

Introduction to Optimization

Lecture 3: Introduction to Continuous Optimization

September 28, 2018

TC2 - Optimisation

Université Paris-Saclay



Dimo Brockhoff
Inria Saclay – Ile-de-France

Course Overview

| | | |
|---|--|--|
| 1 | Mon, 17.9.2018 Thu, 20.9.2018 | Monday's lecture: introduction, example problems, problem types groups defined via wiki everybody went (actively!) through the Getting Started part of github.com/numbbo/coco |
| 2 | Fri, 21.9.2018 | lecture "Benchmarking", final adjustments of groups everybody can run and postprocess the example experiment (~1h for final questions/help during the lecture) |
| 3 | Fri, 28.9.2018 | lecture "Introduction to Continuous Optimization" |
| 4 | Fri, 5.10.2018 | lecture "Gradient-Based Algorithms" |
| 5 | Fri, 12.10.2018 | lecture "Stochastic Algorithms and DFO" |
| 6 | Fri, 19.10.2018 | lecture "Discrete Optimization I: graphs, greedy algos, dyn. progr." deadline for submitting data sets |
| | Wed, 24.10.2018 | deadline for paper submission |
| 7 | Fri, 26.10.2018 | final lecture "Discrete Optimization II: dyn. progr., B&B, heuristics" |
| | 29.10.-2.11.2018 | vacation aka learning for the exams |
| | Thu, 8.11.2018 / Fri, 9.11.2018 | oral presentations (individual time slots) |
| | Fri, 16.11.2018 | written exam |

**All deadlines:
23:59pm Paris time**

Details on Continuous Optimization Lectures

Introduction to Continuous Optimization

- examples (from ML / black-box problems)
- typical difficulties in optimization

Mathematical Tools to Characterize Optima

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

- constraint optimization

Gradient-based Algorithms

- quasi-Newton method (BFGS)
 - [DFO trust-region method]
-

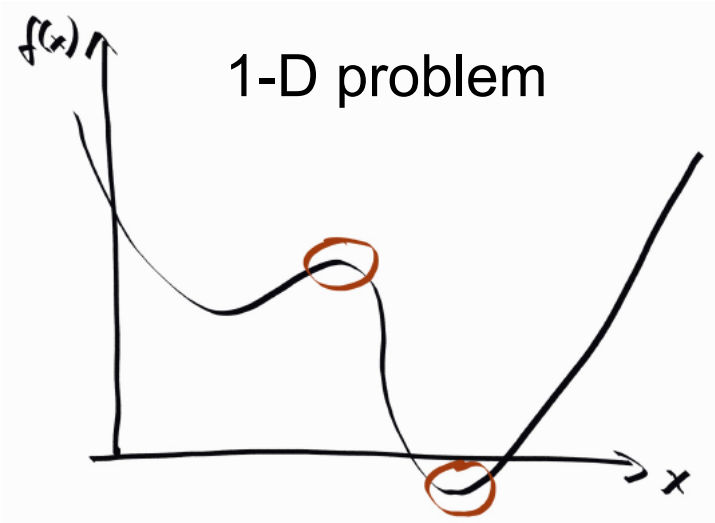
Learning in Optimization / Stochastic Optimization

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

method strongly related to ML / new promising research area
interesting open questions

Continuous Optimization

- Optimize $f: \begin{cases} \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R} \\ x = (x_1, \dots, x_n) \rightarrow f(x_1, \dots, x_n) \\ \quad \quad \quad \swarrow \\ \quad \quad \quad \mathbb{R} \end{cases}$ *unconstrained* optimization
- Search space is continuous, i.e. composed of real vectors $x \in \mathbb{R}^n$
- $n = \begin{cases} \text{dimension of the problem} \\ \text{dimension of the search space } \mathbb{R}^n \text{ (as vector space)} \end{cases}$



Reminder: Different Notions of Optimum

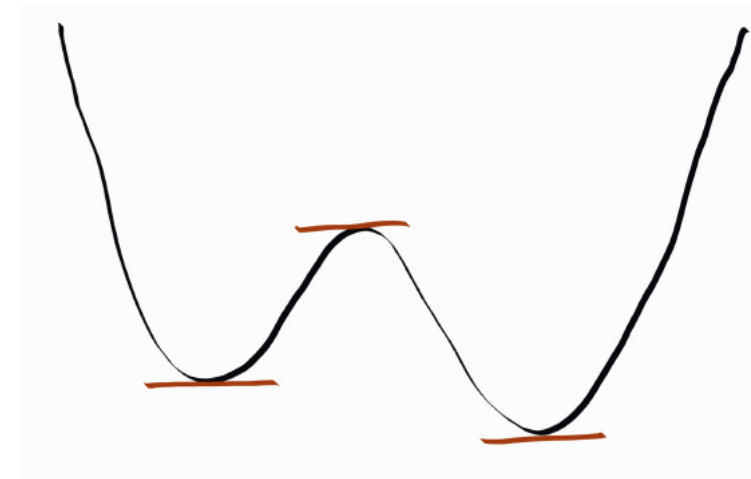
Unconstrained case

- 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

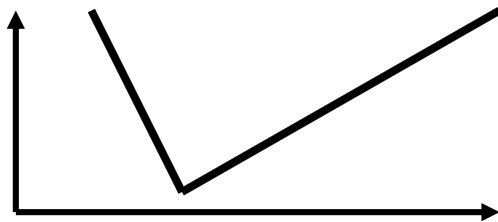
Objective: Derive general characterization of optima

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 derivability

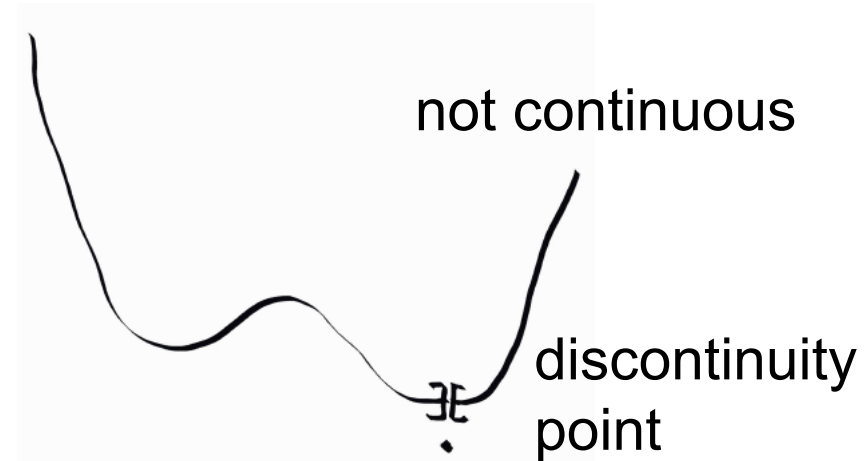
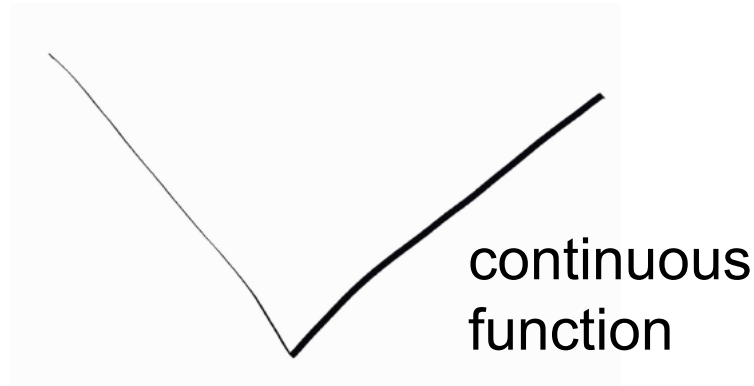


optima of such function can be easily
approached by certain type of methods

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 $\forall y \in V: \|x - y\|_V \leq \eta; \|f(x) - f(y)\|_W \leq \epsilon$



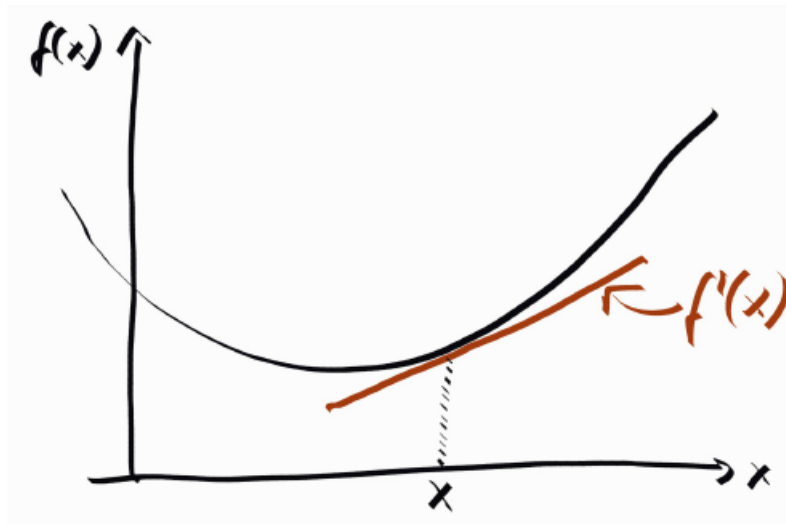
Reminder: Differentiability in 1D ($n=1$)

$f: \mathbb{R} \rightarrow \mathbb{R}$ 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

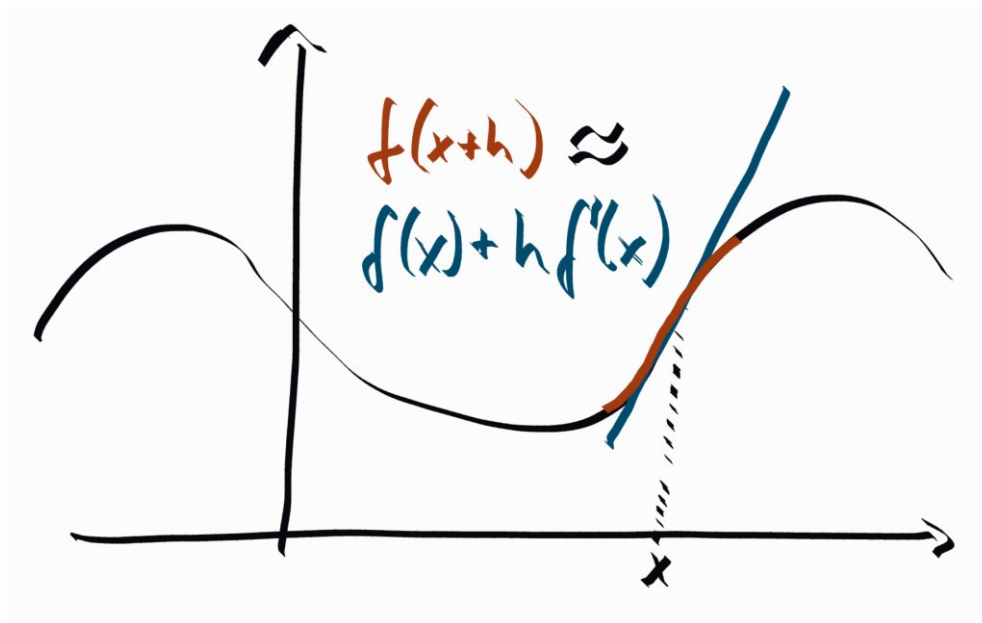
$$f(x + h) = f(x) + f'(x)h + o(\|h\|)$$

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.

How to generalize this to arbitrary dimension?

Gradient Definition Via Partial Derivatives

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

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

- 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: Gradients

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

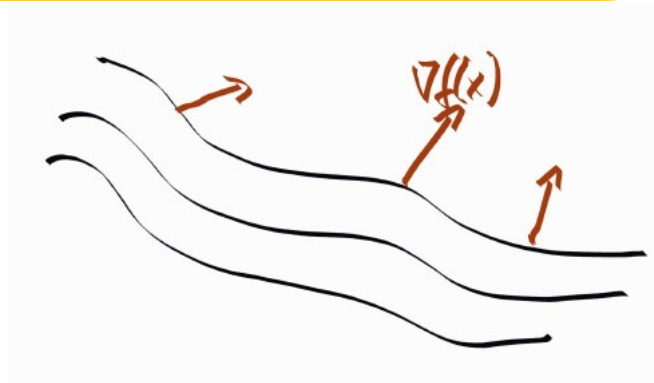
Exercise:

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 $\mathbf{x}_0 \in L_c \neq \emptyset$.

Compute the level sets for $f_1(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$ and $f_2(\mathbf{x}) = \|\mathbf{x}\|^2$ and 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 .

Again: if this seems too difficult, do it for two variables (and a concrete $\mathbf{a} \in \mathbb{R}^2$) and draw the level sets and the gradients.

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



Taylor Formula – Order One

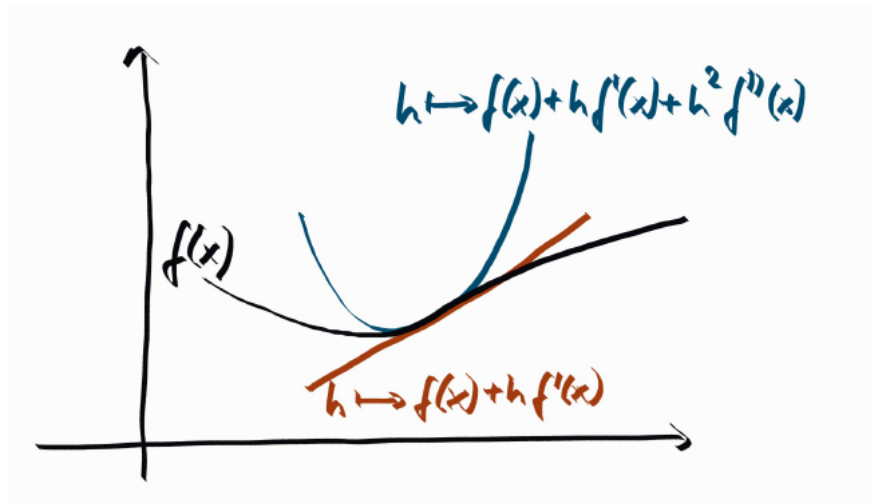
$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + (\nabla f(\mathbf{x}))^T \mathbf{h} + o(\|\mathbf{h}\|)$$

Reminder: Second Order Derivability 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 $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x}$, $\mathbf{x} \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times n}$.

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

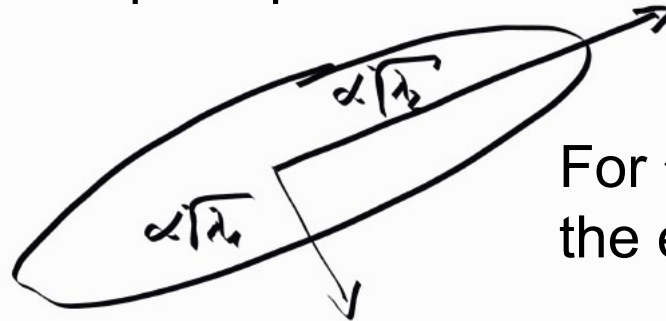
$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + (\nabla f(\mathbf{x}))^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T (\nabla^2 f(\mathbf{x})) \mathbf{h} + o(\|\mathbf{h}\|^2)$$

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.

Optimality Conditions for Unconstrained Problems

Optimality Conditions: First Order Necessary Cond.

For 1-dimensional optimization problems $f: \mathbb{R} \rightarrow \mathbb{R}$

Assume f is differentiable

- \mathbf{x}^* is a local optimum $\Rightarrow f'(\mathbf{x}^*) = 0$

not a sufficient condition: consider $f(x) = x^3$

proof via Taylor formula: $f(\mathbf{x}^ + \mathbf{h}) = f(\mathbf{x}^*) + f'(\mathbf{x}^*)\mathbf{h} + o(\|\mathbf{h}\|)$*

- points \mathbf{y} such that $f'(\mathbf{y}) = 0$ are called **critical** or **stationary** points

Generalization to n -dimensional functions

If $f: U \subset \mathbb{R}^n \mapsto \mathbb{R}$ is differentiable

- necessary condition: If \mathbf{x}^* is a local optimum of f , then $\nabla f(\mathbf{x}^*) = 0$

proof via Taylor formula

Second Order Necessary and Sufficient Opt. Cond.

If f is twice continuously differentiable

- **Necessary condition:** if \mathbf{x}^* is a local minimum, then $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive semi-definite

proof via Taylor formula at order 2

- **Sufficient condition:** if $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite, then \mathbf{x}^* is a strict local minimum

Proof of Sufficient Condition:

- Let $\lambda > 0$ be the smallest eigenvalue of $\nabla^2 f(\mathbf{x}^*)$, using a second order Taylor expansion, we have for all \mathbf{h} :

- $$f(\mathbf{x}^* + \mathbf{h}) - f(\mathbf{x}^*) = \nabla f(\mathbf{x}^*)^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}^*) \mathbf{h} + o(\|\mathbf{h}\|^2)$$
$$> \frac{\lambda}{2} \|\mathbf{h}\|^2 + o(\|\mathbf{h}\|^2) = \left(\frac{\lambda}{2} + \frac{o(\|\mathbf{h}\|^2)}{\|\mathbf{h}\|^2} \right) \|\mathbf{h}\|^2$$

Convex Functions

Let U be a convex open set of \mathbb{R}^n and $f: U \rightarrow \mathbb{R}$. The function f is said to be **convex** if for all $\mathbf{x}, \mathbf{y} \in U$ and for all $t \in [0,1]$

$$f((1-t)\mathbf{x} + t\mathbf{y}) \leq (1-t)f(\mathbf{x}) + tf(\mathbf{y})$$

Theorem

If f is differentiable, then f is convex if and only if for all \mathbf{x}, \mathbf{y}

$$f(\mathbf{y}) - f(\mathbf{x}) \geq (\nabla f(\mathbf{x}))^T (\mathbf{y} - \mathbf{x})$$

if $n = 1$, the curve is on top of the tangent

If f is twice continuously differentiable, then f is convex if and only if $\nabla^2 f(\mathbf{x})$ is positive semi-definite for all \mathbf{x} .

Convex Functions: Why Convexity?

Examples of Convex Functions:

- $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b$
- $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{a}^T \mathbf{x} + b$, A symmetric positive definite
- the negative of the entropy function (i. e. $f(\mathbf{x}) = -\sum_{i=1}^n x_i \ln(x_i)$)

Exercise:

Let $f: U \rightarrow \mathbb{R}$ be a convex and differentiable function on a convex open U .

Show that if $\nabla f(\mathbf{x}^*) = 0$, then \mathbf{x}^* is a global minimum of f

Why is convexity an important concept?