

10-301/601: Introduction to Machine Learning

Lecture 11 – Regularization

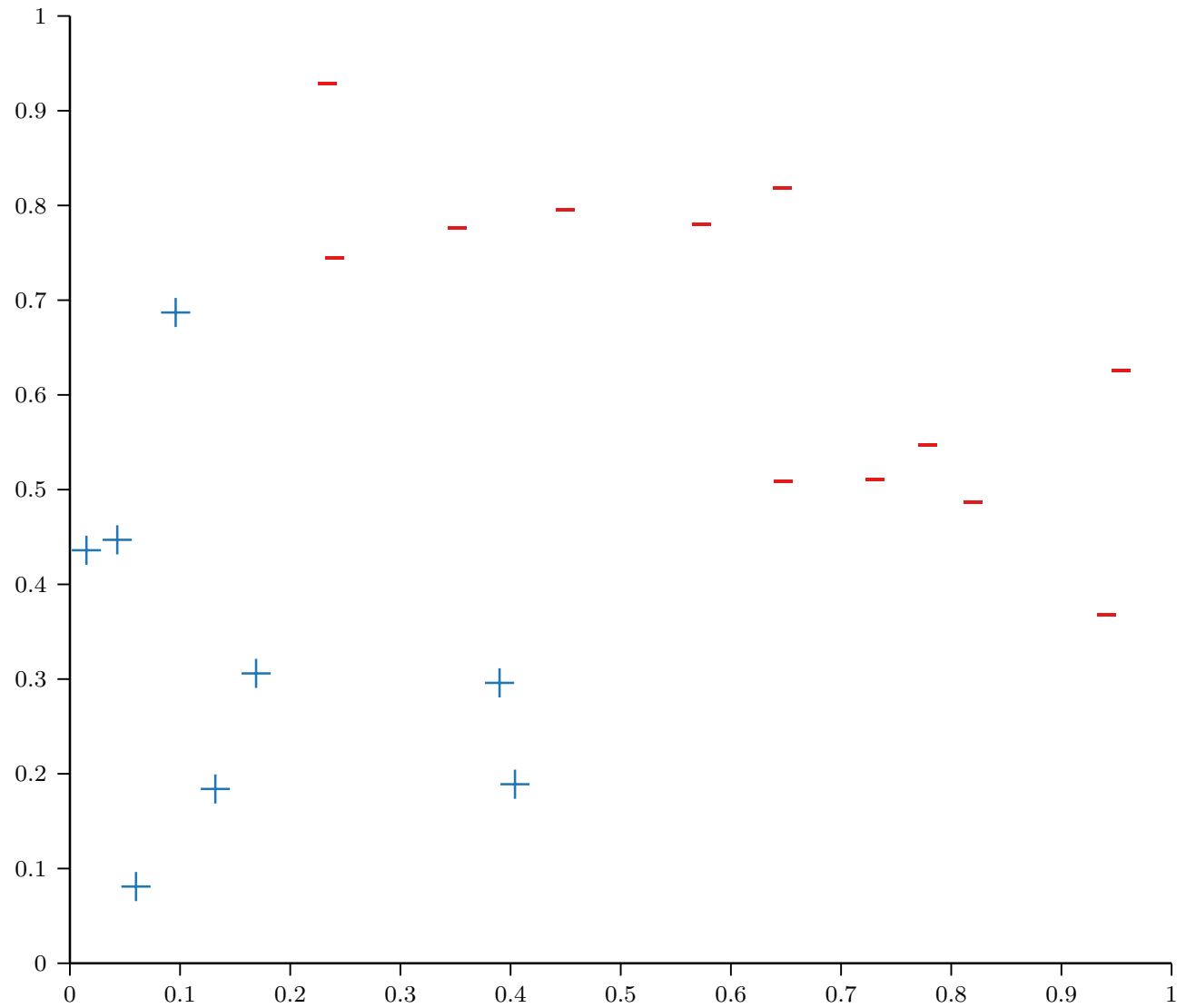
Henry Chai

6/7/23

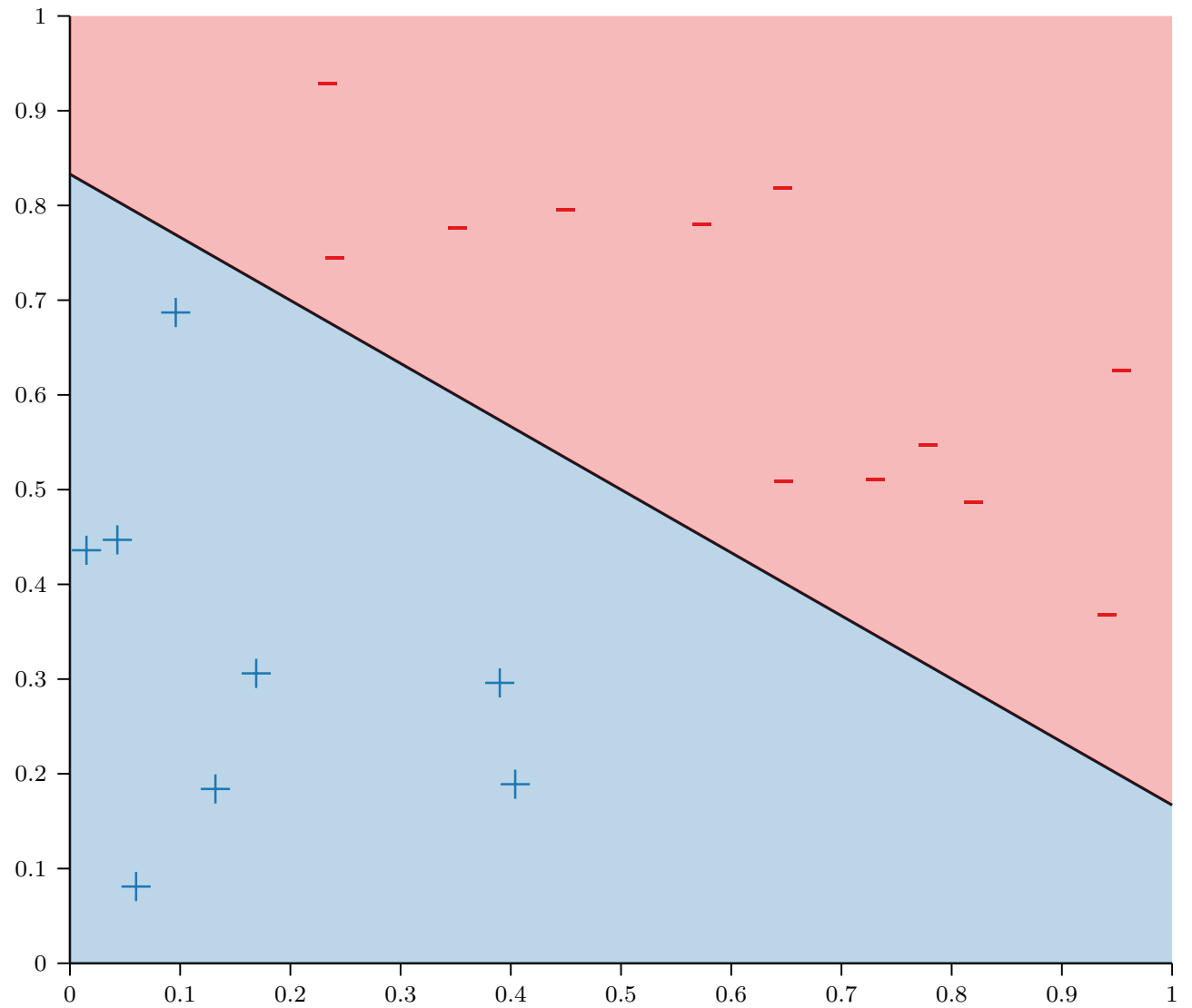
Front Matter

- Announcements:
 - PA3 released 6/8 (tomorrow), due 6/15 at 11:59 PM
 - Midterm on 6/23, two weeks from Friday
 - Practice problems for the Midterm will be posted to the course website on Friday, under [Recitations](#)
- Recommended Readings:
 - Murphy, [Chapter 7.5](#)

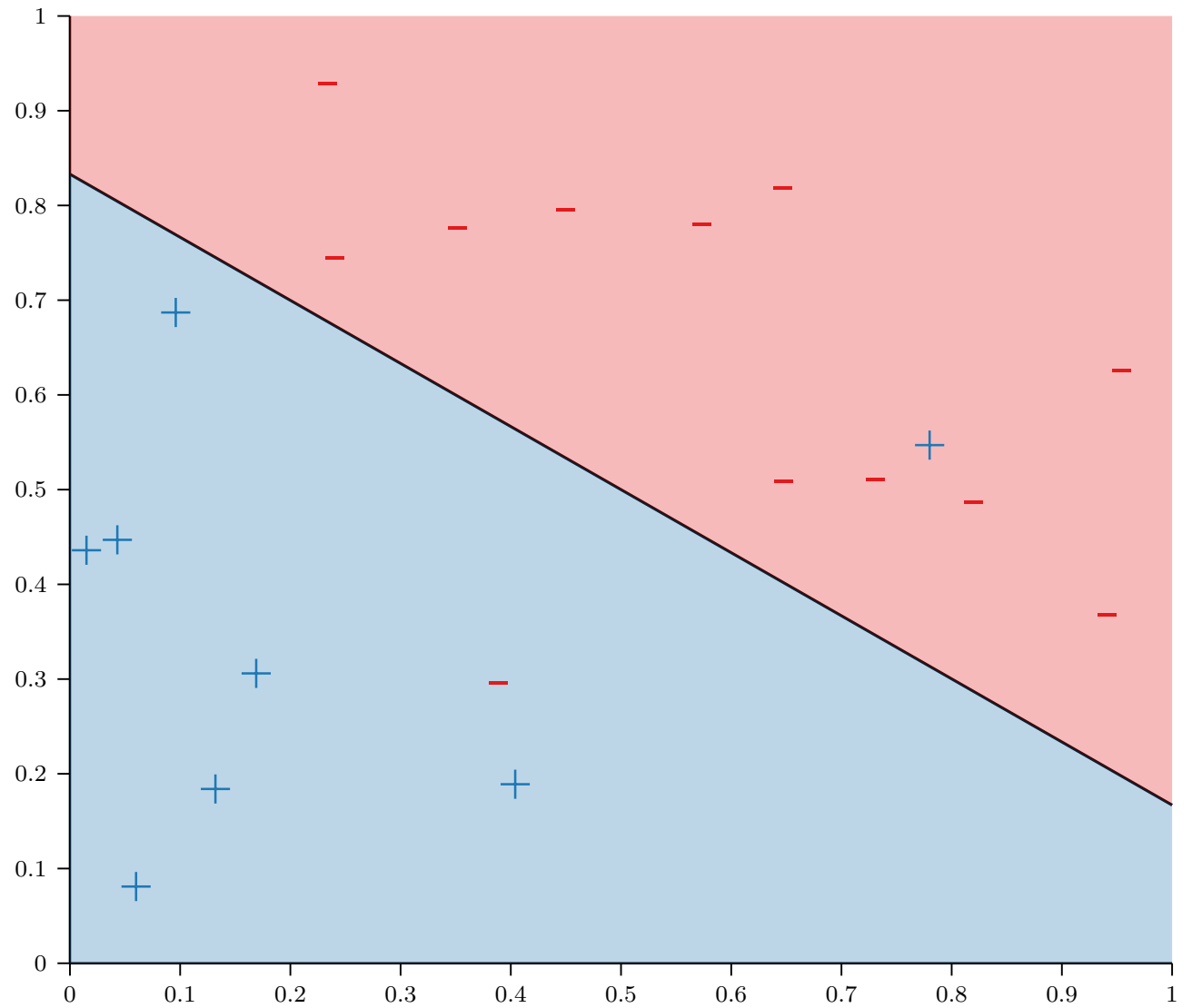
Linear Models



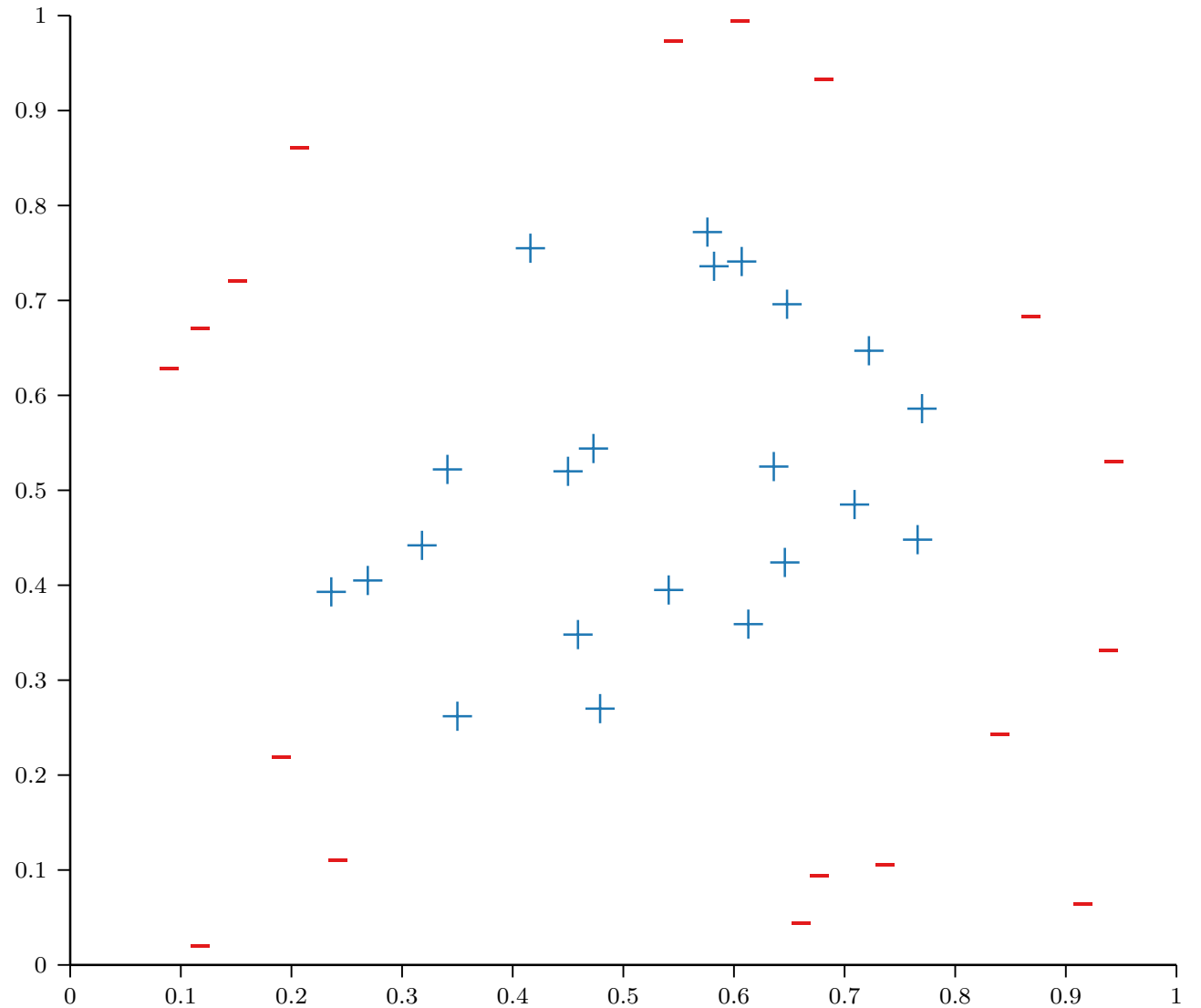
Linear Models



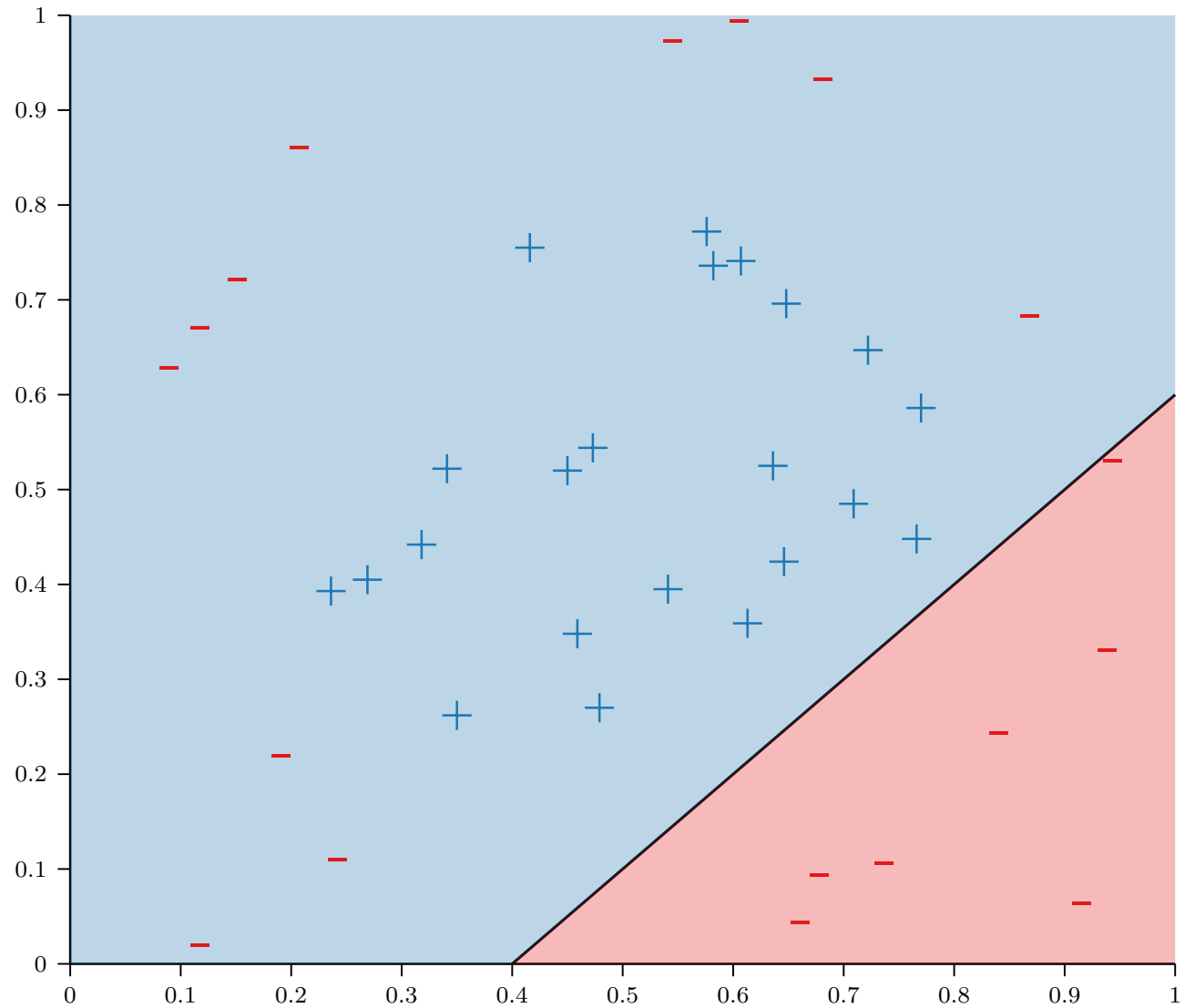
Linear Models



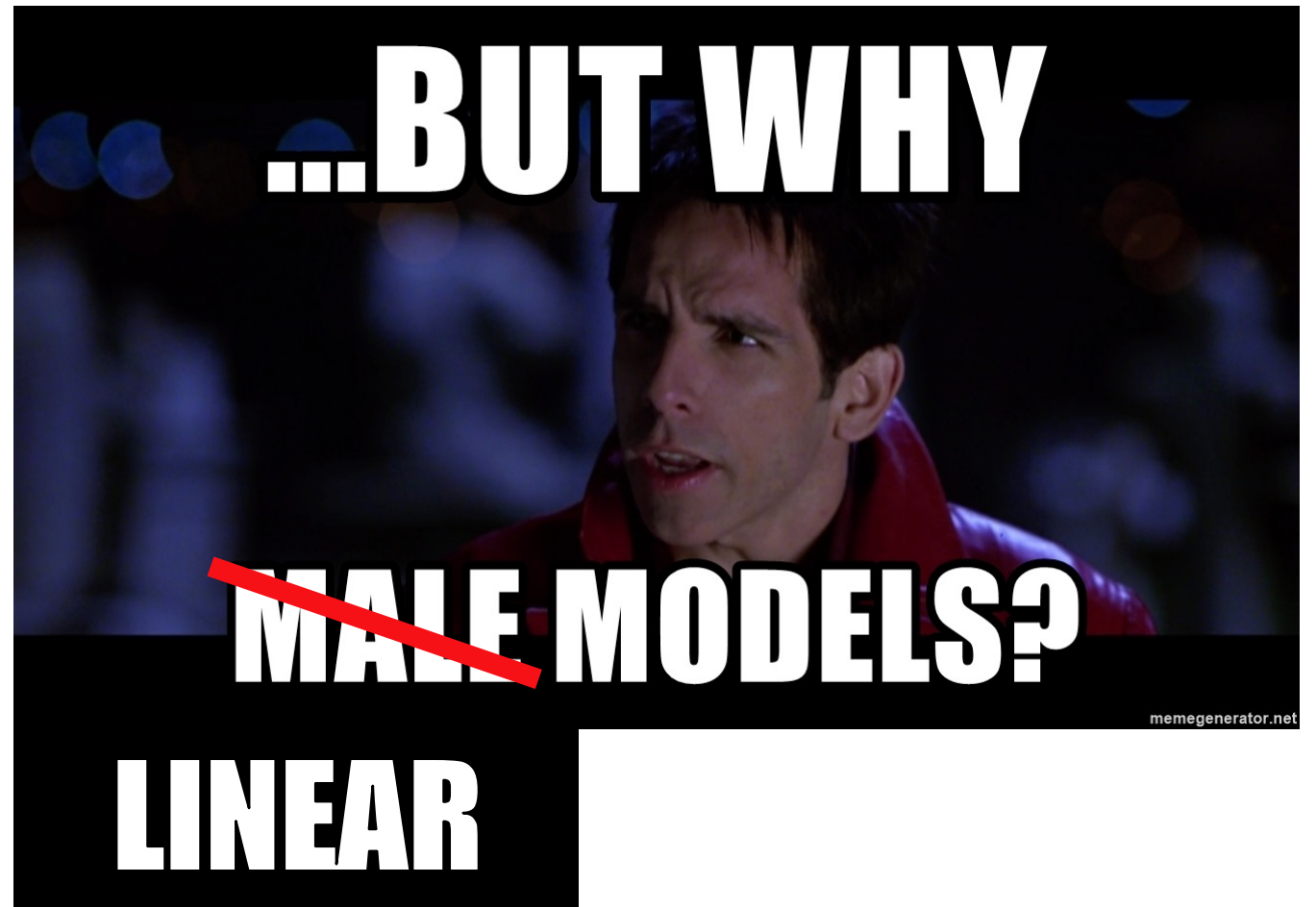
Linear Models?



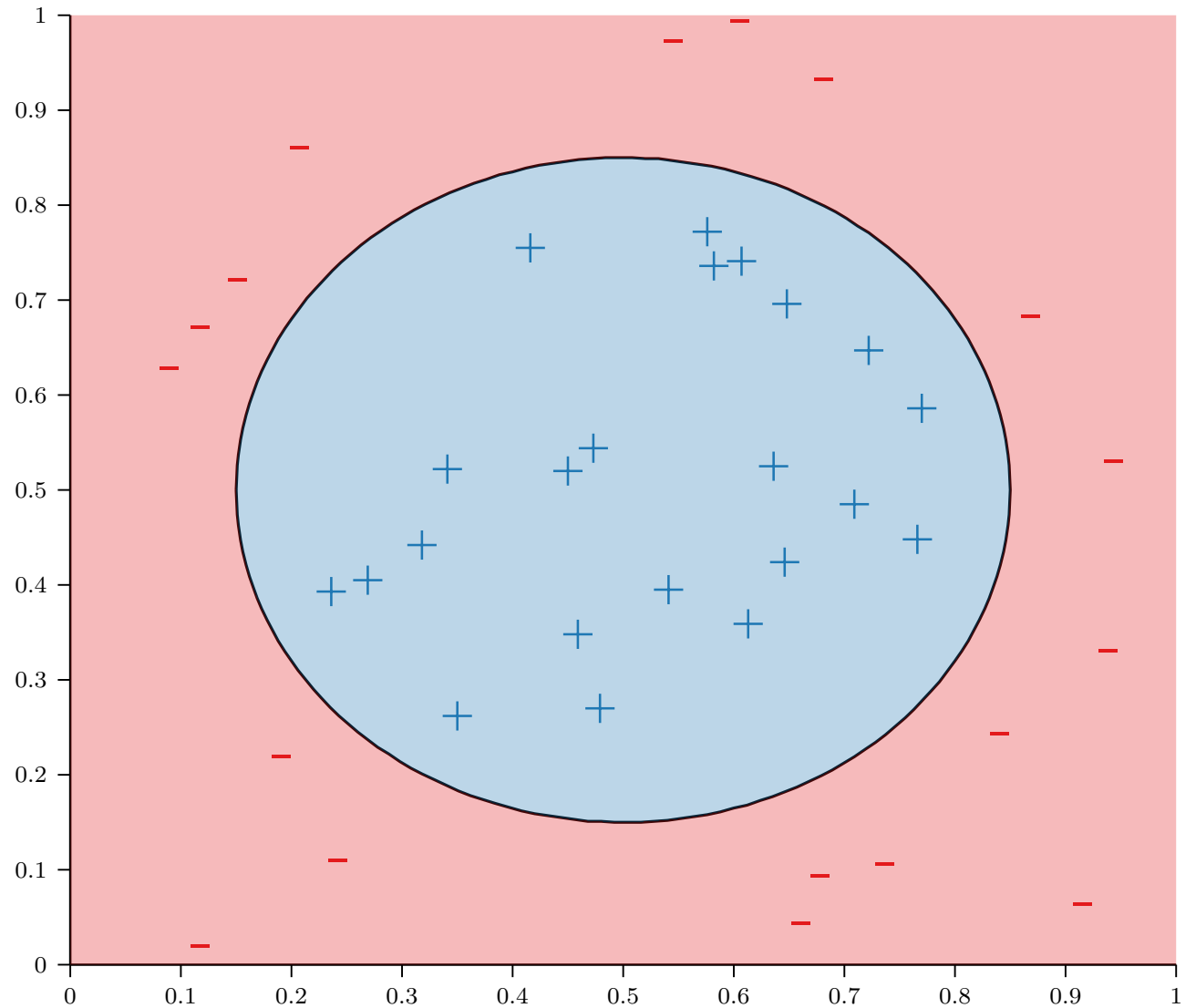
Linear Models?



Linear Models?



Nonlinear Models

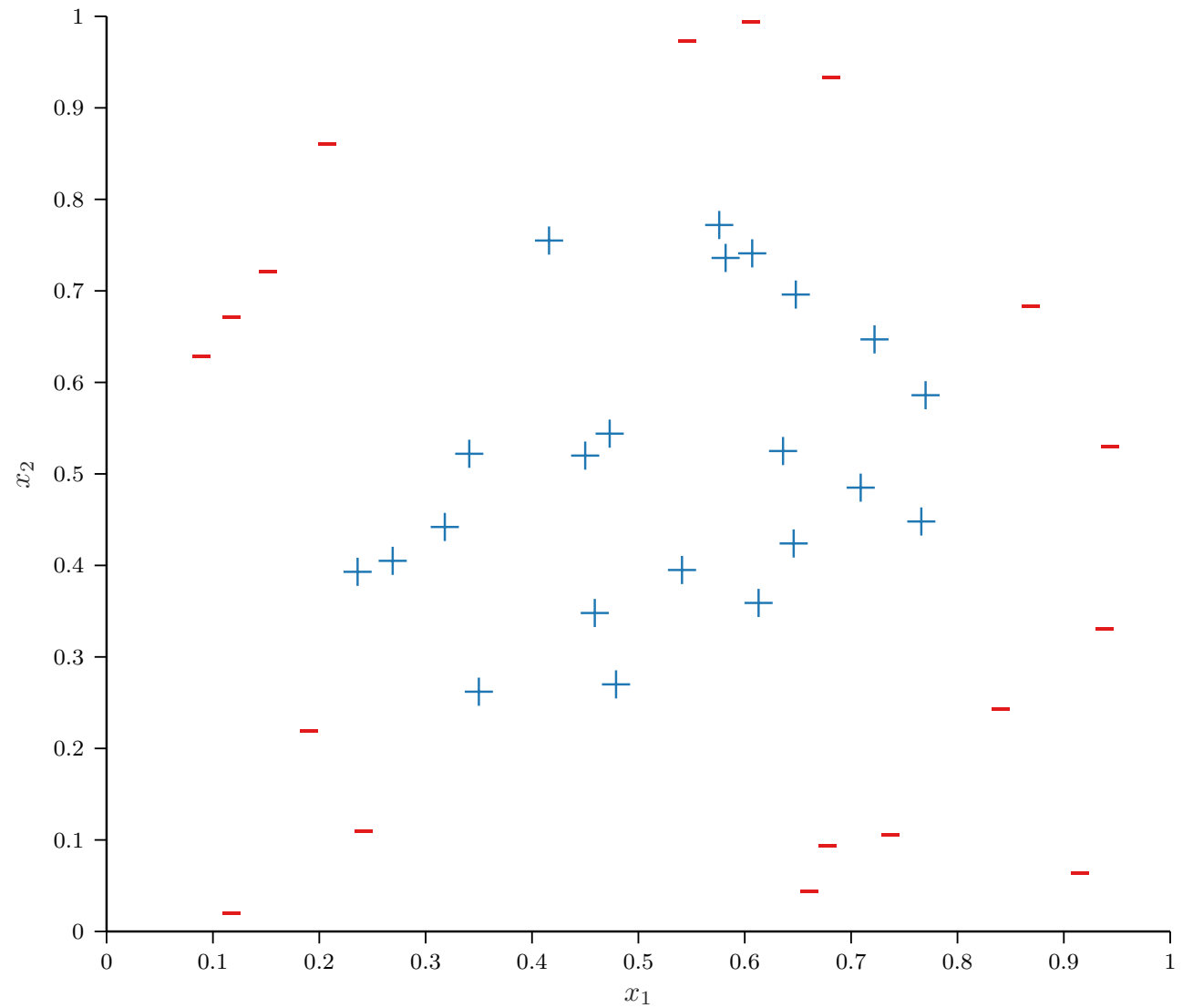


Feature Transforms

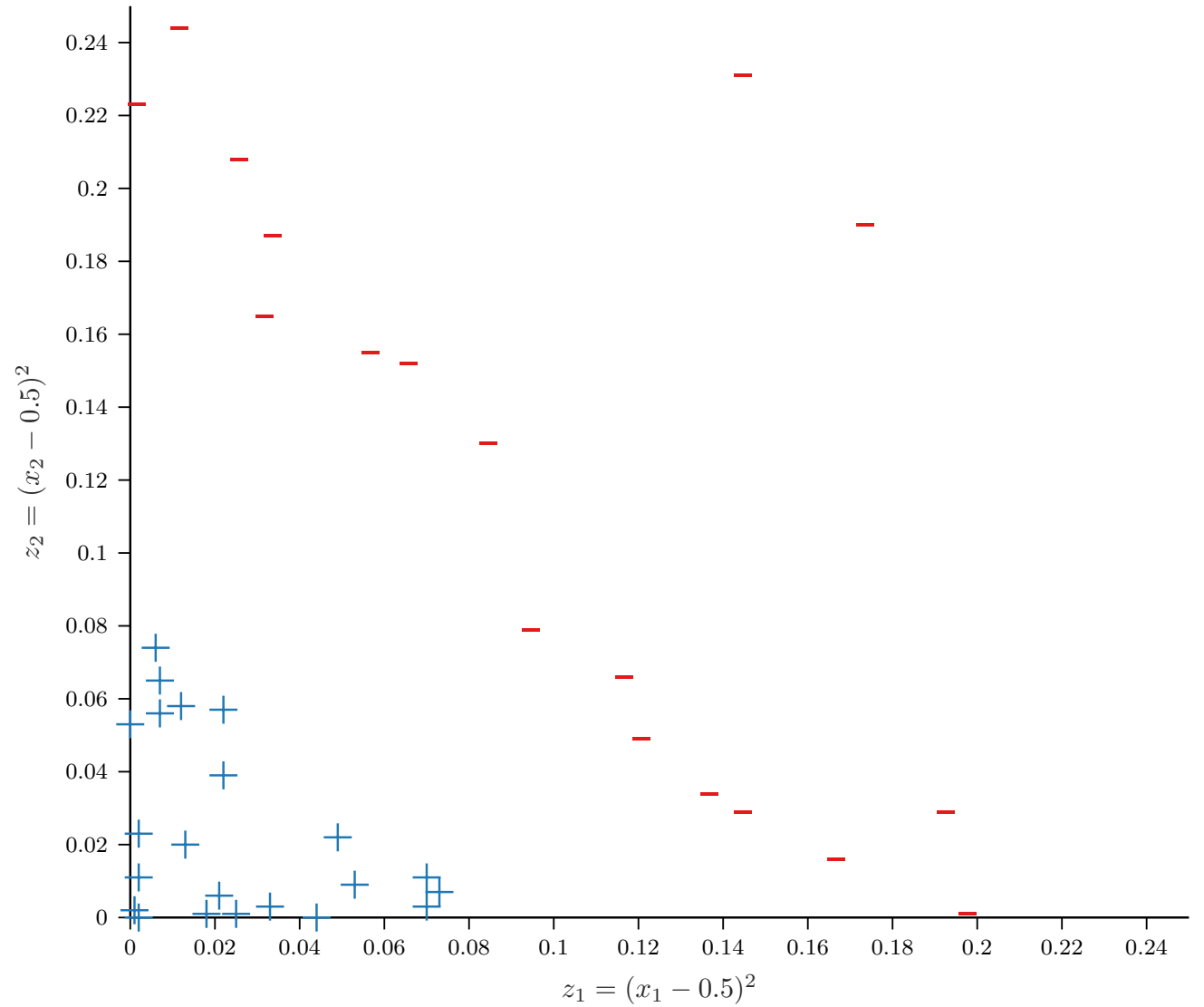
- Given D -dimensional inputs $\mathbf{x} = [x_1, \dots, x_D]$, first compute some transformation of our input, e.g.,

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

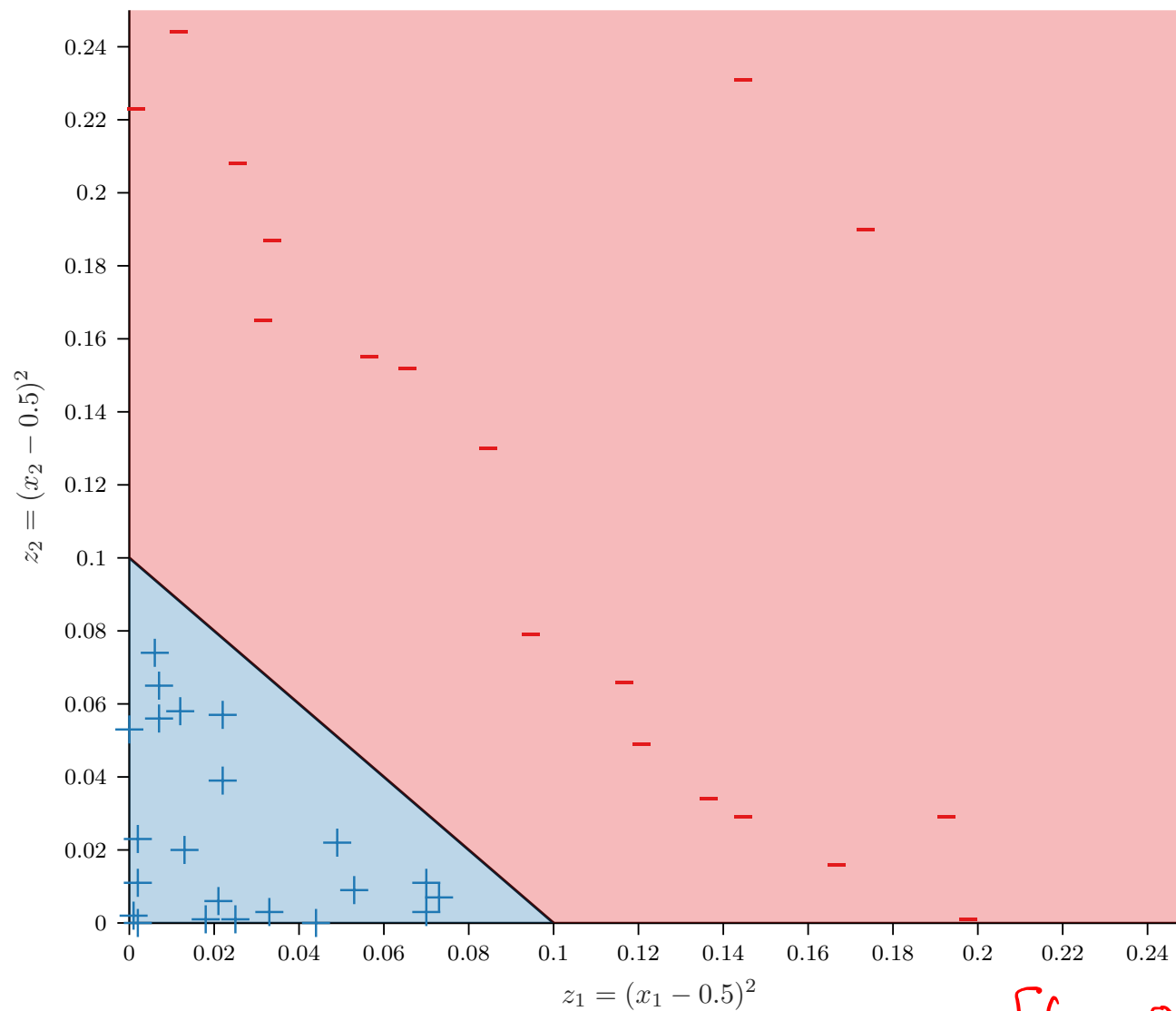
Nonlinear Models



Nonlinear Models

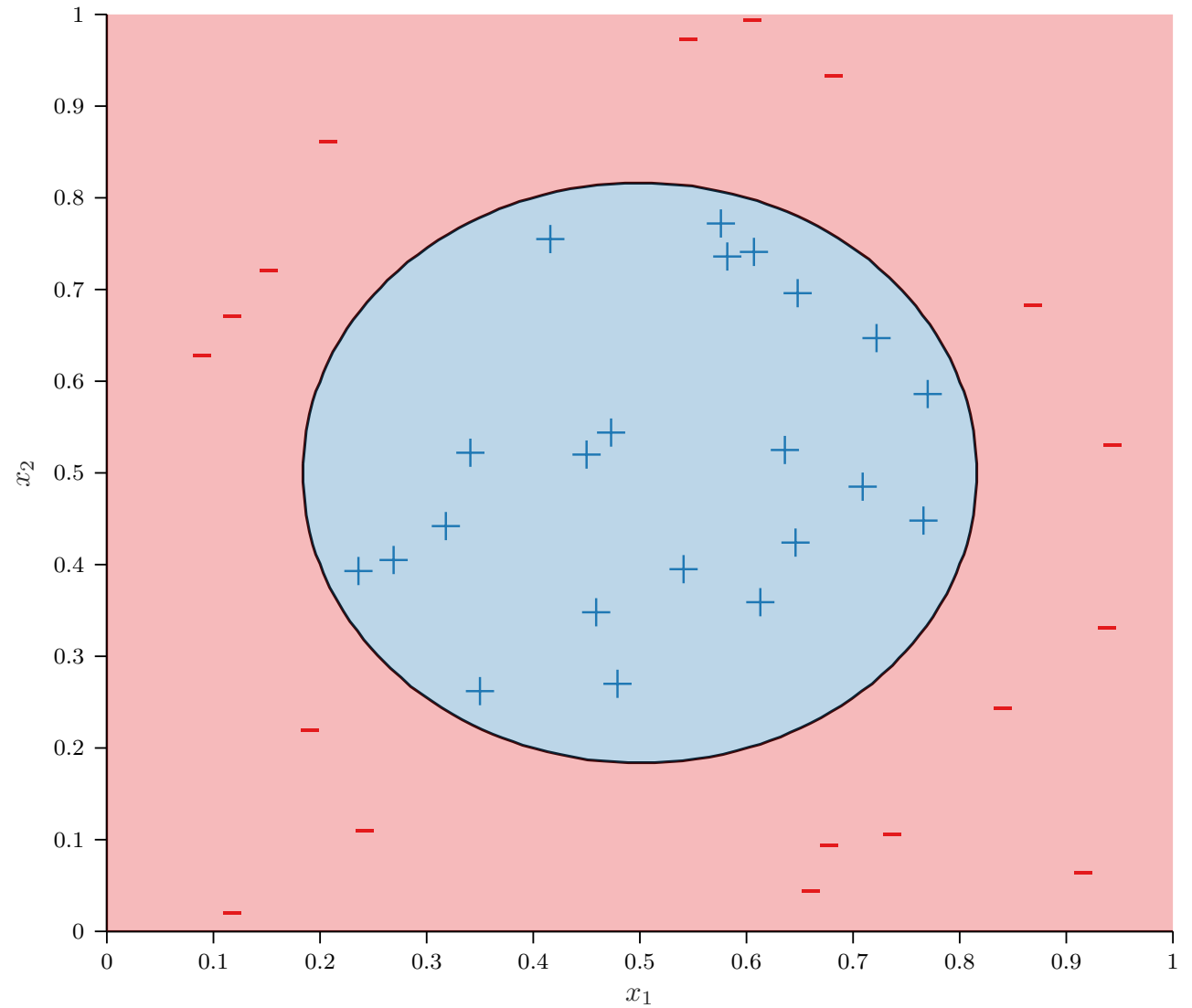


Nonlinear Models



$$h(x) : w^T z = 0 \rightarrow w^T \begin{bmatrix} (x_1 - 0.5)^2 \\ (x_2 - 0.5)^2 \end{bmatrix} = 0$$

Nonlinear Models



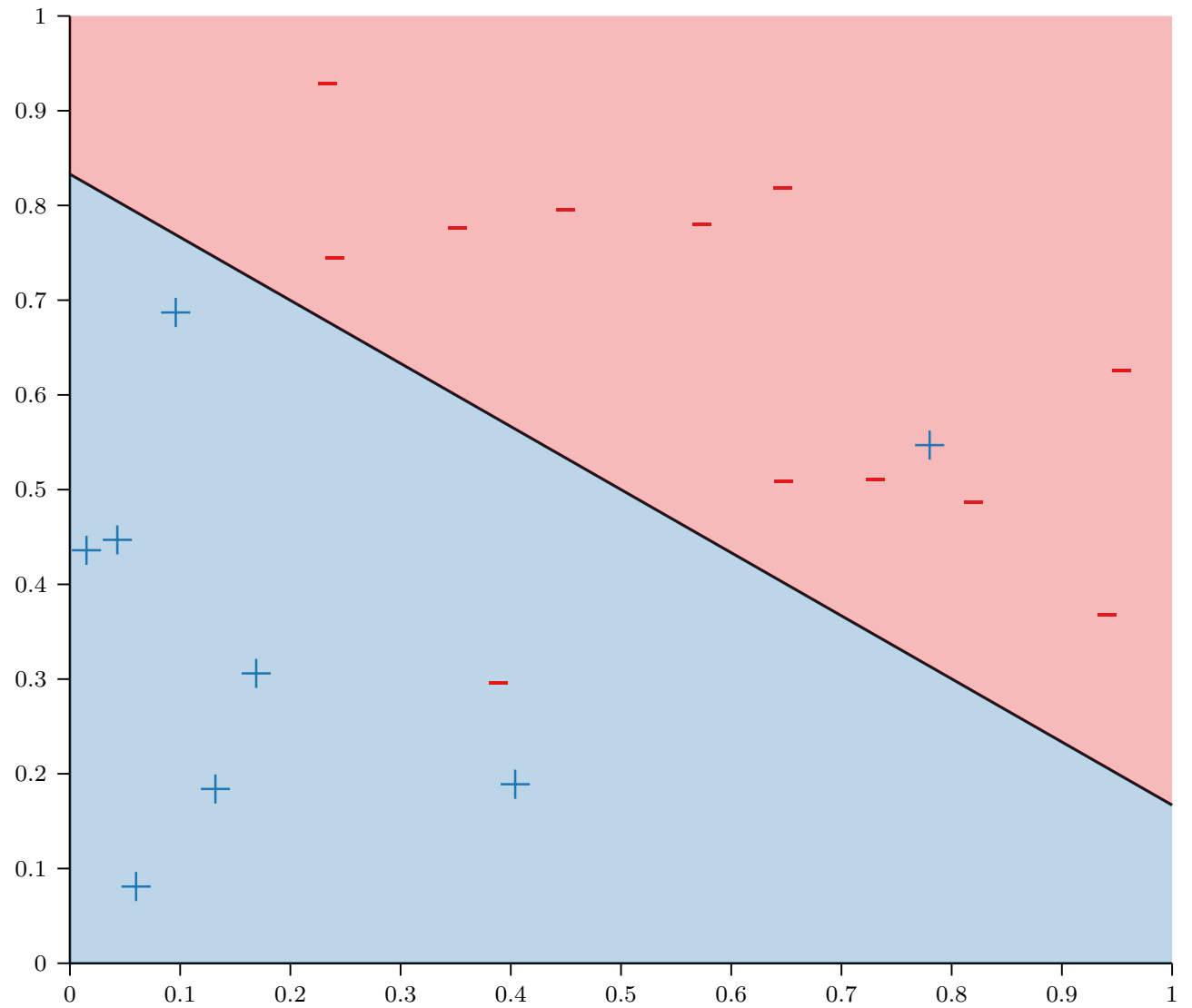
General Q^{th} -order Transforms

two input features
→ 2nd-order transformations

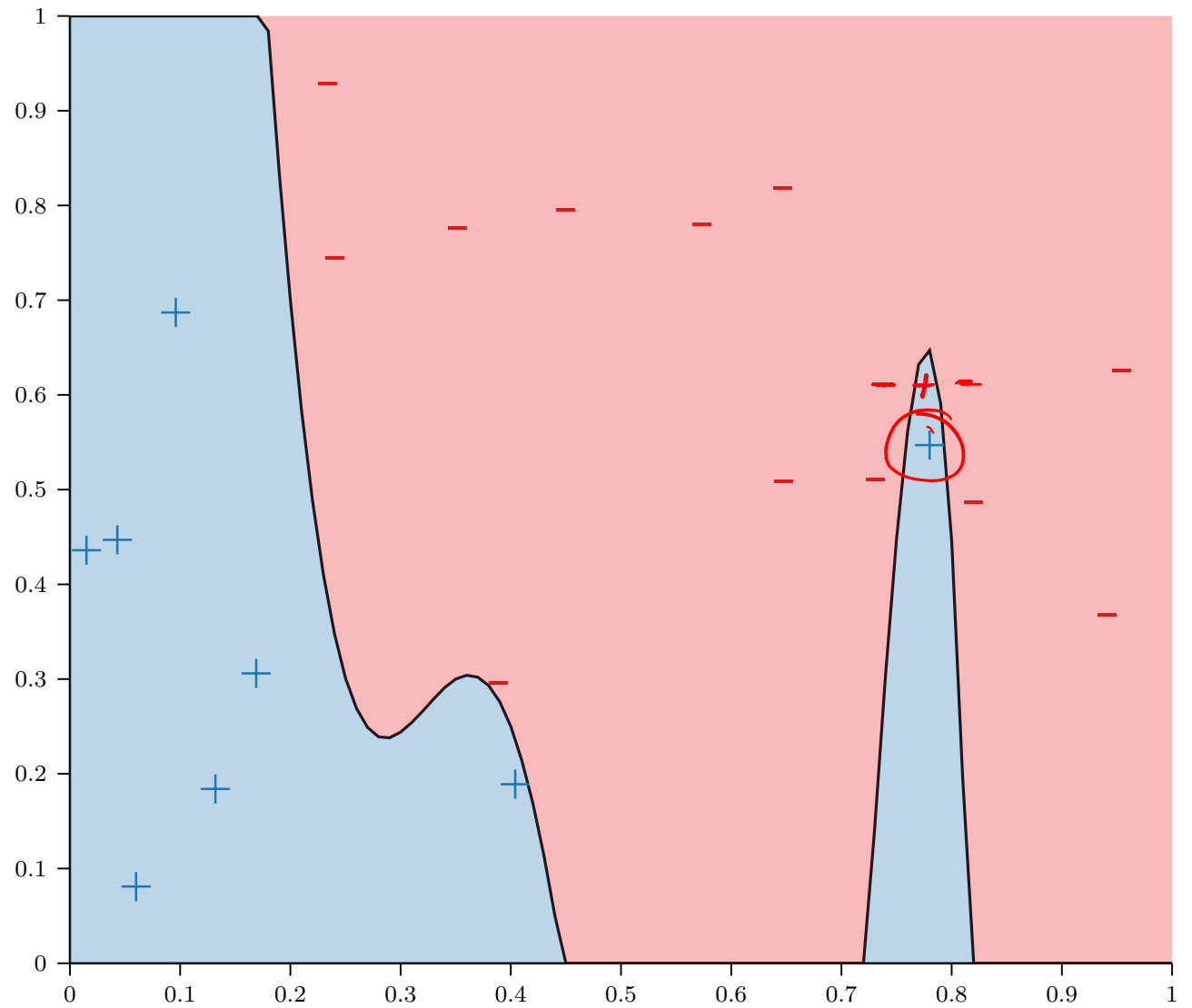
- $\phi_{2,2}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2]$
- $\phi_{2,3}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3]$
- $\phi_{2,4}([x_1, x_2]) = [x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3, x_1^4, x_1^3x_2, x_1^2x_2^2, x_1x_2^3, x_2^4]$
-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

$\phi_{D,Q}$ is a $O(Q^D)$ transformation

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

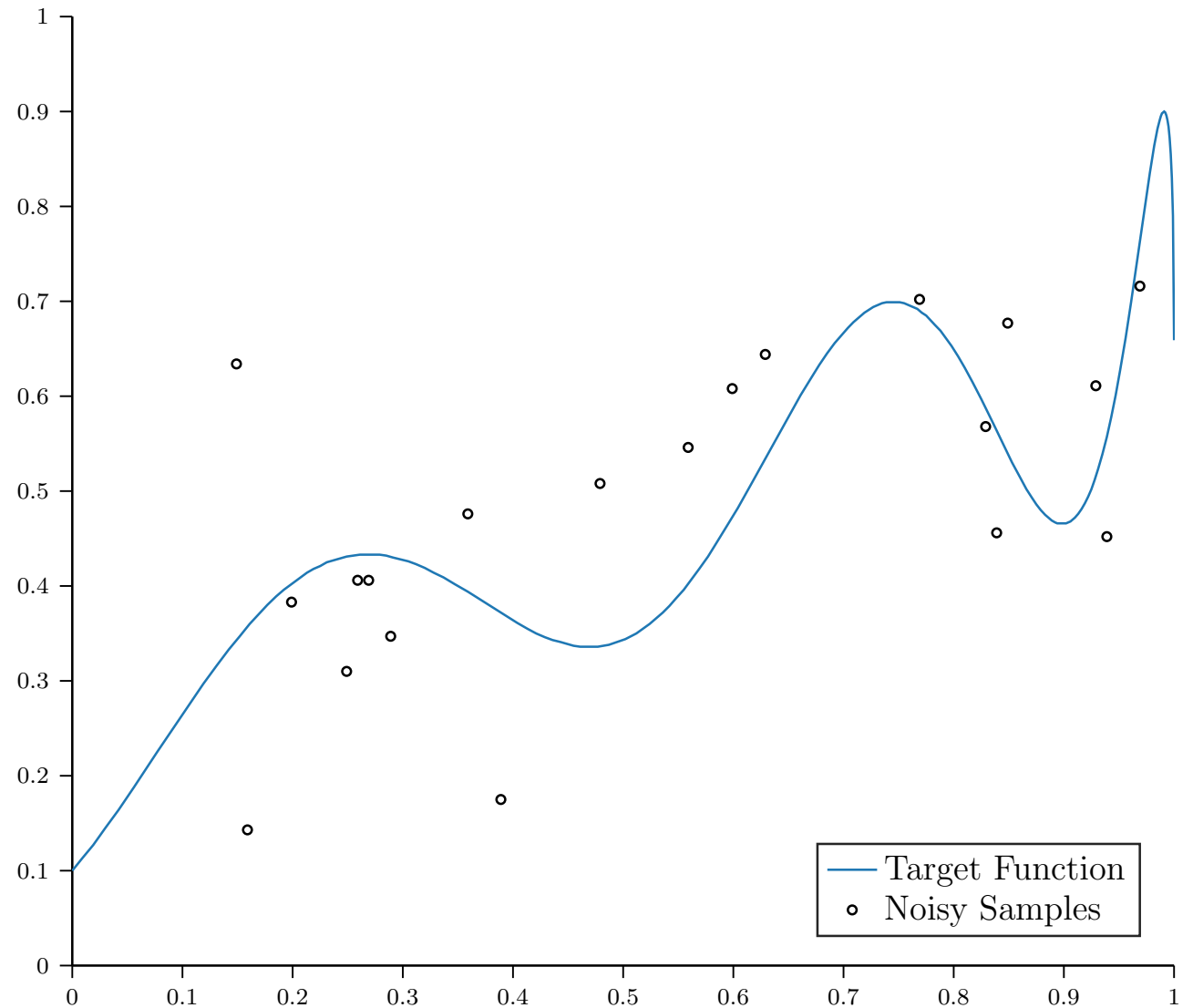
- $x \in \mathbb{R}, y \in \mathbb{R}$ and $N = 20$
- Targets are generated by a 10th-order polynomial in x with additive Gaussian noise:

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

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



Lecture 11 Polls

0 done

 **0 underway**

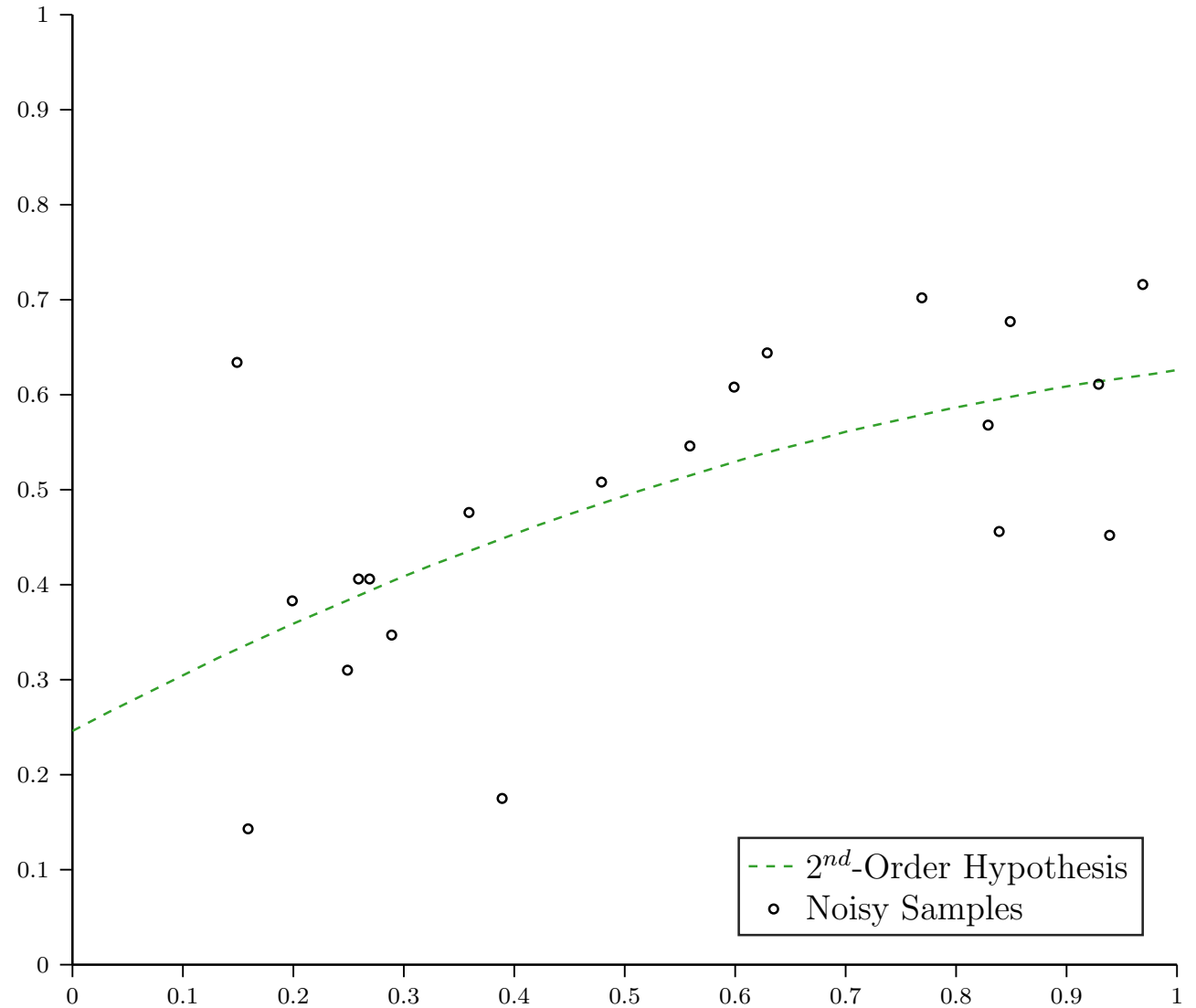
Which model do you think will have a lower true error in this setting (*loading eqn.*)?

\mathcal{H}_2

\mathcal{H}_{10}

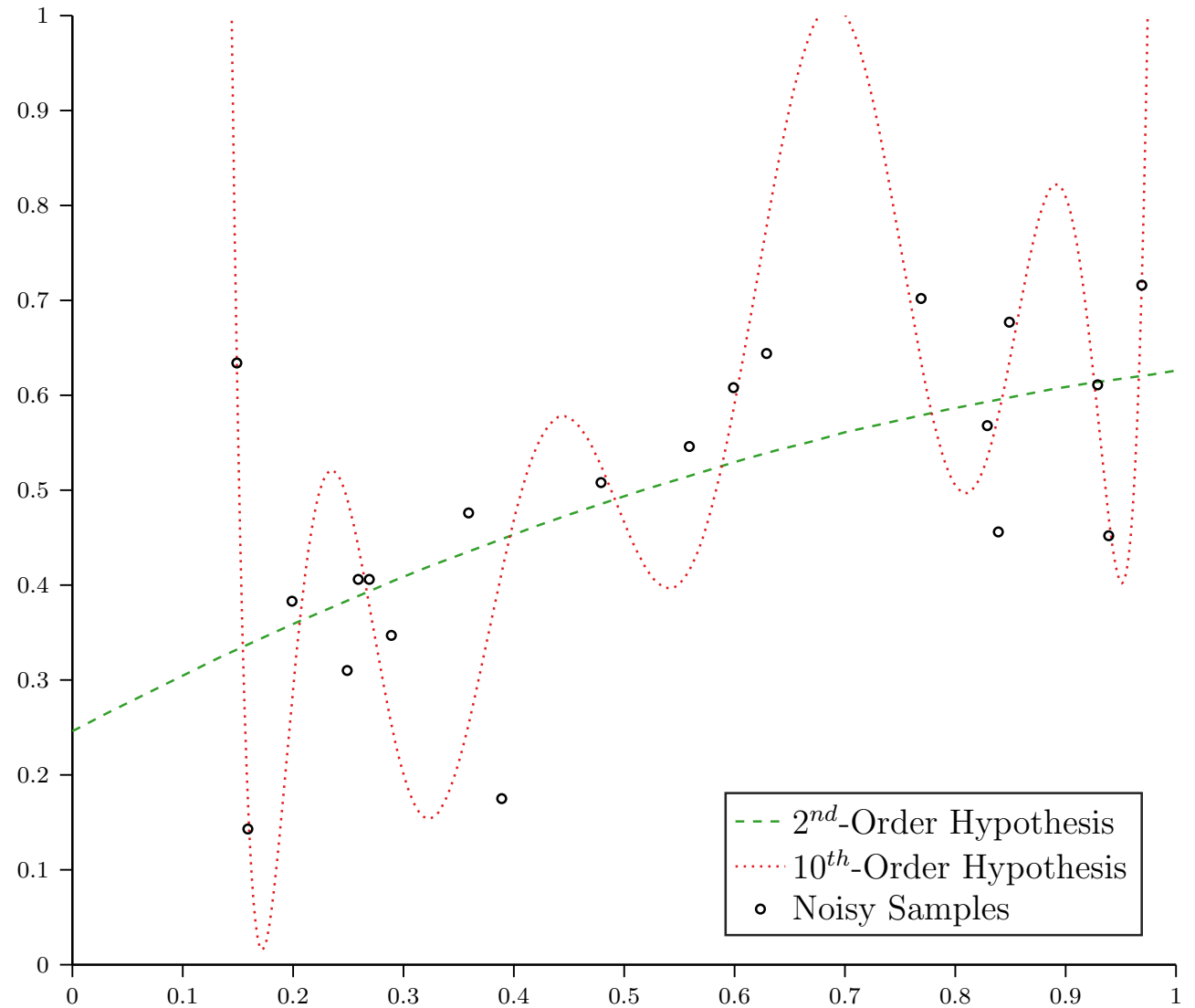
Noisy Targets

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



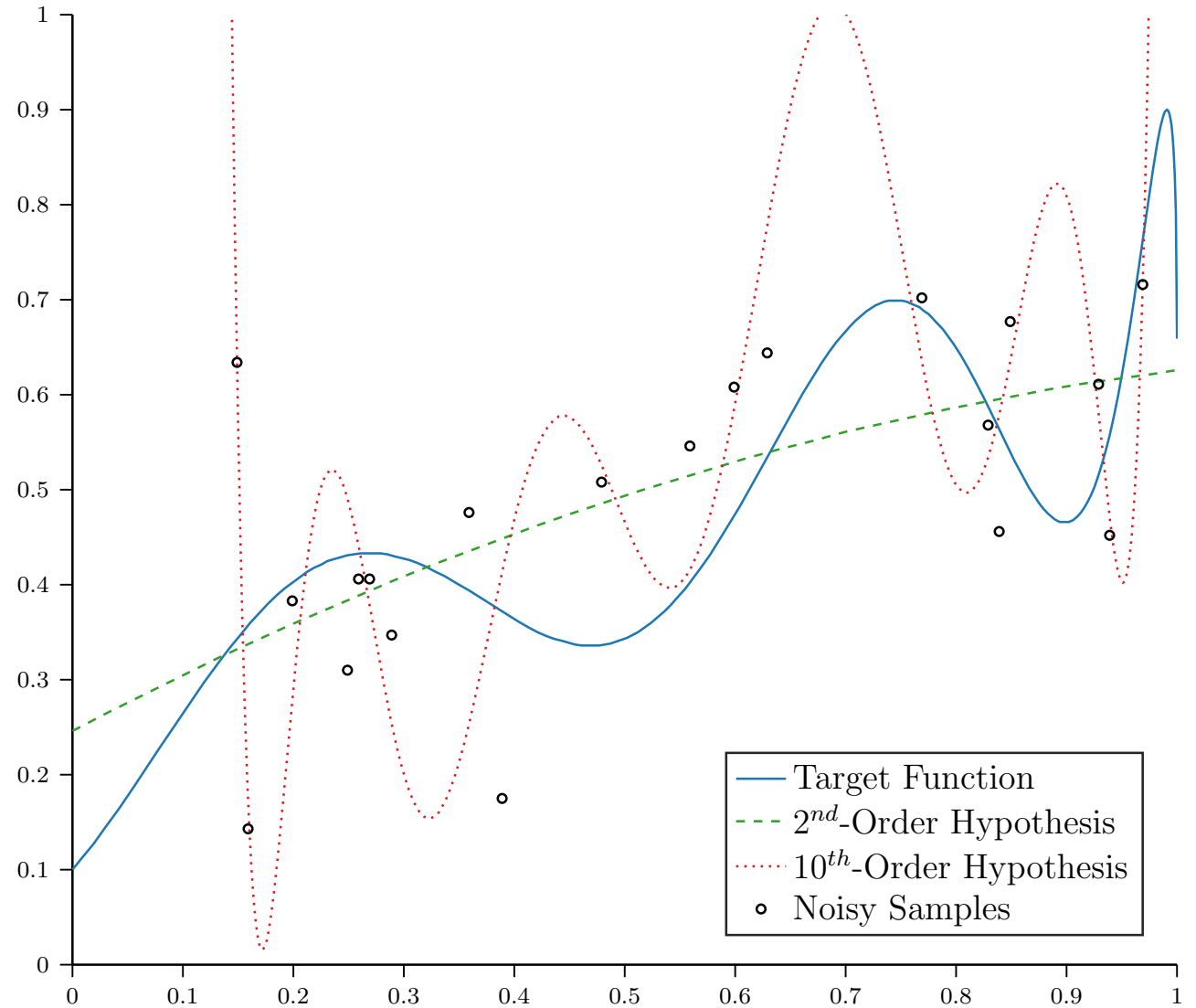
Noisy Targets

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



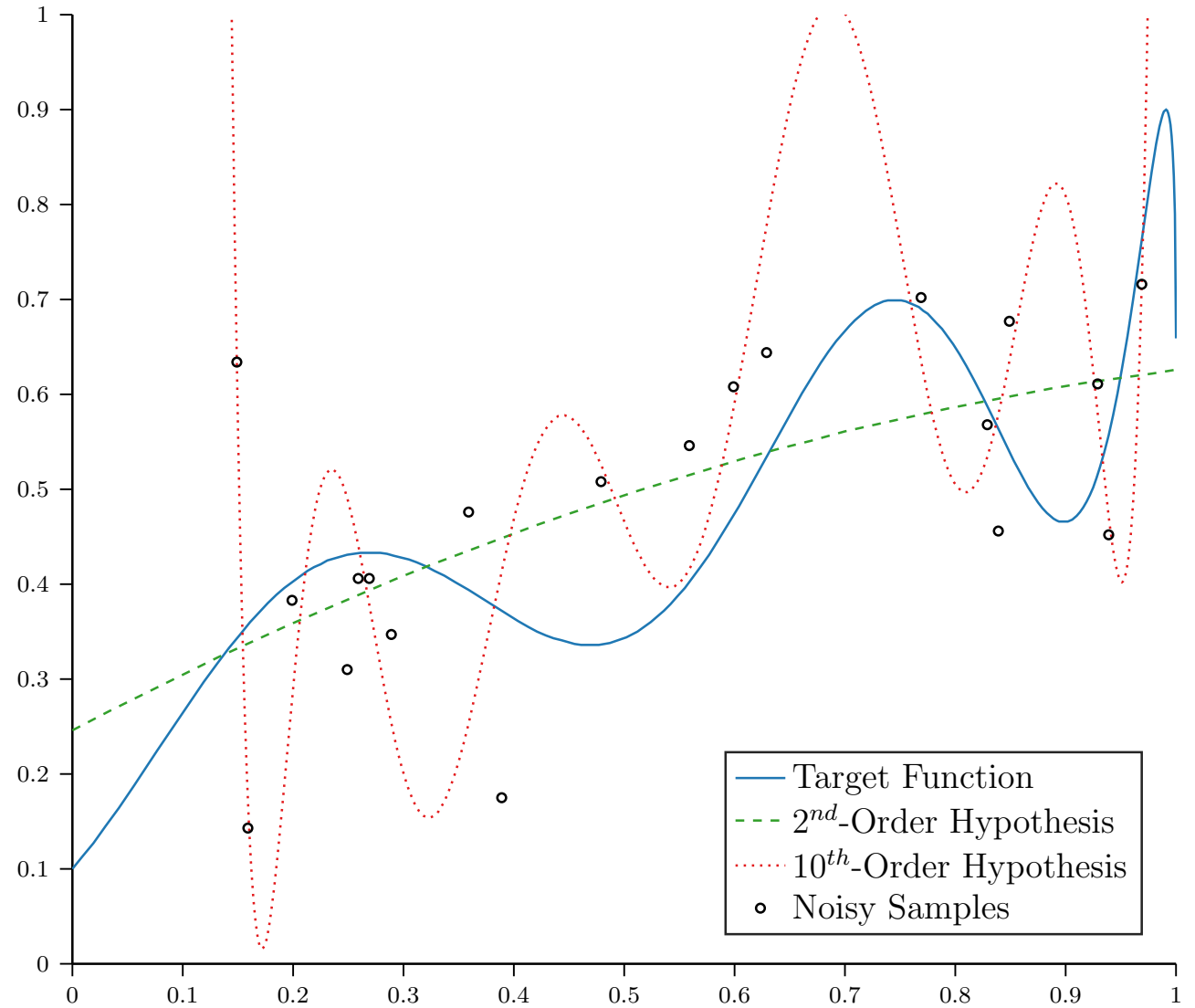
Noisy Targets

- 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

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

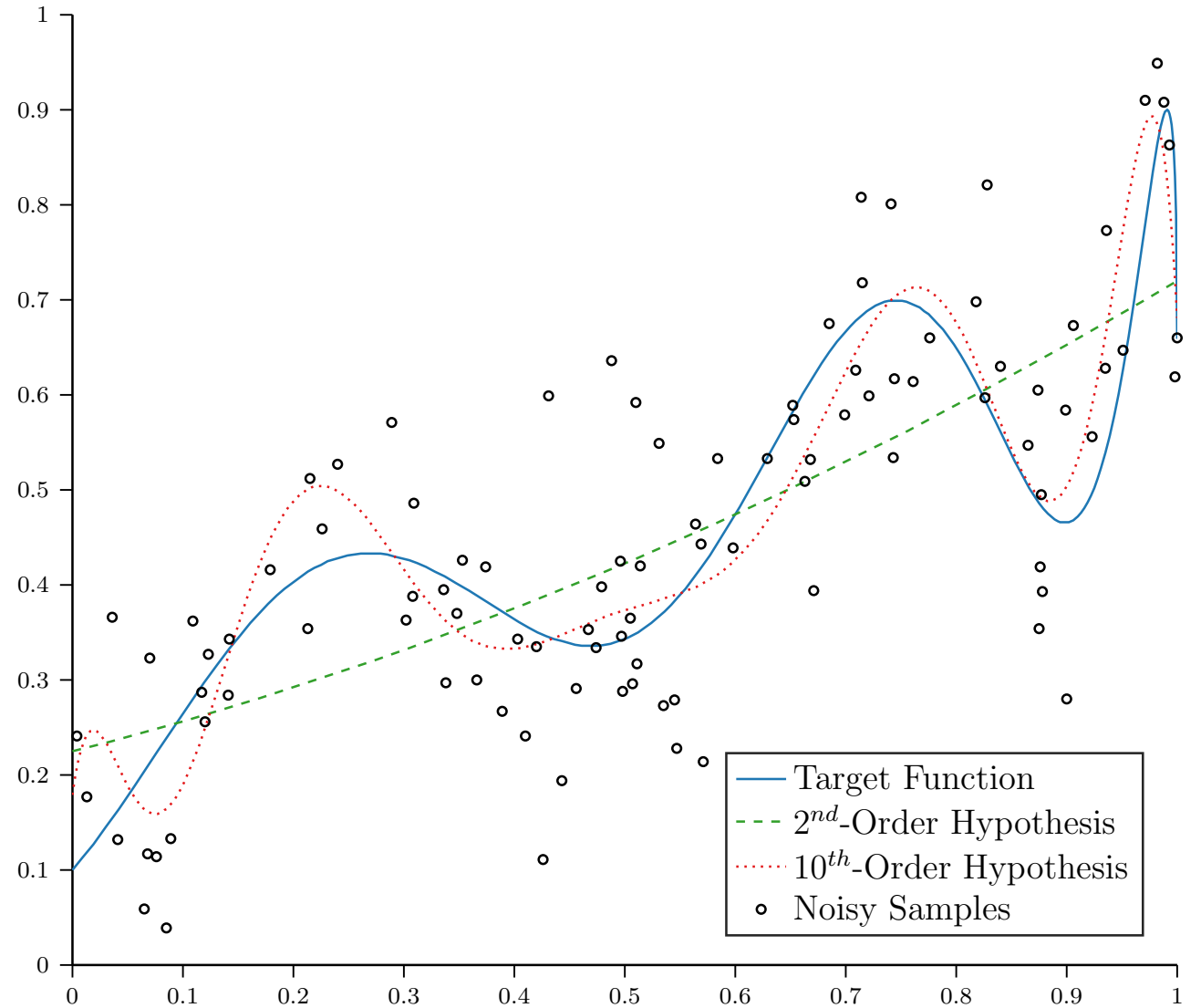
Which model do you think will have a lower true error in this setting ($N = 100$)?

\mathcal{H}_2

\mathcal{H}_{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

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

- 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

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

- Subject to

$$\underbrace{\omega_3 = \omega_4 = \omega_5 = \omega_6 = \omega_7 = \omega_8 = \omega_9 = \omega_{10} = 0}_{\text{Hard Constraints}}$$

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

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

- 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

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

- Subject to nothing!

Soft Constraints

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

- 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

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

- Subject to:

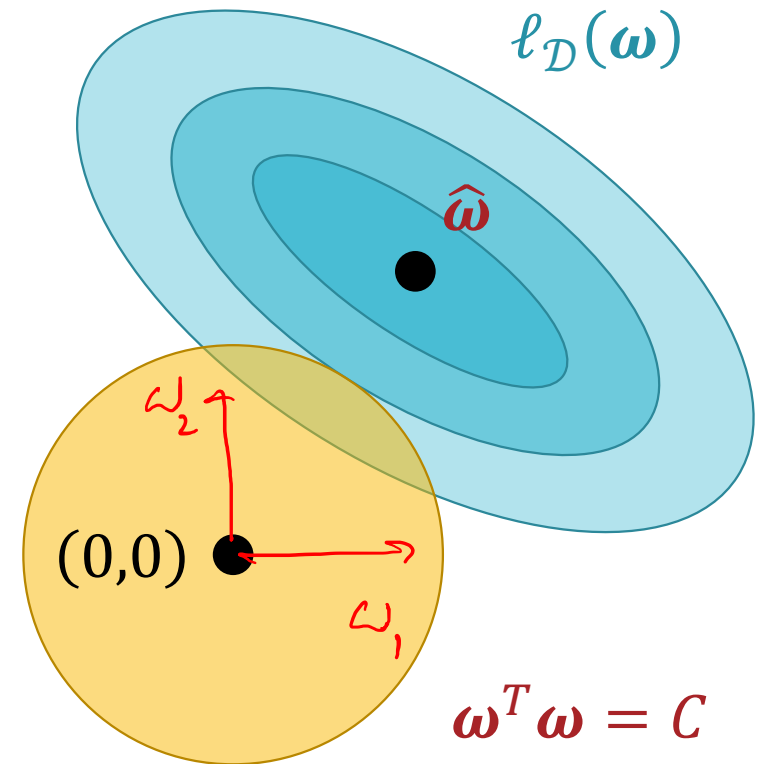
$$\|\boldsymbol{\omega}\|_2^2 = \boldsymbol{\omega}^T \boldsymbol{\omega} = \sum_{d=0}^D \omega_d^2 \leq C$$

Soft Constraints

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

$$\omega_1^2 + \omega_2^2 \leq C$$

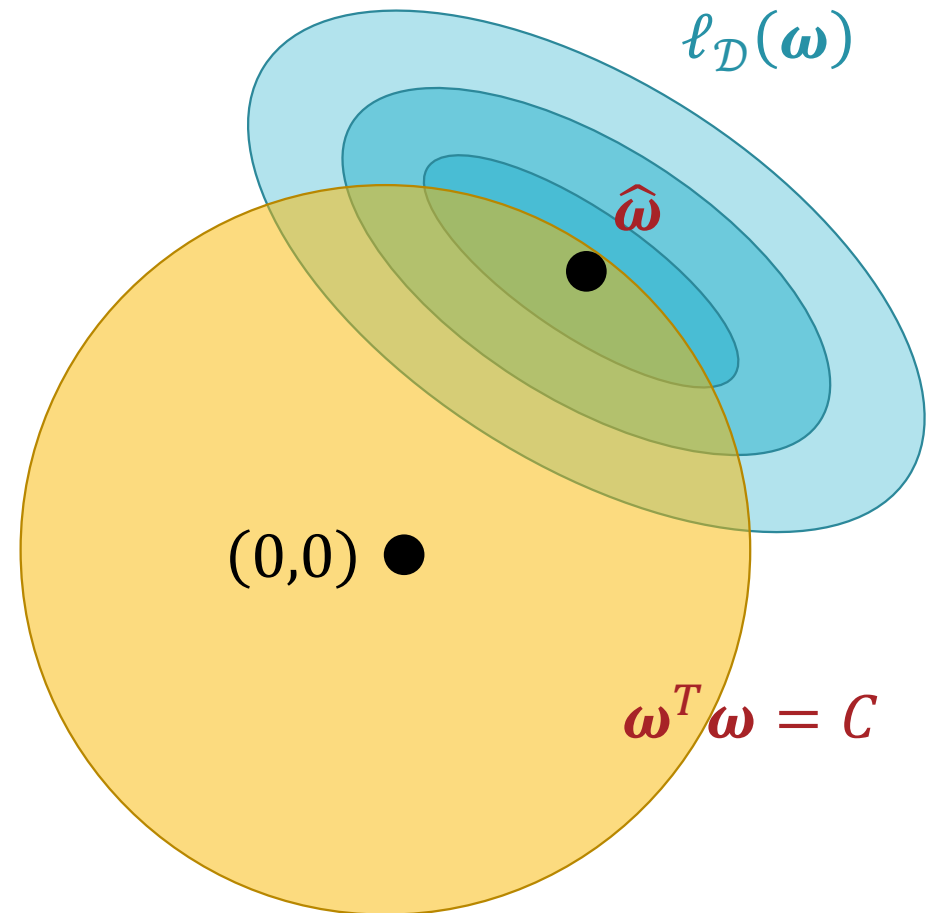
subject to $\boldsymbol{\omega}^T \boldsymbol{\omega} \leq C$



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

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

1.

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

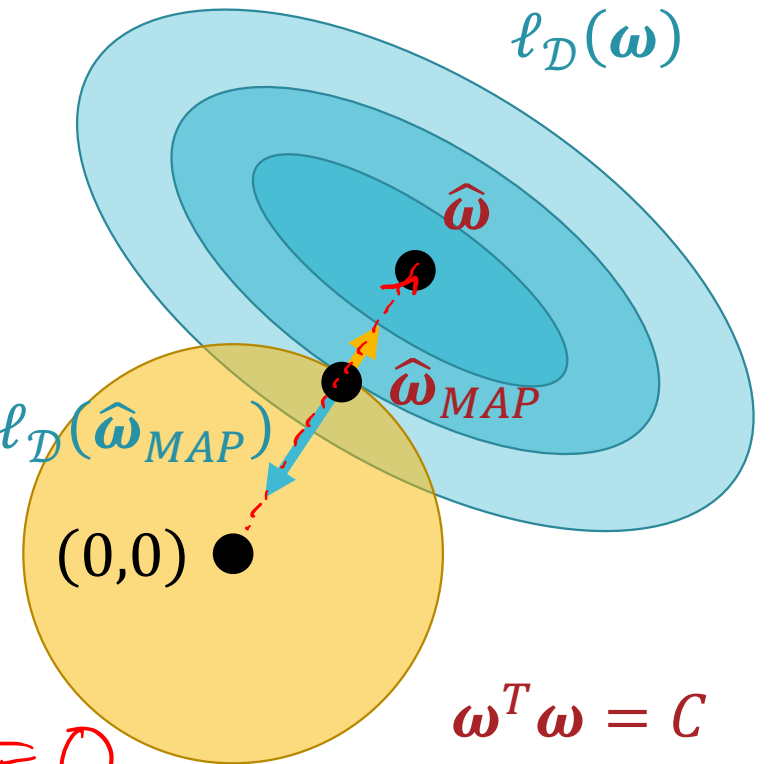
$$2. \quad \nabla_{\boldsymbol{\omega}} \ell_D(\hat{\boldsymbol{\omega}}_{\text{MAP}}) \propto \mathbf{I} - \hat{\boldsymbol{\omega}}_{\text{MAP}}$$

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

$$(\lambda_C > 0)$$

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

$$3. \quad \nabla_{\boldsymbol{\omega}} (\ell_D(\hat{\boldsymbol{\omega}}_{\text{MAP}}) + \lambda_C \hat{\boldsymbol{\omega}}_{\text{MAP}}^T \hat{\boldsymbol{\omega}}_{\text{MAP}}) = \mathbf{0}$$



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

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

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

\Leftrightarrow \rightarrow these problems have the same optimal solution

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

Ridge Regression

$$\text{minimize } \ell_D^{\text{AUG}}(\omega) = \ell_D(\omega) + \underbrace{\lambda_c \omega^T \omega}_{\rightarrow (X\omega - y)^T (X\omega - y)}$$

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

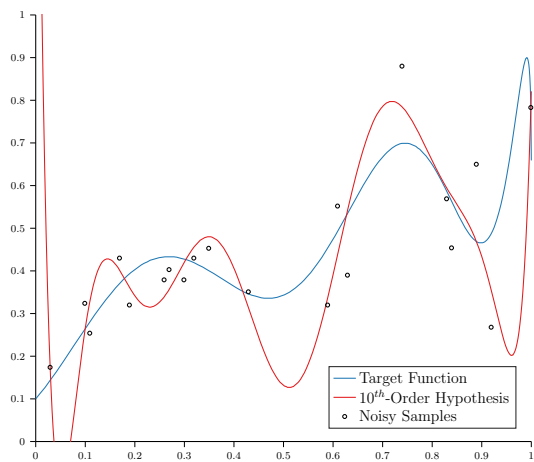
$$2(X^T X \hat{\omega}_{\text{MAP}} - X^T y + \lambda_c \hat{\omega}_{\text{MAP}}) = 0$$

$$X^T X \hat{\omega}_{\text{MAP}} + \lambda_c \hat{\omega}_{\text{MAP}} = X^T y$$

$$\underbrace{(X^T X + \lambda_c I)}_{\text{invertible}} \hat{\omega}_{\text{MAP}} = X^T y$$

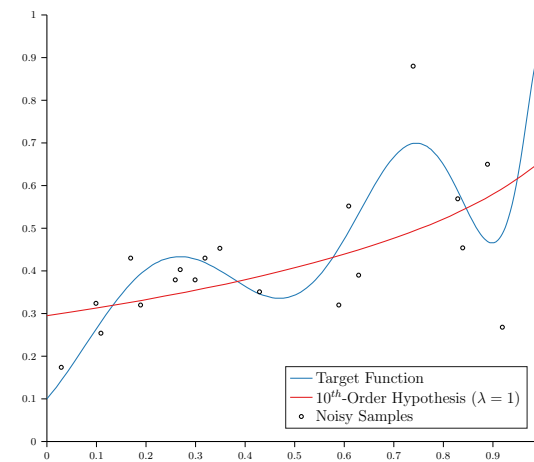
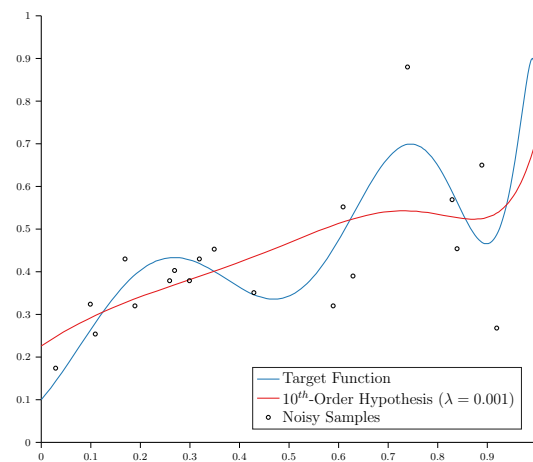
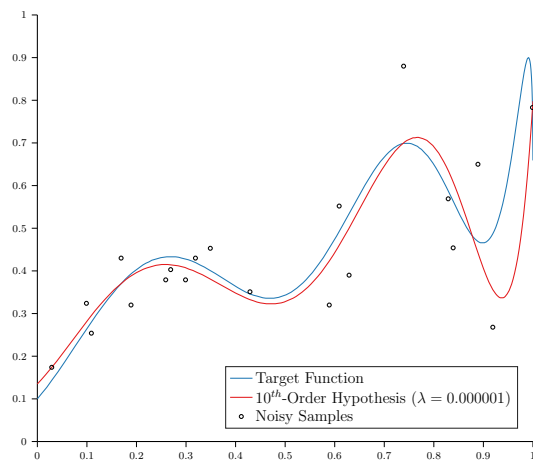
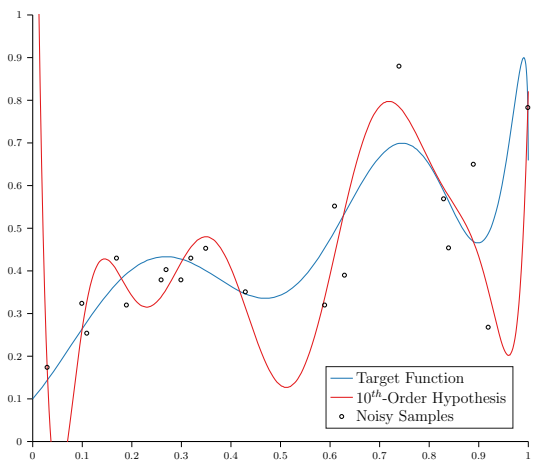
$$\hat{\omega}_{\text{MAP}} = \underbrace{(X^T X + \lambda_c I)^{-1}}_{\text{invertible}} X^T y$$

adding this non-negative term ($\lambda_c \geq 0$) can help invertibility



Ridge Regression

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



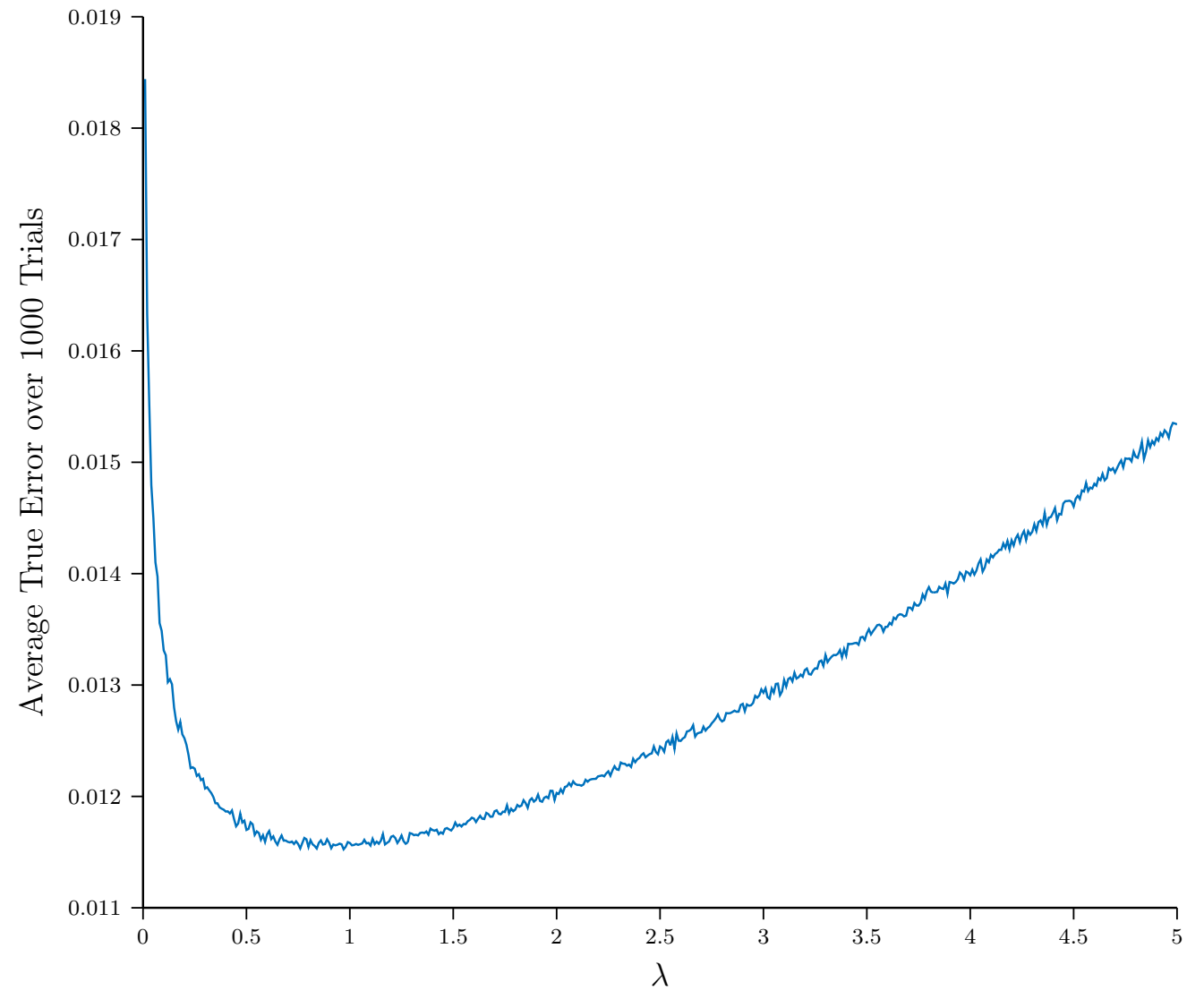
Ridge Regression

$$\lambda_c = 0$$

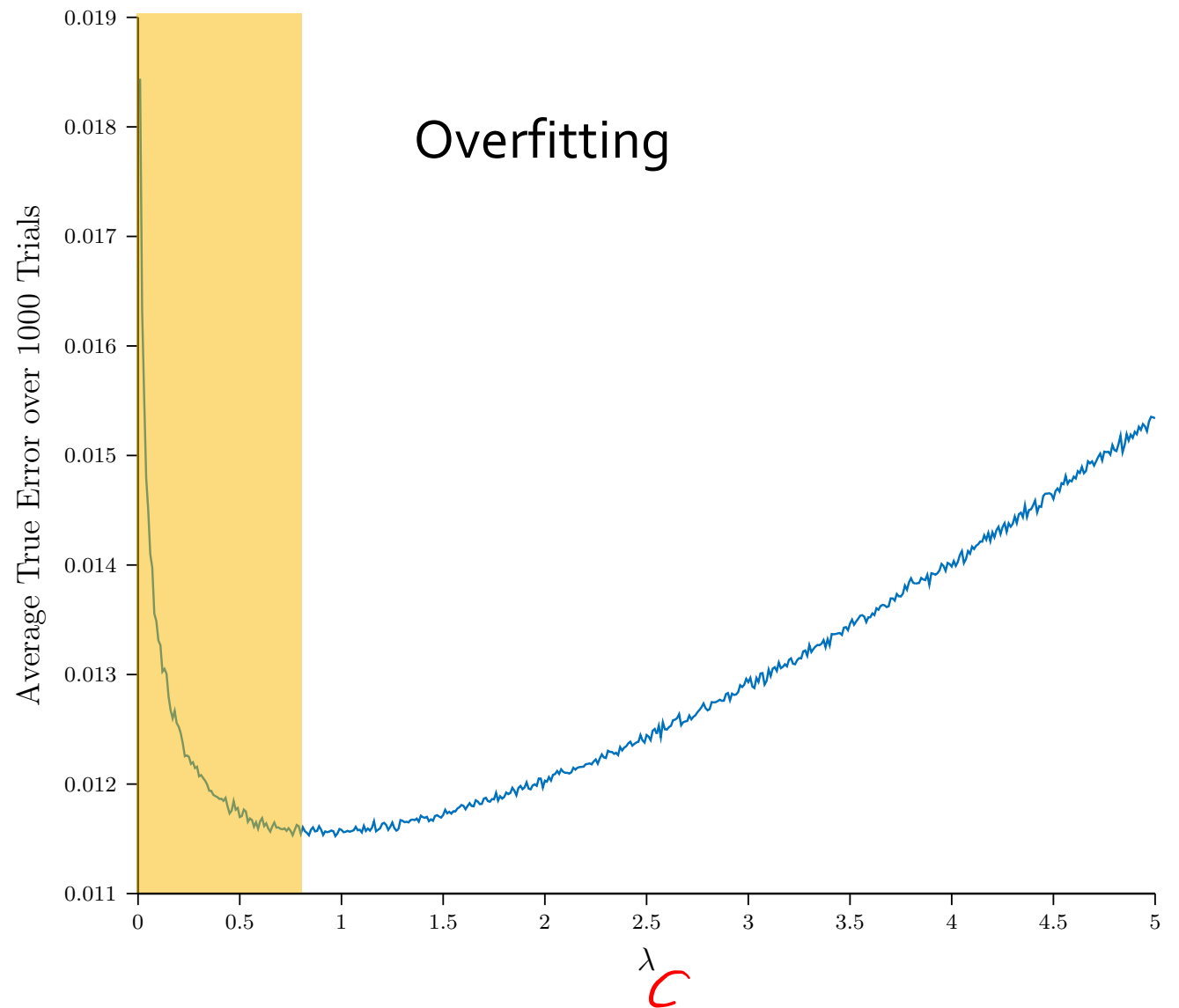
True Error
0.059

Overfit

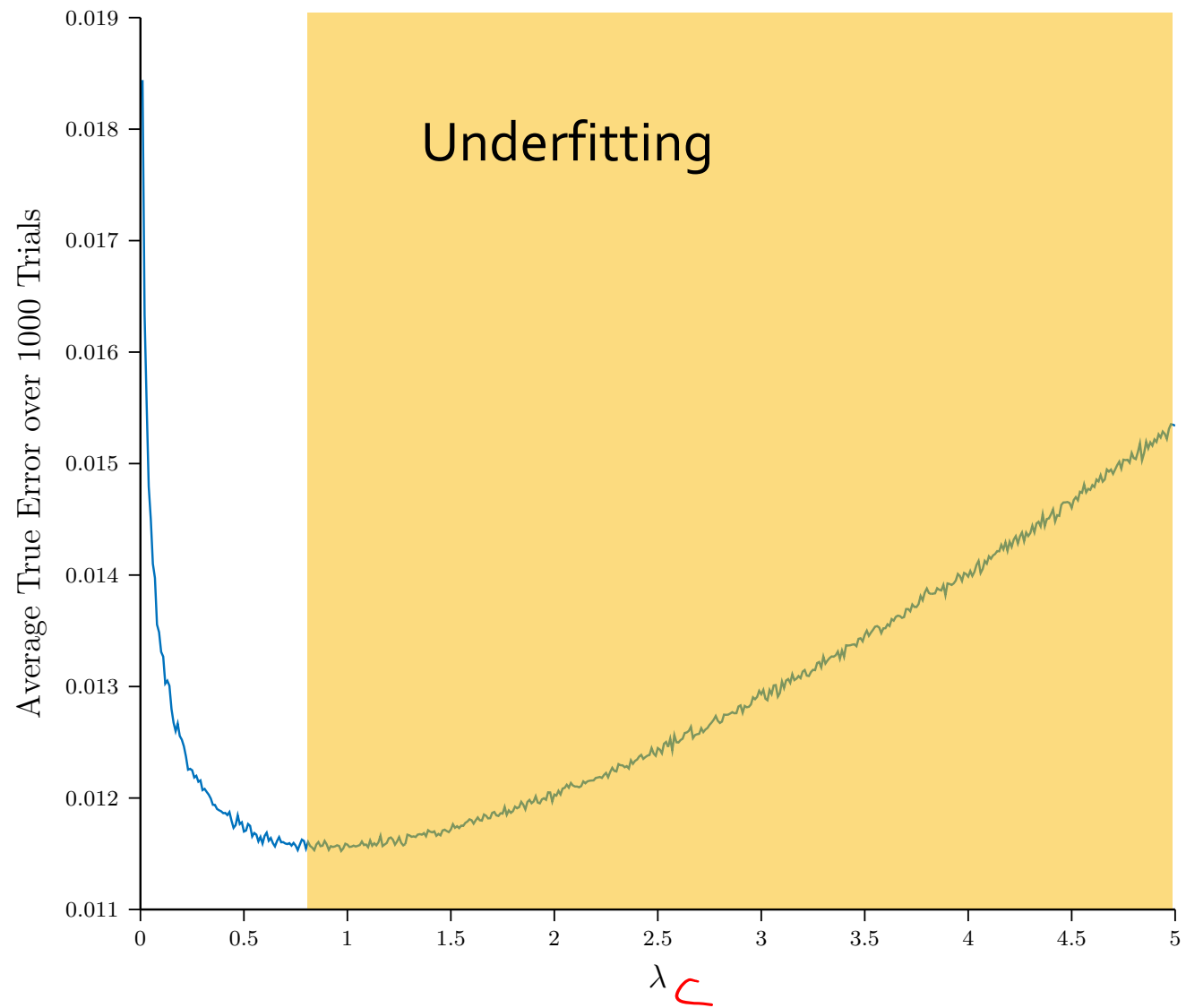
Setting λ



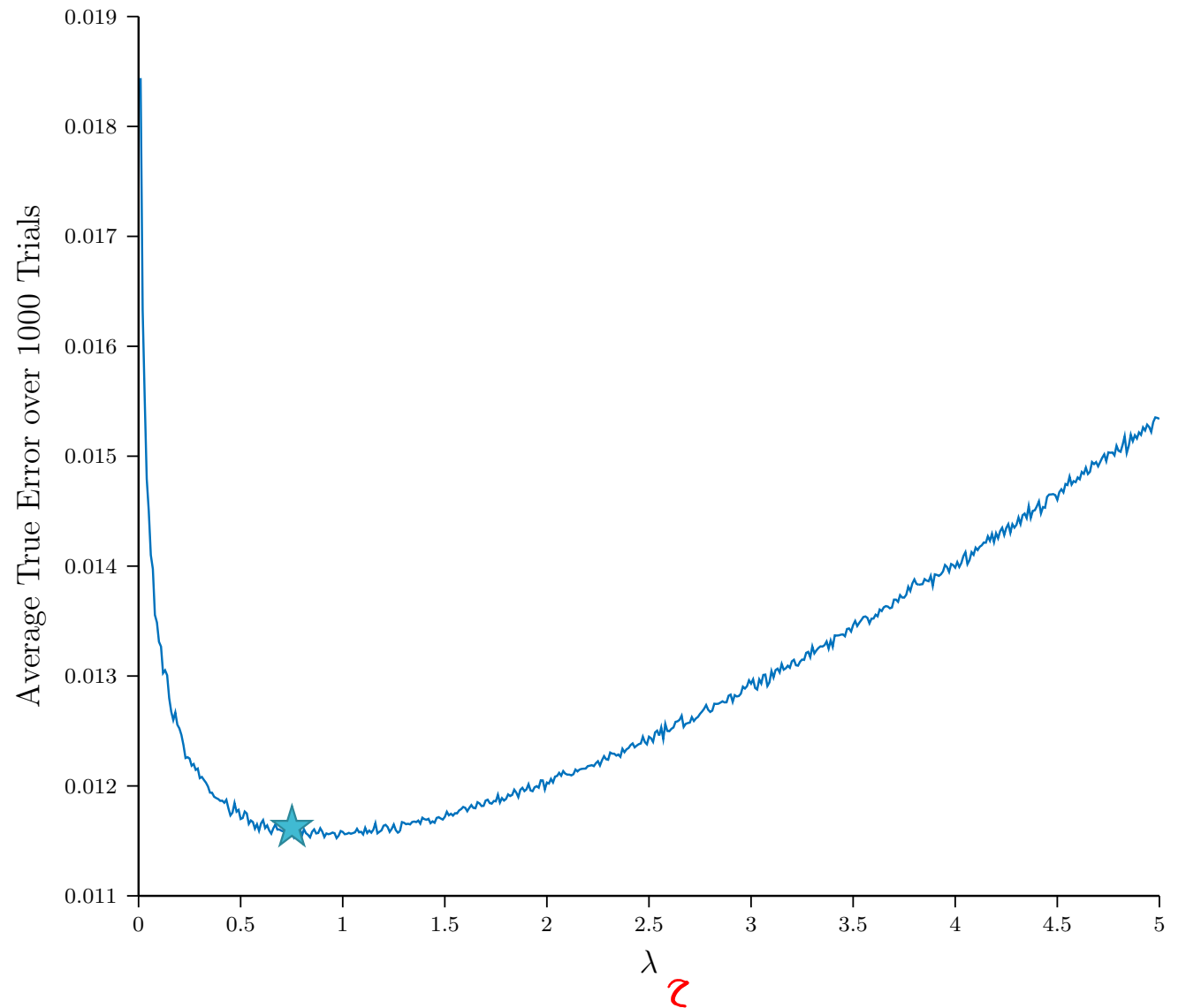
Setting λ



Setting λ



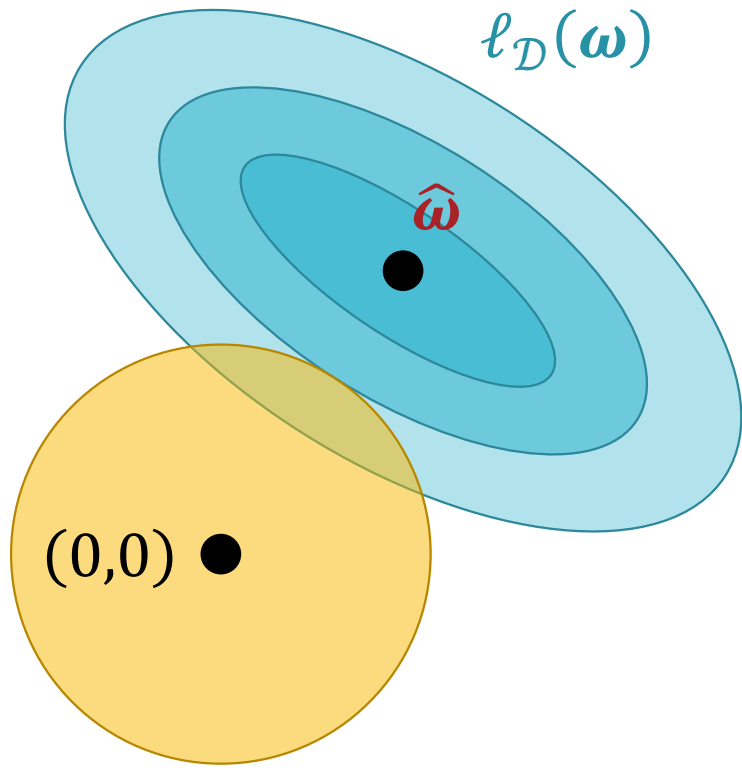
Setting λ



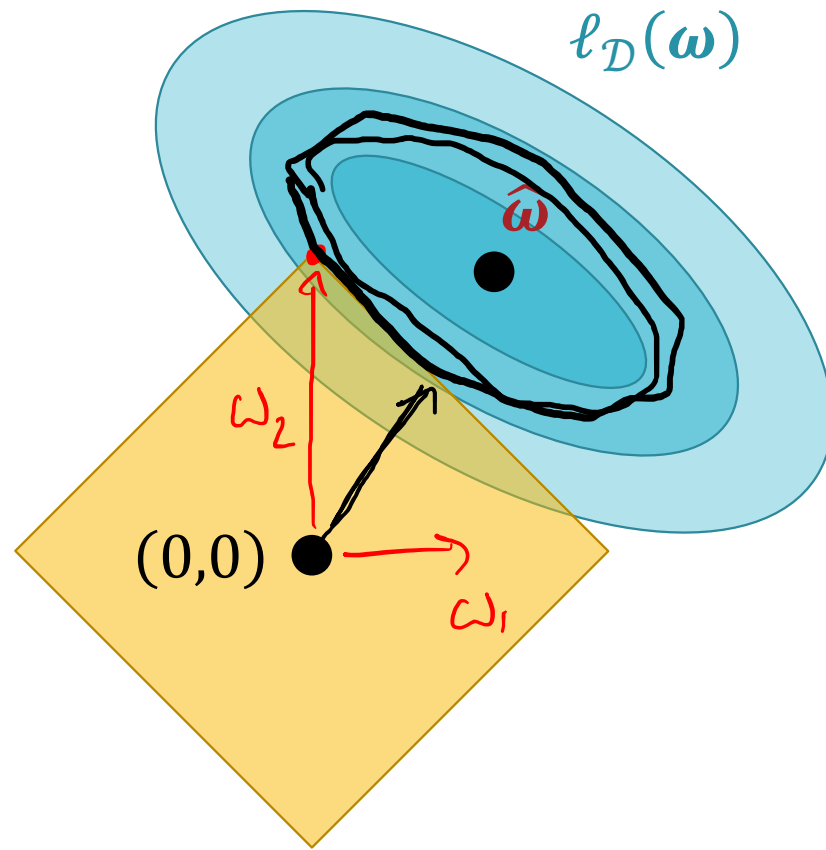
Other Regularizers

$$\ell_{\mathcal{D}}(\boldsymbol{\omega}) + \lambda r(\boldsymbol{\omega})$$

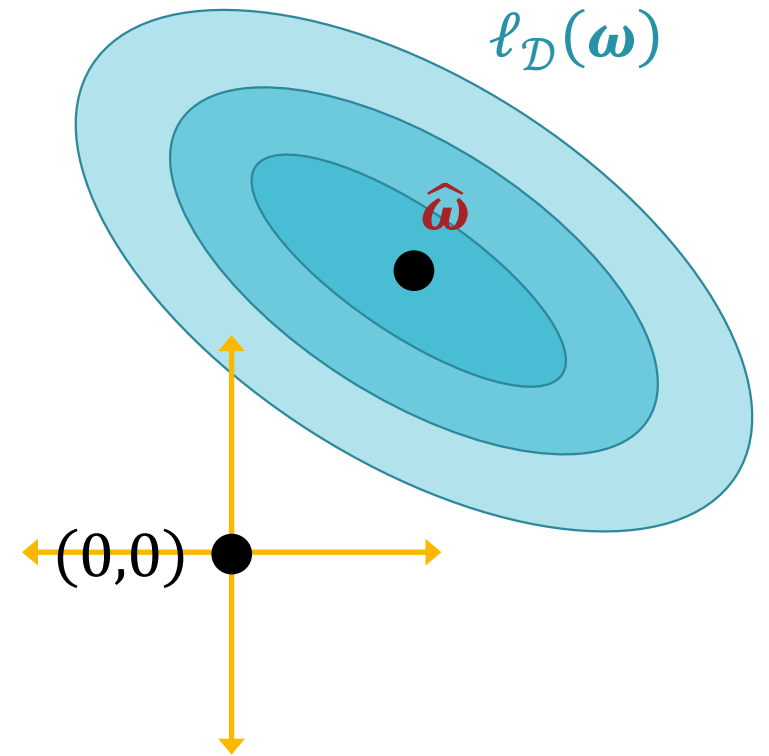
Ridge or $L2$	$r(\boldsymbol{\omega}) = \ \boldsymbol{\omega}\ _2^2 = \sum_{d=0}^D \omega_d^2$	Encourages <u>small weights</u>
<u>Lasso or $L1$</u>	$r(\boldsymbol{\omega}) = \ \boldsymbol{\omega}\ _1 = \sum_{d=0}^D \omega_d $	Encourages <u>sparsity</u>
<u>$L0$</u>	$r(\boldsymbol{\omega}) = \ \boldsymbol{\omega}\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



Ridge or $L2$



Lasso or $L1$



$L0$

Other Regularizers

M(C)LE for Linear Regression

- If we assume a linear model with additive Gaussian noise

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

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

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

$$= \underset{\boldsymbol{\omega}}{\operatorname{argmax}} \log \exp \left(-\frac{1}{2\sigma^2} (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y}) \right)$$

$$= \underset{\boldsymbol{\omega}}{\operatorname{argmin}} (\mathbf{X}\boldsymbol{\omega} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\omega} - \mathbf{y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

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 Gaussian priors on all the weights...

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

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

$$\begin{aligned} \hat{\boldsymbol{\omega}}_{MAP} &= \underset{\boldsymbol{\omega}}{\operatorname{argmin}} (X\boldsymbol{\omega} - \mathbf{y})^T (X\boldsymbol{\omega} - \mathbf{y}) + \lambda \boldsymbol{\omega}^T \boldsymbol{\omega} \\ &= \underbrace{(X^T X + \lambda I_{D+1})^{-1}} X^T \mathbf{y} \end{aligned}$$

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...

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

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

$$\hat{\boldsymbol{\omega}}_{MAP} = \underset{\boldsymbol{\omega}}{\text{argmin}} (X\boldsymbol{\omega} - \mathbf{y})^T (X\boldsymbol{\omega} - \mathbf{y}) + \lambda \|\boldsymbol{\omega}\|_1$$

- No closed form solution but 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 interpreted as MAP for linear regression

Where do features come from?

- More generally, ϕ 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 $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$,

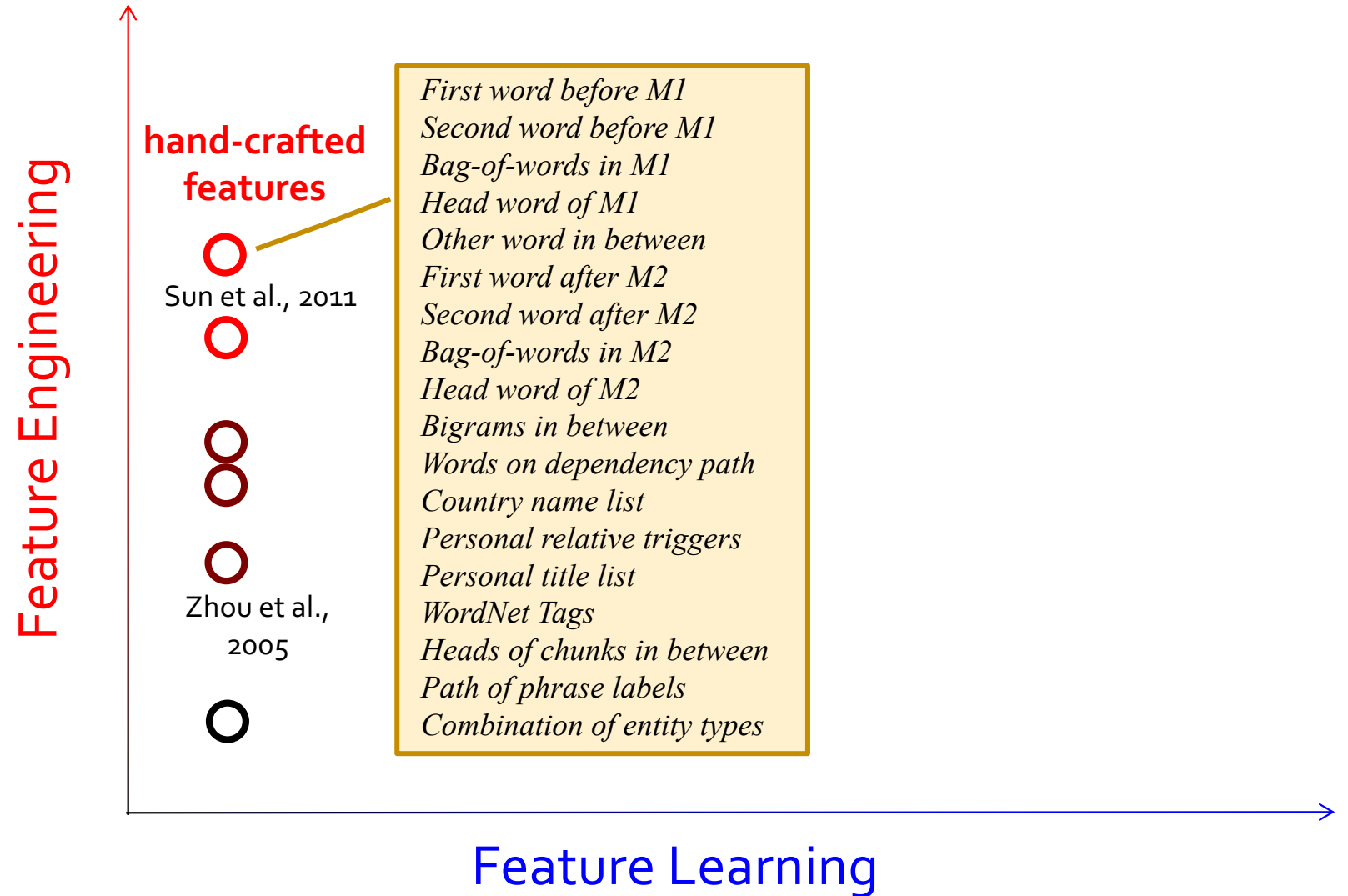
find $\boldsymbol{\omega}$ that minimizes

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

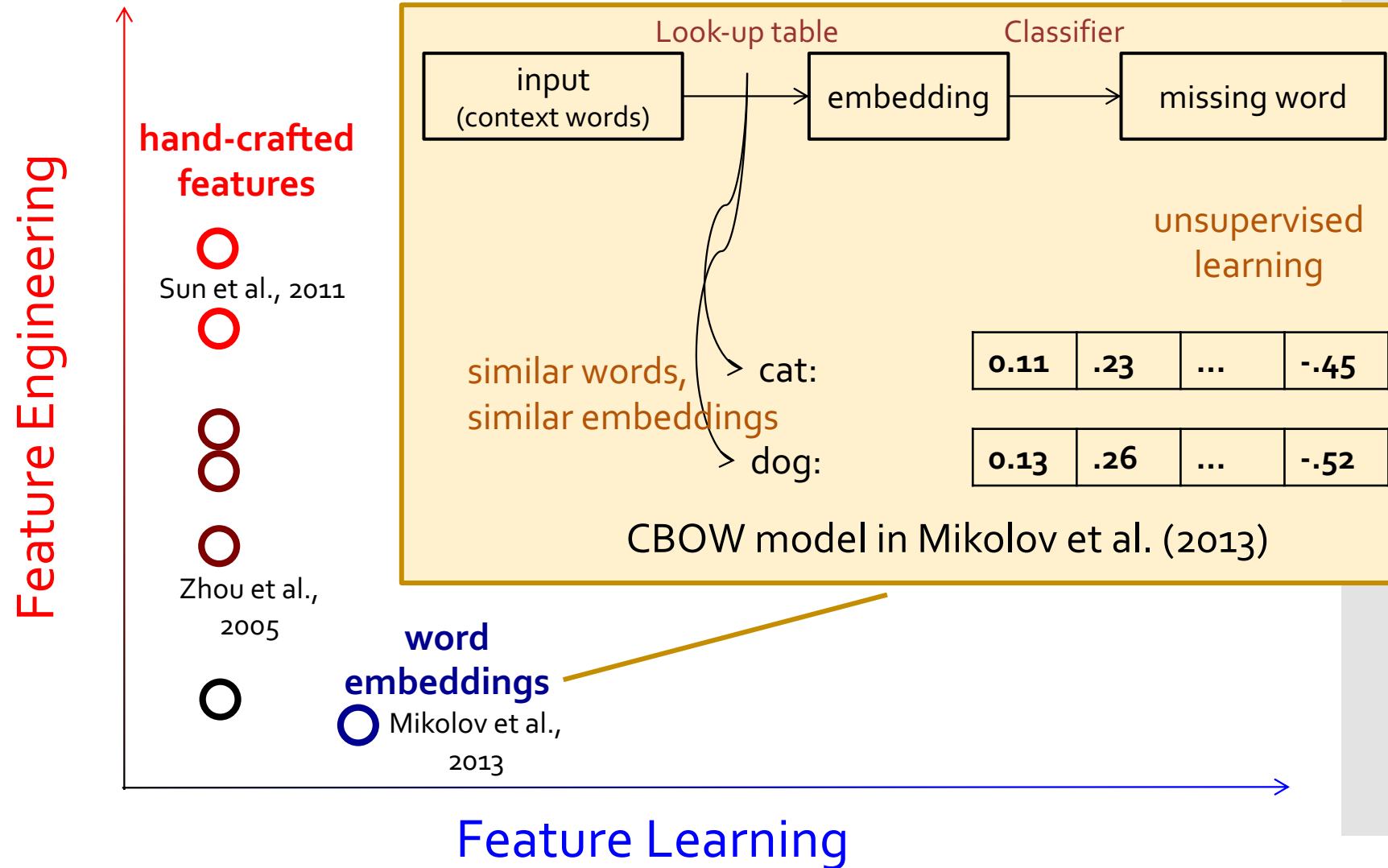
- Subject to:

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

Where do features for text data come from?

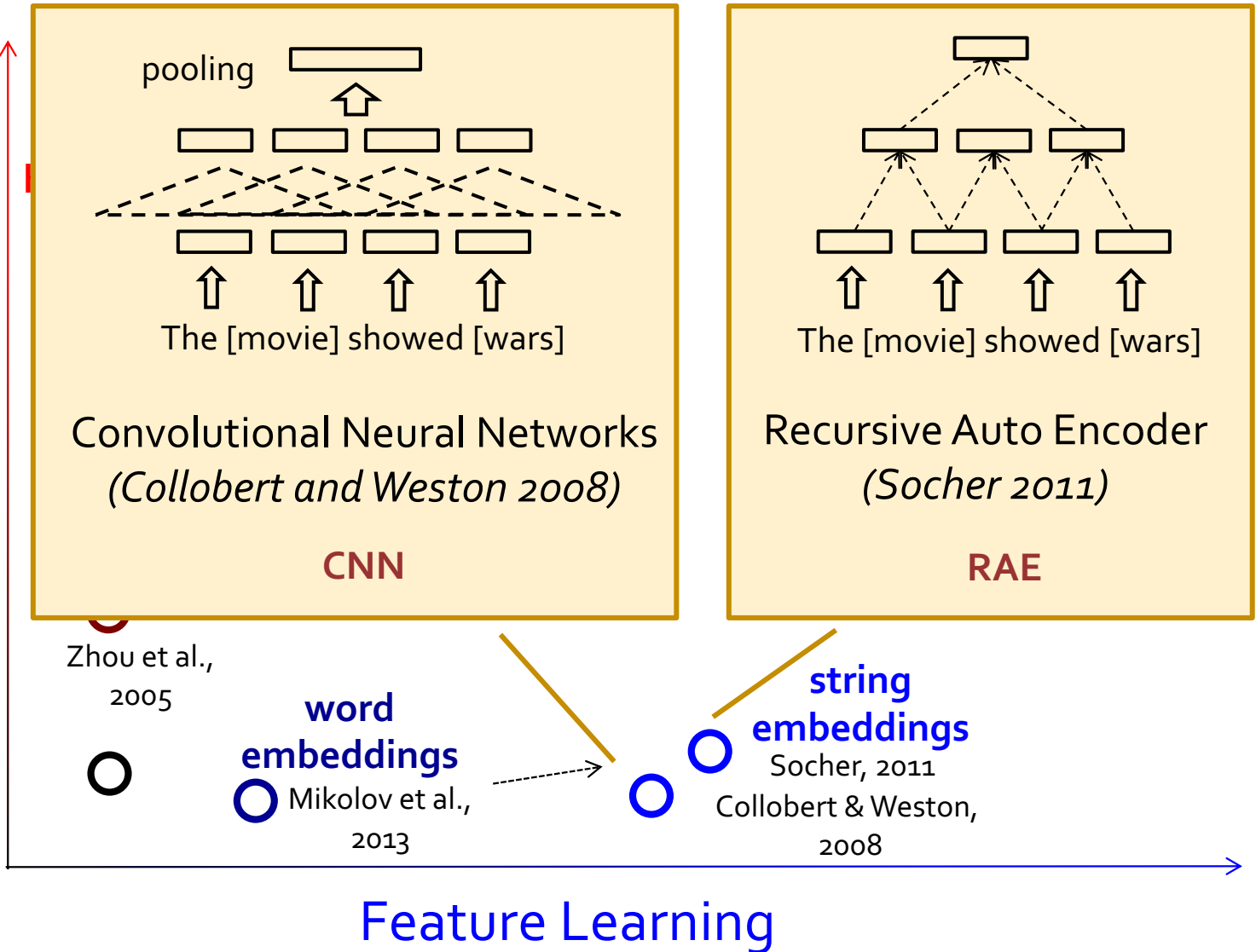


Where do features for text data come from?

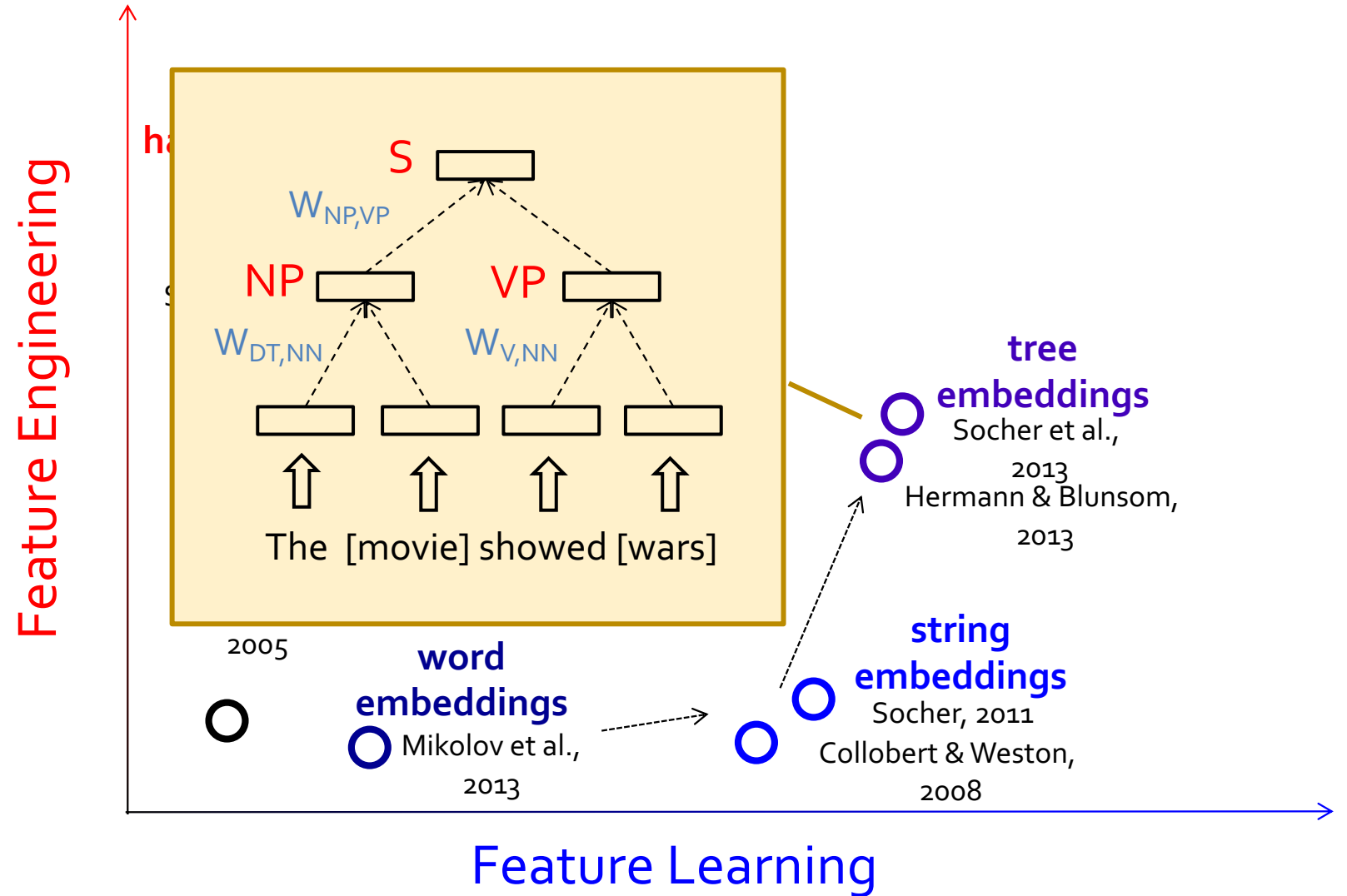


Where do features for text data come from?

Feature Engineering



Where do features for text data come from?



Where do features for text data come from?

