

# \*\*\* Applied Machine Learning Fundamentals \*\*\*

## Regression

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Winter term 2023/2024



Find all slides on [GitHub](#) (DaWe1992/Applied\_ML\_Fundamentals)

# Lecture Overview

Unit I	Machine Learning Introduction
Unit II	Mathematical Foundations
Unit III	Bayesian Decision Theory
<b>Unit IV</b>	<b>Regression</b>
Unit V	Classification I
Unit VI	Evaluation
Unit VII	Classification II
Unit VIII	Clustering
Unit IX	Dimensionality Reduction

# Agenda for this Unit

- 1 Introduction
- 2 Solutions to Regression
- 3 Basis Function Regression
- 4 Regularization Techniques
- 5 Wrap-Up

## Section: Introduction

What is Regression?  
Least Squares Error Function

# Linear Regression Overview

**Type of target variable:**

Continuous

**Type of training information:**

Supervised learning

**Example availability:**

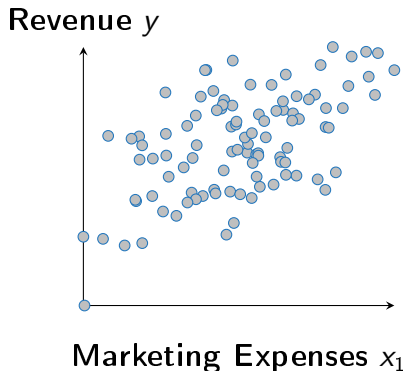
Batch learning

**Algorithm sketch:** Given the training data  $\mathcal{D}$ , the algorithm derives a function (model function / hypothesis) of the type

$$h_{\theta}(\mathbf{x}) := \theta_0 + \theta_1 x_1 + \cdots + \theta_m x_m \quad (\mathbf{x} \in \mathbb{R}^m, \boldsymbol{\theta} \in \mathbb{R}^{m+1}) \quad (1)$$

from the data.  $\boldsymbol{\theta}$  is the parameter vector containing the coefficients to be estimated by the regression algorithm. Once  $\boldsymbol{\theta}$  is learned, it can be used for prediction.

## Example Dataset: Revenues



- Find a linear function of the form:

$$h_{\theta}(\mathbf{x}) := \theta_0 + \theta_1 x_1 + \cdots + \theta_m x_m$$

- Usually we set:  $x_0 := 1$ :

$$\hat{\mathbf{x}} \in \mathbb{R}^{m+1} := [1, \mathbf{x}]^{\top}$$

$$h_{\theta}(\hat{\mathbf{x}}) = \sum_{j=0}^m \theta_j x_j = \boldsymbol{\theta}^{\top} \hat{\mathbf{x}}$$



# Error Function for Regression

- We need an error function  $\mathcal{J}(\boldsymbol{\theta})$  in order to know how well the function fits:

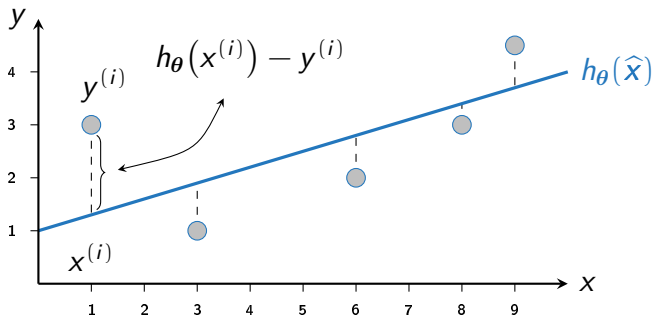
$$\mathcal{J}(\boldsymbol{\theta}) := \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\hat{\mathbf{x}}^{(i)}) - y^{(i)} \right)^2 \quad (2)$$

- We want to minimize  $\mathcal{J}(\boldsymbol{\theta})$  to obtain the optimal model parameters  $\boldsymbol{\theta}^*$ :

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\hat{\mathbf{x}}^{(i)}) - y^{(i)} \right)^2$$

- This is **ordinary least squares (OLS)**

# Error Function Intuition



Question: Why do we consider the square in the error function?



# Origin of Least Squares

**Carl Friedrich Gauss** (1777 – 1855) was a German mathematician, geodesist, and physicist who made significant contributions to many fields in mathematics and science. Gauss ranks among history's most influential mathematicians. Gauss published the second and third complete proofs of the **fundamental theorem of algebra** and made contributions to **number theory**. He was instrumental in the discovery of the dwarf planet Ceres. His work on the motion of planetoids disturbed by large planets led to the introduction of the Gaussian gravitational constant and the method of **least squares**, which is still used in all sciences to minimize measurement error.



## Section: Solutions to Regression

Closed-Form Solutions and Normal Equation  
Geometric Interpretation of Least Squares  
Gradient Descent  
Probabilistic Interpretation of linear Regression

# Closed-Form Solutions

- Usual approach (for two parameters): Calculate  $\theta_0$  and  $\theta_1$  according to (where  $\bar{x} = 1/n \sum_{i=1}^n x^{(i)}$  is the sample set mean)

$$\theta_0^* := \bar{y} - \theta_1^* \bar{x}, \quad \theta_1^* := \frac{\sum_{i=1}^n (x^{(i)} - \bar{x}) \cdot (y^{(i)} - \bar{y})}{\sum_{i=1}^n (x^{(i)} - \bar{x})^2} \quad (3)$$

- ‘Normal equation(s)’ (scales to arbitrary dimensions):

$$\theta^* := (\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top \mathbf{y} \quad (4)$$

$\hat{\mathbf{X}}$  is called ‘**design matrix**’ or ‘**regressor matrix**’

# Design Matrix / Regressor Matrix

- The design matrix  $\hat{\mathbf{X}} \in \mathbb{R}^{n \times (m+1)}$  looks as follows:

$$\hat{\mathbf{X}} := \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_m^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_m^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \cdots & x_m^{(n)} \end{pmatrix} \quad (5)$$

- The label vector  $\mathbf{y} \in \mathbb{R}^n$  is given by:

$$\mathbf{y} := (y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(n)})^\top \quad (6)$$

**Notation:** In the following we use the symbol  $\mathbf{X}$  to denote  $\hat{\mathbf{X}}$  defined in (5)



# Derivation of the Normal Equation

In the derivation of equation (4) we will need the following rules:

**Matrix transposition rules:**

$$(\mathbf{A}^\top)^\top = \mathbf{A} \quad (7)$$

$$(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top \quad (8)$$

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top \quad (9)$$

**Vector derivatives:**

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x} \quad (10)$$

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{a}^\top \mathbf{x} = \mathbf{a} \quad (11)$$

Equation (10) only holds true if  $\mathbf{A}$  is a symmetric matrix!



# Derivation of the Normal Equation (Ctd.)

Step ❶:

- We rewrite the least squares error function (2) in matrix/vector notation:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{2n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2 \stackrel{(\triangle)}{=} \frac{1}{2n} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \quad (12)$$

- **Remarks:**

- $\mathbf{X}\boldsymbol{\theta} \in \mathbb{R}^n$  is the vector containing the model predictions
- $\|\cdot\|$  denotes the **Euclidean norm**:  $\|\mathbf{x}\| := \sqrt{\sum_{j=1}^m x_j^2}$
- In step  $(\triangle)$  we have used  $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x}$



## Derivation of the Normal Equation (Ctd.)

Step ②: Use the matrix transposition rules to further rewrite the error function:

$$\begin{aligned} J(\theta) &= \frac{1}{2n} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) \stackrel{(8)}{=} \frac{1}{2n} ((\mathbf{X}\theta)^\top - \mathbf{y}^\top) (\mathbf{X}\theta - \mathbf{y}) \\ &= \frac{1}{2n} ((\mathbf{X}\theta)^\top \mathbf{X}\theta - (\mathbf{X}\theta)^\top \mathbf{y} - \mathbf{y}^\top (\mathbf{X}\theta) + \mathbf{y}^\top \mathbf{y}) \\ &\stackrel{(9)}{=} \frac{1}{2n} (\theta^\top \mathbf{X}^\top \mathbf{X}\theta - 2(\mathbf{X}\theta)^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}) \\ &\stackrel{(9)}{=} \frac{1}{2n} (\theta^\top \mathbf{X}^\top \mathbf{X}\theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}) \end{aligned} \tag{13}$$



## Derivation of the Normal Equation (Ctd.)

Step ③: We compute the derivative of (13) piece-wise:  
(This is allowed due to the sum rule)

$$\frac{\partial}{\partial \theta} \theta^T X^T X \theta \stackrel{(10)}{=} 2X^T X \theta \quad [X^T X \text{ is symmetric: } (X^T X)^T \stackrel{(9)}{=} X^T (X^T)^T \stackrel{(7)}{=} X^T X]$$

$$\frac{\partial}{\partial \theta} -2\theta^T X^T y \stackrel{(11)}{=} -2X^T y$$

$$\frac{\partial}{\partial \theta} y^T y = 0$$





## Derivation of the Normal Equation (Ctd.)

Step ④: Plug everything together, set the derivative to zero and solve:

$$\frac{\partial}{\partial \theta} \mathcal{J}(\theta) = \frac{1}{2n} (2\mathbf{X}^T \mathbf{X} \theta - 2\mathbf{X}^T \mathbf{y}) = \frac{1}{n} (\mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y}) \stackrel{!}{=} 0$$

$$\iff \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y} = 0$$

$$\iff \mathbf{X}^T \mathbf{X} \theta = \mathbf{X}^T \mathbf{y}$$

[If  $\mathbf{X}^T \mathbf{X}$  is invertible]

$$\iff \boxed{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}} = \theta$$



## Derivation of the Normal Equation (Ctd.)

- We have not yet shown that (4) is a minimum of the error function
- For this we consider the second-order derivative:

$$\frac{\partial^2}{(\partial \boldsymbol{\theta})^2} \mathcal{J}(\boldsymbol{\theta}) = \frac{1}{n} \mathbf{X}^\top \mathbf{X} \quad (14)$$

- We notice that  $\mathbf{X}^\top \mathbf{X} \succeq 0$  (positive semi-definite): Let  $\mathbf{z} \in \mathbb{R}^m \setminus \{\mathbf{0}\}$ , then

$$\mathbf{z}^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{z} \stackrel{(9)}{=} (\mathbf{X} \mathbf{z})^\top (\mathbf{X} \mathbf{z}) = \|\mathbf{X} \mathbf{z}\|^2 \geq 0 \quad (15)$$

- The cost function fulfills the second-order convexity condition

# Geometric Interpretation of Least Squares

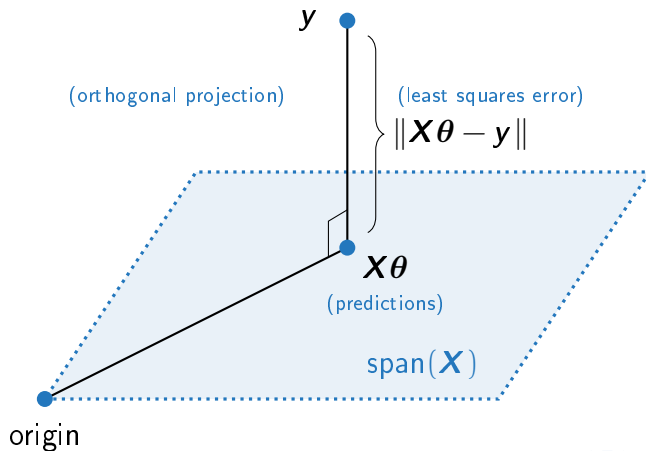
- We want to find a solution to the system of equations

$$\theta_0 \mathbf{x}_0 + \theta_1 \mathbf{x}_1 + \dots + \theta_m \mathbf{x}_m = \mathbf{y} \quad (16)$$

- $\mathbf{y}$  is a **linear combination** of the column vectors  $\mathbf{x}_j$  ( $0 \leq j \leq m$ ) of  $\mathbf{X}$
- However, we will **not find a solution** unless the data is perfectly linear

**Goal:** Find the point in the span of  $\mathbf{X}$  (space spanned by the column vectors of  $\mathbf{X}$  / column space) that is closest to  $\mathbf{y}$ !

# Geometric Interpretation of Least Squares (Ctd.)



## Geometric Interpretation of Least Squares (Ctd.)

- We see that the point closest to  $\mathbf{y}$  is given by the **orthogonal projection** of  $\mathbf{y}$  onto the span of  $\mathbf{X}$
- Any other point on the span of  $\mathbf{X}$  has a larger Euclidean distance to  $\mathbf{y}$  and therefore cannot be optimal
- The geometry of least squares also explains the term 'normal equation':
  - A normal line is a line perpendicular to another line or subspace
  - The normal equation given by (4) computes  $\boldsymbol{\theta}$  such that the residuals are orthogonal to the span of  $\mathbf{X}$



# Problems with Matrix Inversion?

What if  $(\mathbf{X}^\top \mathbf{X})^{-1}$  does not exist? Problems and solutions:

- 1 Linearly dependent (redundant) features or design matrix does not have full rank?  
(E. g. size in  $\text{m}^2$  and size in  $\text{feet}^2$ )  
⇒ **Delete correlated features**
- 2 Too many features ( $m > n$ )?  
⇒ **Delete features (e. g. using PCA) / add training examples**
- 3 Other numerical instabilities?  
⇒ **Add a regularization term** (later)
- 4 Computationally too expensive?  
⇒ **Use gradient descent**



# Gradient Descent

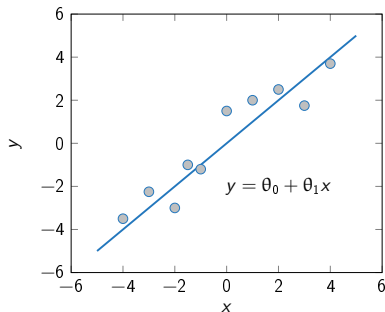
- We want to minimize a continuous function  $\mathcal{J} : \mathbb{R}^{m+1} \rightarrow \mathbb{R}$ :  $\min_{\theta \in \mathbb{R}^{m+1}} \mathcal{J}(\theta)$
- Update the parameters iteratively:

$$\theta_{t+1} \leftarrow \theta_t - \alpha \nabla_{\theta} \mathcal{J}(\theta_t) \quad (17)$$

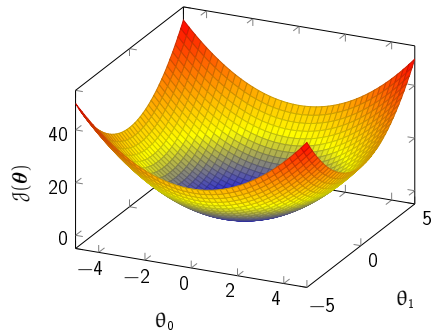
- **Legend:**
  - $\theta_t$  denotes the set of model parameters at timestamp  $t$
  - $\alpha \in (0, 1)$  denotes the **learning rate**
  - $\nabla_{\theta} \mathcal{J}(\theta) := \left( \frac{\partial}{\partial \theta_0} \mathcal{J}(\theta), \frac{\partial}{\partial \theta_1} \mathcal{J}(\theta), \dots, \frac{\partial}{\partial \theta_m} \mathcal{J}(\theta) \right)^{\top}$  is the gradient of  $\mathcal{J}(\theta)$

# Data Input Space vs. Hypothesis Space

## Data input space



## Hypothesis space $\mathcal{H}$





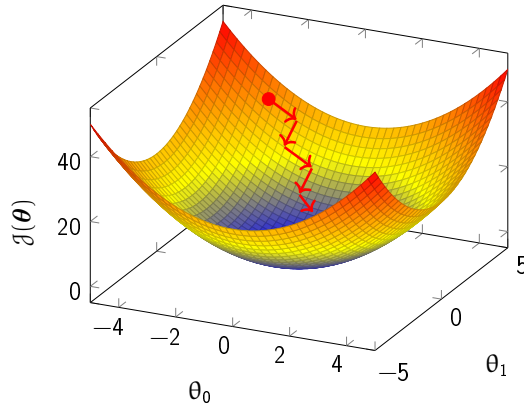
# Data Input Space vs. Hypothesis Space (Ctd.)

- **Data input space**
  - Determined by the  $m + 1$  **attributes** of the dataset  $x_0, x_1, x_2, \dots, x_m$
  - Often high-dimensional
- **Hypothesis space  $\mathcal{H}$** 
  - Determined by the **number of parameters** of the model
  - Each point in the hypothesis space corresponds to a **specific assignment of model parameters**
  - The error function gives information about how good this assignment is
  - **Gradient descent is applied in the hypothesis space  $\mathcal{H}$**



# Data Input Space vs. Hypothesis Space (Ctd.)

# Visualization of Gradient Descent in 3 Dimensions



# Versions of Gradient Descent

- Assume some training data  $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^n$
- The squared error of a **single** example is given by:

$$\ell(y_{pred}, y_{true}) := \frac{1}{2}(y_{pred} - y_{true})^2$$

- Our objective is to minimize the **total error**:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^{m+1}} \mathcal{J}(\boldsymbol{\theta}) = \min_{\boldsymbol{\theta} \in \mathbb{R}^{m+1}} \sum_{i=1}^n \ell(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)})$$

# Batch Gradient Descent

Three versions of gradient descent:

① Batch gradient descent

- Compute the gradient based on ALL  $n$  training data points

$$\theta_{t+1} \leftarrow \theta_t - \alpha \sum_{i=1}^n \nabla \ell(h_{\theta_t}(\mathbf{x}^{(i)}), y^{(i)}) \quad (18)$$

- Most accurate, but may be costly depending on the size of the dataset!

② Stochastic gradient descent

③ Mini-batch gradient descent

# Stochastic Gradient Descent

Three versions of gradient descent:

- 1 Batch gradient descent
- 2 **Stochastic gradient descent**
  - Compute gradient based on a **SINGLE** data point

$$\theta_{t+1} \leftarrow \theta_t - \alpha \nabla \ell(h_{\theta_t}(\mathbf{x}^{(i)}), y^{(i)}) \quad (19)$$

- **Pick training examples randomly and not sequentially!**
  - Efficient, but inaccurate!
- 3 Mini-batch gradient descent

# Mini-Batch Gradient Descent

## Three versions of gradient descent:

- ① Batch gradient descent
- ② Stochastic gradient descent
- ③ **Mini-batch gradient descent**
  - Compute the gradient based on a mini-batch comprising  $n_b$  examples

$$\theta_{t+1} \leftarrow \theta_t - \alpha \sum_{i=1}^{n_b} \nabla \ell(h_{\theta_t}(\mathbf{x}^{(i)}), y^{(i)}) \quad (20)$$

- Good trade-off between batch and stochastic gradient descent

# Stochastic Gradient of the Least Squares Error Function

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \mathcal{J}(\boldsymbol{\theta}) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2 \\&= 2 \cdot \frac{1}{2} (h_{\boldsymbol{\theta}}(\mathbf{x}) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\boldsymbol{\theta}}(\mathbf{x}) - y) \\&= (h_{\boldsymbol{\theta}}(\mathbf{x}) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \cdots + \theta_m x_m - y) \\&= \boxed{(h_{\boldsymbol{\theta}}(\mathbf{x}) - y) x_j}\end{aligned}\tag{21}$$

The vectorized version is given by  $\frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) = (h_{\boldsymbol{\theta}}(\mathbf{x}) - y) \mathbf{x}$





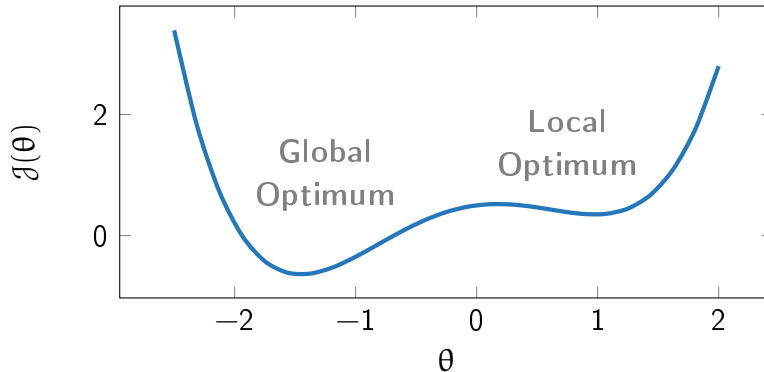
# Batch Gradient of the Least Squares Error Function

- We have computed the batch gradient already when deriving the normal equation (we ignore the factor  $1/n$  here):

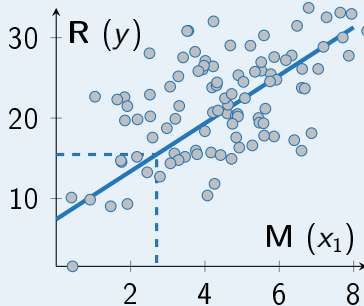
$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) &= \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{X}^\top (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})\end{aligned}\tag{22}$$

- We can use the same formula for the mini-batch gradient by replacing  $\mathbf{X}$  with  $\mathbf{X}_b$  – the design matrix comprising only the mini-batch examples, and  $\mathbf{y}$  with  $\mathbf{y}_b$  – the respective labels

# Disadvantage of Gradient Descent



## Solving the introductory Example



- $\theta_0 \approx 7.4218$
- $\theta_1 \approx 2.9827$
- $\mathcal{J}(\boldsymbol{\theta}) \approx 446.9584$
- $h_{\boldsymbol{\theta}}(\mathbf{x}) = 7.4218 + 2.9827 \cdot x_1$
- $R = h_{\boldsymbol{\theta}}(2.7) = \underline{\underline{15.4750}}$

## A probabilistic View

- **Assumption 1:** The target values  $y$  are generated by an unknown function  $f(\mathbf{x})$  with additive noise  $\varepsilon$ :

$$y = f(\mathbf{x}) + \varepsilon \quad (23)$$

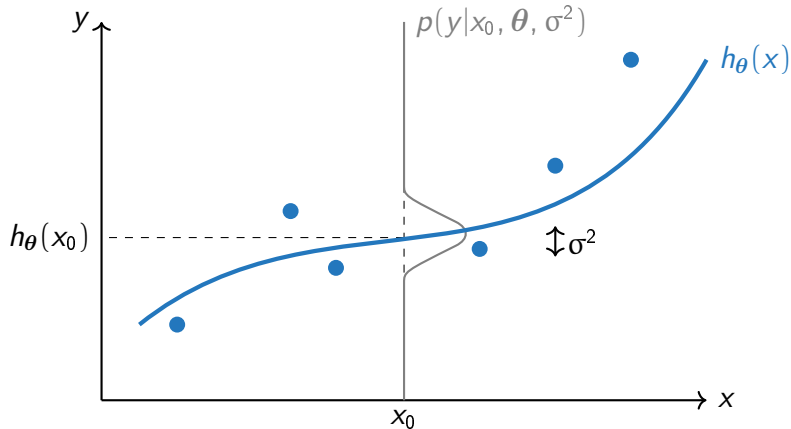
- **Assumption 2:** The noise  $\varepsilon$  is a zero mean Gaussian random variable

$$\varepsilon \sim \mathcal{N}(0, \sigma^2) \quad (24)$$

- We consider  $y$  a random variable which is distributed according to

$$p(y|\mathbf{x}, \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(y|h_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2) \quad (25)$$

## A probabilistic View (Ctd.)



# Likelihood Function for Regression

- We are given a dataset  $\mathcal{D} := \left\{ (\mathbf{x}^{(i)}, y) \right\}_{i=1}^n$
- The (conditional) likelihood is given by:

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma^2) &= \prod_{i=1}^n \mathcal{N}(y^{(i)} | h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), \sigma^2) \\ &= \prod_{i=1}^n \mathcal{N}(y^{(i)} | \boldsymbol{\theta}^\top \mathbf{x}^{(i)}, \sigma^2) \end{aligned} \quad (26)$$



## Likelihood Function for Regression (Ctd.)

- The log-likelihood is given by (we have computed this earlier already):

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{\sigma^2} \cdot \underbrace{\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^\top \mathbf{x}^{(i)})^2}_{\text{least squares error}} \quad (27)$$

- We have to minimize the least squares error to maximize the likelihood!

**When minimizing the squared error we implicitly assume Gaussian noise!**



# The Maximum Likelihood Solution to Regression

- Optimization of the log-likelihood function gives:

$$\boldsymbol{\theta}_{\text{ML}} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (28)$$

$$\sigma_{\text{ML}}^2 := \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}_{\text{ML}}^T \mathbf{x}^{(i)})^2 \quad (29)$$

- The probabilistic interpretation of linear regression allows us to **quantify the (global) uncertainty** of the model

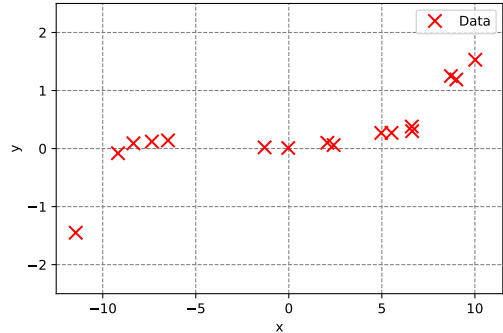


## Section: Basis Function Regression

General Idea  
Polynomial Basis Functions  
Radial Basis Functions

# What if the Data is non-linear?

- So far we have fitted straight lines
- **What if the data is not linear...?**





# Basis Functions

- Remember: ‘When stuck switch to a different perspective’
- We add **higher-order** features by using **basis functions**  $\varphi$ :

$$h_{\theta}(\mathbf{x}) := \sum_{j=0}^p \theta_j \varphi_j(\mathbf{x}) \quad (30)$$

- Commonly used basis functions:
  - **Linear:**  $\varphi_0(\mathbf{x}) = 1$  and  $\varphi_1(\mathbf{x}) = x_1, \dots, \varphi_m(\mathbf{x}) = x_m$
  - **Polynomial** (see below)
  - **Radial basis functions** (see below)

# New Design Matrix

By applying the basis functions to  $\mathbf{X}$  we get a new design matrix  $\Phi$ :

$$\Phi := \begin{pmatrix} \varphi_0(\mathbf{x}^{(1)}) & \varphi_1(\mathbf{x}^{(1)}) & \varphi_2(\mathbf{x}^{(1)}) & \dots & \varphi_p(\mathbf{x}^{(1)}) \\ \varphi_0(\mathbf{x}^{(2)}) & \varphi_1(\mathbf{x}^{(2)}) & \varphi_2(\mathbf{x}^{(2)}) & \dots & \varphi_p(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \varphi_0(\mathbf{x}^{(n)}) & \varphi_1(\mathbf{x}^{(n)}) & \varphi_2(\mathbf{x}^{(n)}) & \dots & \varphi_p(\mathbf{x}^{(n)}) \end{pmatrix} \quad (31)$$

The model is still linear in the parameters, so we can still use the same algorithm as before. **This is still linear regression (!!!)**



# Basis Functions: Polynomial Basis Functions

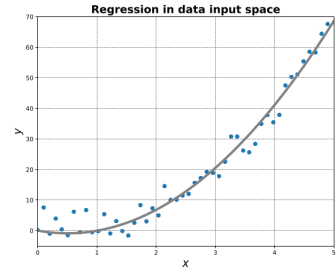
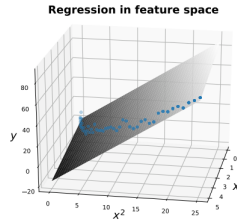
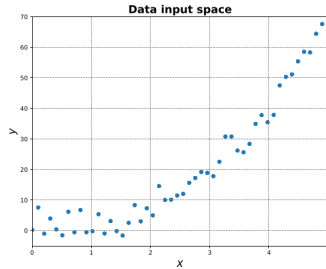
- Let us assume a one-dimensional dataset ( $m = 1$ ) for now
- A quite frequently used basis function is the **polynomial basis**

$$\varphi_0(x) := 1, \varphi_j(x) := x^j$$

$$h_{\theta}(x) := \sum_{j=0}^p \theta_j \varphi_j(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \cdots + \theta_p x^p$$

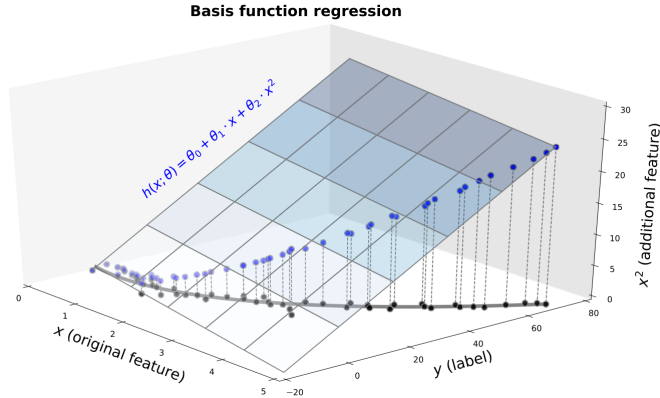
- $p$  is the degree of the polynomial
- For higher dimensional datasets we can also include cross-terms, e. g.  
 $\varphi_j(x) = x_l^2 x_k$

# It is still linear!





# It is still linear! (Ctd.)



# Basis Functions: Radial Basis Functions (RBFs)

- Another possible choice of basis function: **Radial basis functions**

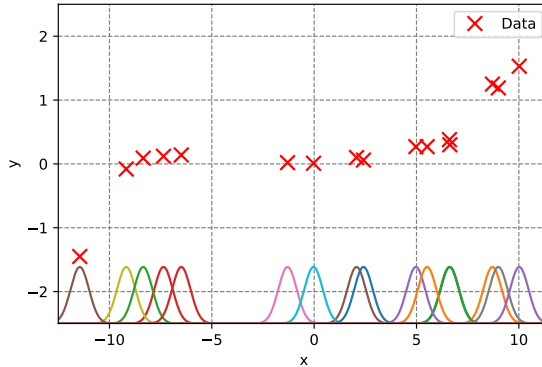
$$\varphi_0(x) := 1 \quad (32)$$

$$\varphi_j(x) := \exp(-1/2 \|x - z_j\|^2 / 2\sigma^2) \quad (33)$$

- $\{z_j\}$  are the centers of the radial basis functions,  $\sigma^2$  is the scale
- $p$  denotes the number of centers / number of radial basis functions
- Often we take each data point as a center, so  $p = n$   
(but in general we are free to put the centers wherever we want)



# Radial Basis Functions (Ctd.)



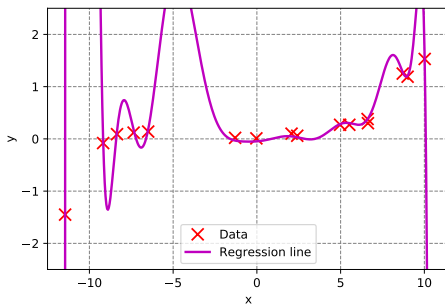
## Section: Regularization Techniques

Underfitting and Overfitting  
L1 and L2 Regularization

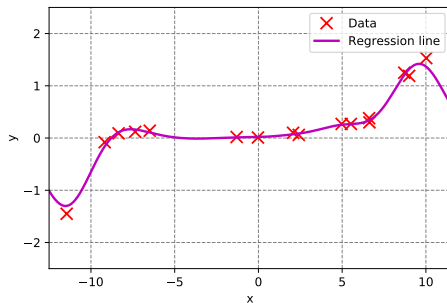


# The Danger of too expressive Models...

Polynomial of degree  $p = 16$   
(💀 **severe overfitting** 💀)



RBF with  $\sigma^2 = 1.00$ ,  $p = n$   
(about right)

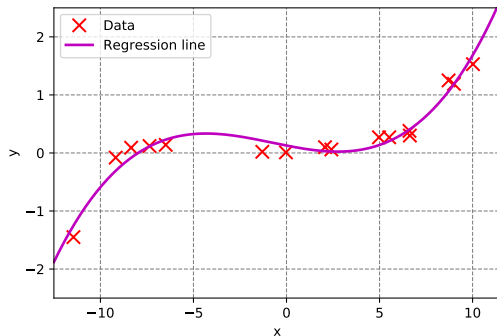


# Overfitting vs. Underfitting

- **Underfitting**
  - The model is not complex enough to fit the data well  $\Rightarrow$  **High bias**
  - Make the model more complex; adding new examples **does not help**
- **Overfitting**
  - The model predicts the training data perfectly
  - But it **fails to generalize** to unseen instances  $\Rightarrow$  **High variance**
  - Decrease the degree of freedom or add more training examples
  - Also: Try **regularization**
- **Bias-Variance trade-off**

# First Solution: Smaller Degree

One solution: Use a **smaller degree** (here:  $p = 3$ )





## Second Solution: Regularization

- Enrich the cost function  $\mathcal{J}(\boldsymbol{\theta})$  with a **regularization term**
- This helps **prevent overfitting** and results in a smoother function

### L1-Regularization:

$$\tilde{\mathcal{J}}(\boldsymbol{\theta}) := \mathcal{J}(\boldsymbol{\theta}) + \lambda|\boldsymbol{\theta}|$$

$$|\boldsymbol{\theta}| = \sum_{j=1}^m |\theta_j|$$

### L2-Regularization:

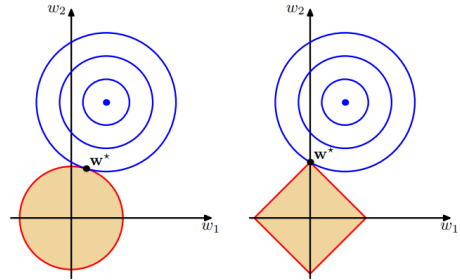
$$\tilde{\mathcal{J}}(\boldsymbol{\theta}) := \mathcal{J}(\boldsymbol{\theta}) + \lambda\|\boldsymbol{\theta}\|^2$$

$$\|\boldsymbol{\theta}\|^2 = \sum_{j=1}^m \theta_j^2$$

( $\lambda \geq 0$  controls the **degree of regularization**)

# Regularization visualized

- Here:  $\mathbf{w} \equiv \boldsymbol{\theta}$
- L1-Regularization  
⇒ **Lasso regression**  
(least abs. shrinkage and select. operator)
- L2-Regularization  
⇒ **Ridge regression**  
(Tikhonov regularization)
- The combination of both is called  
**elastic net**



cf. Bishop.2006, page 146, left: L2, right: L1

# Incorporating Regularization

- Normal equation with regularization: **Ridge regression**

$$\boldsymbol{\theta}^* := (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^\top \mathbf{y} \quad (34)$$

- Regularization also helps overcome numerical issues (matrix inversion!)
- Regularized least squares error gradient:

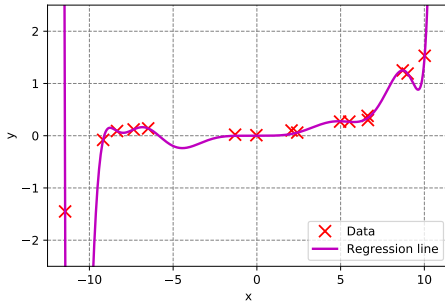
$$\frac{\partial}{\partial \theta_j} \mathcal{J}(\boldsymbol{\theta}) = (h_{\boldsymbol{\theta}}(\boldsymbol{\varphi}(\mathbf{x})) - y) \varphi_j(\mathbf{x}) + \lambda \theta_j \quad (35)$$

where  $\boldsymbol{\varphi}(\mathbf{x}) := (\varphi_0(\mathbf{x}), \dots, \varphi_m(\mathbf{x}))^\top$

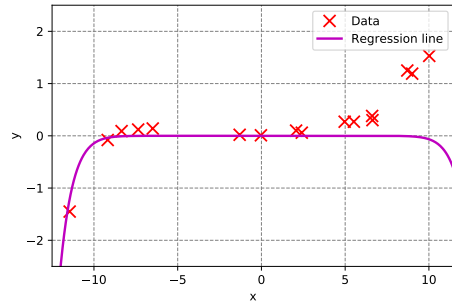


# Polynomial Regression with Regularization

At least better



Way too much regularization



## Section: Wrap-Up

Summary  
Self-Test Questions  
Lecture Outlook

# Summary

- In regression we predict **continuous target variables**
- The algorithm minimizes the **(mean) squared error**
- **Two approaches:**
  - ① Normal equation
  - ② (Batch / stochastic / mini-batch) gradient descent
- Geometric interpretation: Predictions are the orthogonal projection of the label vector onto the span of the design matrix
- Use **basis functions** to fit non-linear regression lines
- **Regularization** is important (especially when using basis functions)



# Self-Test Questions

- 1 What is the goal of regression?
- 2 What can you do if matrix inversion fails for the normal equation?
- 3 What is a suitable cost function for regression? Where does it come from?
- 4 Does gradient descent give the exact solution?
- 5 Which versions of gradient descent do you know?
- 6 What are basis functions? Why use them? State some examples.
- 7 What is overfitting and underfitting?
- 8 What is regularization? Why should you apply it?

# What's next...?

Unit I	Machine Learning Introduction
Unit II	Mathematical Foundations
Unit III	Bayesian Decision Theory
Unit IV	Regression
<b>Unit V</b>	<b>Classification I</b>
Unit VI	Evaluation
Unit VII	Classification II
Unit VIII	Clustering
Unit IX	Dimensionality Reduction

Thank you very much for the attention!

**Topic:** \*\*\* Applied Machine Learning Fundamentals \*\*\* Regression

**Term:** Winter term 2023/2024

**Contact:**

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

[daniel.wehner@sap.com](mailto:daniel.wehner@sap.com)

Do you have any questions?