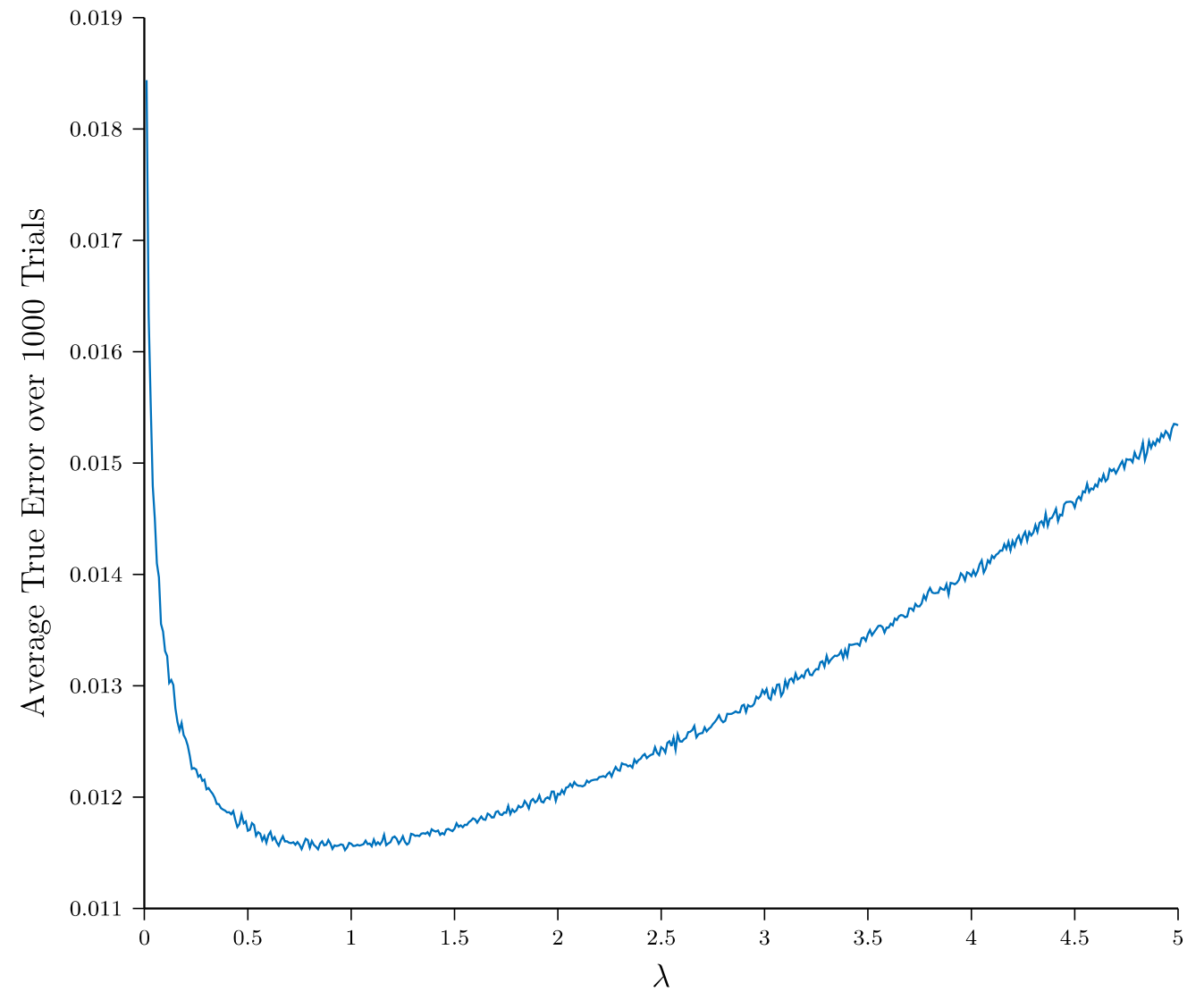
Ridge Regression

$$\lambda_c = 0$$

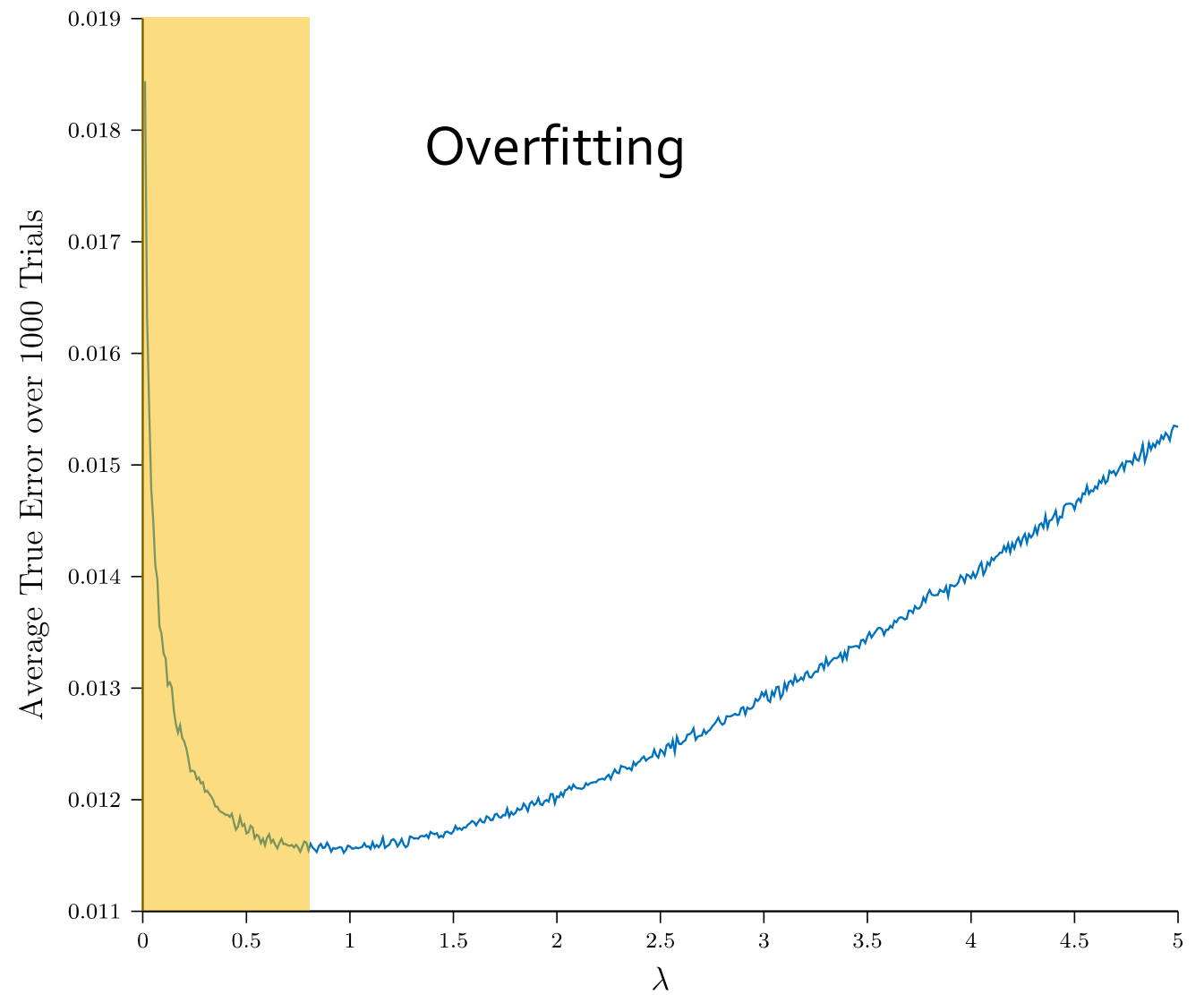
True
Error 0.059

Overfit

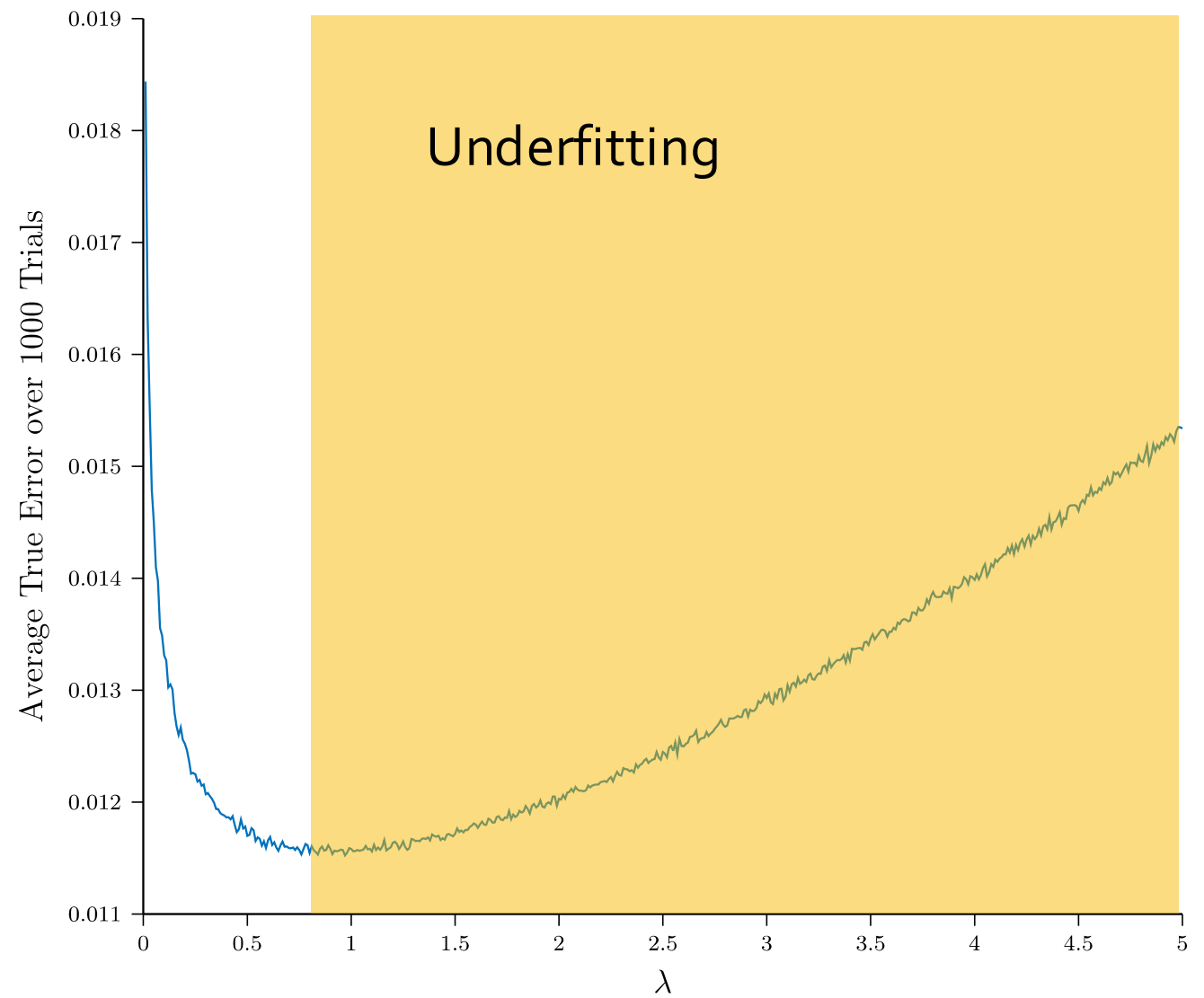
Setting λ



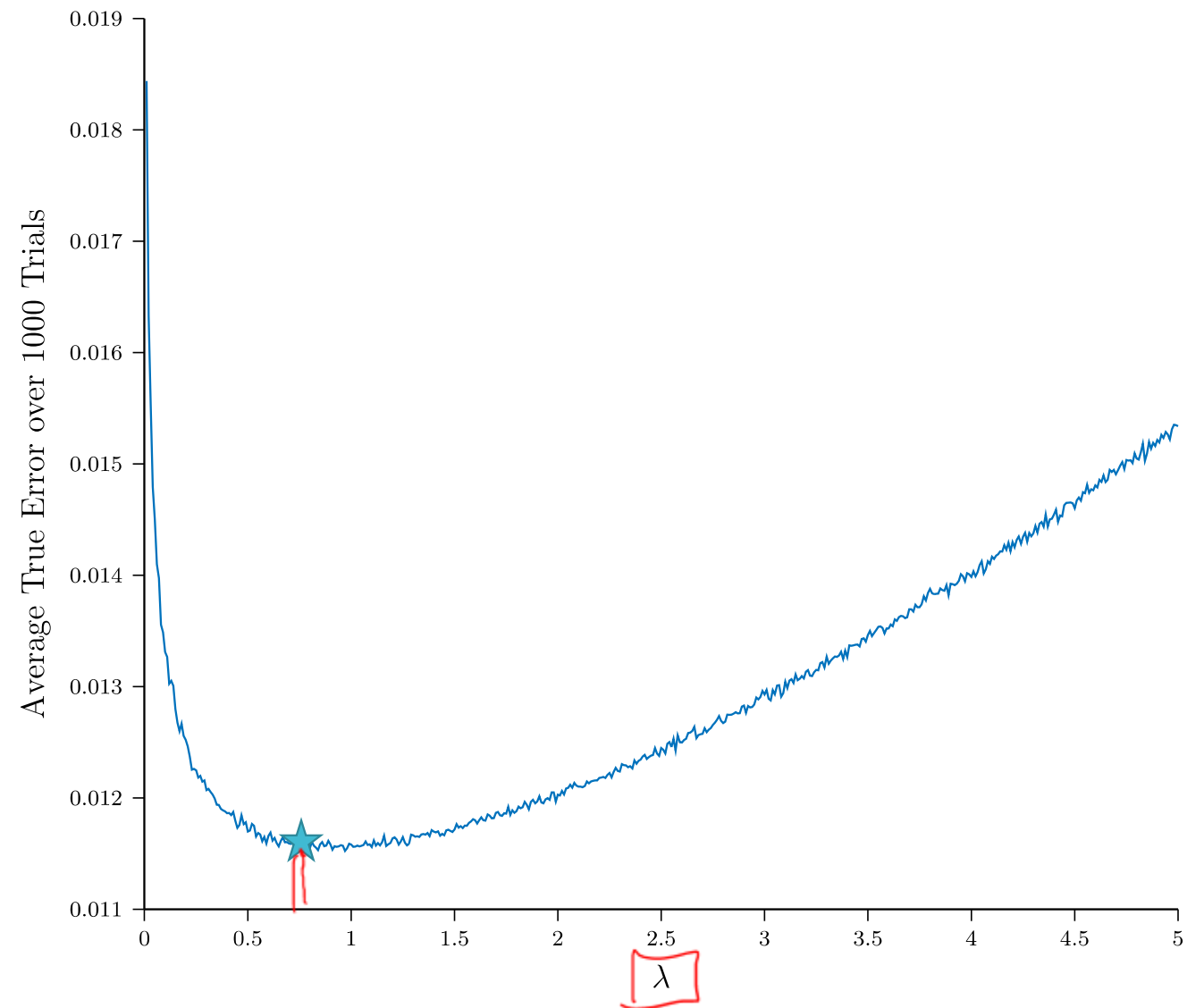
Setting λ



Setting λ

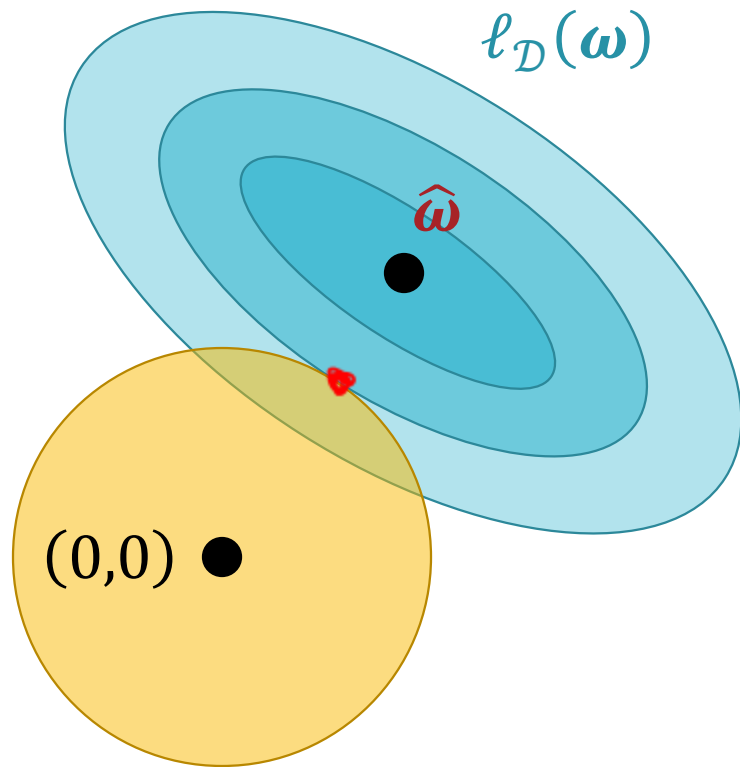


Setting λ

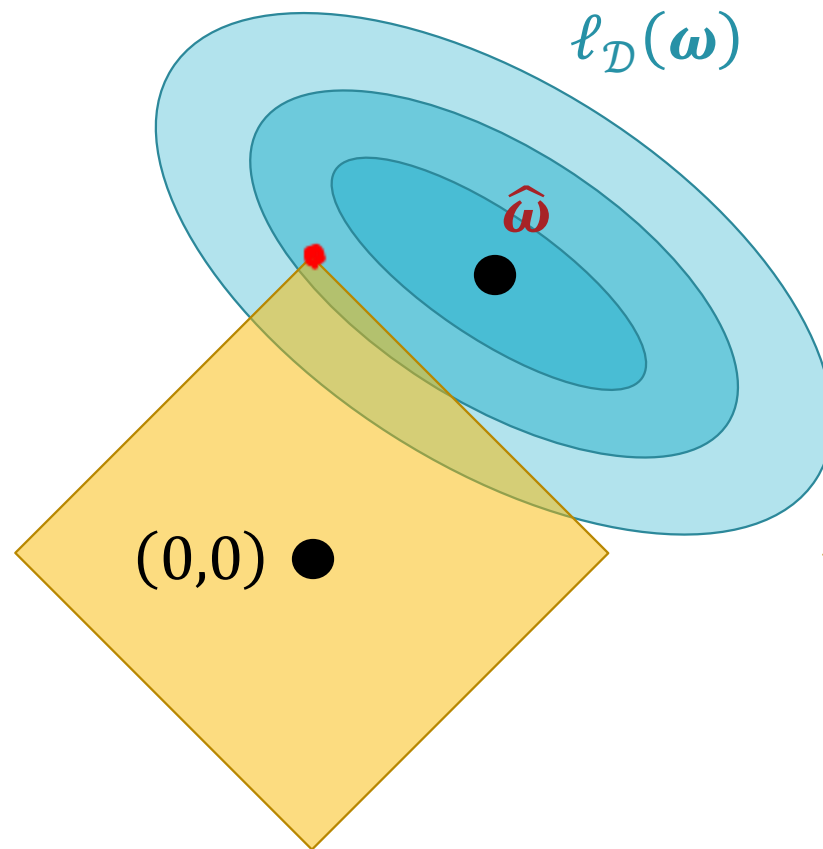


Other Regularizers

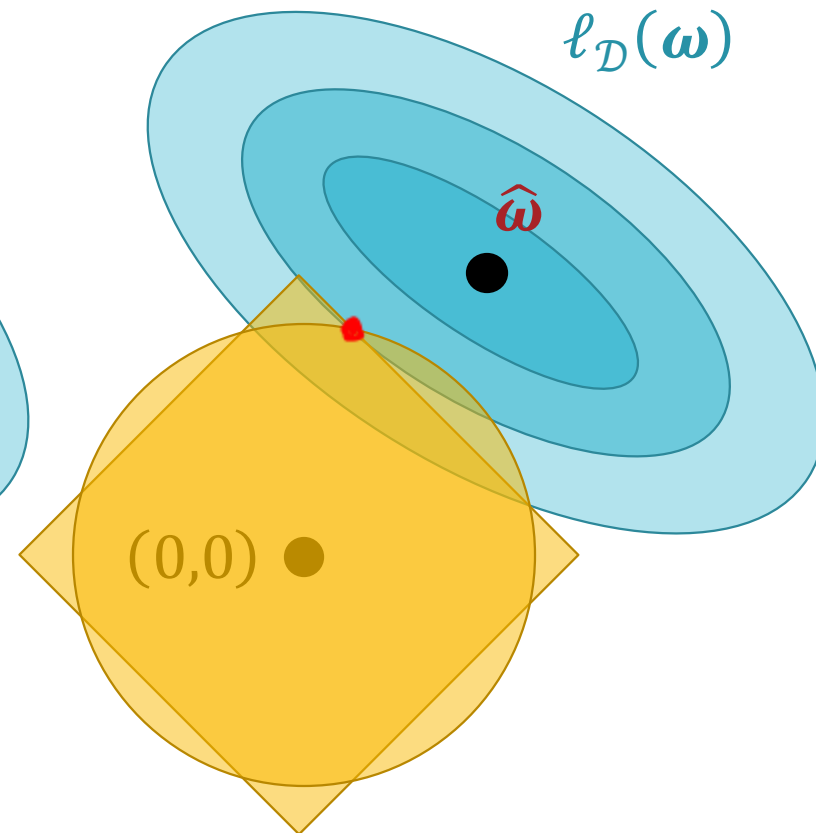
$\ell_D(\omega) + \lambda r(\omega)$		
Ridge or $L2$	$r(\omega) = \underbrace{\ \omega\ _2^2}_{\text{red underline}} = \sum_{d=0}^D \omega_d^2$	Encourages small weights
<u>Lasso or $L1$</u>	$r(\omega) = \underbrace{\ \omega\ _1}_{\text{red underline}} = \underbrace{\sum_{d=0}^D \omega_d }_{\text{red underline}}$	Encourages sparsity
<u>Elastic Net $L1 + L2$</u>	$r(\omega) = \lambda \ \omega\ _1 + \lambda \ \omega\ _2^2$	Encourages sparsity & shrinkage



Ridge or L_2



Lasso or L_1



Elastic Net

Other Regularizers

M(C)LE for Linear Regression

- If we assume a linear model with additive Gaussian noise

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

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

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

\vdots

$$= \underbrace{(X^T X)^{-1} X^T}_{\downarrow} \mathbf{y}$$

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...

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

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

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

\vdots

$$= \boxed{(X^T X + \lambda_c I_{D+1})^{-1} X^T \mathbf{y}} \equiv \text{ridge reg solution}$$

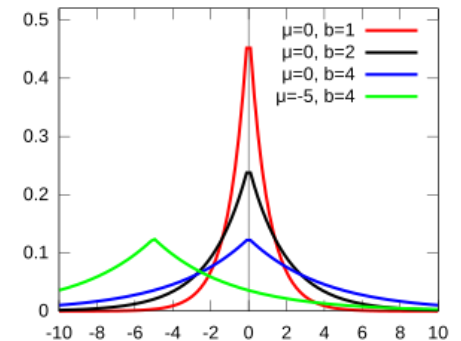
MAP for Linear Regression

- 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 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

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