



03. Non-linear transformations and regularization

Winter 2020/2021

MACHINE LEARNING / BURKHARDT FUNK



Problem set 1

Given a linearly separable training set S and learning rate $\eta \in \mathbb{R}^+$

$w_0 \leftarrow 0; b_0 \leftarrow 0; k \leftarrow 0$

$R \leftarrow \max_{1 \leq i \leq l} \|x_i\|$

repeat

 for $i = 1$ to l

 if $y_i(\langle w_k \cdot x_i \rangle + b_k) \leq 0$ then

$w_{k+1} \leftarrow w_k + \eta y_i x_i$

$b_{k+1} \leftarrow b_k + \eta y_i R^2$

$k \leftarrow k + 1$

 end if

 end for

Until there are no mistakes within the *for* loop

Return the list (w_k, b_k)



Agenda

- Introduction
- Learning problem & linear classification
- Linear models: regression & logistic regression
- **Non-linear transformation, overfitting & regularization**
- Support Vector Machines and kernel learning
- Neural Networks: shallow [and deep]
- Theoretical foundation of supervised learning
- Unsupervised learning



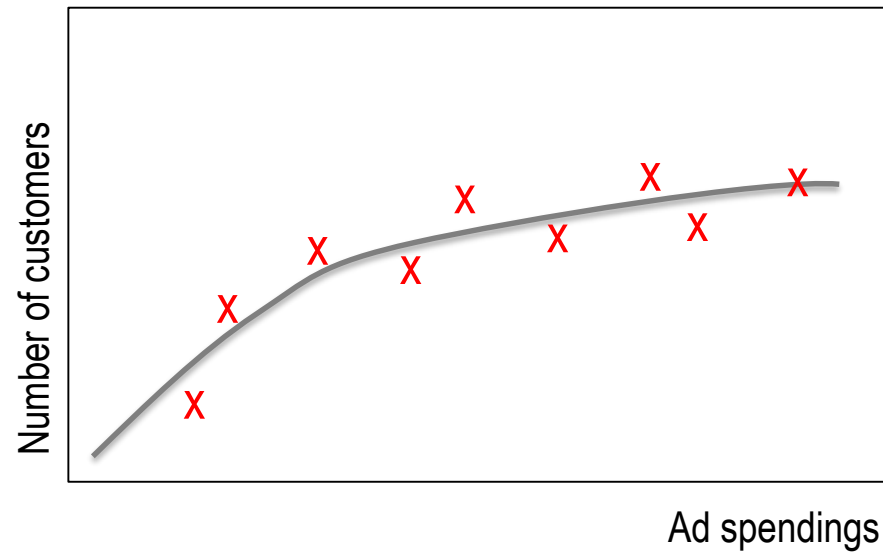
Today's agenda

- **Non-linear transformations**
- Regularization: restricting solutions



Linear is limited

Often, the target does not linearly depend on the variables, but they are related





Again: linear in what?

Linear regression implements

$$\sum_{i=0}^d \mathbf{w}_i x_i$$

Linear classification implements

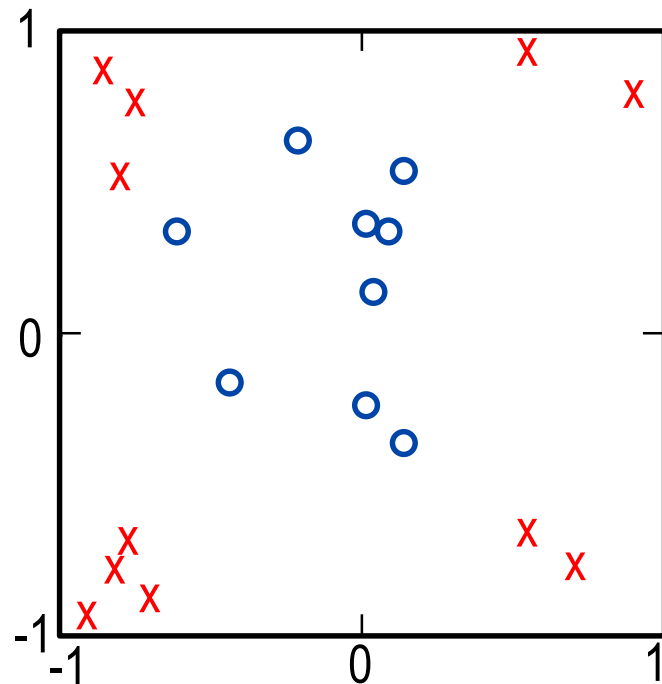
$$\text{sign}(\sum_{i=0}^d \mathbf{w}_i x_i)$$

Algorithms work because of **linearity in the weights**

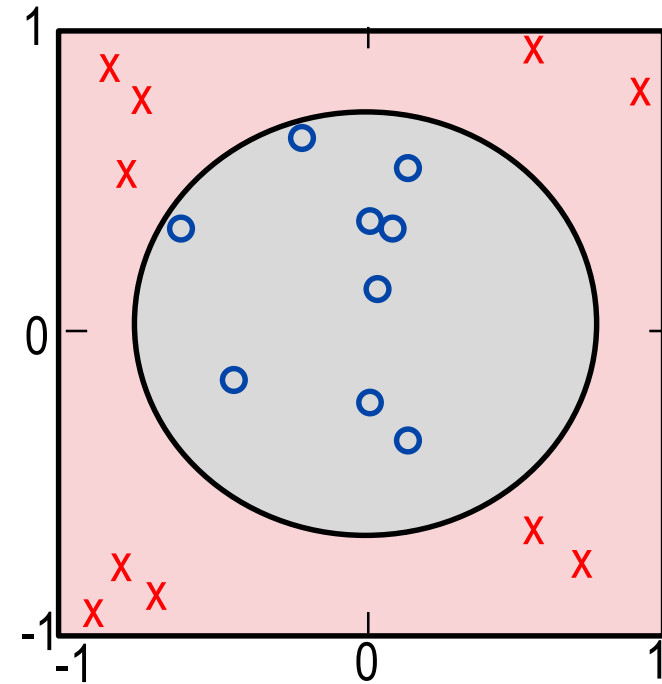


What we would like to have

Data:



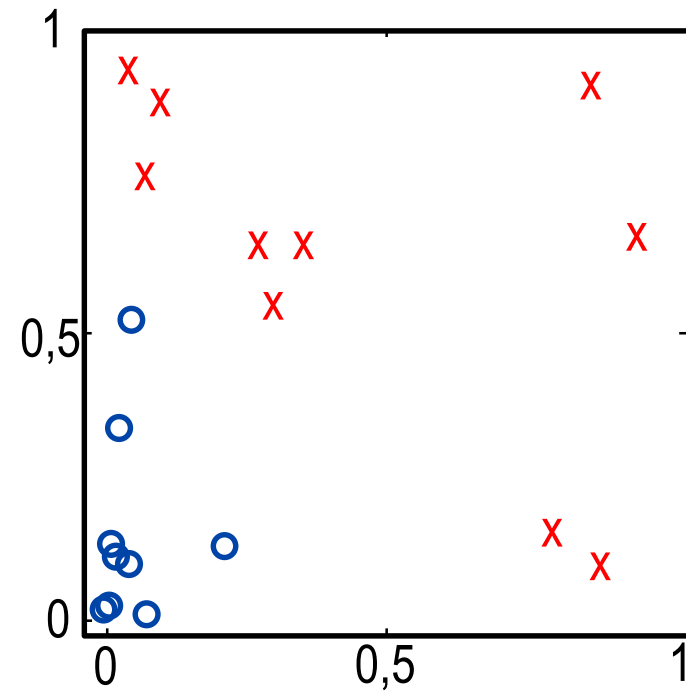
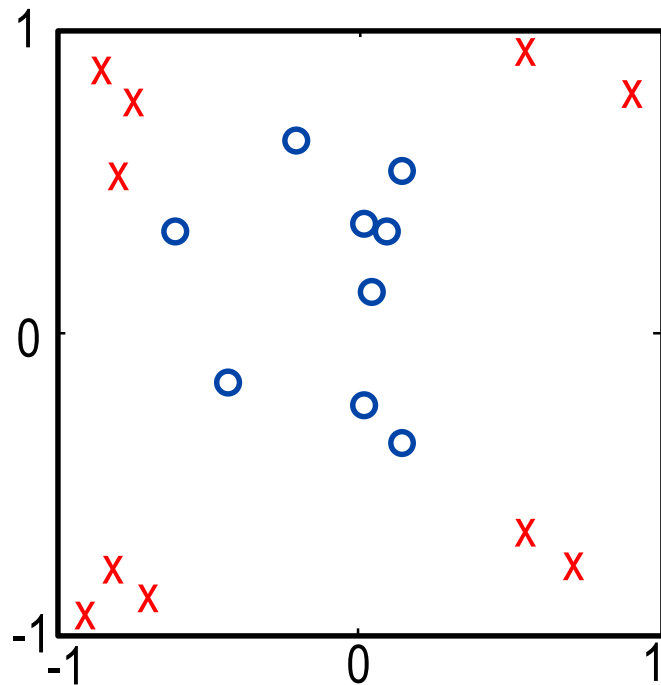
Hypothesis:





Transform the data nonlinearly

$$(x_1, x_2) \xrightarrow{\Phi} (x_1^2, x_2^2)$$





Nonlinear transforms

$$\mathbf{x} = (x_0, x_1, \dots, x_d) \xrightarrow{\Phi} \mathbf{z} = (z_0, z_1, \dots, z_{\tilde{d}})$$

$$\text{Each } z_i = \phi_i(\mathbf{x}) \qquad \mathbf{z} = \Phi(\mathbf{x})$$

$$\text{Example: } \mathbf{z} = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$$

Final hypothesis $g(\mathbf{x})$ operates on \mathcal{X} :

$$\text{sign}(\tilde{\mathbf{w}}^\top \Phi(\mathbf{x})) \quad \text{or} \quad \tilde{\mathbf{w}}^\top \Phi(\mathbf{x})$$



Quiz

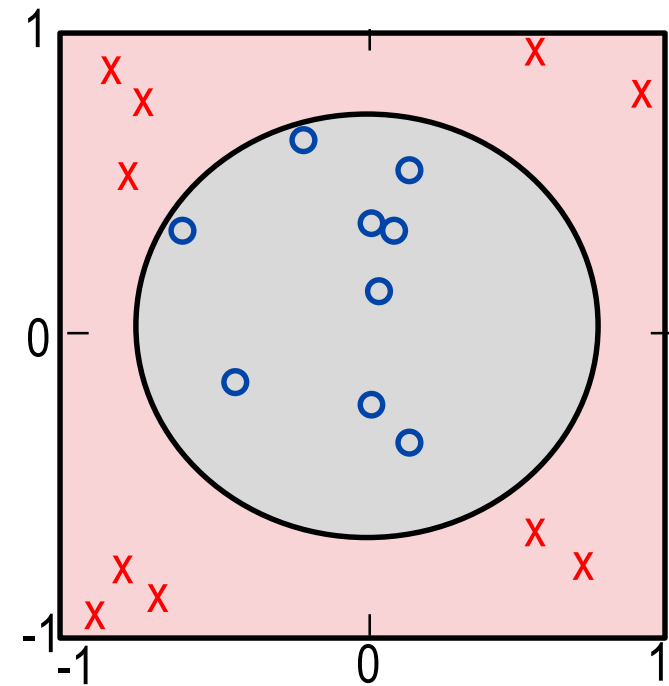
- Let the hypothesis set be $\tilde{h} = \text{sign}(\tilde{\mathbf{w}}^\top \Phi(\mathbf{x}))$ with $\Phi(\mathbf{x}) = \mathbf{z} = (1, x_1^2, x_2^2)$
- Can you tell what geometric forms $\tilde{\mathbf{w}} = (\tilde{w}_0, \tilde{w}_1, \tilde{w}_2)$ corresponds to?

$$\tilde{\mathbf{w}} = (1, -1, -1)$$

$$\tilde{\mathbf{w}} = (-1, 1, 1)$$

$$\tilde{\mathbf{w}} = (1, -1, -2)$$

$$\tilde{\mathbf{w}} = (1, 1, -1)$$





*We will explore this
in a problem set*

A flexible way to transform raw data: radial basis functions

— The following transformation uses radial basis functions

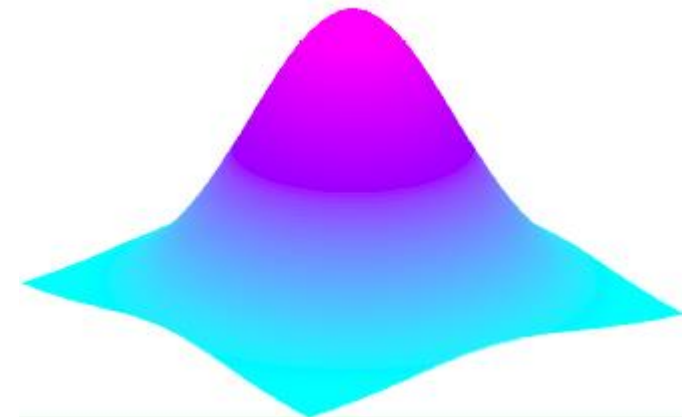
$$z = \Phi(x) = \begin{pmatrix} 1 \\ B(x, \mu_1, \gamma) \\ \vdots \\ B(x, \mu_d, \gamma) \end{pmatrix} \text{ where } B(x, \mu, \gamma) = \exp(-\gamma \|x - \mu\|^2)$$



Basic RBF model (using training data points)

Standard form (where $(x_n, y_n) \in D$ is our training data)

$$h(x) = \sum_{n=1}^N w_n \exp(-\gamma \|x - x_n\|^2)$$





Today's agenda

- Non-linear transformations
- **Regularization: restricting solutions**



Review

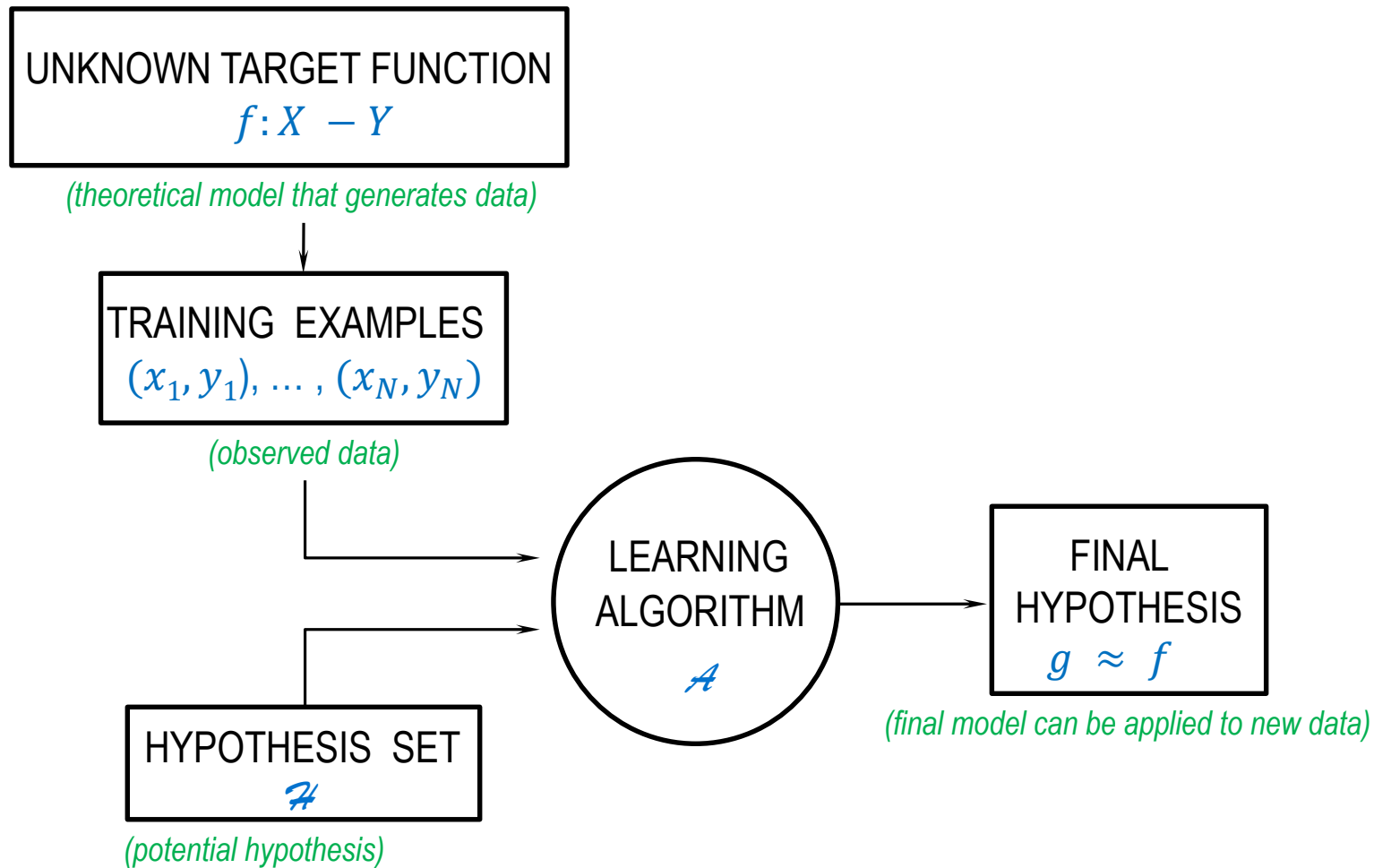




Illustration of overfitting

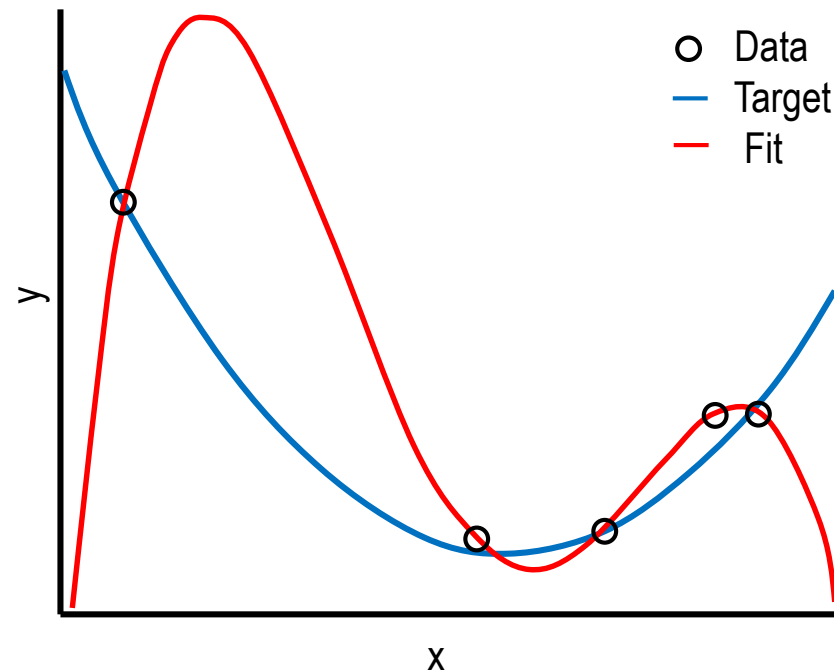
Simple target function

5 data points include a bit **noise**

4th-order polynomial fit

$E_{in} = 0$ (in-sample error)

E_{out} is huge (out-of-sample error)





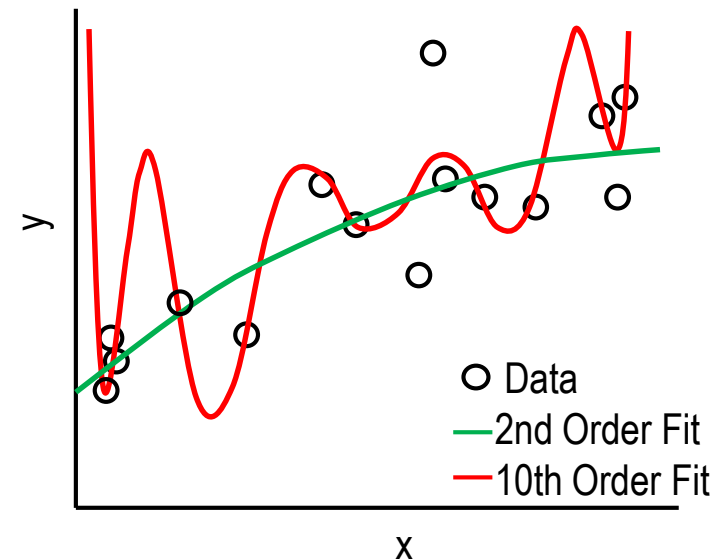
What is the optimal hypothesis set?

Two learners O (10th order polynomial)
and R (2nd order polynomial)

We know the target is 10th order

O chooses \mathcal{H}_{10}

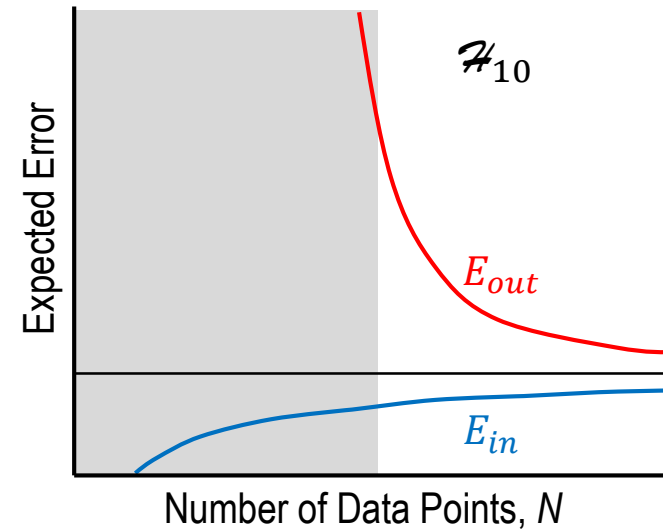
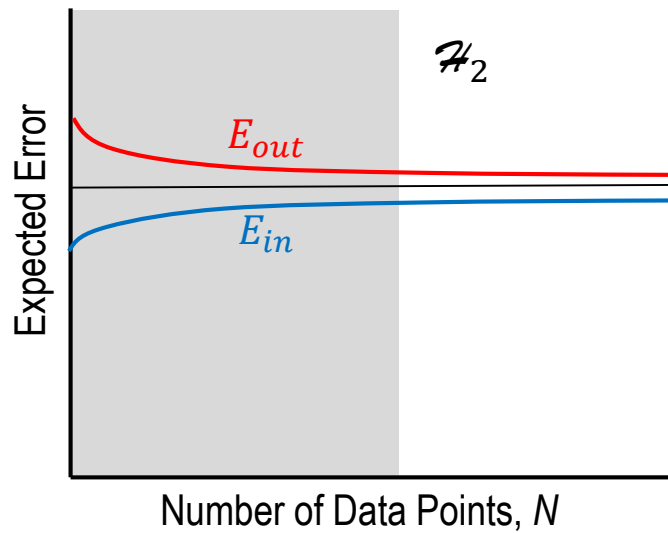
R chooses \mathcal{H}_2



Learning a 10th order target

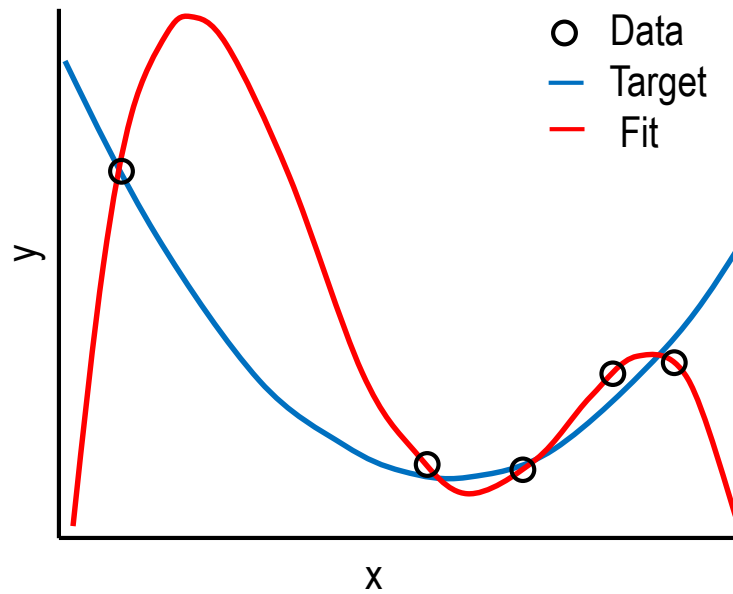


Learning curves for both hypothesis sets

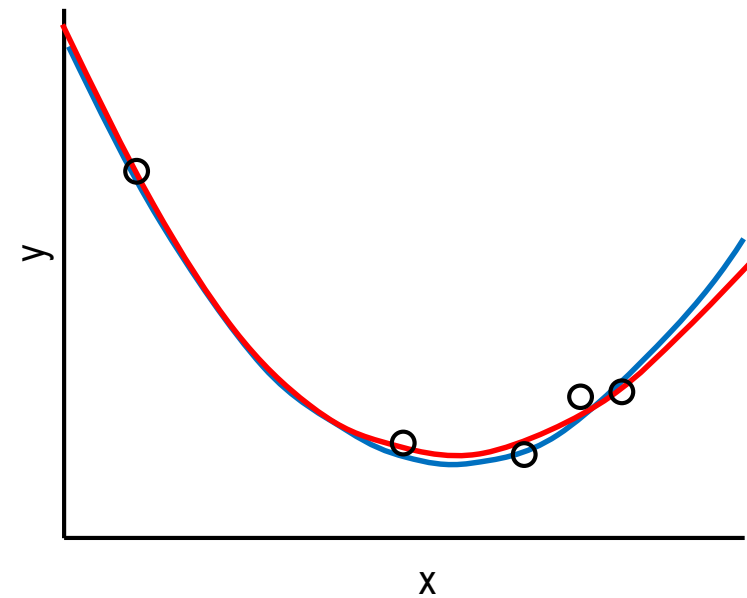




The goal of **regularization**



free fit



restrained fit



Unconstrained solution

Given $(x_1, y_1), \dots, (x_N, y_N) \rightarrow (\mathbf{z}_1, y_1), \dots, (\mathbf{z}_N, y_N)$

$$\text{Minimize } E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{z}_n - y_n)^2$$

$$\text{Minimize } \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^\top (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

$$\mathbf{w}_{lin} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{y}$$



Constraining the weights

Hard constraint: \mathcal{H}_2 is constrained version of \mathcal{H}_{10} with $w_q = 0$ for $q > 2$

Softer version: $\sum_{q=0}^Q w_q^2 \leq \mathcal{C}$ “soft-order” constraint

Minimize $\frac{1}{N} (Z\mathbf{w} - \mathbf{y})^\top (Z\mathbf{w} - \mathbf{y})$

subject to: $\mathbf{w}^\top \mathbf{w} \leq \mathcal{C}$

Solution: \mathbf{w}_{reg} instead of \mathbf{w}_{lin}



Solving for \mathbf{w}_{reg}

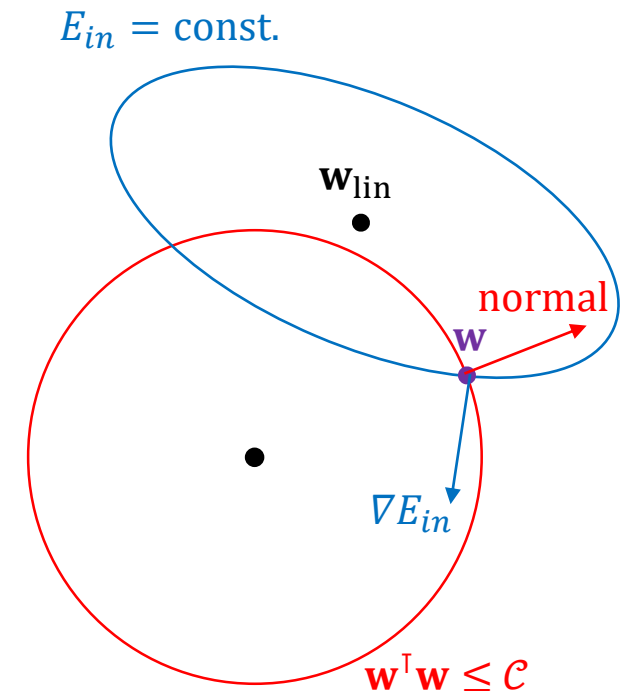
$$\text{Minimize } E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^\top (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

$$\text{subject to } \mathbf{w}^\top \mathbf{w} \leq \mathcal{C}$$

$$\nabla E_{in}(\mathbf{w}_{\text{reg}}) \propto -\mathbf{w}_{\text{reg}} =: -2 \frac{\lambda}{N} \mathbf{w}_{\text{reg}}$$

$$\nabla E_{in}(\mathbf{w}_{\text{reg}}) + 2 \frac{\lambda}{N} \mathbf{w}_{\text{reg}} = \mathbf{0}$$

$$\text{Minimize } E_{in}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^\top \mathbf{w}$$





Augmented error

Minimizing $E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$

$$= \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} \text{ unconditionally}$$

— solves —

Minimizing $E_{\text{in}}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$

subject to: $\mathbf{w}^T \mathbf{w} \leq C$



The solution

Minimize $E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$

$$= \frac{1}{N} ((Z\mathbf{w} - \mathbf{y})^T (Z\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w})$$

$$\nabla E_{\text{aug}}(\mathbf{w}) = 0 \quad \Rightarrow \quad Z^T(Z\mathbf{w} - \mathbf{y}) + 2\lambda\mathbf{w} = 0 \text{ (typically we drop the 2)}$$

$$\boxed{\mathbf{w}_{\text{reg}} = (Z^T Z + \lambda I)^{-1} Z^T \mathbf{y}} \quad \text{(with regularization)}$$

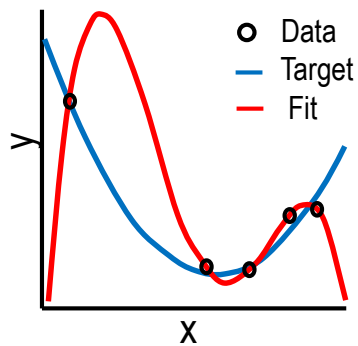
as opposed to $\mathbf{w}_{\text{lin}} = (Z^T Z)^{-1} Z^T \mathbf{y}$ (without regularization)



The result

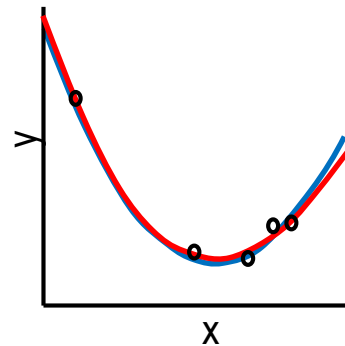
Minimizing $E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$ for different λ 's:

$\lambda = 0$

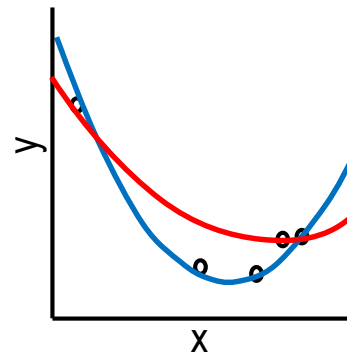


overfitting

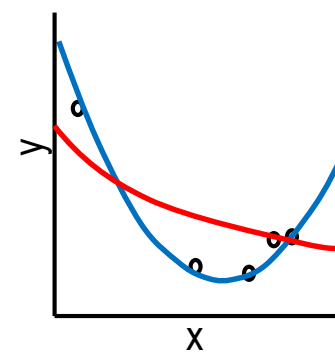
$\lambda = 0.0001$



$\lambda = 0.01$



$\lambda = 1$



underfitting



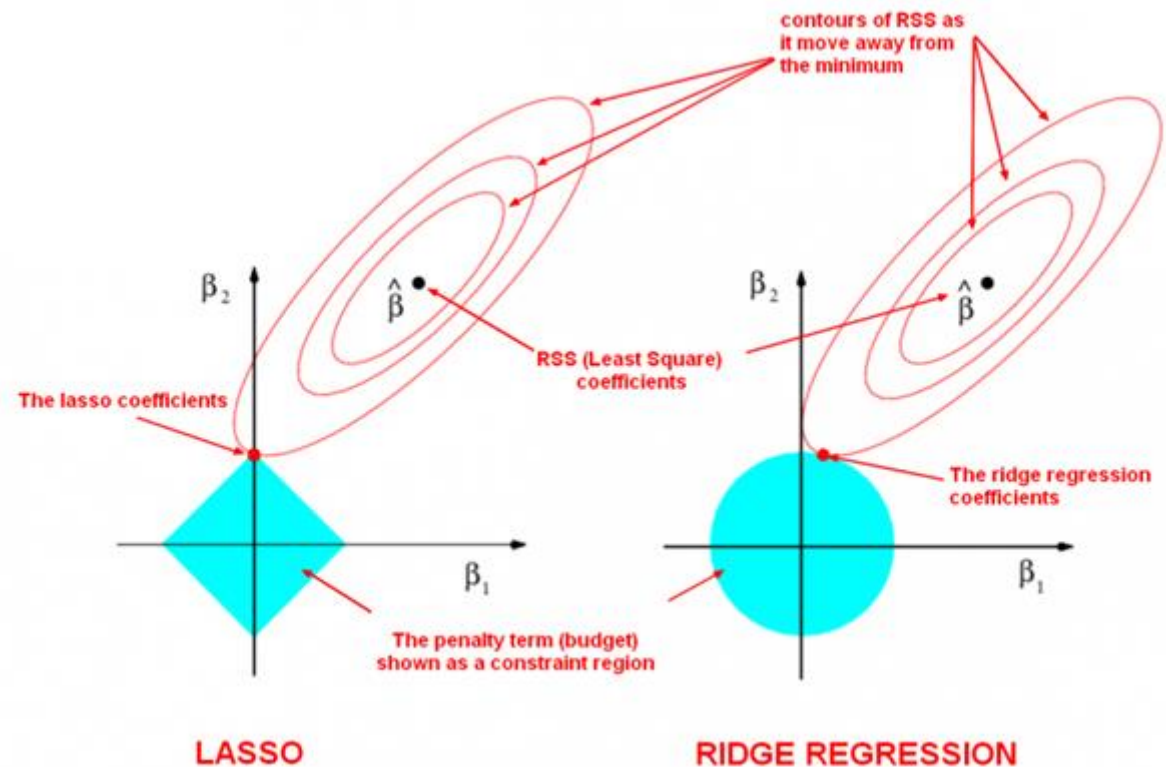


Another exciting regularizer

—In LASSO regression we use

$$E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} |\mathbf{w}|$$

—LASSO regression can be used for variable selection

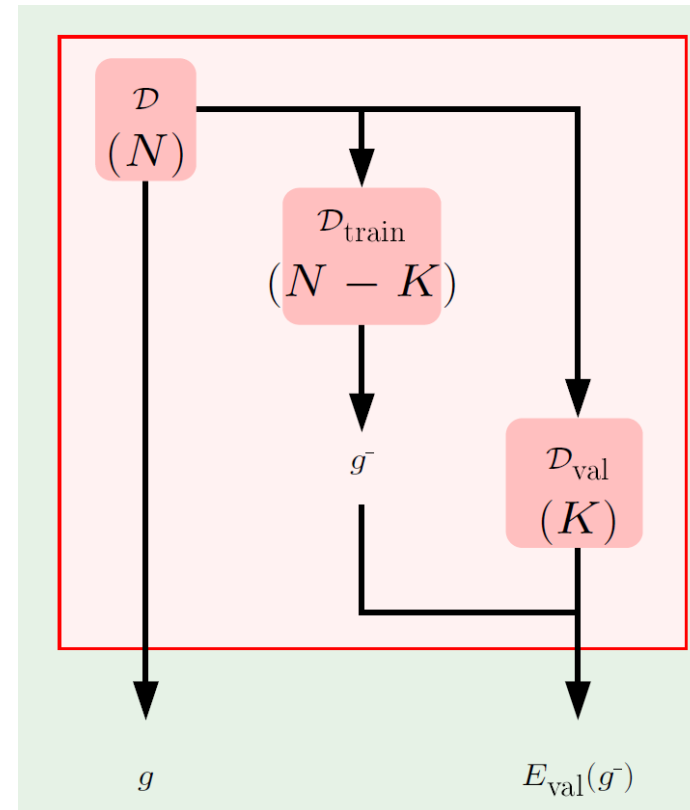


Source: <https://www.quora.com/How-would-you-describe-the-difference-between-linear-regression-lasso-regression-and-ridge-regression>



Finding the best λ

- We want to estimate/minimize E_{out} - we can do so by using a separate dataset that has not been used for training
- This approach is called validation
- Now we can minimize E_{val} by finding the best λ





Additional material



The learning algorithm

—Finding w_1, \dots, w_n for

$$h(x) = \sum_{n=1}^N w_n \exp(-\gamma \|x - x_n\|^2)$$

—Based on $D = (x_1, y_1), \dots, (x_n, y_n)$

—Can we choose w_1, \dots, w_n such that $E_{in} = 0$ or $h(x_n) = y_n$?

—We need to solve $y_n = \sum_{m=1}^N w_m \exp(-\gamma \|x_n - x_m\|^2)$



Solution

— $y_n = \sum_{m=1}^N w_m \exp(-\gamma \|x_n - x_m\|^2)$ gives us N equations with N unknowns

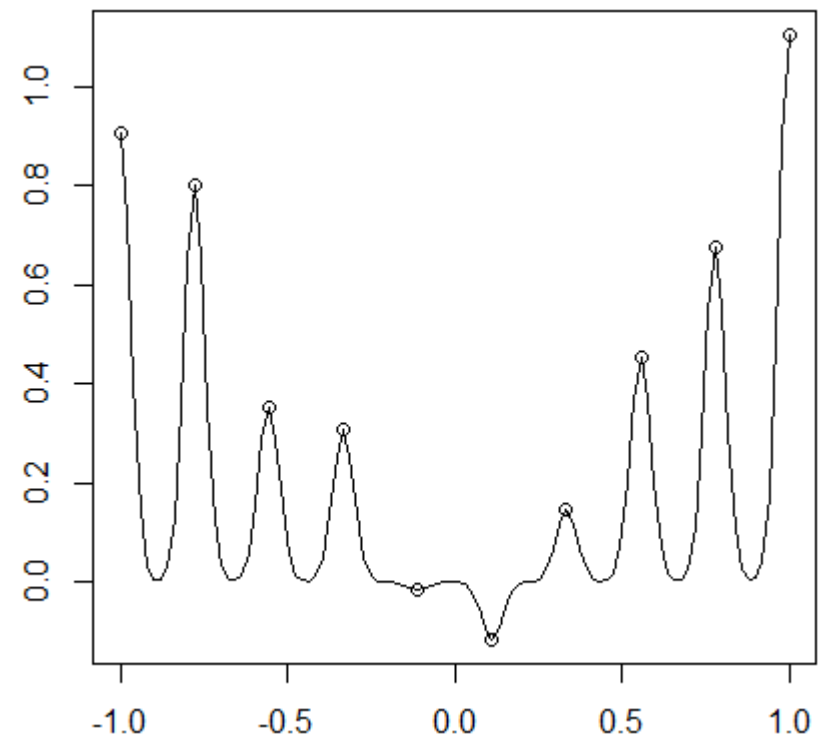
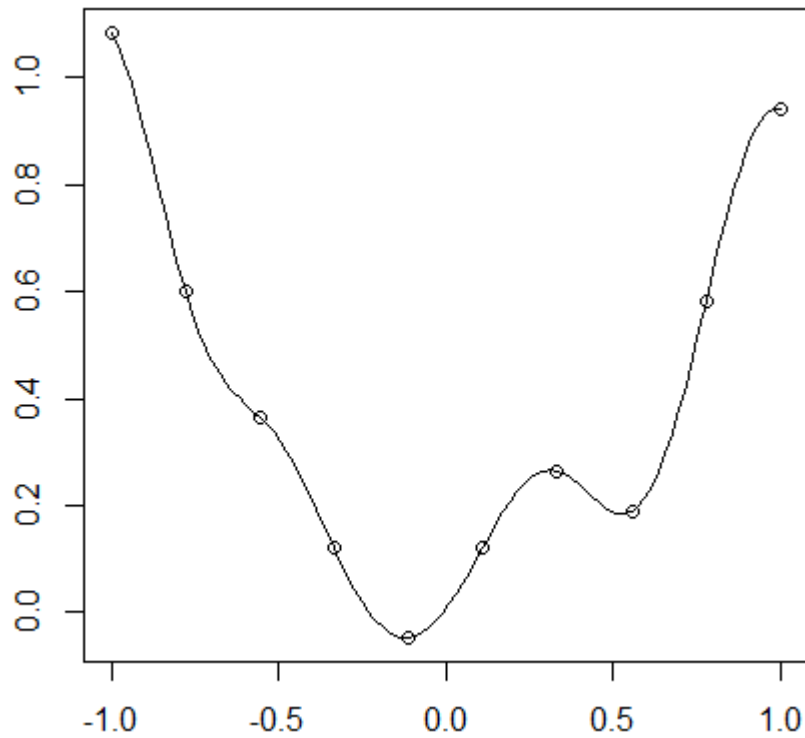
$$\underbrace{\begin{bmatrix} \exp(-\gamma \|x_1 - x_1\|^2) & \cdots & \exp(-\gamma \|x_1 - x_N\|^2) \\ \exp(-\gamma \|x_2 - x_1\|^2) & \cdots & \exp(-\gamma \|x_2 - x_N\|^2) \\ \vdots & & \vdots \\ \exp(-\gamma \|x_N - x_1\|^2) & & \exp(-\gamma \|x_N - x_N\|^2) \end{bmatrix}}_{\Phi} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

— If Φ is invertible, then $\mathbf{w} = \Phi^{-1} \mathbf{y}$



The impact of γ

— Large or small γ ?





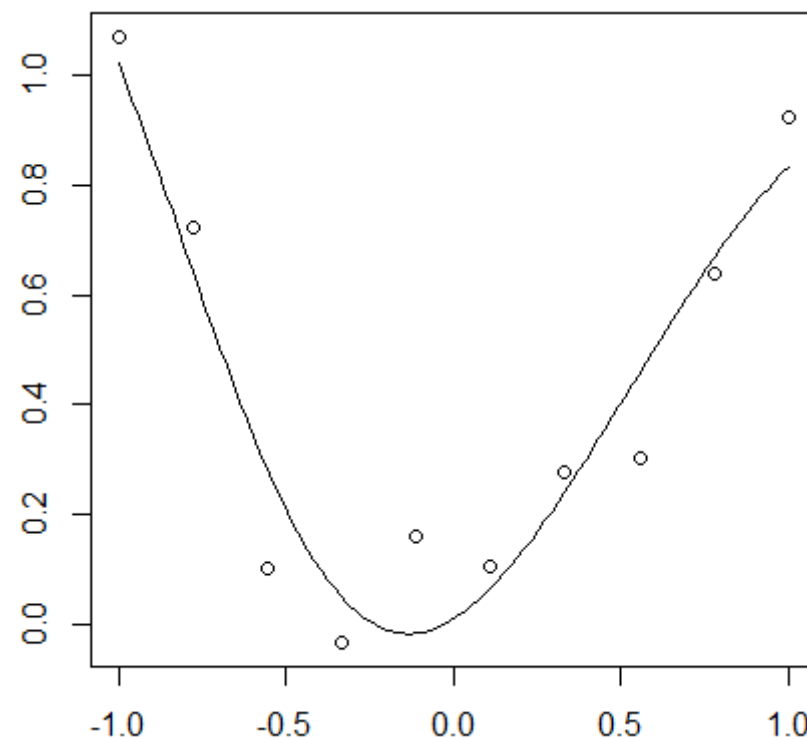
What about generalization?

— So far, we have fitted the training dataset perfectly → hypothesis is not likely to generalize

— Solution: Regularization – add penalty term:

$$\mathbf{w}_{ridge} = (\Phi^T \Phi + \delta^2 \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

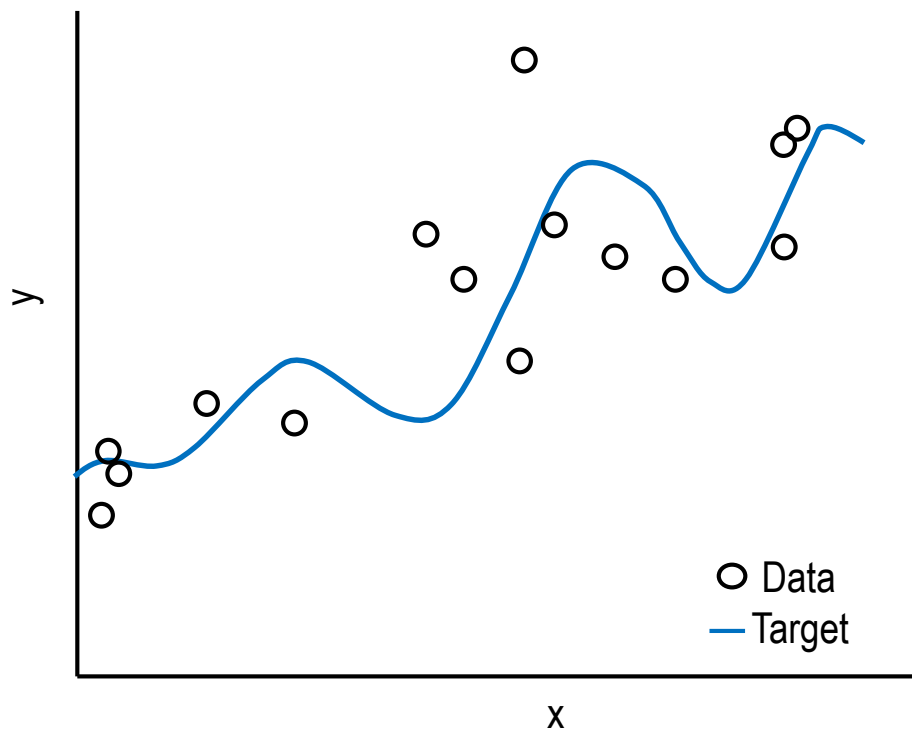
— This will be the topic of the lecture on overfitting and regularization





A detailed experiment to understand overfitting

Impact of noise level and target complexity



$$y = f(x) + \underbrace{\epsilon(x)}_{\sigma^2} = \sum_{q=0}^{Q_f} a_q x^q + \epsilon(x)$$

noise level: σ^2

target complexity: Q_f

Data set size: N



Legendre polynomials

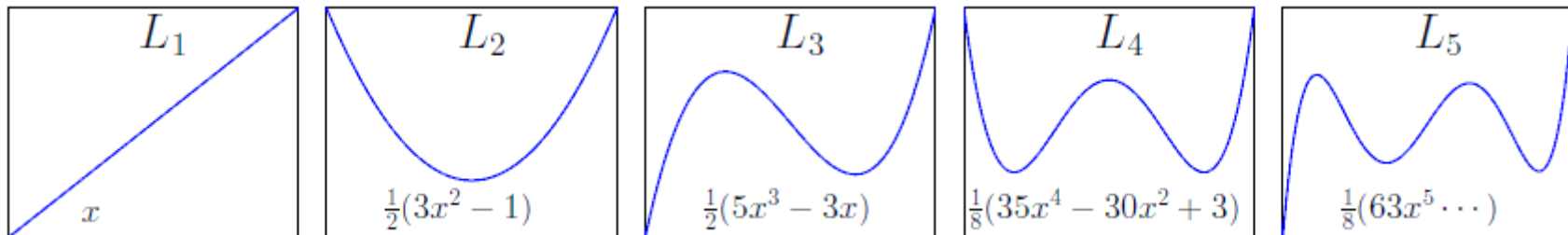
\mathcal{H}_Q : polynomials of Order Q

linear regression in \mathcal{Z} space

$$z = \begin{bmatrix} 1 \\ L_1(x) \\ \vdots \\ L_Q(x) \end{bmatrix}$$

$$\mathcal{H}_Q = \left\{ \sum_{q=0}^Q w_q L_q(x) \right\}$$

Legendre polynomials:



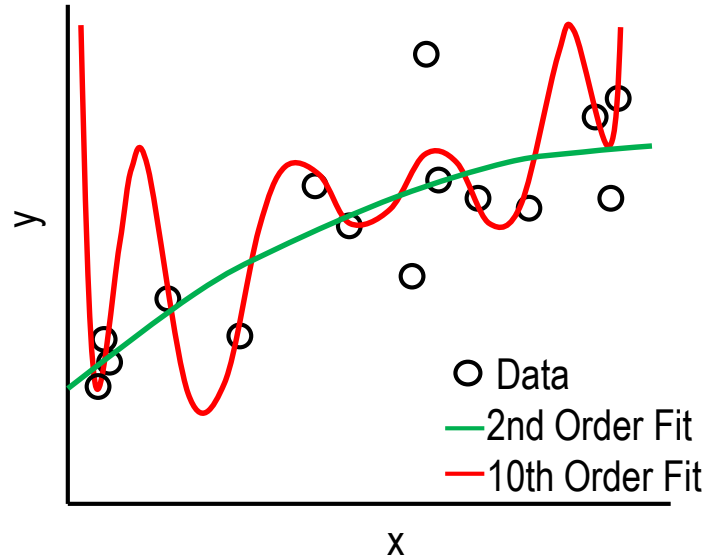


The overfit measure

We fit the data set $(x_1, y_1), \dots, (x_N, y_N)$ using our two models:

\mathcal{H}_2 : 2nd-order polynomials

\mathcal{H}_{10} : 10th-order polynomials



Compare out-of-sample errors of

$g_2 \in \mathcal{H}_2$ and $g_{10} \in \mathcal{H}_{10}$

overfit measure: $E_{out}(g_{10}) - E_{out}(g_2)$



Weight ‘decay’

Minimizing $E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$ is called *weight decay*. Why?

Gradient descent: $\mathbf{w}(t + 1) = \mathbf{w}(t) - \eta \nabla E_{\text{in}}(\mathbf{w}(t)) - 2\eta \frac{\lambda}{N} \mathbf{w}(t)$

$$= \mathbf{w}(t) \left(1 - \frac{2\eta\lambda}{N}\right) - \eta \nabla E_{\text{in}}(\mathbf{w}(t))$$