

Optimization Techniques

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Summer term 2025



Find all slides on [GitHub](#) (DaWe1992/Applied_ML_Fundamentals)

Lecture Overview

- **I** Machine Learning Introduction
- II** Optimization Techniques
- III** Bayesian Decision Theory
- IV** Non-parametric Density Estimation
- V** Probabilistic Graphical Models
- VI** Linear Regression
- VII** Logistic Regression
- VIII** Deep Learning
- IX** Evaluation
- X** Decision Trees
- XI** Support Vector Machines
- XII** Clustering
- XIII** Principal Component Analysis
- XIV** Reinforcement Learning
- XV** Advanced Regression

Agenda for this Unit

① Introduction and Motivation

② Important Concepts and Definitions

③ Unconstrained and constrained
Optimization

④ Numerical Optimization

⑤ Wrap-Up

Section:
Introduction and Motivation

What Learning is all about
A motivational Example

Learning equals Optimization

- In every machine learning problem, we have:
 - ① An **objective function** we want to optimize
 - ② **Data** we want to learn from
 - ③ **Parameters** which need to be learned
 - ④ Assumptions about the problem and the data
- We would like to have general solutions to the problem of learning
- Different algorithms embody different objective functions and assumptions

Every machine learning problem is an optimization problem!

Example: Linear Regression

- Suppose we want to fit an **affine-linear** function to a dataset with N data points:

$$h_{\theta}(x) = \theta_1 x + \theta_0$$

- To do so, we have to find parameters $\theta_0 \in \mathbb{R}$ (y -intercept) and $\theta_1 \in \mathbb{R}$ (slope) such that the function is a good fit
- Question: How to measure the goodness of fit?**
- Answer:** Use a **cost function** (objective function) which encodes the desired behavior, i. e. prefer functions which are close to the data points

Example: Linear Regression (Ctd.)

- In linear regression we usually use the **squared error** cost function

$$\mathfrak{J}(\theta_0, \theta_1) := \sum_{n=1}^N (\theta_1 x_n + \theta_0 - y_n)^2 \quad (1)$$

- Choose θ_0, θ_1 so as to minimize the cost function

$$\theta_0^*, \theta_1^* = \arg \min_{\theta_0, \theta_1} \mathfrak{J}(\theta_0, \theta_1)$$

It is the task of the optimization algorithm to find the best solution!

Section: **Important Concepts and Definitions**

Partial Derivatives and Gradients
HESSE Matrix
TAYLOR Expansion
Convex Functions and convex Sets

Partial Derivatives of Functions $\mathbb{R}^D \rightarrow \mathbb{R}$

Definition: Let $S \subset \mathbb{R}^D$ be an open set and $f : S \rightarrow \mathbb{R}$. We say that f is **partially differentiable** in $\mathbf{x}_0 \in S$ with respect to the i -th input variable x_i , if the limit

$$\lim_{h \rightarrow 0} \frac{f(\mathbf{x}_0 + h\mathbf{e}^i) - f(\mathbf{x}_0)}{h} \quad (2)$$

exists.

- \mathbf{e}^i is the i -th canonical unit vector
- The limit (2) – if existent – is denoted by $\frac{\partial}{\partial x_i} f(\mathbf{x}_0)$

Gradients ∇

- The **gradient** of the function $f : S \rightarrow \mathbb{R}$ is obtained by stacking all D partial derivatives of f into a column vector:

$$\nabla f(\mathbf{x}_0) := \begin{pmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x}_0) \\ \vdots \\ \frac{\partial}{\partial x_D} f(\mathbf{x}_0) \end{pmatrix} \quad (3)$$

- We will often use **transposition** to save space: $\left(\frac{\partial}{\partial x_1} f(\mathbf{x}_0), \dots, \frac{\partial}{\partial x_D} f(\mathbf{x}_0) \right)^\top$
- The symbol ∇ is pronounced ‘*nabla*’ or ‘*gradient of*’ (*it is **not** a capital delta Δ !*)

Computation of the Gradient

How can we compute the gradient?

- Luckily, we can apply the same rules of differentiation we know from one-dimensional calculus:
 - **Chain rule,**
 - **product rule,**
 - **Quotient rule,**
 - etc.
- Differentiate with respect to one variable and hold the others fixed

Example: Computation of the Gradient

- Let the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}, (x_1, x_2) \mapsto \log(x_1 x_2) + \sin(x_2)$ be given
- Partial derivatives:

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \frac{1}{x_1} \quad \text{and} \quad \frac{\partial}{\partial x_2} f(x_1, x_2) = \frac{1}{x_2} + \cos(x_2)$$

- Thus, the gradient is:

$$\nabla f(x_1, x_2) = \begin{pmatrix} \frac{1}{x_1} \\ \frac{1}{x_2} + \cos(x_2) \end{pmatrix}$$

Direction of steepest Ascent

Lemma: The gradient ∇f points into the direction of steepest **ascent**.

Proof: Let $\mathbf{v} \in \mathbb{R}^D$ be a unit vector, i.e. $\|\mathbf{v}\| = 1$. The directional derivative of f in the direction of \mathbf{v} is given by the scalar product $\mathbf{v}^\top \nabla f$. We aim to find the vector \mathbf{v} which maximizes the directional derivative. This vector will then point into the direction of steepest ascent. We can rewrite the scalar product using its geometrical definition and receive $\mathbf{v}^\top \nabla f = \|\mathbf{v}\| \cdot \|\nabla f\| \cdot \cos(\alpha)$.

Since \mathbf{v} has length 1, the equation simplifies to $\mathbf{v}^\top \nabla f = \|\nabla f\| \cdot \cos(\alpha)$. Only the angle α can be altered to maximize this quantity (since the gradient is given). The cosine reaches a maximum value for $\alpha = 0$. This means ∇f is parallel to \mathbf{v} , i.e. ∇f points into the direction of steepest ascent. ■

Direction of steepest Descent

The following observation will be especially useful when minimizing functions!

Lemma: The negative gradient $-\nabla f$ points into the direction of steepest **descent**.

Proof: Analogous to the proof above. ■

HESSE Matrix ∇^2

The HESSE matrix $\nabla^2 f$ (also referred to as **Hessian**) contains all second-order partial derivatives of f :

$$\nabla^2 f(\mathbf{x}_0) := \begin{pmatrix} \frac{\partial^2}{\partial x_1 \partial x_1} f(\mathbf{x}_0) & \frac{\partial^2}{\partial x_1 \partial x_2} f(\mathbf{x}_0) & \dots & \frac{\partial^2}{\partial x_1 \partial x_D} f(\mathbf{x}_0) \\ \frac{\partial^2}{\partial x_2 \partial x_1} f(\mathbf{x}_0) & \frac{\partial^2}{\partial x_2 \partial x_2} f(\mathbf{x}_0) & \dots & \frac{\partial^2}{\partial x_2 \partial x_D} f(\mathbf{x}_0) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_D \partial x_1} f(\mathbf{x}_0) & \frac{\partial^2}{\partial x_D \partial x_2} f(\mathbf{x}_0) & \dots & \frac{\partial^2}{\partial x_D \partial x_D} f(\mathbf{x}_0) \end{pmatrix} \in \mathbb{R}^{D \times D} \quad (4)$$

SCHWARZ's Theorem

SCHWARZ's theorem: If f has continuous second-order derivatives, then

$$\frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} f(\mathbf{x}_0) = \frac{\partial^2}{\partial \mathbf{x}_j \partial \mathbf{x}_i} f(\mathbf{x}_0) \quad \forall i, j = 1, 2, \dots, D$$

- In other words: If f has continuous second-order derivatives, then the order in which the derivatives are taken is not important
- **Corollary:** If f has continuous second-order derivatives, then (4) is **symmetric**

Example: Computation of the HESSE Matrix

- Consider again the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, (x_1, x_2) \mapsto \log(x_1 x_2) + \sin(x_2)$$

- The partial **first-order** derivatives are given by:

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \frac{1}{x_1} \quad \text{and} \quad \frac{\partial}{\partial x_2} f(x_1, x_2) = \frac{1}{x_2} + \cos(x_2)$$

- We compute the partial **second-order** derivatives on the next slide by taking the derivatives of the first-order derivatives by x_1 and x_2

Example: Computation of the HESSE Matrix (Ctd.)

$$\frac{\partial^2}{\partial x_1 \partial x_1} f(x_1, x_2) = -\frac{1}{x_1^2}$$

$$\frac{\partial^2}{\partial x_1 \partial x_2} f(x_1, x_2) = 0$$

$$\frac{\partial^2}{\partial x_2 \partial x_1} f(x_1, x_2) = 0$$

$$\frac{\partial^2}{\partial x_2 \partial x_2} f(x_1, x_2) = -\frac{1}{x_2^2} - \sin(x_2)$$

Example: Computation of the HESSE Matrix (Ctd.)

- Therefore, we have:

$$\nabla^2 f(x_1, x_2) = \begin{pmatrix} -\frac{1}{x_1^2} & 0 \\ 0 & -\frac{1}{x_2^2} - \sin(x_2) \end{pmatrix}$$

- Please note that SCHWARZ's theorem holds true in this case as f has continuous second-order derivatives, i. e. $\nabla^2 f(x_1, x_2)$ is a symmetric matrix

TAYLOR Expansion for Functions $\mathbb{R} \rightarrow \mathbb{R}$

Lemma: Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be n -times differentiable at $x_0 \in \mathbb{R}$. We can approximate f using a **TAYLOR polynomial** of degree n given by

[Expansion in x_0]

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n \quad (5)$$

[Expansion in x]

$$f(x + \delta_x) \approx f(x) + f'(x)\delta_x + \frac{f''(x)}{2!}\delta_x^2 + \dots + \frac{f^{(n)}(x)}{n!}\delta_x^n \quad (6)$$

Example: TAYLOR Polynomial for $\sin(x)$

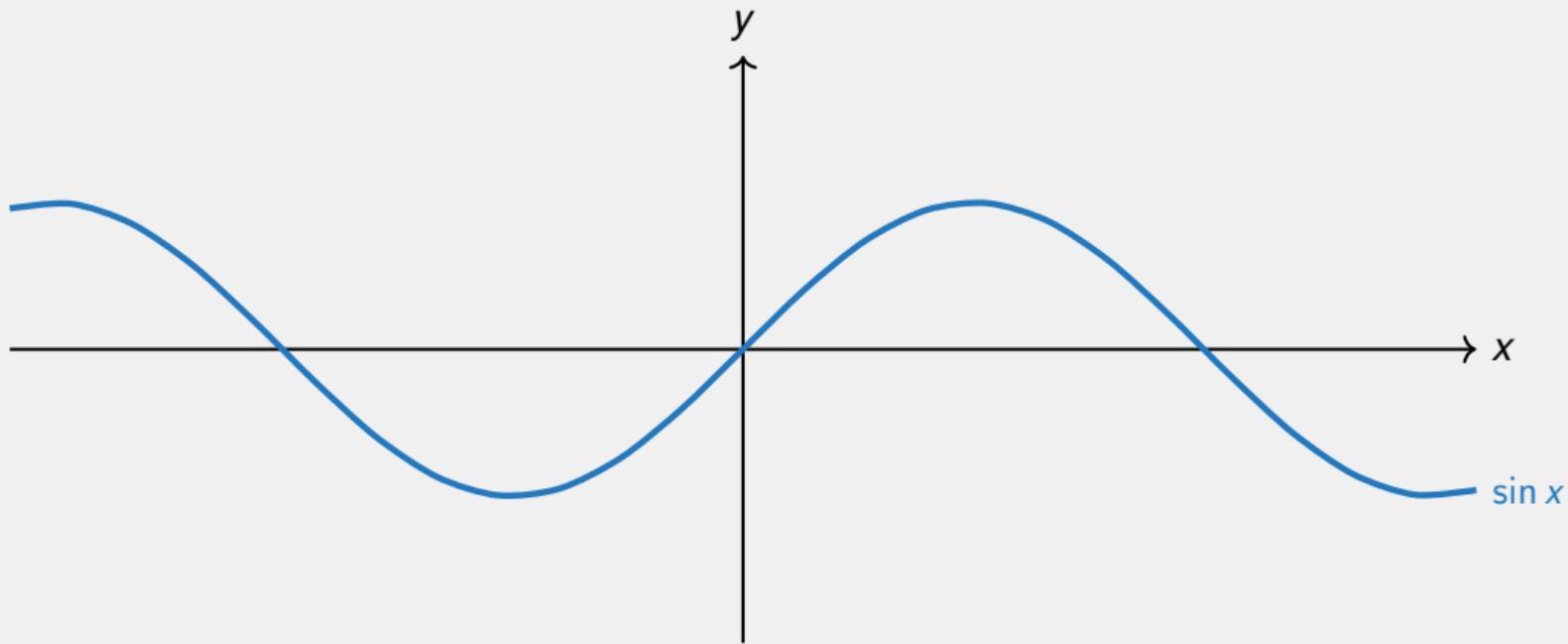
- Let us compute a degree 5 TAYLOR expansion of $\sin(x)$ at $x_0 = 0$:

$$\begin{aligned}\sin(x) &\approx \sin(0) + x \cos(0) - \frac{x^2}{2} \sin(0) - \frac{x^3}{6} \cos(0) + \frac{x^4}{24} \sin(0) + \frac{x^5}{120} \cos(0) \\ &= x - \frac{x^3}{6} + \frac{x^5}{120}\end{aligned}$$

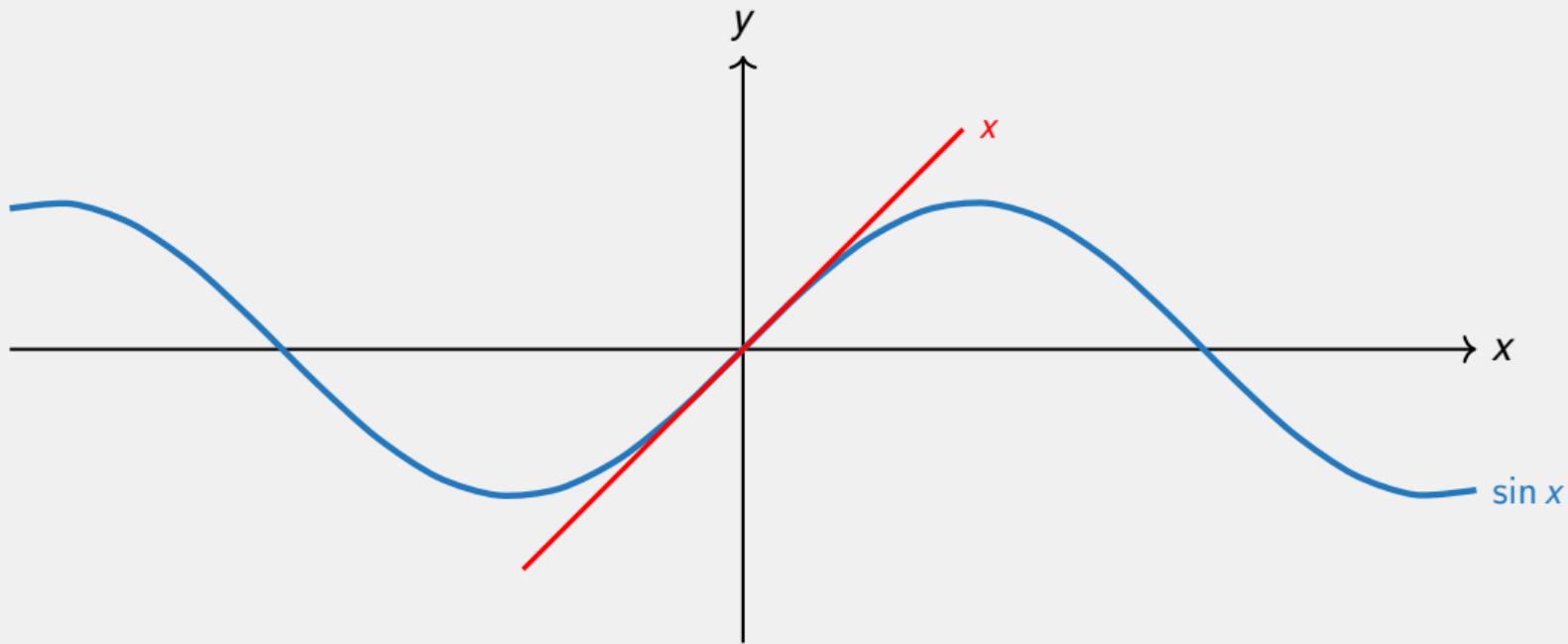
- We could continue the expansion forever because $\sin(x)$ is infinitely differentiable:

$$\sin(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}$$

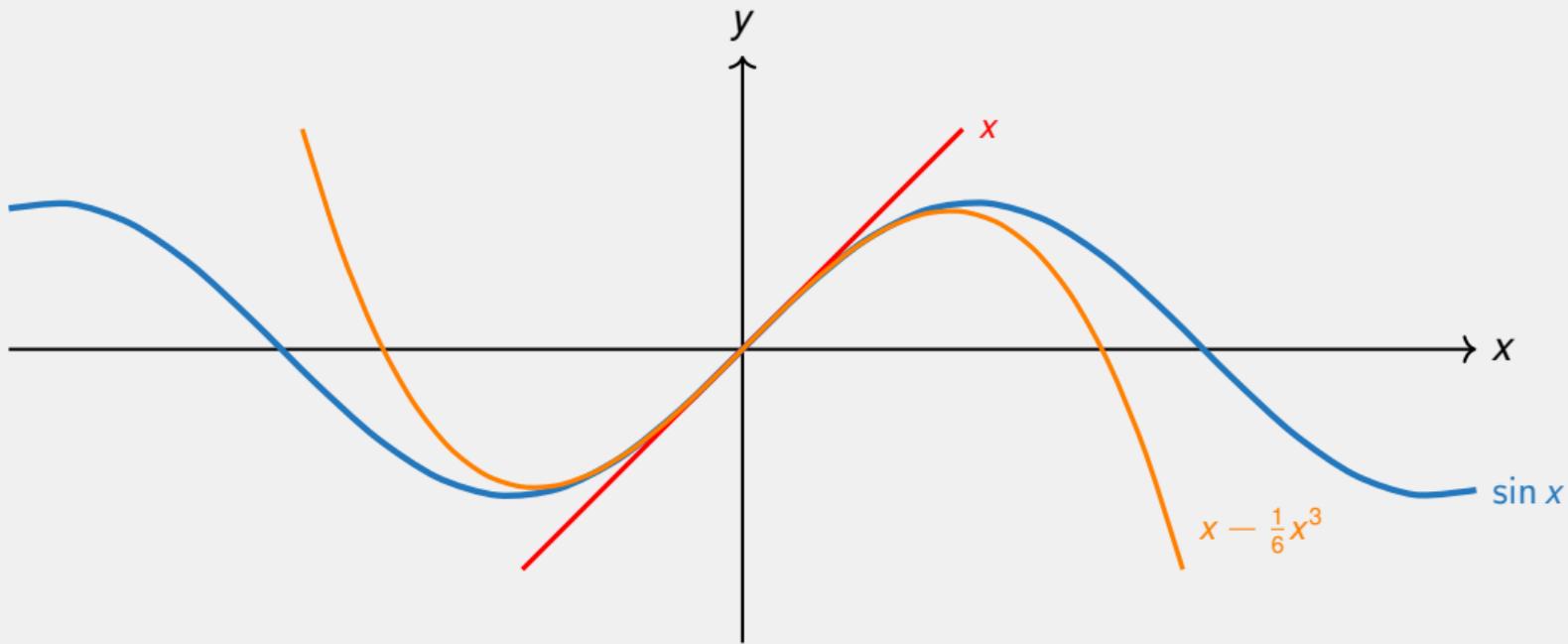
Example: TAYLOR Polynomial for $\sin(x)$ (Ctd.)



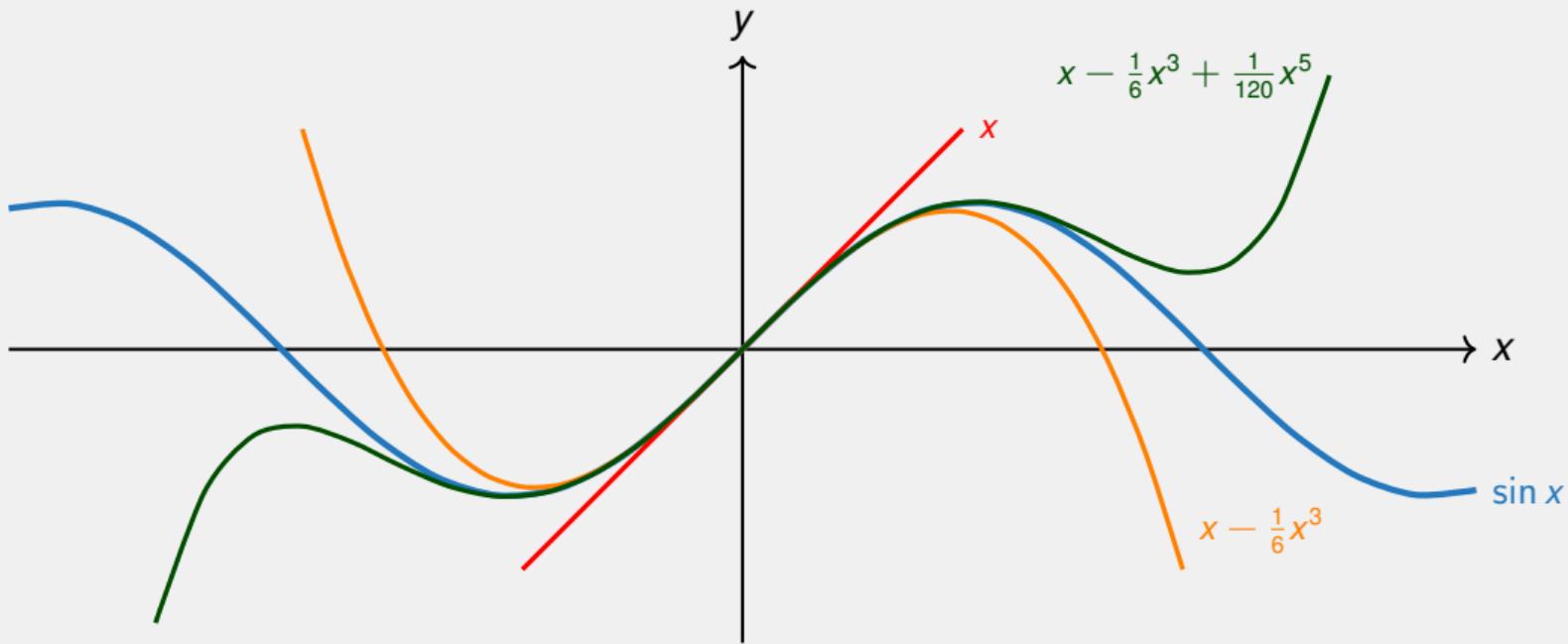
Example: TAYLOR Polynomial for $\sin(x)$ (Ctd.)



Example: TAYLOR Polynomial for $\sin(x)$ (Ctd.)



Example: TAYLOR Polynomial for $\sin(x)$ (Ctd.)



TAYLOR Expansion for Functions $\mathbb{R}^D \rightarrow \mathbb{R}$

We generalize the concept of TAYLOR expansions to multivariate functions:

Let $f : \mathbb{R}^D \rightarrow \mathbb{R}$ be 2-times differentiable in $\mathbf{x} \in \mathbb{R}^D$. We can approximate f using a **TAYLOR polynomial** of degree 2 given by

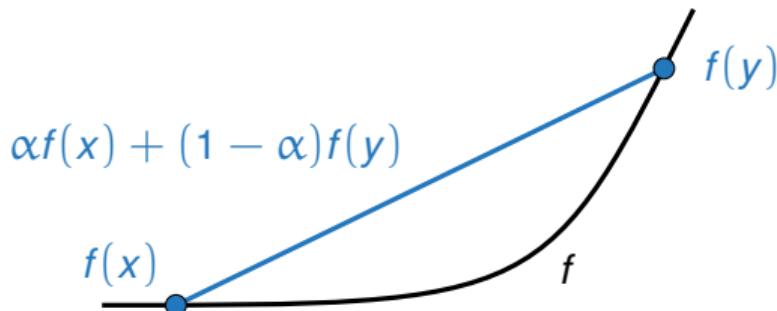
$$f(\mathbf{x} + \boldsymbol{\delta}_{\mathbf{x}}) \approx f(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top \boldsymbol{\delta}_{\mathbf{x}} + \frac{1}{2} \boldsymbol{\delta}_{\mathbf{x}}^\top \mathbf{H}(\mathbf{x}) \boldsymbol{\delta}_{\mathbf{x}} \quad (7)$$

- \mathbf{g} is the gradient of f , and \mathbf{H} the HESSE matrix of f
- A TAYLOR polynomial of degree 2 is sufficient in our case. Of course we could consider polynomials of degree $n > 2$ as well

Convex Functions

Definition: A function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is said to be **convex**, if $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and $\forall \alpha \in [0, 1]$:

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) \quad (8)$$



Strictly convex Functions

- We call a function **strictly convex**, if $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and $\forall \alpha \in (0, 1)$:

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) < \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y})$$

- Remark:** A function which is strictly convex is also convex (*but not vice versa!*)
- Check the definiteness of the HESSE matrix to test a function f on convexity:
 - If $\nabla^2 f(\mathbf{x}) \succeq 0$ (positive semi-definite), then f is convex
 - If $\nabla^2 f(\mathbf{x}) \succ 0$ (positive definite), then f is strictly convex

Examples: Convex Functions

Examples:

- **Convex, but not strictly convex:** (Affine-)linear functions ($\mathbf{a} \in \mathbb{R}^D, b \in \mathbb{R}$):

$$f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + b$$

- **Strictly convex:** Quadratic functions ($\mathbf{A} \in \mathbb{R}^{D \times D}, \mathbf{b} \in \mathbb{R}^D, c \in \mathbb{R}$):

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c$$

- **Strictly convex:** Exponential function

Convex and Non-convex Functions

Convex function



Non-convex function



Advantages of Convexity

Question: Why are convex cost functions so appealing?

Advantages of convex functions:

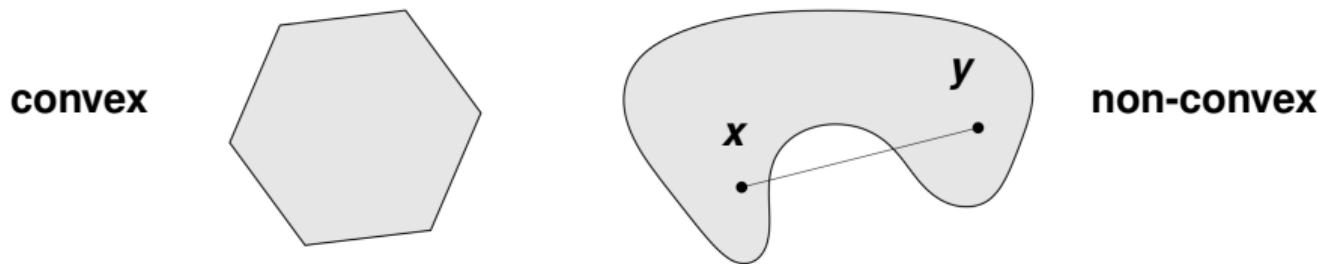
- Local solutions are global optima
- Numeric optimizers do not get stuck in local optima
- Efficient implementations of optimizers are available

Convex Sets

Definition: A set $C \subset \mathbb{R}^D$ is called **convex**, if $\forall \mathbf{x}, \mathbf{y} \in C$ and $\forall \alpha \in [0, 1]$:

$$\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in C \quad (9)$$

(this is the equation of the line segment between \mathbf{x} and \mathbf{y})



Section: **Unconstrained and constrained Optimization**

- Unconstrained Optimization
- Constrained Optimization and KKT conditions
- Optimization Problems with Equality Constraints
- Optimization Problems with Inequality Constraints
- Duality (Primal and Dual Optimization Problem)

Optimization of Functions $\mathbb{R} \rightarrow \mathbb{R}$

- Let a function $f : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto f(x)$ be given
- How do we find minimums and maximums of f ?**
- Necessary condition:** A point $x_0 \in \mathbb{R}$ is a potential candidate (**stationary point**), if $f'(x_0) = 0$
- We then check the sign of the second-order derivative

$$f''(x_0) > 0 \implies \text{minimum}$$

$$f''(x_0) < 0 \implies \text{maximum}$$

Optimization of Functions $\mathbb{R}^D \rightarrow \mathbb{R}$

- Let a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$ be given
- How do we find minimums and maximums of f ?**
- Necessary condition:** A point $\mathbf{x}_0 \in \mathbb{R}^D$ is a stationary point, if

$$\nabla f(\mathbf{x}_0) = \mathbf{0} \tag{10}$$

- We then check if the HESSE matrix $\nabla^2 f(\mathbf{x}_0)$ is
 - positive definite** (*minimum*) or
 - negative definite** (*maximum*)

How to check the Definiteness of Matrices?

For simplicity we consider (2×2) -matrices:

Lemma: Let S be a non-empty subset of \mathbb{R}^2 and $f : S \rightarrow \mathbb{R}$ be a 2-times differentiable function with continuous derivatives.

We define (determinant of the HESSE matrix):

$$\Delta(\mathbf{x}_0) := \frac{\partial f(\mathbf{x}_0)}{\partial x_1 \partial x_1} \cdot \frac{\partial f(\mathbf{x}_0)}{\partial x_2 \partial x_2} - \left(\frac{\partial f(\mathbf{x}_0)}{\partial x_1 \partial x_2} \right)^2 \quad (11)$$

Let $\mathbf{x}_0 \in S$ be a candidate solution, i. e. $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

How to check the Definiteness of Matrices? (Ctd.)

- ① If $\Delta(\mathbf{x}_0) > 0$, then \mathbf{x}_0 is a
$$\begin{cases} \textbf{local maximum} & \text{if } \frac{\partial f(\mathbf{x}_0)}{\partial x_1 \partial x_1} < 0. \\ \textbf{local minimum} & \text{if } \frac{\partial f(\mathbf{x}_0)}{\partial x_1 \partial x_1} > 0. \end{cases}$$
- ② If $\Delta(\mathbf{x}_0) < 0$, then \mathbf{x}_0 is a **saddle point** (German: *Sattelpunkt*) of f .
- ③ If $\Delta(\mathbf{x}_0) = 0$, we cannot decide.

Example: Finding Minimums of Functions $\mathbb{R}^2 \rightarrow \mathbb{R}$

- Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}, (x_1, x_2) \mapsto x_1^2 + x_2^2$
- The gradient is given by $\nabla f(x_1, x_2) = (2x_1, 2x_2)^\top$ and therefore $\mathbf{x}_0 = (0, 0)^\top$ the only stationary point
- The HESSE matrix is:

$$\nabla^2 f(x_1, x_2) = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$

- We compute the determinant $\Delta(\mathbf{x}_0) = 2 \cdot 2 - 0 = 4 > 0$
- Since $\frac{\partial f(\mathbf{x}_0)}{\partial x_1 \partial x_1} = 2 > 0$, we conclude that \mathbf{x}_0 is a local minimum

Constrained Optimization

- In many practical applications we have to consider **additional constraints** in the optimization process
- In the following we shall consider the following optimization problem (P):

$$\text{minimize} \quad f(\mathbf{x})$$

$$\text{subject to} \quad h_k(\mathbf{x}) = 0 \quad (k = 1, \dots, K)$$

$$g_\ell(\mathbf{x}) \leq 0 \quad (\ell = 1, \dots, L)$$

Constrained Optimization (Ctd.)

- The differentiable function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is referred to as **objective function**
- The functions h_k ($k = 1, \dots, K$) represent K **equality constraints**
- The functions g_ℓ ($\ell = 1, \dots, L$) represent L **inequality constraints**
- We define the set of **feasible points**

$$\mathcal{Z}_P := \left\{ \mathbf{x} \in \mathbb{R}^D \mid h_k(\mathbf{x}) = 0 \ (k = 1, \dots, K), g_\ell(\mathbf{x}) \leq 0 \ (\ell = 1, \dots, L) \right\}$$

- For any solution \mathbf{x}^* to (P) we have $\mathbf{x}^* \in \mathcal{Z}_P$
- In the unconstrained case all points are feasible, i. e. we have $\mathcal{Z}_P := \mathbb{R}^D$

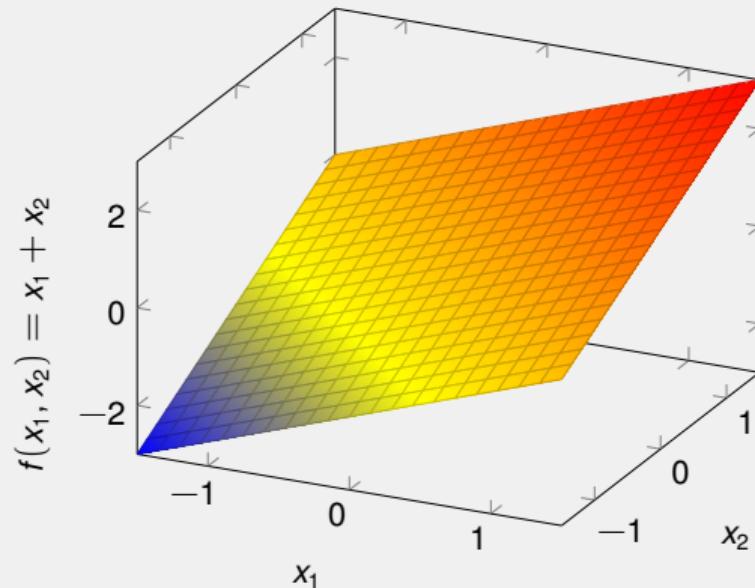
Example: Constrained Optimization

- Minimize

$$f(x_1, x_2) := x_1 + x_2$$

- Subject to the constraint

$$h(x_1, x_2) := x_1^2 + x_2^2 - 1 = 0$$



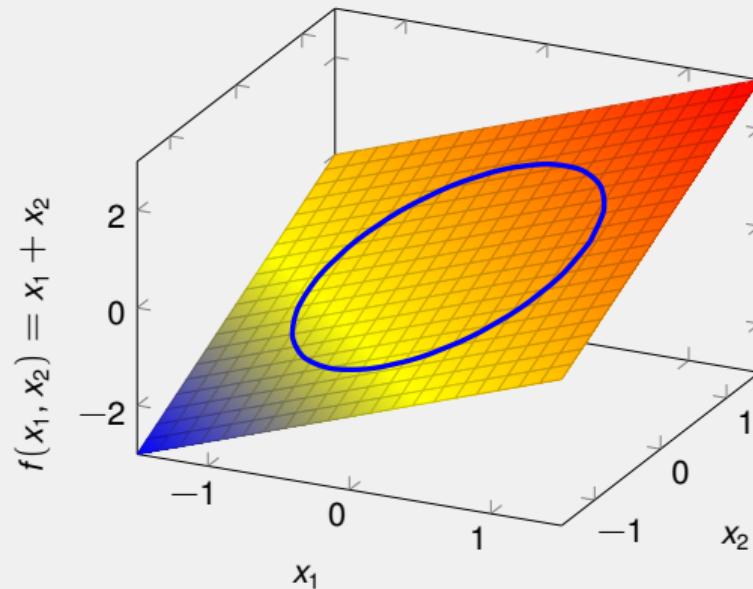
Example: Constrained Optimization

- Minimize

$$f(x_1, x_2) := x_1 + x_2$$

- Subject to the constraint

$$h(x_1, x_2) := x_1^2 + x_2^2 - 1 = 0$$



Necessary Conditions for Solutions

- For unconstrained optimization problems we search for points \mathbf{x} which satisfy

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

- In the following we shall introduce the so-called

KARUSH-KUHN-TUCKER (KKT) conditions

which represent first-order necessary conditions for a solution to problem (P) , provided that some **regularity conditions/constraint qualifications** are satisfied

- Please note:** In general, these conditions are not sufficient for a solution!

KKT Conditions

The following set of equalities and inequalities with the variables $(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \in \mathbb{R}^{D+K+L}$ are called KKT conditions for problem (P) :

$$h_k(\mathbf{x}) = 0 \quad (k = 1, \dots, K) \quad (12)$$

$$g_\ell(\mathbf{x}) \leq 0 \quad (\ell = 1, \dots, L) \quad (13)$$

$$\nabla f(\mathbf{x}) + \sum_{k=1}^K \lambda_k \nabla h_k(\mathbf{x}) + \sum_{\ell=1}^L \mu_\ell \nabla g_\ell(\mathbf{x}) = \mathbf{0} \quad (14)$$

$$\mu_\ell g_\ell(\mathbf{x}) = 0 \quad (\ell = 1, \dots, L) \quad (15)$$

$$\mu_\ell \geq 0 \quad (\ell = 1, \dots, L) \quad (16)$$

KKT Conditions Remarks

- A point $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ satisfying the KKT conditions is called **KKT point**
- The conditions (12) and (13) ensure that $\mathbf{x}^* \in \mathcal{Z}_P$
- The equality (14) requires that the gradient of f be a linear combination of the gradients of the constraints, where $\boldsymbol{\lambda} \in \mathbb{R}^K$ and $\boldsymbol{\mu} \in \mathbb{R}^L$ are called **(LAGRANGE multipliers)**
- The conditions in (15) are called **complementary slackness conditions** and require for a KKT point \mathbf{x}^* that $\mu_\ell = 0$ or $g_\ell(\mathbf{x}^*) = 0$ holds true
- **Remark:** In the unconstrained case, the KKT conditions reduce to $\nabla f(\mathbf{x}) = \mathbf{0}$

KKT Conditions Remarks (Ctd.)

- A local minimum \mathbf{x}^* of (P) is a KKT point if all constraints are (affine-)linear functions, i. e. the KKT conditions are necessary in this case
- In general, the KKT conditions are necessary conditions if the set of feasible points satisfies additional **constraint qualifications** (*not covered here*)

Lemma: Let f be a convex differentiable function with a continuous derivative and let h_k and g_ℓ be affine-linear constraints. Then \mathbf{x}^* is a solution to problem (P) if and only if \mathbf{x}^* is a KKT point for (P) .

Optimization Problems

Question: What should an ideal optimization problem look like?

minimize $f(\mathbf{x})$ \longleftarrow cost function / objective function

subject to $h_k(\mathbf{x}) = 0 \quad (k = 1, \dots, K)$ \longleftarrow equality constraints

$g_\ell(\mathbf{x}) \leq 0 \quad (\ell = 1, \dots, L)$ \longleftarrow inequality constraints

Optimization Problems (Ctd.)

Answer:

$$\text{minimize} \quad f(\mathbf{x}) \qquad \longleftarrow \text{convex function}$$

$$\text{subject to} \quad h_k(\mathbf{x}) = 0 \quad (k = 1, \dots, K) \qquad \longleftarrow \text{affine-linear functions}$$

$$g_\ell(\mathbf{x}) \leq 0 \quad (\ell = 1, \dots, L) \qquad \longleftarrow \text{convex set}$$

LAGRANGE Function

- We define the **LAGRANGE function** (*also called Lagrangian*):

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) := f(\mathbf{x}) + \sum_{k=1}^K \lambda_k h_k(\mathbf{x}) + \sum_{\ell=1}^L \mu_\ell g_\ell(\mathbf{x}) \quad (17)$$

$$= f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x}) \quad (18)$$

- We can write condition (14) in terms of the LAGRANGE function:

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0} \quad (19)$$

Optimization Problems with Equality Constraints

- In the following we consider problems **without inequality constraints**
- Thus, problem (P) reduces to the special case

$$\text{minimize} \quad f(\mathbf{x})$$

$$\text{subject to} \quad h_k(\mathbf{x}) = 0 \quad (k = 1, \dots, K)$$

- We call this problem (P_{Eq}) as the solution is constrained exclusively by equality constraints

KKT Conditions for Problems with Equality Constraints

- The KKT conditions for problem (P_{Eq}) reduce to a system of $D + K$ equations in $D + K$ variables:

$$h_k(\mathbf{x}) = 0 \quad (k = 1, \dots, K)$$

$$\nabla f(\mathbf{x}) + \sum_{k=1}^K \lambda_k \nabla h_k(\mathbf{x}) = \mathbf{0}$$

- These conditions can be rewritten in terms of the LAGRANGE function (17):

$$\nabla_{\mathbf{x}} \mathfrak{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0} \quad \text{and} \quad \nabla_{\boldsymbol{\lambda}} \mathfrak{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0} \quad (20)$$

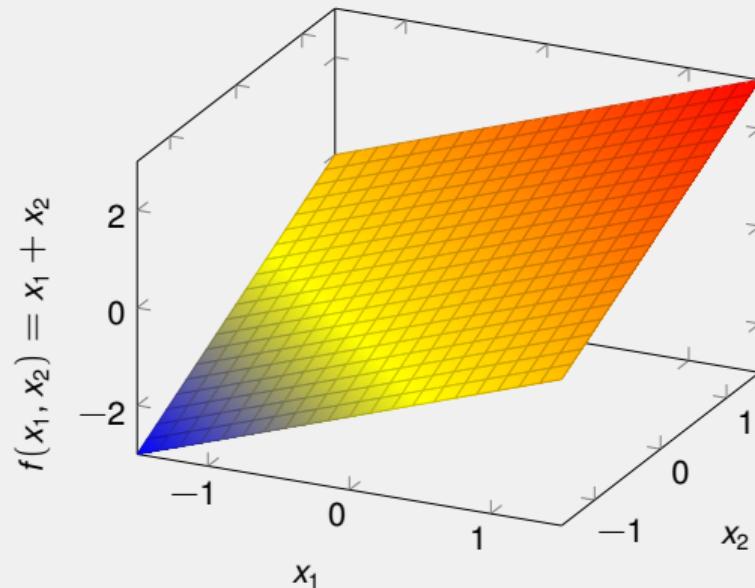
Example Revisited

- Let us revisit the example from above:
- Minimize

$$f(x_1, x_2) := x_1 + x_2$$

- Subject to the constraint

$$h(x_1, x_2) := x_1^2 + x_2^2 - 1 = 0$$



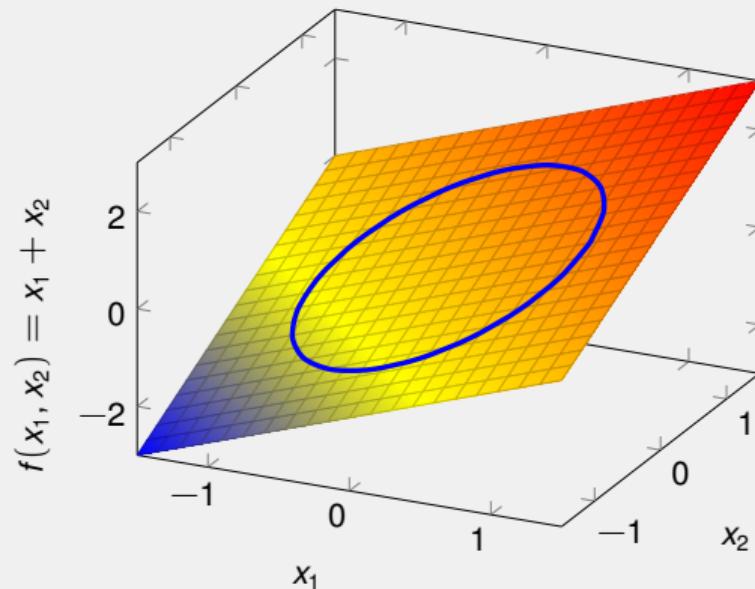
Example Revisited

- Let us revisit the example from above:
- Minimize

$$f(x_1, x_2) := x_1 + x_2$$

- Subject to the constraint

$$h(x_1, x_2) := x_1^2 + x_2^2 - 1 = 0$$



Example Revisited (Ctd.)

- **Step 1)** We set up the LAGRANGE function for the problem:

$$\mathcal{L}(\mathbf{x}, \lambda) := x_1 + x_2 + \lambda(x_1^2 + x_2^2 - 1)$$

- **Step 2)** We compute the partial derivatives and set them to zero:

$$(i) \quad \partial \mathcal{L} / \partial x_1 = 1 + 2\lambda x_1 \stackrel{!}{=} 0$$

$$(ii) \quad \partial \mathcal{L} / \partial x_2 = 1 + 2\lambda x_2 \stackrel{!}{=} 0$$

$$(iii) \quad \partial \mathcal{L} / \partial \lambda = x_1^2 + x_2^2 - 1 \stackrel{!}{=} 0$$

Example Revisited (Ctd.)

- **Step 3)** We solve the system of equations from step 2):
 - Assuming $\lambda \neq 0$ we can solve (i) for x_1 and (ii) for x_2
 - We obtain $x_1 = -\frac{1}{2\lambda} = x_2$
 - Plugging this result into (iii), we get

$$2\lambda^2 = 1 \iff \lambda = \pm \frac{1}{\sqrt{2}}$$

- Thus, the candidate solutions are:

$$\underline{\mathbf{x}} = \left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2} \right) \quad \text{and} \quad \bar{\mathbf{x}} = \left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right)$$



Bordered HESSE Matrix

- We can use the **bordered HESSE matrix** (German: *geränderte HESSE Matrix*) to decide if a candidate solution is a minimum or a maximum
- In the example above the bordered HESSE matrix has the following structure:

$$\tilde{\nabla}^2 f(\mathbf{x}, \lambda) := \begin{pmatrix} 0 & \frac{\partial h}{\partial x_1} & \frac{\partial h}{\partial x_2} \\ \frac{\partial h}{\partial x_1} & \frac{\partial^2 \mathfrak{L}}{\partial x_1 \partial x_1} & \frac{\partial^2 \mathfrak{L}}{\partial x_1 \partial x_2} \\ \frac{\partial h}{\partial x_2} & \frac{\partial^2 \mathfrak{L}}{\partial x_2 \partial x_1} & \frac{\partial^2 \mathfrak{L}}{\partial x_2 \partial x_2} \end{pmatrix} \quad (21)$$



Bordered HESSE Matrix (Ctd.)

The sign of the determinant tells us if we have a local minimum or maximum:

Lemma:

- $\det(\tilde{\nabla}^2 f(\mathbf{x}, \lambda)) > 0 \implies \mathbf{x}$ is a **local maximum**
- $\det(\tilde{\nabla}^2 f(\mathbf{x}, \lambda)) < 0 \implies \mathbf{x}$ is a **local minimum**
- $\det(\tilde{\nabla}^2 f(\mathbf{x}, \lambda)) = 0 \implies \text{We cannot decide}$

Attention: We have to check the sign of the determinant of the bordered HESSE matrix, not its definiteness!



Example Revisited (Ctd.)

- **Step 4)** We decide for minimum or maximum:
 - We compute the bordered HESSE matrix according to (21):
$$\tilde{\nabla}^2 f(\mathbf{x}, \lambda) = \begin{pmatrix} 0 & 2x_1 & 2x_2 \\ 2x_1 & 2\lambda & 0 \\ 2x_2 & 0 & 2\lambda \end{pmatrix}$$
 - Now we plug the two candidates $\underline{\mathbf{x}}$ and $\bar{\mathbf{x}}$ into the bordered HESSE matrix and compute the determinants



Example Revisited (Ctd.)

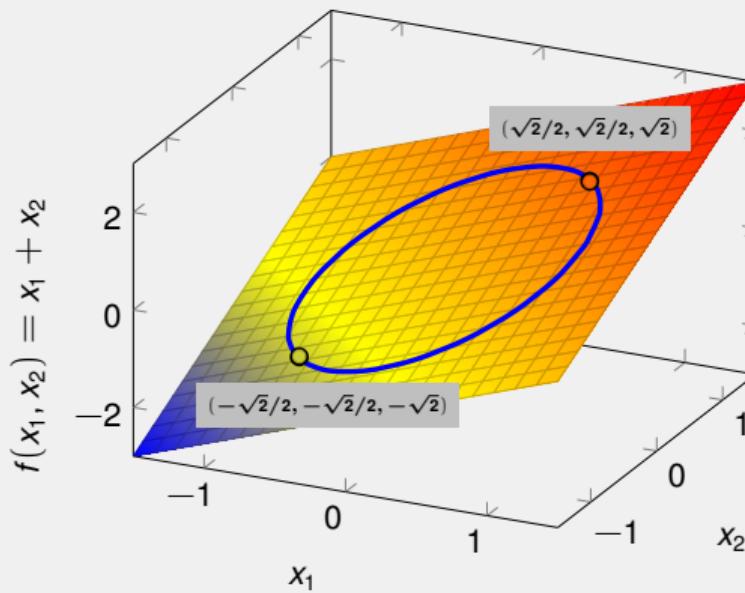
For $\underline{\mathbf{x}}$ we get a **minimum**:

$$\det(\tilde{\nabla}^2 f(\underline{\mathbf{x}}, 1/\sqrt{2})) = \det \begin{pmatrix} 0 & -\sqrt{2} & -\sqrt{2} \\ -\sqrt{2} & 2/\sqrt{2} & 0 \\ -\sqrt{2} & 0 & 2/\sqrt{2} \end{pmatrix} = -4\sqrt{2} < 0$$

For $\bar{\mathbf{x}}$ we get a **maximum**:

$$\det(\tilde{\nabla}^2 f(\bar{\mathbf{x}}, -1/\sqrt{2})) = \det \begin{pmatrix} 0 & \sqrt{2} & \sqrt{2} \\ \sqrt{2} & -2/\sqrt{2} & 0 \\ \sqrt{2} & 0 & -2/\sqrt{2} \end{pmatrix} = 4\sqrt{2} > 0$$

Solution to the Example



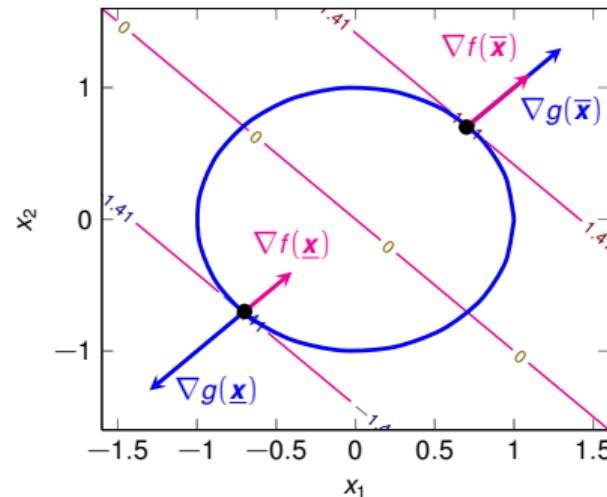
Geometry of LAGRANGE Optimization

- This is a **contour plot** of f
- Both points, $\underline{\mathbf{x}}$ and $\bar{\mathbf{x}}$, fulfill

$$\nabla f(\underline{\mathbf{x}}) = \lambda \nabla h(\underline{\mathbf{x}}) \quad (22)$$

$$h(\underline{\mathbf{x}}) = 0 \quad (23)$$

- ∇f and ∇h are **collinear**
- λ is a stretching/compression factor
- This leads to the same system of equations we have solved above!

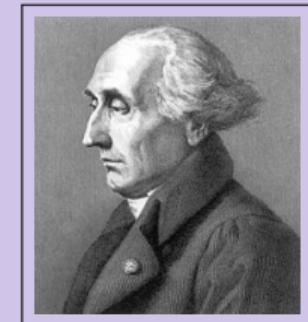




Portrait: JOSEPH-LOUIS LAGRANGE

JOSEPH-LOUIS LAGRANGE (1736 – 1813) was an Italian mathematician, physicist and astronomer. He made significant contributions to the fields of analysis, number theory, and both classical and celestial mechanics.

LAGRANGE was one of the creators of the calculus of variations, deriving the EULER–LAGRANGE equations for extrema of functionals. He extended the method **to include possible constraints**, arriving at the method of **LAGRANGE multipliers**. LAGRANGE invented the method of solving differential equations known as variation of constants and applied differential calculus to the theory of probabilities.



(Wikipedia)

Exercise

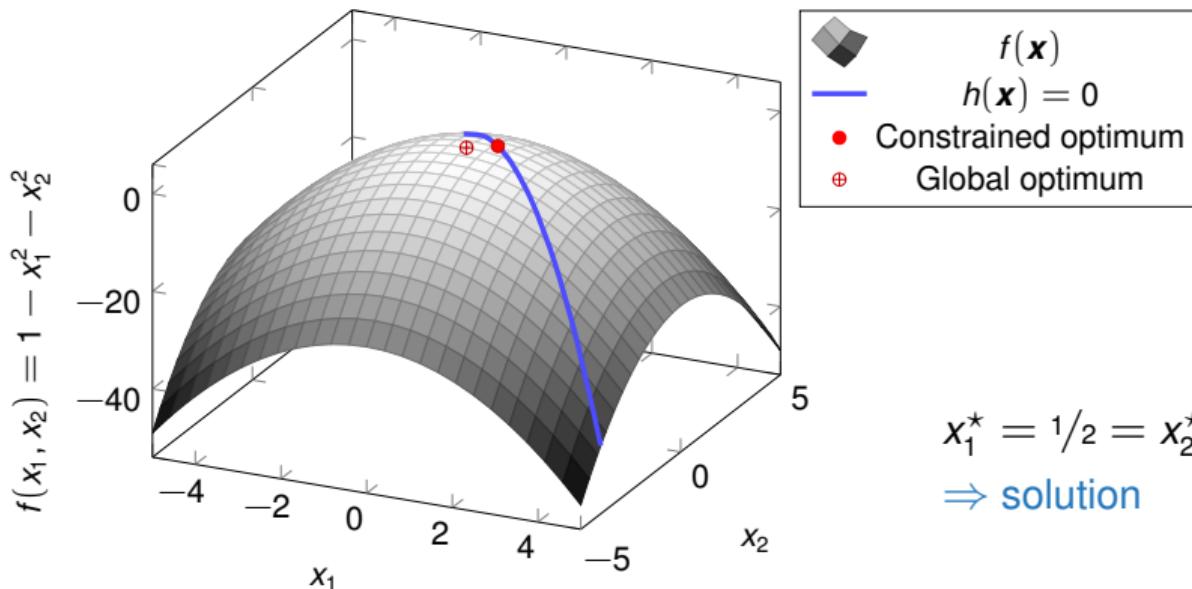
Solve the following optimization problem:

$$\text{maximize} \quad f(x_1, x_2) := 1 - x_1^2 - x_2^2$$

$$\text{subject to} \quad h(x_1, x_2) := x_1 + x_2 - 1 = 0$$

(This is a concave function with linear constraints, therefore there is only one global maximum which you don't have to show!)

Exercise Solution



Dealing with additional Inequality Constraints

- For simple problems we can try to find a solution by directly solving the KKT system given by (12) - (16)
- We usually start with the complementary slackness conditions in (15) and consider all combinations for which they are satisfied
- We then plug these values into the other conditions and check if all of them hold true, if yes we have found a KKT point
- This approach can become tedious for large problems
- Optimization libraries implement efficient algorithms to solve such problems numerically

Dealing with additional Inequality Constraints (Ctd.)

Let the following optimization problem (P) be given:

$$\text{minimize} \quad q(x_1, x_2) = (x_1 - 1)^2 + (x_2 - 2.5)^2$$

$$\text{subject to} \quad -x_1 + 2x_2 - 2 \leqslant 0 \quad 1$$

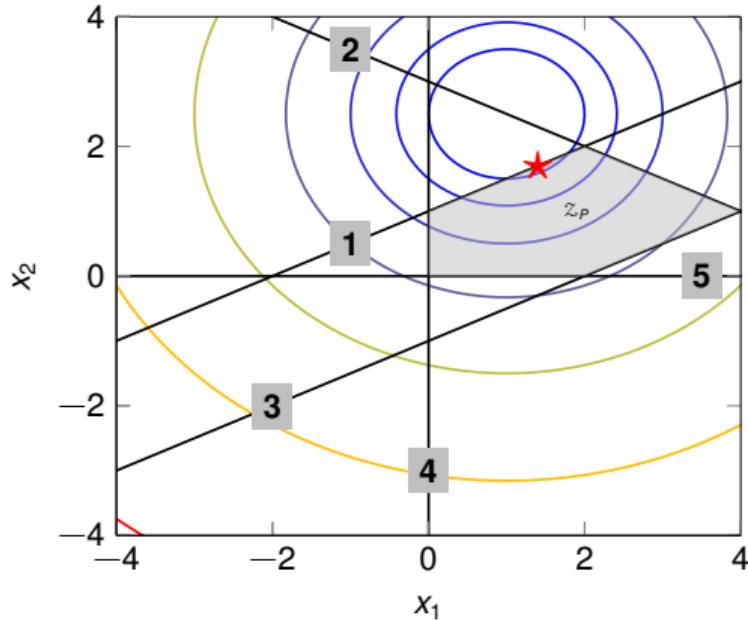
$$x_1 + 2x_2 - 6 \leqslant 0 \quad 2$$

$$x_1 - 2x_2 - 2 \leqslant 0 \quad 3$$

$$x_1 \geqslant 0 \quad 4$$

$$x_2 \geqslant 0 \quad 5$$

Dealing with additional Inequality Constraints (Ctd.)



LAGRANGE Duality

- **Duality** in optimization is the idea of converting an optimization problem (**primal problem**) from one set of variables x (**primal variables**) into a different set of variables μ (**dual variables**) which gives rise to the **dual problem**
- The dual problem might be easier to solve than the primal problem
- The primal and dual optimization problems have the same solution, but the value of the objective function may differ (**duality gap**)
- For convex optimization problems the duality gap is zero if **Slater's condition** holds
- There are different approaches to duality
- We shall now introduce the concept of **LAGRANGE duality**

Primal and dual Optimization Problems

- Consider the **primal optimization problem**

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && g_\ell(\mathbf{x}) \leq 0 \quad (\ell = 1, \dots, L). \end{aligned}$$

- The **dual optimization problem** is then given by

$$\begin{aligned} & \text{maximize (!)} && \mathfrak{D}(\boldsymbol{\mu}) := \inf_{\mathbf{x} \in \mathbb{R}^D} \mathfrak{L}(\mathbf{x}, \boldsymbol{\mu}) \\ & \text{subject to} && \boldsymbol{\mu} \geq \mathbf{0}. \end{aligned}$$

The dual optimization Problem

- \mathfrak{D} is a **concave function** (*which is why we have to maximize*)
- We can turn the dual optimization problem into a minimization problem by flipping the sign of \mathfrak{D}
- $-\mathfrak{D}$ is then a **convex function**
- We assume f and g_ℓ to be differentiable functions
- We can then compute $\inf_{\mathbf{x} \in \mathbb{R}^D} \mathfrak{L}(\mathbf{x}, \boldsymbol{\mu})$ for a given value of $\boldsymbol{\mu}$ by differentiating \mathfrak{L} with respect to \mathbf{x} , setting the differential to zero, and solving for the optimal value
- The next slides will demonstrate this for a quadratic objective function

Example: Duality for quadratic Programs

- Consider the case of a quadratic objective function with affine constraints, e.g.

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A} \mathbf{x} - \mathbf{b} \leq \mathbf{0}, \end{aligned}$$

where $\mathbf{A} \in \mathbb{R}^{L \times D}$, $\mathbf{b} \in \mathbb{R}^L$, and $\mathbf{c} \in \mathbb{R}^D$

- We assume that $\mathbf{Q} \in \mathbb{R}^{D \times D}$ is symmetric and positive definite, and therefore, the objective function is a convex function
- This problem is known as a **quadratic program**

Example: Duality for quadratic Programs (Ctd.)

- The Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\mu})$ for this problem is given by

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x} + \boldsymbol{\mu}^\top (\mathbf{A} \mathbf{x} - \mathbf{b}) \\ &= \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu})^\top \mathbf{x} - \boldsymbol{\mu}^\top \mathbf{b}.\end{aligned}$$

- We take the derivative of $\mathcal{L}(\mathbf{x}, \boldsymbol{\mu})$ with respect to \mathbf{x} and set it to zero:

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\mu}) = \mathbf{Q} \mathbf{x} + (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu}) = \mathbf{0} \tag{24}$$

Example: Duality for quadratic Programs (Ctd.)

- \mathbf{Q} is invertible (because we assumed it to be positive definite)
- We can therefore rearrange equation (24) for \mathbf{x} and obtain

$$\mathbf{x} = -\mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu}).$$

- We plug this result into $\mathfrak{L}(\mathbf{x}, \boldsymbol{\mu})$ which yields the dual objective function

$$\mathfrak{D}(\boldsymbol{\mu}) = -\frac{1}{2}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu})^\top \mathbf{A} \mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu}) - \boldsymbol{\mu}^\top \mathbf{b}.$$

Example: Duality for quadratic Programs (Ctd.)

- The dual optimization problem is therefore given by

$$\begin{aligned} & \text{maximize} && -\frac{1}{2}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu})^\top \mathbf{A} \mathbf{Q}^{-1} (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\mu}) - \boldsymbol{\mu}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\mu} \geq \mathbf{0}. \end{aligned}$$

- A quadratic optimization problem subject to constraints has to be solved when training **support vector machines**

Section: Numerical Optimization

Introduction to numerical Optimization Techniques

First-Order Method: Gradient Descent

First-Order Method: Gradient Descent with Momentum

Second-Order Method: NEWTON's Method

Numerical Optimization

Problem:

Some optimization problems are too complex to solve them analytically!

Solution: Use numerical optimization techniques!

There is a plethora of techniques which differ in the **number of iterations** required, the **computational cost**, the **convergence guarantees**, the **robustness** with noisy cost functions, and their **memory usage**

Numerical Optimization (Ctd.)

- In the following we shall assume a cost function – consider e.g. (1)

$$\mathfrak{J} : \mathbb{R}^M \rightarrow \mathbb{R}, \boldsymbol{\theta} \mapsto \mathfrak{J}(\boldsymbol{\theta}),$$

where $\boldsymbol{\theta} \in \mathbb{R}^M$ is the vector of model parameters

- General idea:** Incrementally update an initial estimate of the parameters $\boldsymbol{\theta}$:

$$\boldsymbol{\theta}^{t+1} \leftarrow \boldsymbol{\theta}^t + \delta_{\boldsymbol{\theta}}^t \quad (25)$$

- Choose the **step** $\delta_{\boldsymbol{\theta}}^t$ such that $\mathfrak{J}(\boldsymbol{\theta}^{t+1}) < \mathfrak{J}(\boldsymbol{\theta}^t)$
- Optimization algorithms differ in how they determine $\delta_{\boldsymbol{\theta}}^t$

Numerical Optimization Algorithms (some Examples)

Gradient based methods:

- **Gradient descent**
- Conjugate gradient descent
- **NEWTON's method**
- (L-)BFGS

(BROYDEN-FLETCHER-GOLDFARB-SHANNO)

Non-gradient based methods:

- Genetic algorithms
- Particle swarm
- Non-linear simplex
- NELDER-MEAD

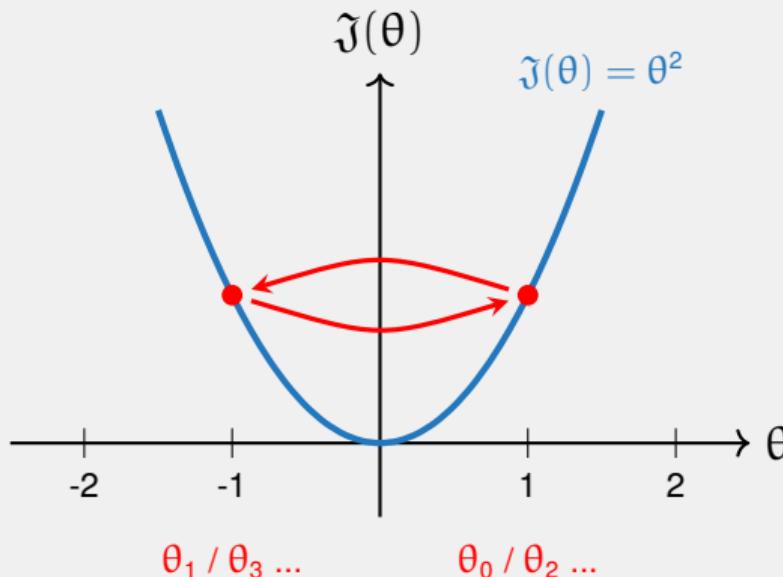
Numerical optimization techniques may not find the global optimum!

Gradient Descent

- **General idea:**
 - ① Start with a random (or informed) guess of the parameters θ^0
 - ② Then go into the direction of the **steepest descent**
 - ③ Iterate until the solution is '*good enough*'
- We have already seen that the negative gradient points into the direction of steepest descent, therefore:
- **First idea:**

$$\delta_{\theta}^t \leftarrow -\nabla_{\theta} \mathfrak{J}(\theta^t) \quad \theta^{t+1} \leftarrow \theta^t + \delta_{\theta}^t \quad (26)$$

Example: Gradient Descent



- Consider $J(\theta) = \theta^2$
- $J'(\theta) = 2\theta$
- Start with $\theta_0 = 1$

$$\theta_1 \leftarrow 1 - 2 = -1$$

$$\theta_2 \leftarrow -1 - (-2) = 1$$

⋮

Gradient Descent: Update Rule

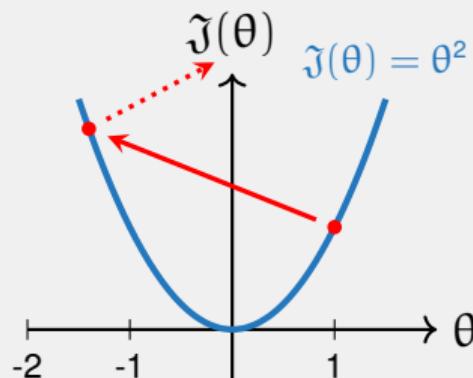
- It turns out that (26) is numerically unstable
- Therefore, we introduce the **learning rate** $\alpha \in (0, 1)$ which stabilizes the learning process
- **Better idea:**

$$\delta_{\theta}^t \leftarrow -\alpha \nabla_{\theta} \mathfrak{J}(\theta^t) \quad \theta^{t+1} \leftarrow \theta^t + \delta_{\theta}^t \quad (27)$$

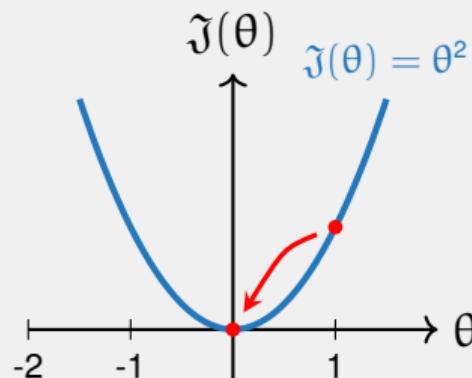
Gradient descent is a basic, but most commonly used optimization algorithm in machine learning!

Example revisited: Gradient Descent

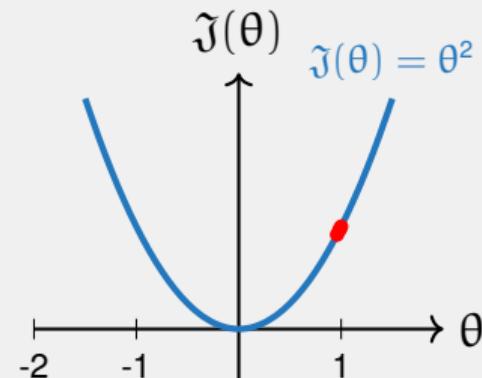
$\alpha = 1.2$ (too high)



$\alpha = 0.5$ (about right)



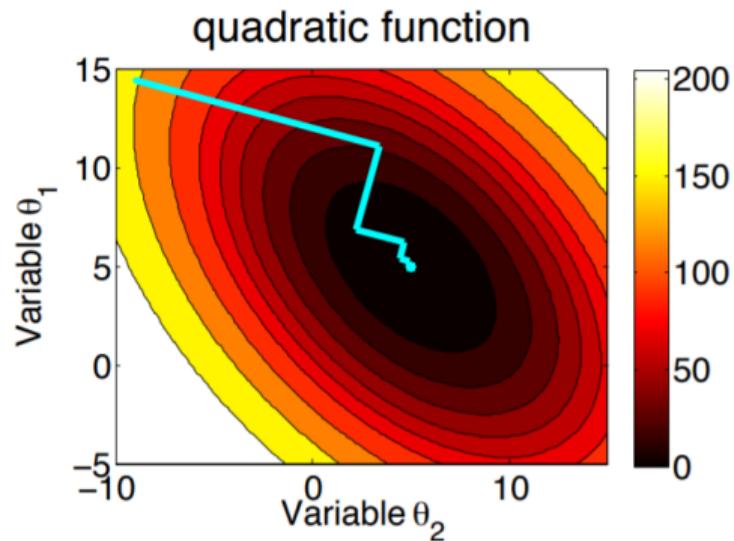
$\alpha = 0.001$ (too low)



Caution: The value of the learning rate depends on the optimization problem!

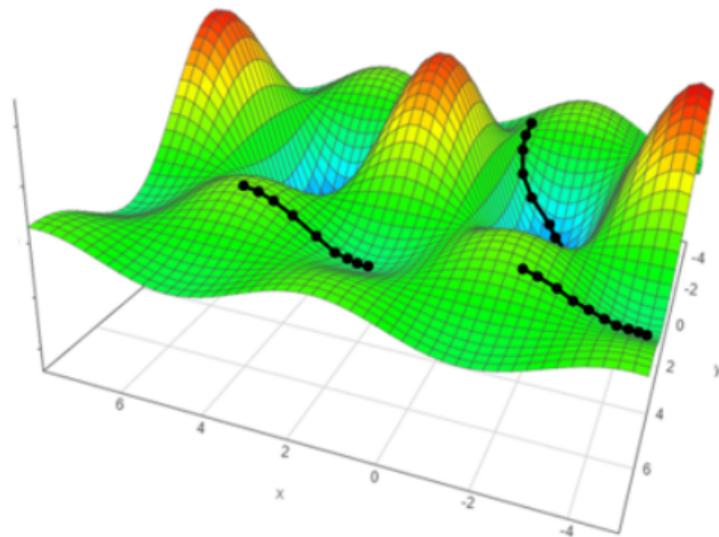
Problems of Gradient Descent

The parameter updates tend to 'zig-zag' down the valley (as the **curvature** of the function is not considered):



Problems of Gradient Descent (Ctd.)

Initialization also matters...



Gradient Descent with Momentum

- If the function is not convex, we run into the risk of not finding the optimal solution
- We adjust the algorithm to make it less likely to **get stuck in local solutions** and simultaneously **accelerate its convergence**
- Gradient descent with **momentum**:

$$\delta_{\theta}^t \leftarrow \mu \delta_{\theta}^{t-1} + (1 - \mu) \nabla_{\theta} \mathfrak{J}(\theta^t) \quad (28)$$

$$\theta^{t+1} \leftarrow \theta^t - \alpha \delta_{\theta}^t \quad (29)$$

- $\mu \in [0, 1]$ is the momentum parameter ($\mu = 0$ recovers gradient descent)

NEWTON's Method

- Unlike gradient descent, NEWTON's method is a **second-order method**
- As a second-order method it also takes the **curvature** of the function into account
- Another name is NEWTON-RAPHSON algorithm

Lemma: The NEWTON step size δ_{θ} to optimize a function $\mathfrak{J} : \mathbb{R}^M \rightarrow \mathbb{R}$ is given by the expression

$$\delta_{\theta} := -\mathbf{H}^{-1}(\boldsymbol{\theta})\mathbf{g}(\boldsymbol{\theta}), \quad (30)$$

where $\mathbf{g}(\boldsymbol{\theta})$ is the gradient, and $\mathbf{H}(\boldsymbol{\theta})$ the HESSE matrix of \mathfrak{J} evaluated at $\boldsymbol{\theta}$



Proof of NEWTON's Method

Proof: We consider the second-order TAYLOR expansion of \mathfrak{J} at the point θ which is given by

$$\mathfrak{J}(\theta + \delta_\theta) \approx \mathfrak{J}(\theta) + \mathbf{g}(\theta)^\top \delta_\theta + \frac{1}{2} \delta_\theta^\top \mathbf{H}(\theta) \delta_\theta, \quad (31)$$

where $\mathbf{g}(\theta)$ is the gradient, and $\mathbf{H}(\theta)$ the HESSE matrix of \mathfrak{J} evaluated at θ . Subsequently, we differentiate (31) with respect to δ_θ and set the derivative to zero:

$$\nabla_{\delta_\theta} \mathfrak{J}(\theta + \delta_\theta) = \mathbf{g}(\theta) + \mathbf{H}(\theta) \delta_\theta \stackrel{!}{=} \mathbf{0} \quad (32)$$

Rearranging for δ_θ yields

$$\delta_\theta = -\mathbf{H}^{-1}(\theta) \mathbf{g}(\theta) \quad (33)$$

which concludes the proof. ■

NEWTON's Method: Update Rule

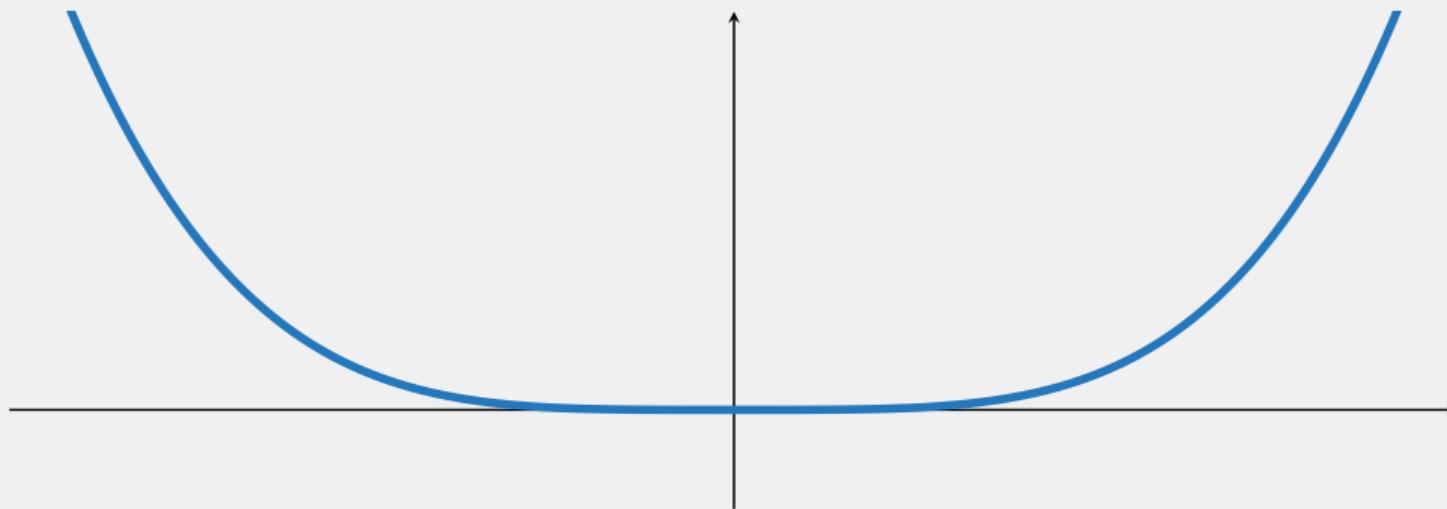
- The NEWTON update rule therefore takes the form

$$\delta_{\theta}^t \leftarrow -H^{-1}(\theta^t)g(\theta^t) \quad (34)$$

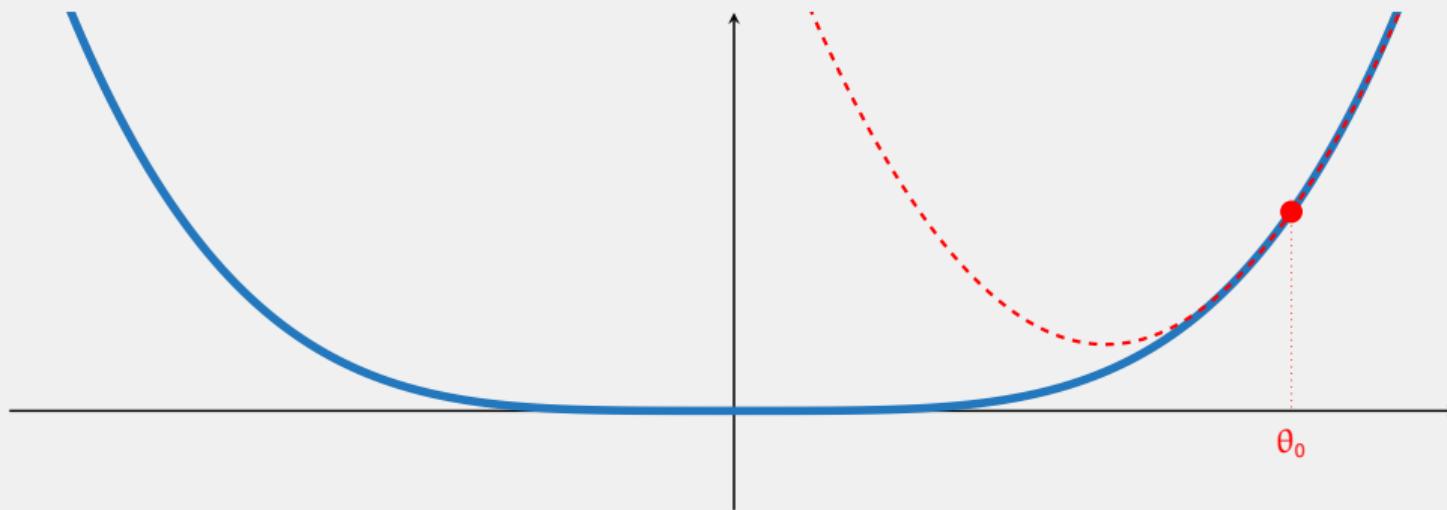
$$\theta^{t+1} \leftarrow \theta^t + \delta_{\theta}^t \quad (35)$$

- The learning rate α has been replaced by the inverse of the Hessian matrix
- Newton's method **converges faster** (quadratic convergence) than gradient descent, but the inverse of the HESSE matrix is **expensive to compute** $\approx \mathcal{O}(n^3)$

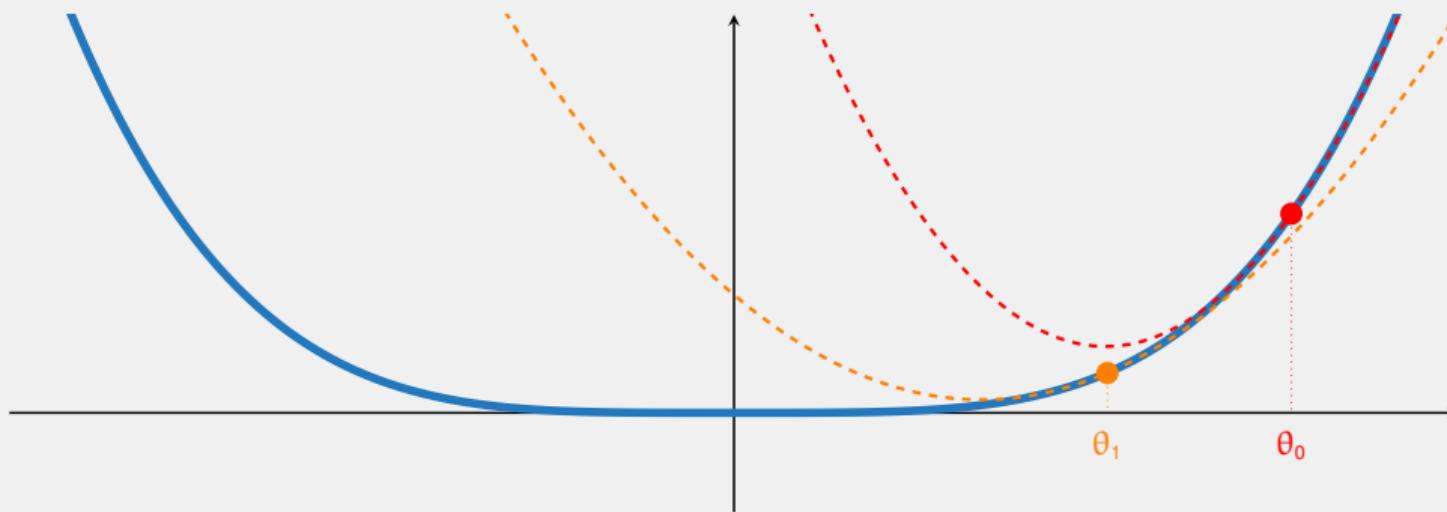
Visualization: NEWTON's Method



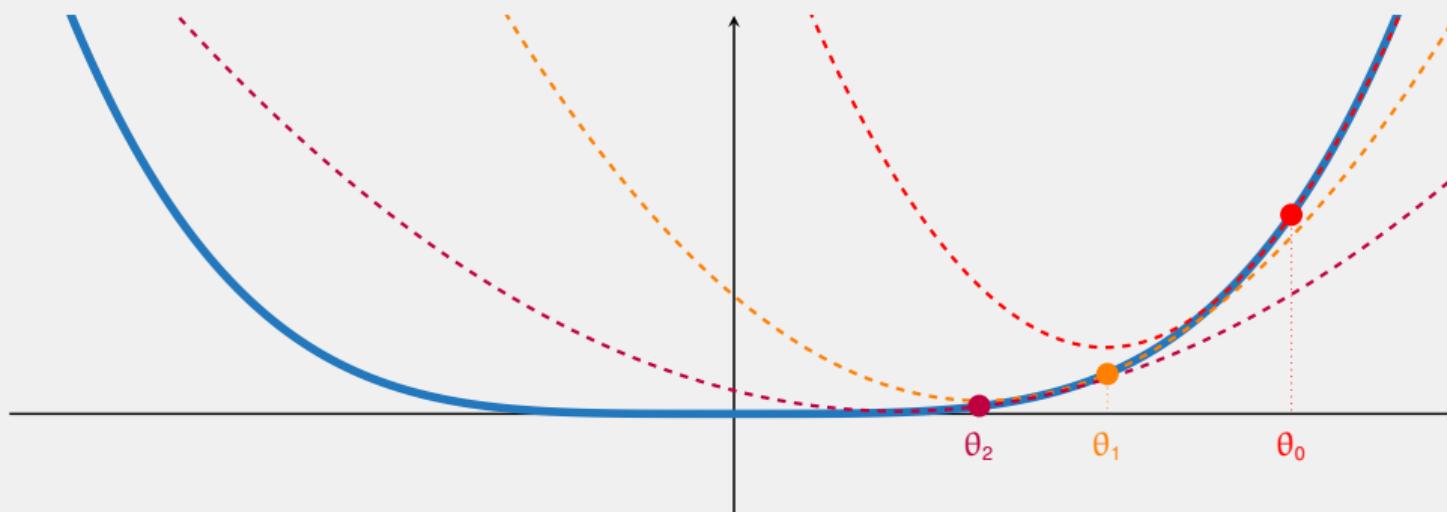
Visualization: NEWTON's Method



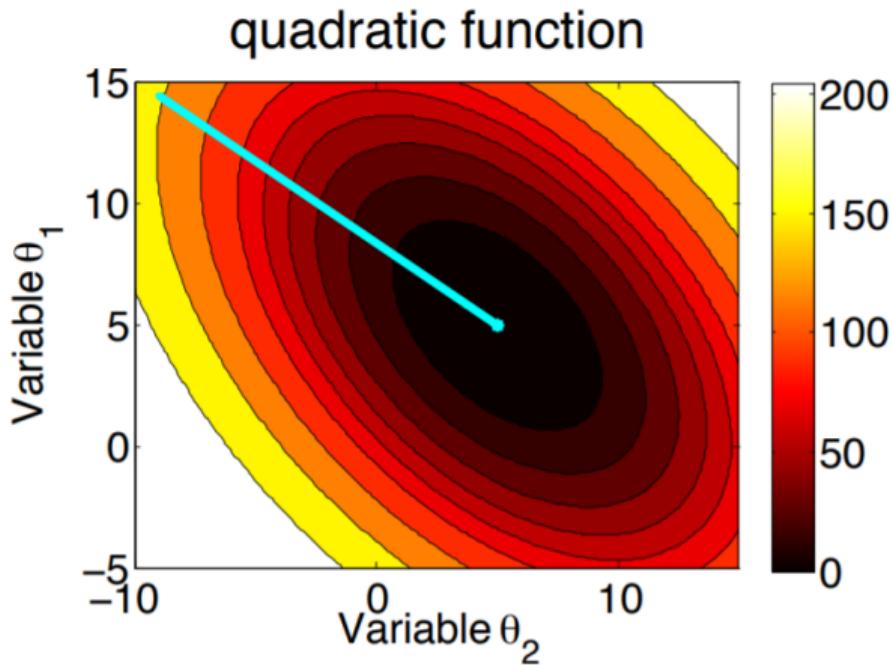
Visualization: NEWTON's Method



Visualization: NEWTON's Method

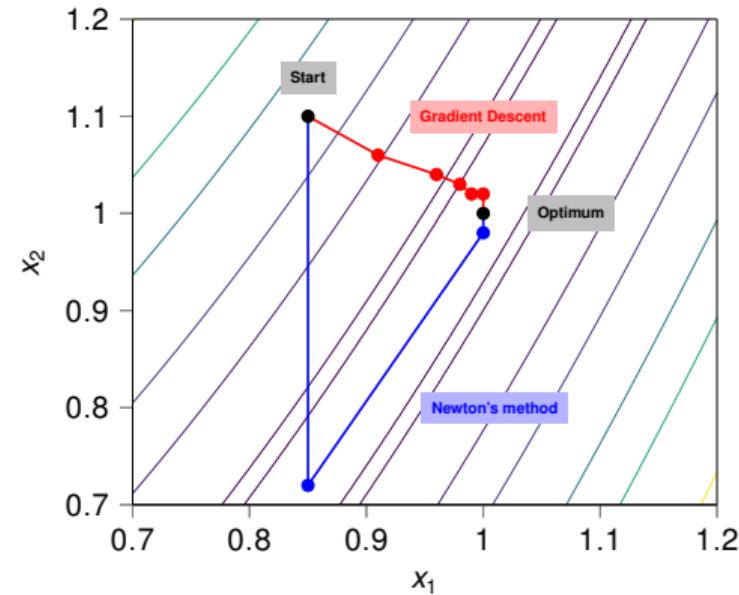


Visualization: NEWTON's Method (Ctd.)



NEWTON's Method in Comparison to Gradient Descent

- NEWTON's methods (blue) needs fewer iterations than gradient descent (red)
- **However, matrix inversion is expensive!**



Section: Wrap-Up

- Summary
- Recommended Literature
- Self-Test Questions
- Lecture Outlook

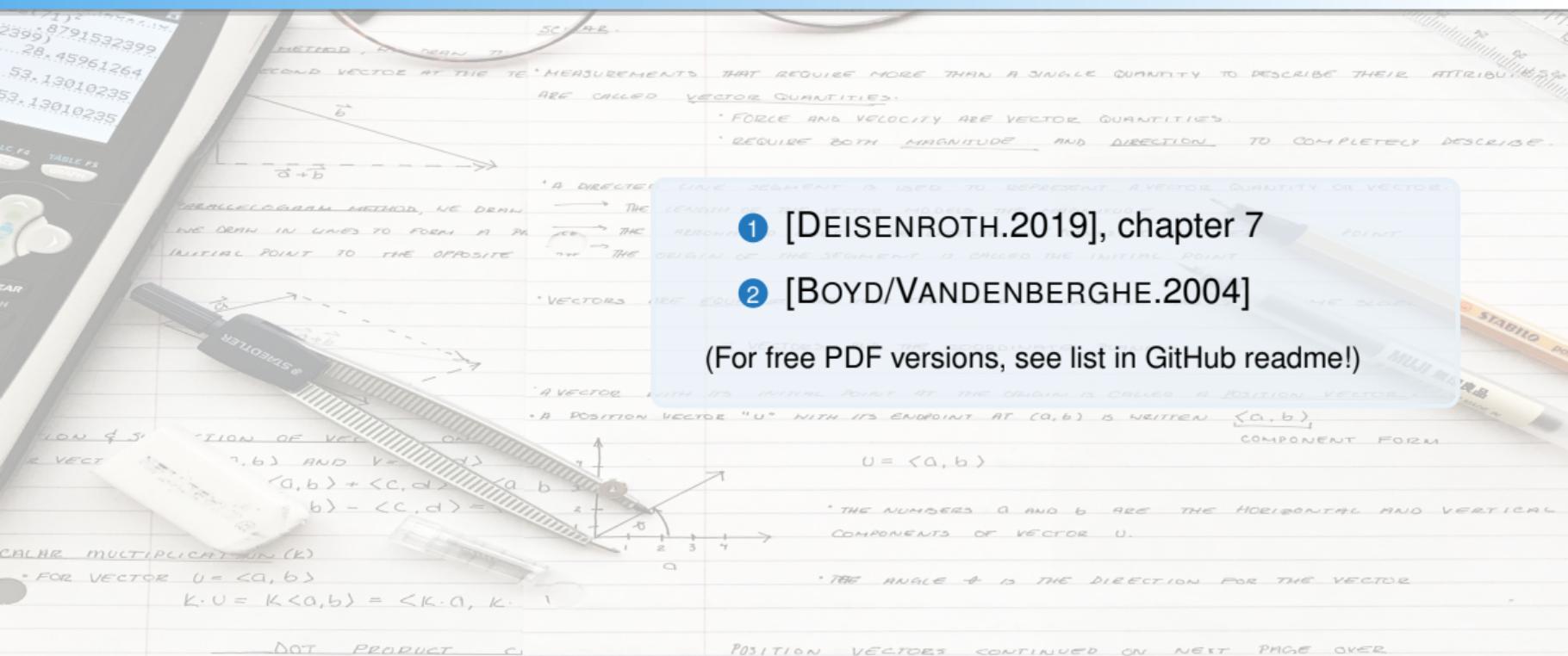
Summary

- **Every machine learning problem is an optimization problem!**
- We have learned about gradients, the HESSE matrix, and the TAYLOR expansion
- Good cost functions are **convex**
- You should be familiar with unconstrained and constrained optimization, especially with the concept of LAGRANGE optimization
- Closed-form solutions are often not feasible. In such cases we resort to **numerical optimization techniques**
- The most prominent numerical technique is called **gradient descent**
- NEWTON's method requires less iterations, but is computationally more demanding

Recommended Literature

- 1 [DEISENROTH.2019], chapter 7
- 2 [BOYD/VANDENBERGHE.2004]

(For free PDF versions, see list in GitHub readme!)



Self-Test Questions

- ① What is a gradient and how is it used in optimization?
- ② Which property of gradients is exploited in gradient descent?
- ③ What is convexity? Why should cost functions be convex?
- ④ How does LAGRANGE optimization work? What is the geometric interpretation?
- ⑤ Describe the gradient descent algorithm! Which issues does it have?
- ⑥ Explain the role of the learning rate in the gradient descent algorithm!
- ⑦ What is momentum?
- ⑧ How does NEWTON's method differ from gradient descent?

What's next...?

- I Machine Learning Introduction
- II Optimization Techniques
- III Bayesian Decision Theory
- IV Non-parametric Density Estimation
- V Probabilistic Graphical Models
- VI Linear Regression
- VII Logistic Regression
- VIII Deep Learning
- IX Evaluation
- X Decision Trees
- XI Support Vector Machines
- XII Clustering
- XIII Principal Component Analysis
- XIV Reinforcement Learning
- XV Advanced Regression

Thank you very much for the attention!

* * * Artificial Intelligence and Machine Learning * * *

Topic: Optimization Techniques

Term: Summer term 2025

Contact:

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

daniel.wehner@sap.com

Do you have any questions?