

Introduction to Optimization

Approximation Algorithms and Heuristics

November 9, 2015

École Centrale Paris, Châtenay-Malabry, France



Dimo Brockhoff
INRIA Lille – Nord Europe

Course Overview

Date		Topic
Mon, 21.9.2015		Introduction
Mon, 28.9.2015	D	Basic Flavors of Complexity Theory
Mon, 5.10.2015	D	Greedy algorithms
Mon, 12.10.2015	D	Branch and bound (switched w/ dynamic programming)
Mon, 2.11.2015	D	Dynamic programming [<i>salle Proto</i>]
Fri, 6.11.2015	D	Approximation algorithms and heuristics [<i>S205/S207</i>]
Mon, 9.11.2015	C	Introduction to Continuous Optimization I [<i>S118</i>]
Fri, 13.11.2015	C	Introduction to Continuous Optimization II <i>[from here onwards always: S205/S207]</i>
Fri, 20.11.2015	C	Gradient-based Algorithms
Fri, 27.11.2015	C	End of Gradient-based Algorithms + Linear Programming <i>Stochastic Optimization and Derivative Free Optimization I</i>
Fri, 4.12.2015	C	Stochastic Optimization and Derivative Free Optimization II
Tue, 15.12.2015		Exam

Overview of Today's Lecture

Introduction to Continuous Optimization

- examples (from ML / black-box problems)
- typical difficulties in optimization (e.g. constraints)

Mathematical Tools to Characterize Optima

- reminders about differentiability, gradient, Hessian matrix

Further Details on Remaining Lectures

Introduction to Continuous Optimization

- examples (from ML / black-box problems)
- typical difficulties in optimization (e.g. constraints)

Mathematical Tools to Characterize Optima

- reminders about differentiability, gradient, Hessian matrix
- unconstrained optimization
 - first and second order conditions
 - convexity
- constrained optimization

Gradient-based Algorithms

- quasi-Newton method (BFGS)

Learning in Optimization / Stochastic Optimization

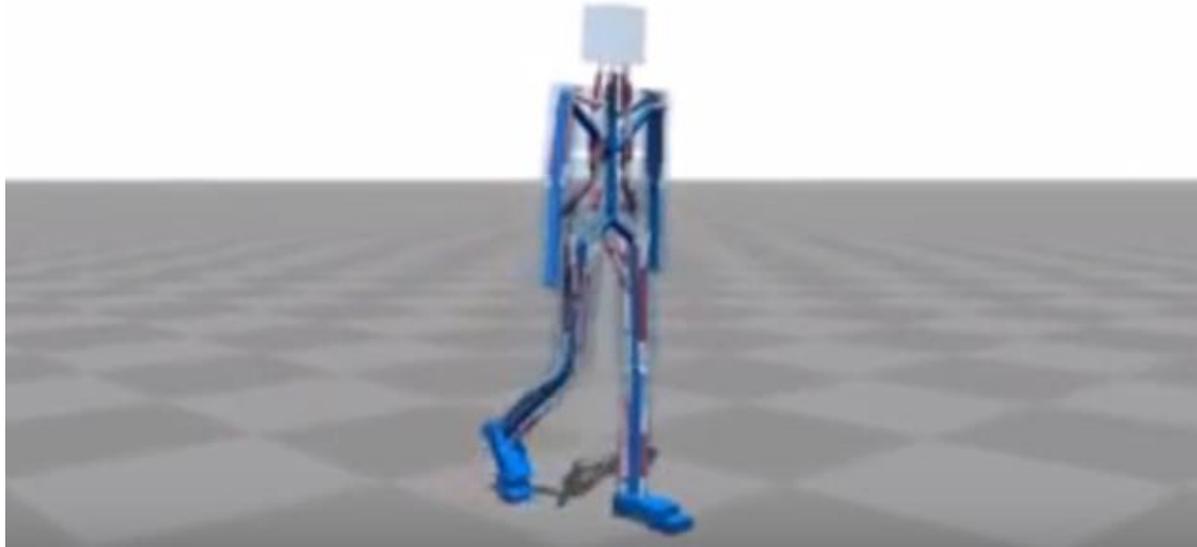
- CMA-ES (adaptive algorithms / Information Geometry)
- PhD thesis possible on this topic

strongly related to ML, new promising research area, interesting open questions

First Example of a Continuous Optimization Problem

Computer simulation teaches itself to walk upright (virtual robots (of different shapes) learning to walk, through stochastic optimization (CMA-ES)), by Utrecht University:

We present a control system based on 3D muscle actuation



<https://www.youtube.com/watch?v=yci5Fu1ovk>

T. Geitjenbeek, M. Van de Panne, F. Van der Stappen: "Flexible Muscle-Based Locomotion for Bipedal Creatures", SIGGRAPH Asia, 2013.

Unconstrained vs. Constrained Optimization

Unconstrained optimization

$$\inf \{f(x) \mid x \in \mathbb{R}^n\}$$

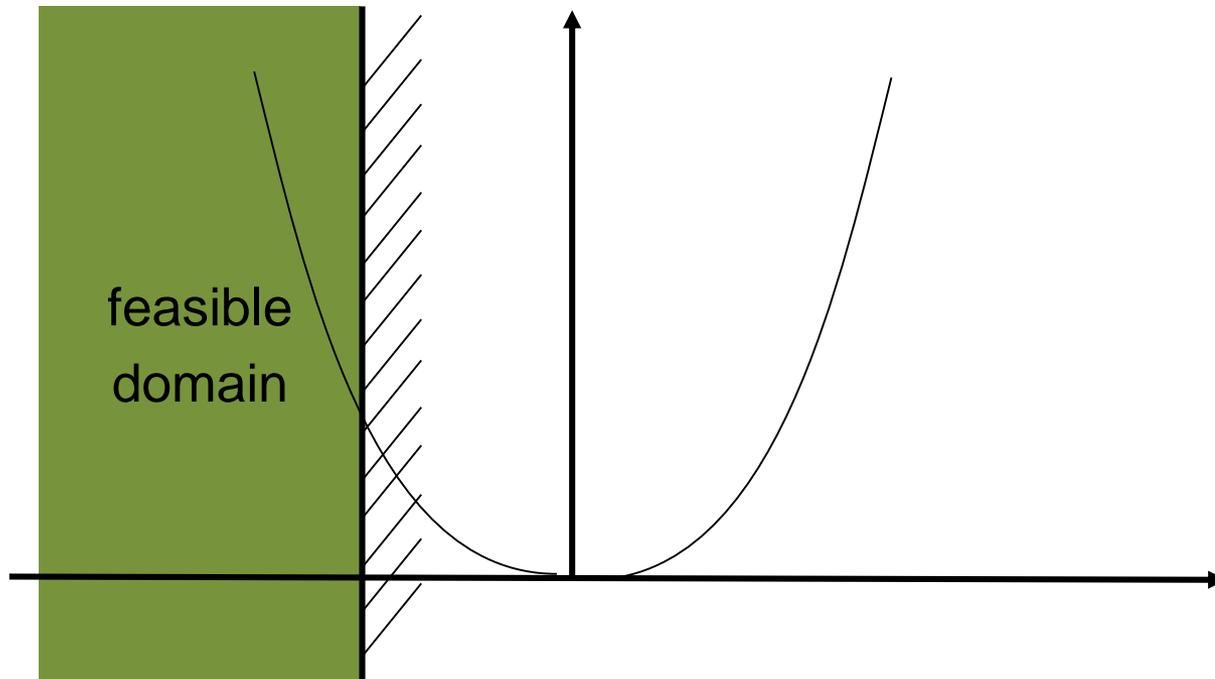
Constrained optimization

- Equality constraints: $\inf \{f(x) \mid x \in \mathbb{R}^n, g_k(x) = 0, 1 \leq k \leq p\}$
- Inequality constraints: $\inf \{f(x) \mid x \in \mathbb{R}^n, g_k(x) \leq 0, 1 \leq k \leq p\}$

where always $g_k: \mathbb{R}^n \rightarrow \mathbb{R}$

Example of a Constraint

$$\min_{x \in \mathbb{R}} f(x) = x^2 \text{ such that } x \leq -1$$



Analytical Functions

Example: 1-D

$$f_1(x) = a(x - x_0)^2 + b$$

where $x, x_0, b \in \mathbb{R}, a \in \mathbb{R}$

Generalization:

convex quadratic function

$$f_2(x) = (x - x_0)^T A (x - x_0) + b$$

where $x, x_0, b \in \mathbb{R}^n, A \in \mathbb{R}^{\{n \times n\}}$
and A symmetric positive definite (SPD)

Exercise:

What is the minimum of $f_2(x)$?

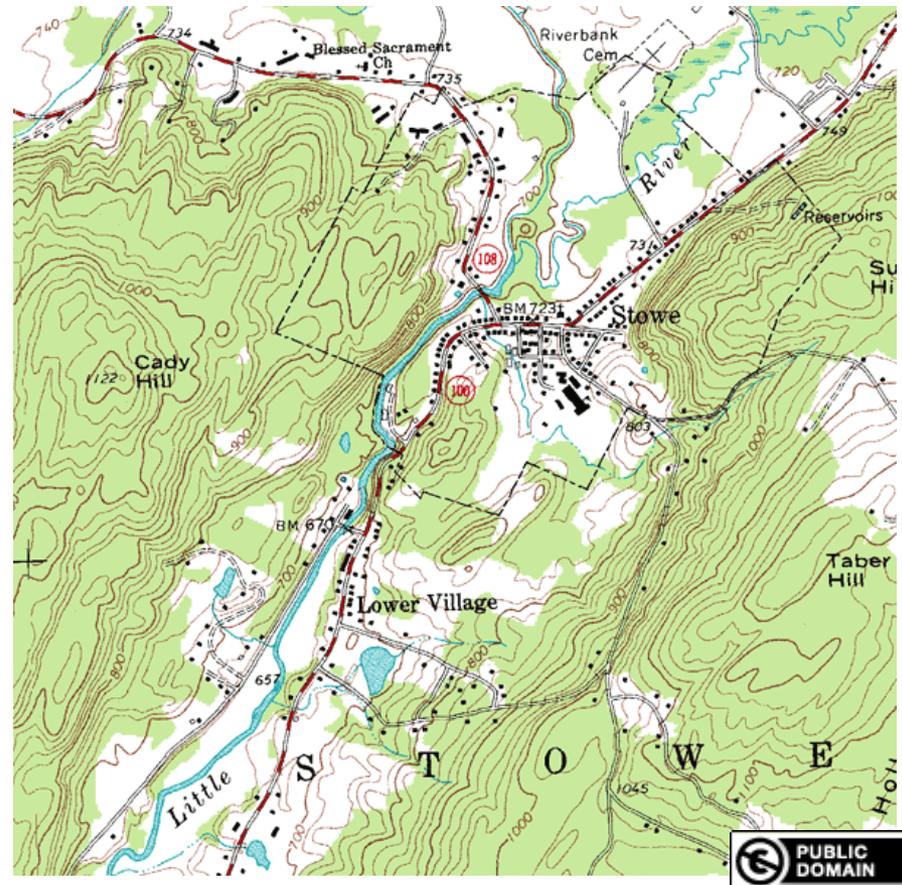
Levels Sets of Convex Quadratic Functions

Continuation of exercise:
What are the level sets of f_2 ?

Reminder: level sets of a function

$$L_c = \{x \in \mathbb{R}^n \mid f(x) = c\}$$

(similar to topography lines /
level sets on a map)



Continuation of exercise:

What are the level sets of f_2 ?

- Probably too complicated in general, thus an example here
- Consider $A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$, $b = 0$, $n = 2$
 - a) Compute $f_2(x)$.
 - b) Plot the level sets of $f_2(x)$.
 - c) Optional: More generally, for $n = 2$, if A is SPD with eigenvalues $\lambda_1 = 9$ and $\lambda_2 = 1$, what are the level sets of $f_2(x)$?

Example Problems

Data Fitting – Data Calibration

Objective

- Given a sequence of data points $(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}, i = 1, \dots, N$, find a model " $y = f(\mathbf{x})$ " that explains the data
experimental measurements in biology, chemistry, ...
- In general, choice of a parametric model or family of functions $(f_\theta)_{\theta \in \mathbb{R}^n}$
use of expertise for choosing model or simple models only affordable (linear, quadratic)
- Try to find the parameter $\theta \in \mathbb{R}^n$ fitting best to the data

Fitting best to the data

Minimize the quadratic error:

$$\min_{\theta \in \mathbb{R}^n} \sum_{i=1}^N |f_\theta(\mathbf{x}_i) - y_i|^2$$

Supervised Learning:

Predict $y \in \mathcal{Y}$ from $x \in \mathcal{X}$, given a set of observations (examples)

$$\{y_i, \mathbf{x}_i\}_{i=1, \dots, N}$$

(Simple) Linear regression

Given a set of data: $\{y_i, \underbrace{x_i^1, \dots, x_i^p}_{\mathbf{x}_i^T}\}_{i=1 \dots N}$

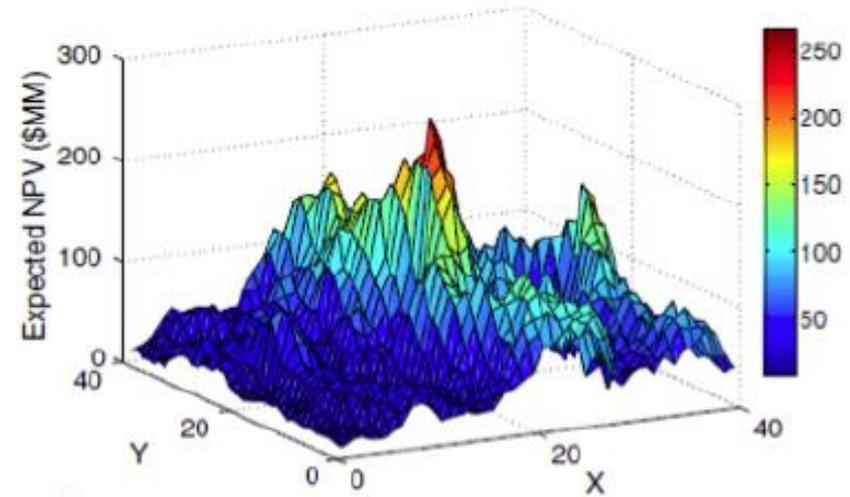
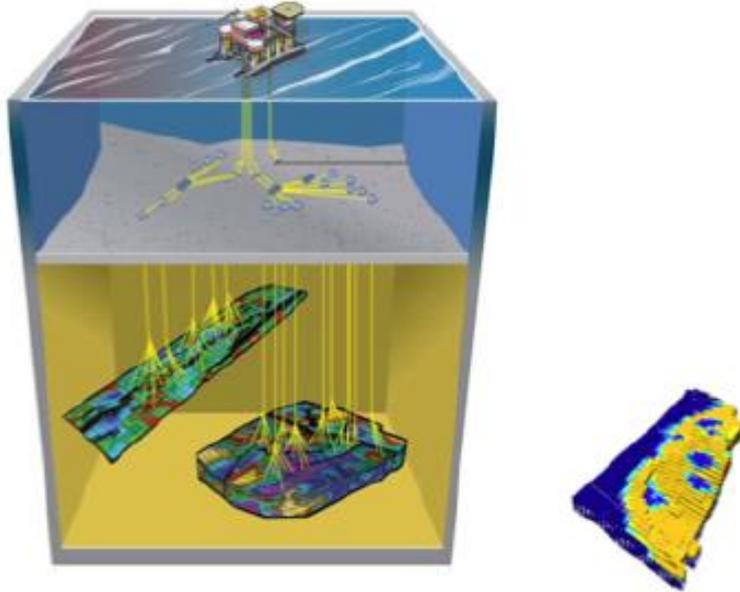
$$\min_{\mathbf{w} \in \mathbb{R}^p, \beta \in \mathbb{R}} \sum_{i=1}^N \underbrace{|\mathbf{w}^T \mathbf{x}_i + \beta - y_i|^2}_{\|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|^2}$$

$$\tilde{\mathbf{X}} \in \mathbb{R}^{N \times (p+1)}, \tilde{\mathbf{w}} \in \mathbb{R}^{p+1}$$

same as data fitting with linear model, i.e. $f_{(\mathbf{w}, \beta)}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \beta$,
 $\theta \in \mathbb{R}^{p+1}$

A Real-World Problem in Petroleum Engineering

Well Placement Problem



Onwunalu & Durlofsky (2010)

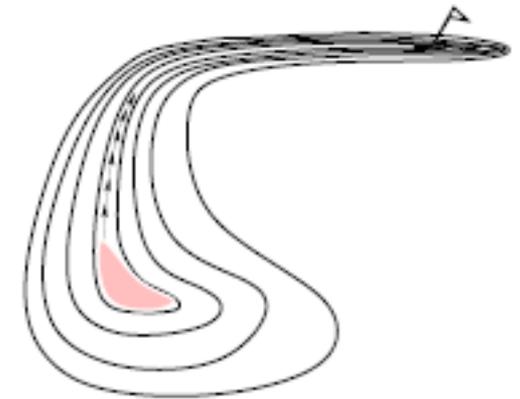


Fluid flow simulation

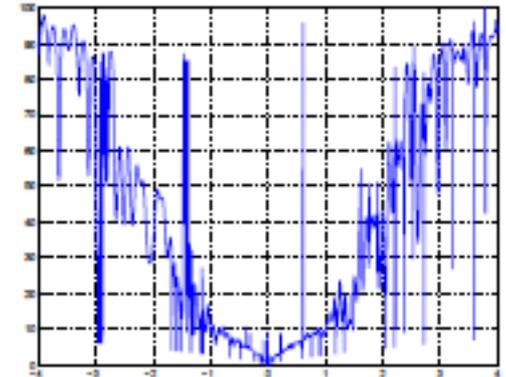
Function Difficulties

What Makes a Function Difficult to Solve?

- dimensionality
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function



a narrow ridge



cut from 3D example,
solvable with an
evolution strategy

Curse of Dimensionality

- The term *Curse of dimensionality* (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding extra dimensions to a (mathematical) space.
- Example: Consider placing 100 points onto a real interval, say $[0,1]$. To get **similar coverage**, in terms of distance between adjacent points, of the 10-dimensional space $[0,1]^{10}$ would require $100^{10} = 10^{20}$ points. The original 100 points appear now as isolated points in a vast empty space.
- Consequently, a **search policy** (e.g. exhaustive search) that is valuable in small dimensions **might be useless** in moderate or large dimensional search spaces.

Separable Problems

Definition (Separable Problem)

A function f is separable if

$$\operatorname{argmin}_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left(\operatorname{argmin}_{x_1} f(x_1, \dots), \dots, \operatorname{argmin}_{x_n} f(\dots, x_n) \right)$$

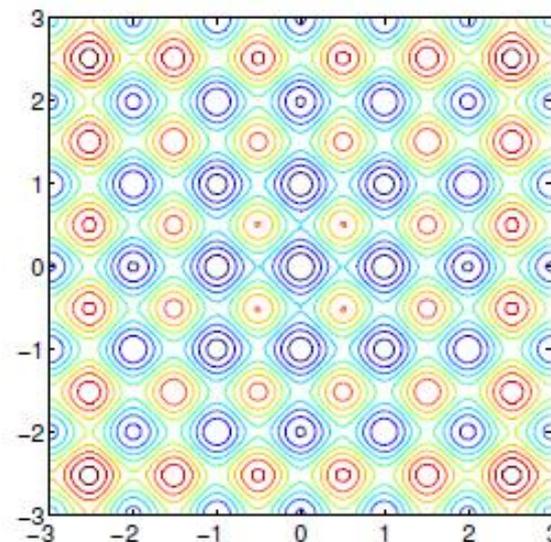
\Rightarrow it follows that f can be optimized in a sequence of n independent 1-D optimization processes

Example:

Additively decomposable functions

$$f(x_1, \dots, x_n) = \sum_{i=1}^n f_i(x_i)$$

Rastrigin function



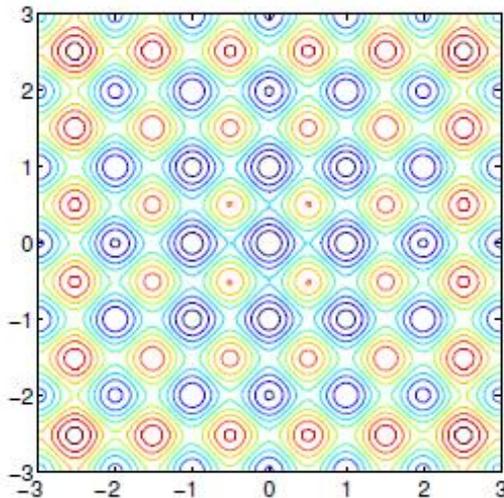
Non-Separable Problems

Building a non-separable problem from a separable one [1,2]

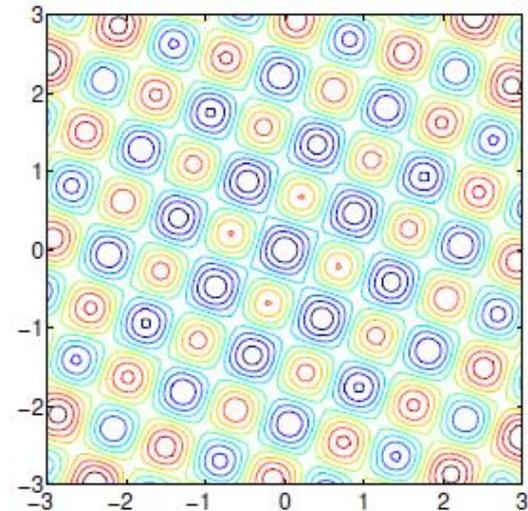
Rotating the coordinate system

- $f: \mathbf{x} \mapsto f(\mathbf{x})$ separable
- $f: \mathbf{x} \mapsto f(R\mathbf{x})$ non-separable

R rotation matrix



R
→



[1] N. Hansen, A. Ostermeier, A. Gawelczyk (1995). "On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation". Sixth ICGA, pp. 57-64, Morgan Kaufmann

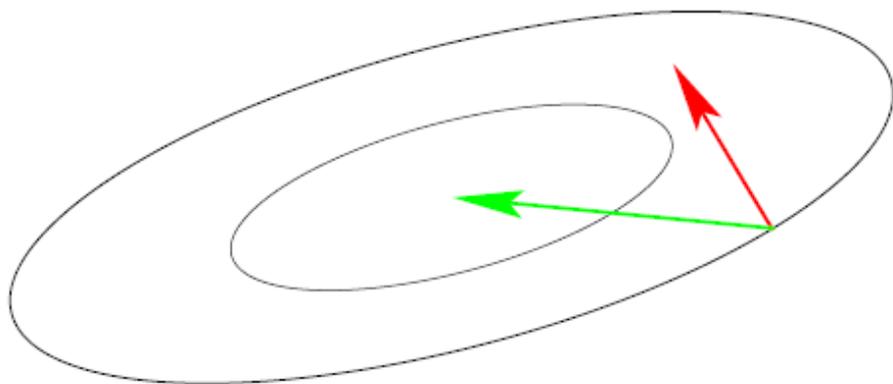
[2] R. Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions; A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278

Ill-Conditioned Problems: Curvature of Level Sets

Consider the convex-quadratic function

$$f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T H(\mathbf{x} - \mathbf{x}^*) = \frac{1}{2} \sum_i h_{i,i} x_i^2 + \frac{1}{2} \sum_{i,j} h_{i,j} x_i x_j$$

H is Hessian matrix of f and symmetric positive definite



gradient direction $-f'(\mathbf{x})^T$

Newton direction $-H^{-1}f'(\mathbf{x})^T$

*Ill-conditioning means **squeezed level sets** (high curvature).*

Condition number of SPD matrix A = ratio between largest and smallest eigenvalue

Condition number equals nine here (kind of well-conditioned). Condition numbers up to 10^{10} are not unusual in real-world problems.

Mathematical Tools to Characterize Optima

Different Notions of Optimum

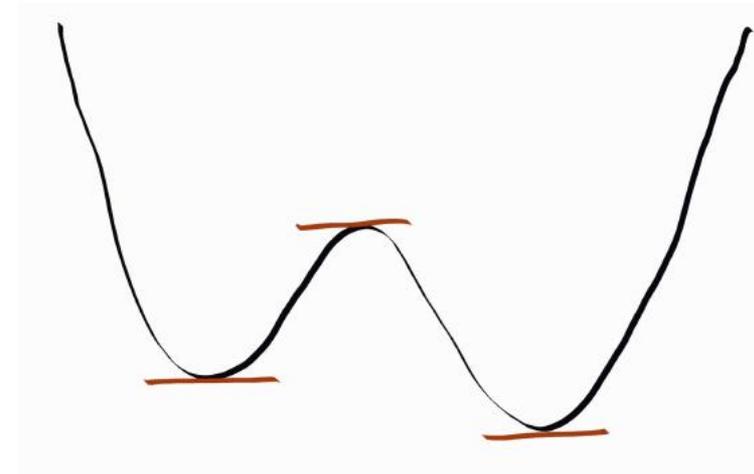
Unconstrained case

- local vs. global
 - local minimum \mathbf{x}^* : \exists a neighborhood V of \mathbf{x}^* such that
$$\forall \mathbf{x} \in V: f(\mathbf{x}) \geq f(\mathbf{x}^*)$$
 - global minimum: $\forall \mathbf{x} \in \Omega: f(\mathbf{x}) \geq f(\mathbf{x}^*)$
- strict local minimum if the inequality is strict

Mathematical Characterization of Optima

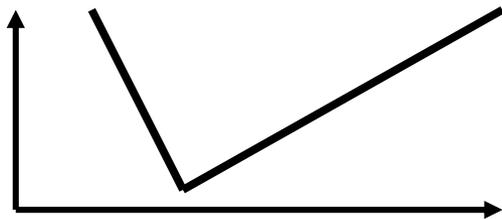
Objective: Derive general characterization of optima

Example: if $f: \mathbb{R} \rightarrow \mathbb{R}$ differentiable,
 $f'(x) = 0$ at optimal points



- generalization to $f: \mathbb{R}^n \rightarrow \mathbb{R}$?
- generalization to constrained problems?

Remark: notion of optimum independent of notion of differentiability

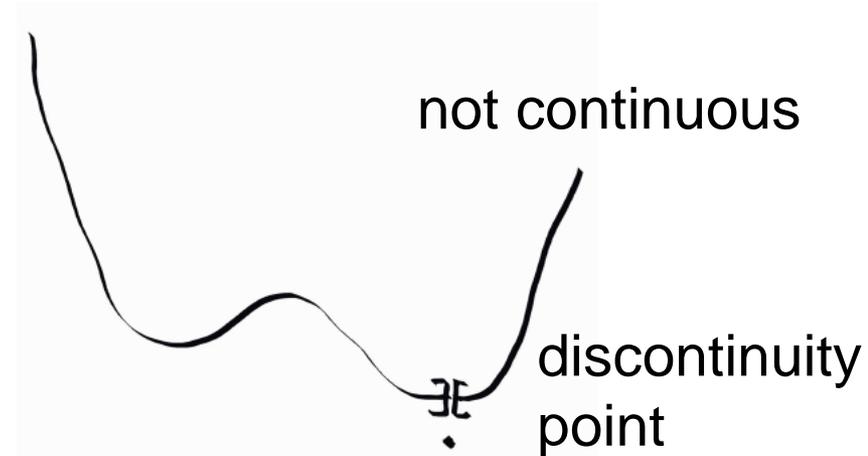
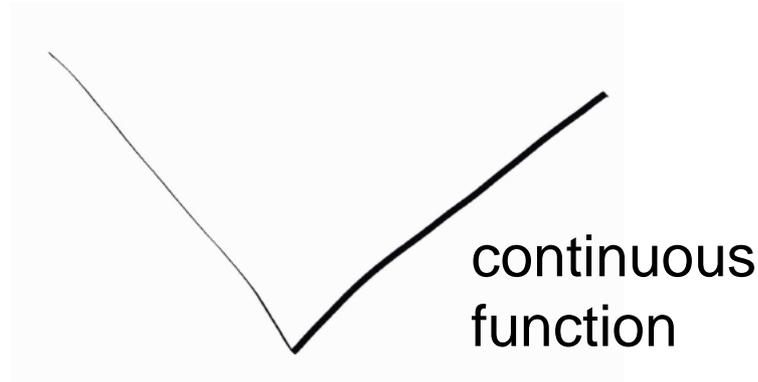


optima of such function can be easily
approached by certain type of methods

Reminder: Continuity of a Function

$f: (V, \| \cdot \|_V) \rightarrow (W, \| \cdot \|_W)$ is continuous in $x \in V$ if

$\forall \epsilon > 0, \exists \eta > 0$ such that $\forall y \in V: \|x - y\|_V \leq \eta; \|f(x) - f(y)\|_W \leq \epsilon$



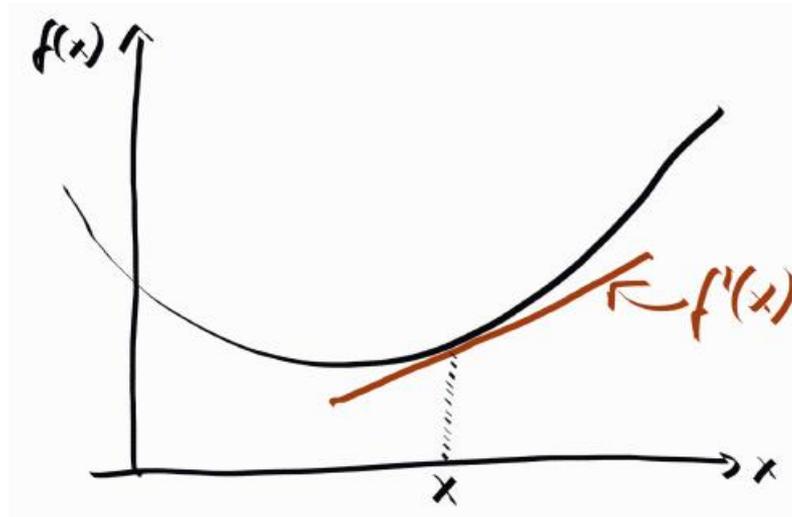
Reminder: Differentiability in 1D (n=1)

$f: \mathbb{R} \rightarrow \mathbb{R}$ is differentiable in $x \in \mathbb{R}$ if

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \text{ exists, } h \in \mathbb{R}$$

Notation:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$



The derivative corresponds to the slope of the tangent in x .

Reminder: Differentiability in 1D ($n=1$)

Taylor Formula (Order 1)

If f is differentiable in x then

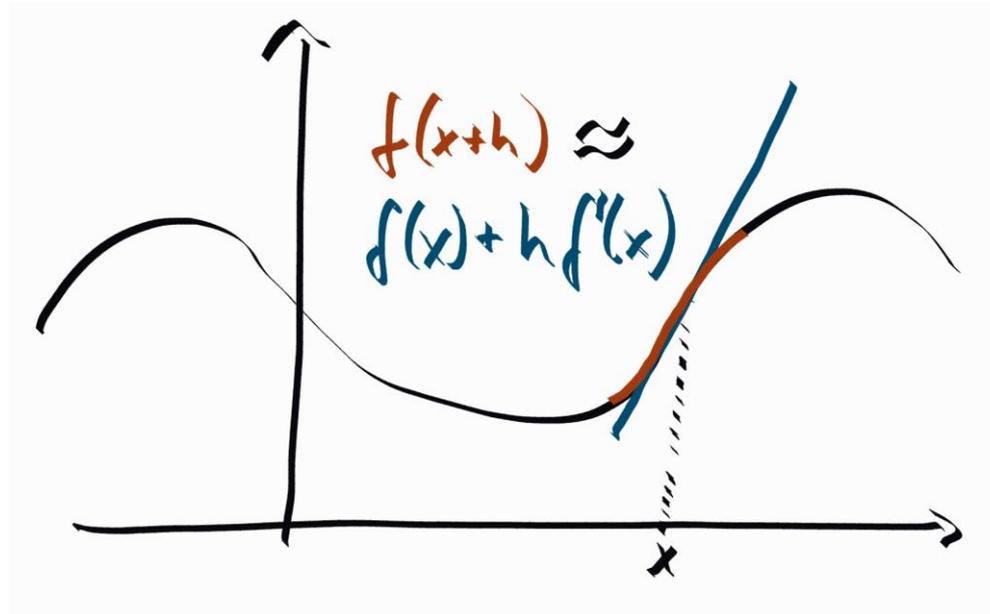
$$f(x + h) = f(x) + f'(x)h + o(\|h\|)$$

i.e. for h small enough, $h \mapsto f(x + h)$ is approximated by $h \mapsto f(x) + f'(x)h$

$h \mapsto f(x) + f'(x)h$ is called a **first order approximation** of $f(x + h)$

Reminder: Differentiability in 1D ($n=1$)

Geometrically:



The notion of derivative of a function defined on \mathbb{R}^n is generalized via this idea of a linear approximation of $f(x + h)$ for h small enough.

Gradient Definition Via Partial Derivatives

- In $(\mathbb{R}^n, \|\cdot\|_2)$ where $\|x\|_2 = \sqrt{\langle x, x \rangle}$ is the Euclidean norm deriving from the scalar product $\langle x, y \rangle = x^T y$

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

- Reminder: partial derivative in x_0

$$f_i: y \rightarrow f(x_0^1, \dots, x_0^{i-1}, y, x_0^{i+1}, \dots, x_0^n)$$

$$\frac{\partial f}{\partial x_i}(x_0) = f_i'(x_0)$$

Exercise:

Compute the gradients of

a) $f(x) = x_1$ with $x \in \mathbb{R}^n$

b) $f(x) = a^T x$ with $a, x \in \mathbb{R}^n$

c) $f(x) = x^T x (= \|x\|^2)$ with $x \in \mathbb{R}^n$

Exercise:

Compute the gradients of

- a) $f(x) = x_1$ with $x \in \mathbb{R}^n$
- b) $f(x) = a^T x$ with $a, x \in \mathbb{R}^n$
- c) $f(x) = x^T x (= \|x\|^2)$ with $x \in \mathbb{R}^n$

Some more examples:

- in \mathbb{R}^n , if $f(x) = x^T A x$, then $\nabla f(x) = (A + A^T)x$
- in \mathbb{R} , $\nabla f(x) = f'(x)$

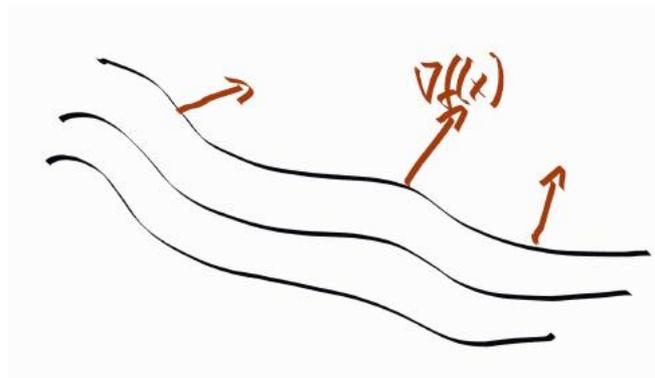
Gradient: Geometrical Interpretation

Exercise:

Let $L_c = \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) = c\}$ be again a level set of a function $f(\mathbf{x})$.
Let $\mathbf{x}_0 \in L_c \neq \emptyset$.

Plot the level sets for $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$ and $f(\mathbf{x}) = \|\mathbf{x}\|^2$, compute the gradient in a chosen point \mathbf{x}_0 and observe that $\nabla f(\mathbf{x}_0)$ is **orthogonal** to the level set in \mathbf{x}_0 .

More generally, the gradient of a differentiable function is orthogonal to its level sets.



Taylor Formula – Order One

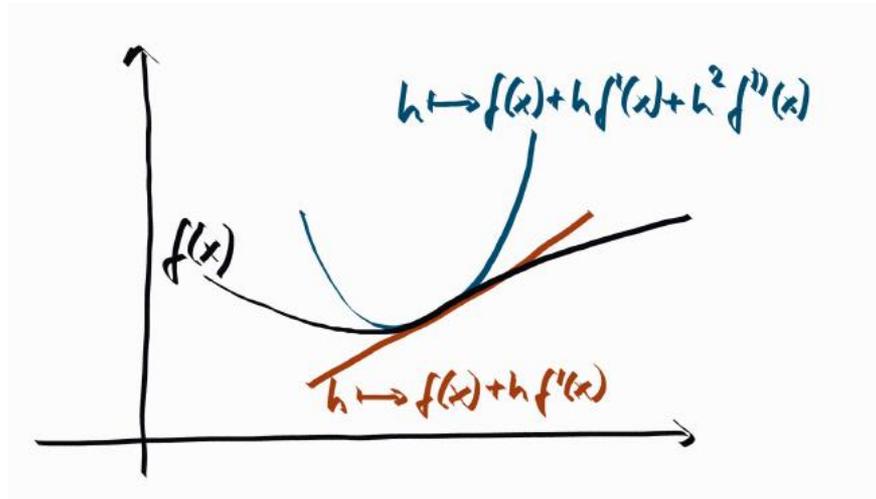
$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + (\nabla f(\mathbf{x}))^T \mathbf{h} + o(\|\mathbf{h}\|)$$

Reminder: Second Order Differentiability in 1D

- Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be a differentiable function and let $f': x \rightarrow f'(x)$ be its derivative.
- If f' is differentiable in x , then we denote its derivative as $f''(x)$
- $f''(x)$ is called the *second order derivative* of f .

Taylor Formula: Second Order Derivative

- If $f: \mathbb{R} \rightarrow \mathbb{R}$ is two times differentiable then
$$f(x+h) = f(x) + f'(x)h + f''(x)h^2 + o(\|h\|^2)$$
i.e. for h small enough, $h \rightarrow f(x) + hf'(x) + h^2f''(x)$ approximates $h + f(x+h)$
- $h \rightarrow f(x) + hf'(x) + h^2f''(x)$ is a quadratic approximation (or order 2) of f in a neighborhood of x



- The second derivative of $f: \mathbb{R} \rightarrow \mathbb{R}$ generalizes naturally to larger dimension.

Hessian Matrix

In $(\mathbb{R}^n, \langle x, y \rangle = x^T y)$, $\nabla^2 f(x)$ is represented by a symmetric matrix called the Hessian matrix. It can be computed as

$$\nabla^2(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Exercise on Hessian Matrix

Exercise:

Let $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x}$, $\mathbf{x} \in \mathbb{R}^n$, and $A \in \mathbb{R}^{n \times n}$ symmetric.

Compute the Hessian matrix of f .

If it is too complex, consider $f: \begin{cases} \mathbb{R}^2 \rightarrow \mathbb{R} \\ \mathbf{x} \rightarrow \frac{1}{2} \mathbf{x}^T A \mathbf{x} \end{cases}$ with $A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$

Taylor Formula – Order Two

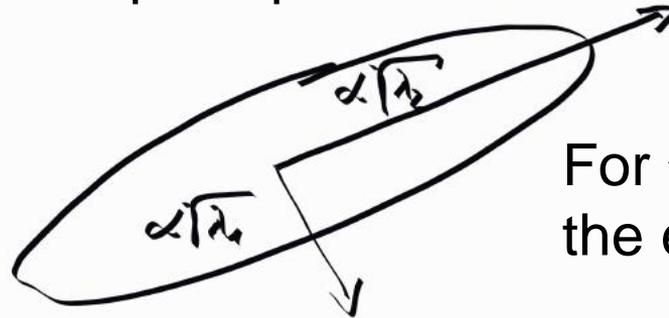
$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + (\nabla f(\mathbf{x}))^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T (\nabla^2 f(\mathbf{x})) \mathbf{h} + o(\|\mathbf{h}\|^2)$$

Back to Ill-Conditioned Problems

We have seen that for a convex quadratic function

$$f(x) = \frac{1}{2}(x - x_0)^T A(x - x_0) + b \text{ of } x \in \mathbb{R}^n, A \in \mathbb{R}^{n \times n}, A \text{ SPD}, b \in \mathbb{R}^n:$$

- 1) The level sets are ellipsoids. The eigenvalues of A determine the lengths of the principle axes of the ellipsoid.



For $n = 2$, let λ_1, λ_2 be the eigenvalues of A .

- 2) The Hessian matrix of f equals to A .

Ill-conditioned convex quadratic problems are problems with large ratio between largest and smallest eigenvalue of A which means large ratio between longest and shortest axis of ellipsoid.

This corresponds to having an ill-conditioned Hessian matrix.

Gradient Direction Vs. Newton Direction

Gradient direction: $\nabla f(\mathbf{x})$

Newton direction: $(H(\mathbf{x}))^{-1} \cdot \nabla f(\mathbf{x})$

with $H(\mathbf{x}) = \nabla^2 f(\mathbf{x})$ being the Hessian at \mathbf{x}

Exercise:

Let again $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x}$, $\mathbf{x} \in \mathbb{R}^2$, $A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$.

Plot the gradient and Newton direction of f in a point $\mathbf{x} \in \mathbb{R}^n$ of your choice (which should not be on a coordinate axis) into the same plot with the level sets, we created before.

Conclusions

I hope it became clear...

...what are the difficulties to cope with when solving numerical optimization problems

in particular dimensionality, non-separability and ill-conditioning

...what are **gradient** and **Hessian**

...what is the difference between **gradient** and **Newton direction**