

Advanced Optimization

Lecture 3: Continuous Optimization II

October 27, 2021

CentraleSupélec / ESSEC Business School

dimo.brockhoff@inria.fr



Dimo Brockhoff
Inria Saclay – Ile-de-France



Course Overview

		Topic
Wed, 13.10.2021	PM	Introduction, examples of problems, problem types
Wed, 20.10.2021	PM	Continuous (unconstrained) optimization: convexity, gradients, Hessian, ... [technical test Evalmee]
Wed, 27.10.2021	PM	Continuous optimization II: gradient descent, Newton direction, quasi-Newton (BFGS) [1 st mini-exam] Constrained optimization: Lagrangian, optimality conditions
Wed, 03.11.2021	PM	Linear programming: duality, maxflow/mincut, simplex algo
Wed, 10.11.2021	PM	Gradient-based and derivative-free stochastic algorithms: SGD and CMA-ES
Wed, 17.11.2021	PM	Other blackbox optimizers: Nelder-Mead, Bayesian optimization [2 nd mini-exam]
Wed, 24.11.2021	PM	Benchmarking solvers: runtime distributions, performance profiles
Tue, 30.11.2021	23:59	Deadline open source project (PDF sent by email)
Wed, 01.12.2021	PM	Discrete optimization: branch and bound, branch and cut, k-means clustering
Wed, 15.12.2021	PM	Exam

Overlap with Other Classes

- I was not aware that gradients/Hessian/optimalty conditions/... were already taught in some M1 course(s)
- I will try my best to adapt
- What I propose for today and next week:
 - lecture with basics needed for later on (for those who do not have the basics)
 - additional exercise (for those who are bored otherwise 😊)

Don't forget the group project:

<https://docs.google.com/spreadsheets/d/1WV8yfl1T0rYqtdoPYzOu7ORVx9qKC1kwvOSFE6MVaX4/edit?usp=sharing>

Deadline for the report: November 30, 2021

Advanced Exercise

- Shall be done quietly (to not disturb the class)
- I will be there during the normal exercise times to help
- Jupyter notebook available on Edunao

Idea:

- Self-teaching of concepts, discussed later in the class via learning-by-doing
- practical python/programming exercise
- Part I: Invariances (today)
- Part II: Benchmarking (next week)

back to lecture

Details on Continuous Optimization Lectures

Introduction to Continuous Optimization

- examples and 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
 - Lagrangian, optimality conditions

Gradient-based Algorithms

- gradient descent
- quasi-Newton method (BFGS) and invariances

Linear programming, duality

Learning in Optimization / Optimization in Machine Learning

- Stochastic gradient descent (SGD) + Adam
- CMA-ES (adaptive algorithms / Information Geometry)
- Other derivative-free algorithms: Nelder-Mead, Bayesian opt.

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 $f(\mathbf{x}) = \mathbf{x}^3$

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:

- $f(\mathbf{x}) = a^T \mathbf{x} + b$
- $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + a^T \mathbf{x} + b$, 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?

Mini-Exam in Evalmee

Constrained Optimization

Equality Constraint

Objective:

Generalize the necessary condition of $\nabla f(x) = 0$ at the optima of f
when f is in \mathcal{C}^1 , i.e. is differentiable and its differential is continuous

Theorem:

Be U an open set of $(E, || \cdot ||)$, and $f: U \rightarrow \mathbb{R}$, $g: U \rightarrow \mathbb{R}$ in \mathcal{C}^1 .

Let $a \in E$ satisfy

$$\begin{cases} f(a) = \inf \{f(x) \mid x \in \mathbb{R}^n, g(x) = 0\} \\ g(a) = 0 \end{cases}$$

i.e. a is optimum of the problem

If $\nabla g(a) \neq 0$, then there exists a constant $\lambda \in \mathbb{R}$ called *Lagrange multiplier*, such that

$$\underbrace{\nabla f(a) + \lambda \nabla g(a)} = 0 \quad \text{Euler – Lagrange equation}$$

i.e. gradients of f and g in a are colinear

Geometrical Interpretation Using an Example

Exercise:

Consider the problem

$$\inf \{ f(x, y) \mid (x, y) \in \mathbb{R}^2, g(x, y) = 0 \}$$

$$f(x, y) = y - x^2 \quad g(x, y) = x^2 + y^2 - 1 = 0$$

- 1) Plot the level sets of f , plot $g = 0$
- 2) Compute ∇f and ∇g
- 3) Find the solutions with $\nabla f + \lambda \nabla g = 0$
equation solving with 3 unknowns (x, y, λ)
- 4) Plot the solutions of 3) on top of the level set graph of 1)

Interpretation of Euler-Lagrange Equation

Intuitive way to retrieve the Euler-Lagrange equation:

- In a local minimum a of a constrained problem, the hypersurfaces (or level sets) $f = f(a)$ and $g = 0$ are necessarily tangent (otherwise we could decrease f by moving along $g = 0$).
- Since the gradients $\nabla f(a)$ and $\nabla g(a)$ are orthogonal to the level sets $f = f(a)$ and $g = 0$, it follows that $\nabla f(a)$ and $\nabla g(a)$ are colinear.

Generalization to More than One Constraint

Theorem

- Assume $f: U \rightarrow \mathbb{R}$ and $g_k: U \rightarrow \mathbb{R}$ ($1 \leq k \leq p$) are \mathcal{C}^1 .
- Let a be such that
$$\begin{cases} f(a) = \inf \{f(x) \mid x \in \mathbb{R}^n, & g_k(x) = 0, & 1 \leq k \leq p\} \\ g_k(a) = 0 \text{ for all } 1 \leq k \leq p \end{cases}$$
- If $(\nabla g_k(a))_{1 \leq k \leq p}$ are linearly independent, then there exist p real constants $(\lambda_k)_{1 \leq k \leq p}$ such that

$$\nabla f(a) + \sum_{k=1}^p \lambda_k \nabla g_k(a) = 0$$



Lagrange multiplier

again: a does not need to be global but local minimum

The Lagrangian

- Define the Lagrangian on $\mathbb{R}^n \times \mathbb{R}^p$ as

$$\mathcal{L}(x, \{\lambda_k\}) = f(x) + \sum_{k=1}^p \lambda_k g_k(x)$$

- To find optimal solutions, we can solve the optimality system

$$\left\{ \begin{array}{l} \text{Find } (x, \{\lambda_k\}) \in \mathbb{R}^n \times \mathbb{R}^p \text{ such that } \nabla f(x) + \sum_{k=1}^p \lambda_k \nabla g_k(x) = 0 \\ g_k(x) = 0 \text{ for all } 1 \leq k \leq p \end{array} \right.$$

$$\Leftrightarrow \left\{ \begin{array}{l} \text{Find } (x, \{\lambda_k\}) \in \mathbb{R}^n \times \mathbb{R}^p \text{ such that } \nabla_x \mathcal{L}(x, \{\lambda_k\}) = 0 \\ \nabla_{\lambda_k} \mathcal{L}(x, \{\lambda_k\})(x) = 0 \text{ for all } 1 \leq k \leq p \end{array} \right.$$

Inequality Constraint: Definitions

Let $\mathcal{U} = \{x \in \mathbb{R}^n \mid g_k(x) = 0 \text{ (for } k \in E), g_k(x) \leq 0 \text{ (for } k \in I)\}$.

Definition:

The points in \mathbb{R}^n that satisfy the constraints are also called *feasible* points.

Definition:

Let $a \in \mathcal{U}$, we say that the constraint $g_k(x) \leq 0$ (for $k \in I$) is *active* in a if $g_k(a) = 0$.

Inequality Constraint: Karush-Kuhn-Tucker Theorem

Theorem (Karush-Kuhn-Tucker, KKT):

Let U be an open set of $(\mathbb{R}^n, || \cdot ||)$ and $f: U \rightarrow \mathbb{R}$, $g_k: U \rightarrow \mathbb{R}$, all \mathcal{C}^1

Furthermore, let $a \in U$ satisfy

$$\begin{cases} f(a) = \inf\{f(x) \mid x \in \mathbb{R}^n, g_k(x) = 0 \text{ (for } k \in E), g_k(x) \leq 0 \text{ (for } k \in I)\} \\ g_k(a) = 0 \text{ (for } k \in E) \\ g_k(a) \leq 0 \text{ (for } k \in I) \end{cases}$$

also works again for a being a local minimum

Let I_a^0 be the set of constraints that are active in a . Assume that $(\nabla g_k(a))_{k \in E \cup I_a^0}$ are linearly independent.

Then there exist $(\lambda_k)_{1 \leq k \leq p}$ that satisfy

$$\begin{cases} \nabla f(a) + \sum_{k=1}^p \lambda_k \nabla g_k(a) = 0 \\ g_k(a) = 0 \text{ (for } k \in E) \\ g_k(a) \leq 0 \text{ (for } k \in I) \\ \lambda_k \geq 0 \text{ (for } k \in I_a^0) \\ \lambda_k g_k(a) = 0 \text{ (for } k \in E \cup I) \end{cases}$$

Inequality Constraint: Karush-Kuhn-Tucker Theorem

Theorem (Karush-Kuhn-Tucker, KKT):

Let U be an open set of $(\mathbb{R}^n, || \cdot ||)$ and $f: U \rightarrow \mathbb{R}$, $g_k: U \rightarrow \mathbb{R}$, all \mathcal{C}^1

Furthermore, let $a \in U$ satisfy

$$\begin{cases} f(a) = \inf\{f(x) \mid x \in \mathbb{R}^n, g_k(x) = 0 \text{ (for } k \in E), g_k(x) \leq 0 \text{ (for } k \in I)\} \\ g_k(a) = 0 \text{ (for } k \in E) \\ g_k(a) \leq 0 \text{ (for } k \in I) \end{cases}$$

Let I_a^0 be the set of constraints that are active in a . Assume that $(\nabla g_k(a))_{k \in E \cup I_a^0}$ are linearly independent.

Then there exist $(\lambda_k)_{1 \leq k \leq p}$ that satisfy

$$\begin{cases} \nabla f(a) + \sum_{k=1}^p \lambda_k \nabla g_k(a) = 0 \\ g_k(a) = 0 \text{ (for } k \in E) \\ g_k(a) \leq 0 \text{ (for } k \in I) \\ \lambda_k \geq 0 \text{ (for } k \in I_a^0) \\ \lambda_k g_k(a) = 0 \text{ (for } k \in E \cup I) \end{cases}$$

either active constraint
or $\lambda_k = 0$

Descent Methods

Descent Methods

General principle

- ① choose an initial point x_0 , set $t = 0$
- ② while not happy
 - choose a descent direction $d_t \neq 0$
 - line search:
 - choose a step size $\sigma_t > 0$
 - set $x_{t+1} = x_t + \sigma_t d_t$
 - set $t = t + 1$

Remaining questions

- how to choose d_t ?
- how to choose σ_t ?

Gradient Descent

Rationale: $\mathbf{d}_t = -\nabla f(\mathbf{x}_t)$ is a descent direction
indeed for f differentiable

$$\begin{aligned} f(x - \sigma \nabla f(x)) &= f(x) - \sigma \|\nabla f(x)\|^2 + o(\sigma \|\nabla f(x)\|) \\ &< f(x) \text{ for } \sigma \text{ small enough} \end{aligned}$$

Step-size

- optimal step-size: $\sigma_t = \underset{\sigma}{\operatorname{argmin}} f(\mathbf{x}_t - \sigma \nabla f(\mathbf{x}_t))$
- **Line Search:** **total** or partial optimization w.r.t. σ
Total is however often too "expensive" (needs to be performed at each iteration step)
Partial optimization: execute a limited number of trial steps until a loose approximation of the optimum is found. Typical rule for partial optimization: **Armijo rule** (see next slides)

Typical stopping criterium:

norm of gradient smaller than ϵ

The Armijo-Goldstein Rule

Choosing the step size:

- Only to decrease f -value not enough to converge (quickly)
- Want to have a reasonably large decrease in f

Armijo-Goldstein rule:

- also known as backtracking line search
- starts with a (too) large estimate of σ and reduces it until f is reduced enough
- what is enough?
 - assuming a linear f e.g. $m_k(x) = f(x_k) + \nabla f(x_k)^T (x - x_k)$
 - expected decrease if step of σ_k is done in direction \mathbf{d} :
 $\sigma_k \nabla f(x_k)^T \mathbf{d}$
 - actual decrease: $f(x_k) - f(x_k + \sigma_k \mathbf{d})$
 - stop if actual decrease is at least constant times expected decrease (constant typically chosen in $[0, 1]$)

The Armijo-Goldstein Rule

The Actual Algorithm:

Input: descent direction \mathbf{d} , point \mathbf{x} , objective function $f(\mathbf{x})$ and its gradient $\nabla f(\mathbf{x})$, parameters $\sigma_0 = 10$, $\theta \in [0, 1]$ and $\beta \in (0, 1)$

Output: step-size σ

Initialize σ : $\sigma \leftarrow \sigma_0$

while $f(\mathbf{x} + \sigma\mathbf{d}) > f(\mathbf{x}) + \theta\sigma\nabla f(\mathbf{x})^T\mathbf{d}$ **do**

$\sigma \leftarrow \beta\sigma$

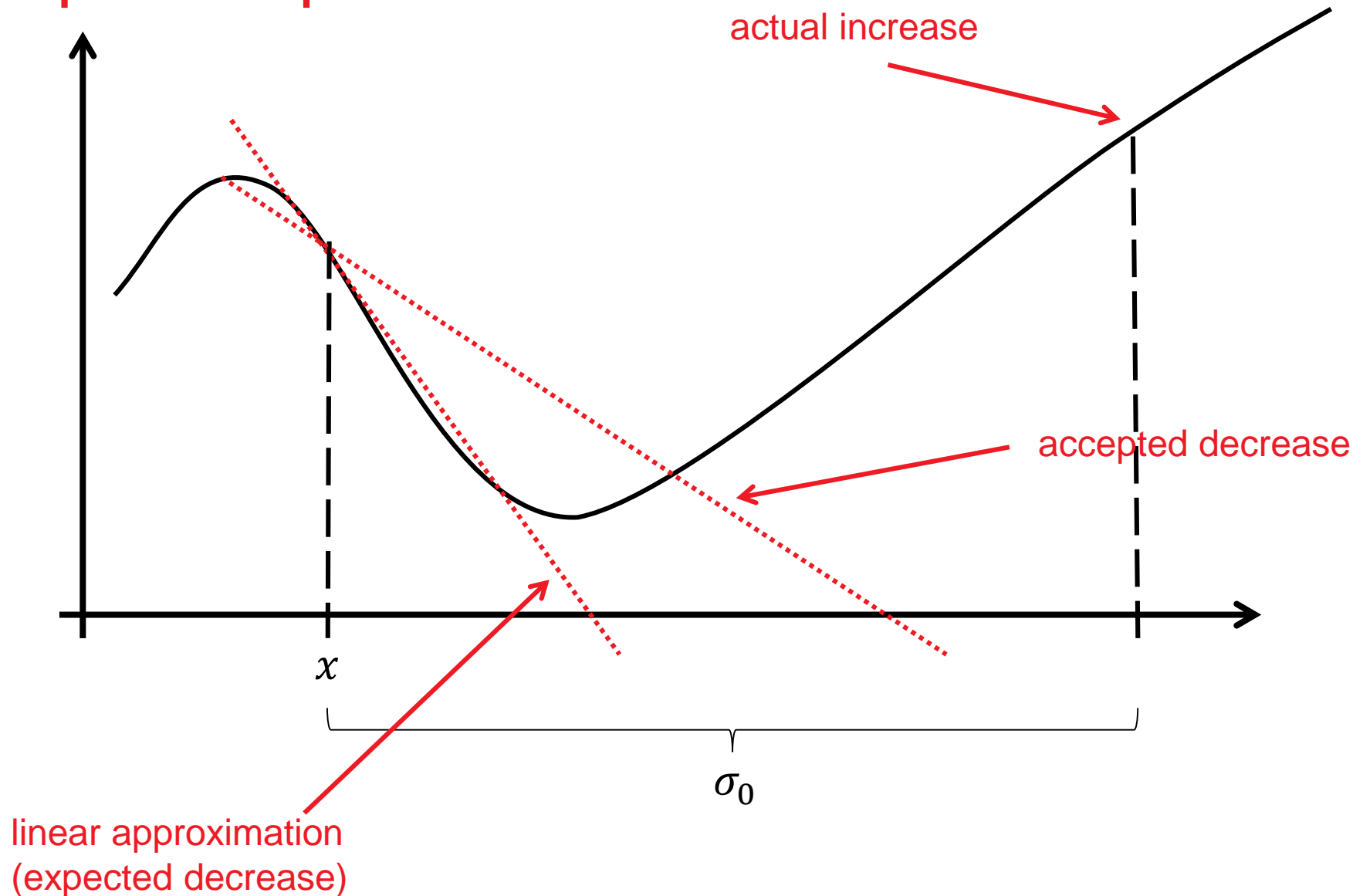
end while

Armijo, in his original publication chose $\beta = \theta = 0.5$.

Choosing $\theta = 0$ means the algorithm accepts any decrease.

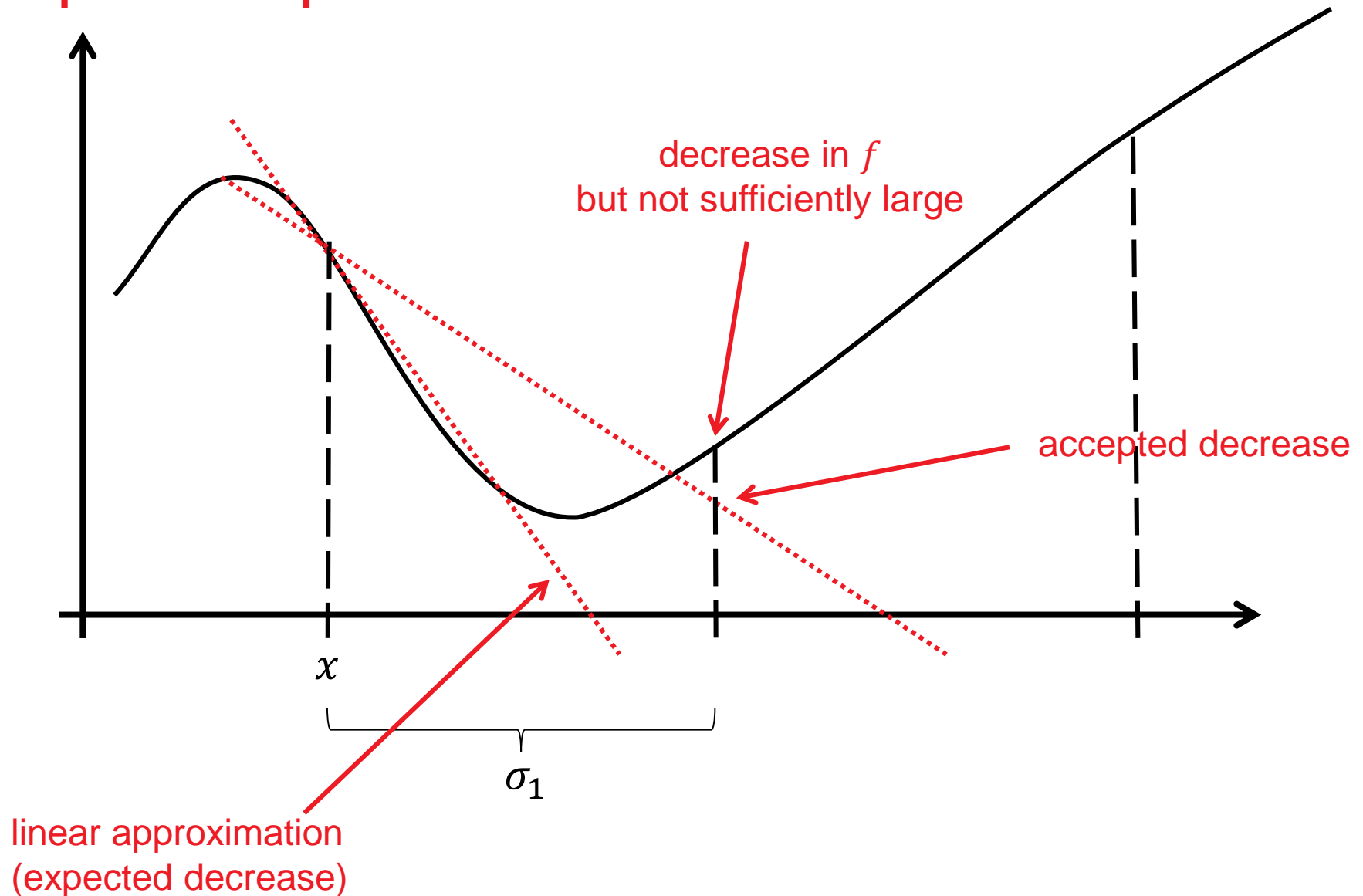
The Armijo-Goldstein Rule

Graphical Interpretation



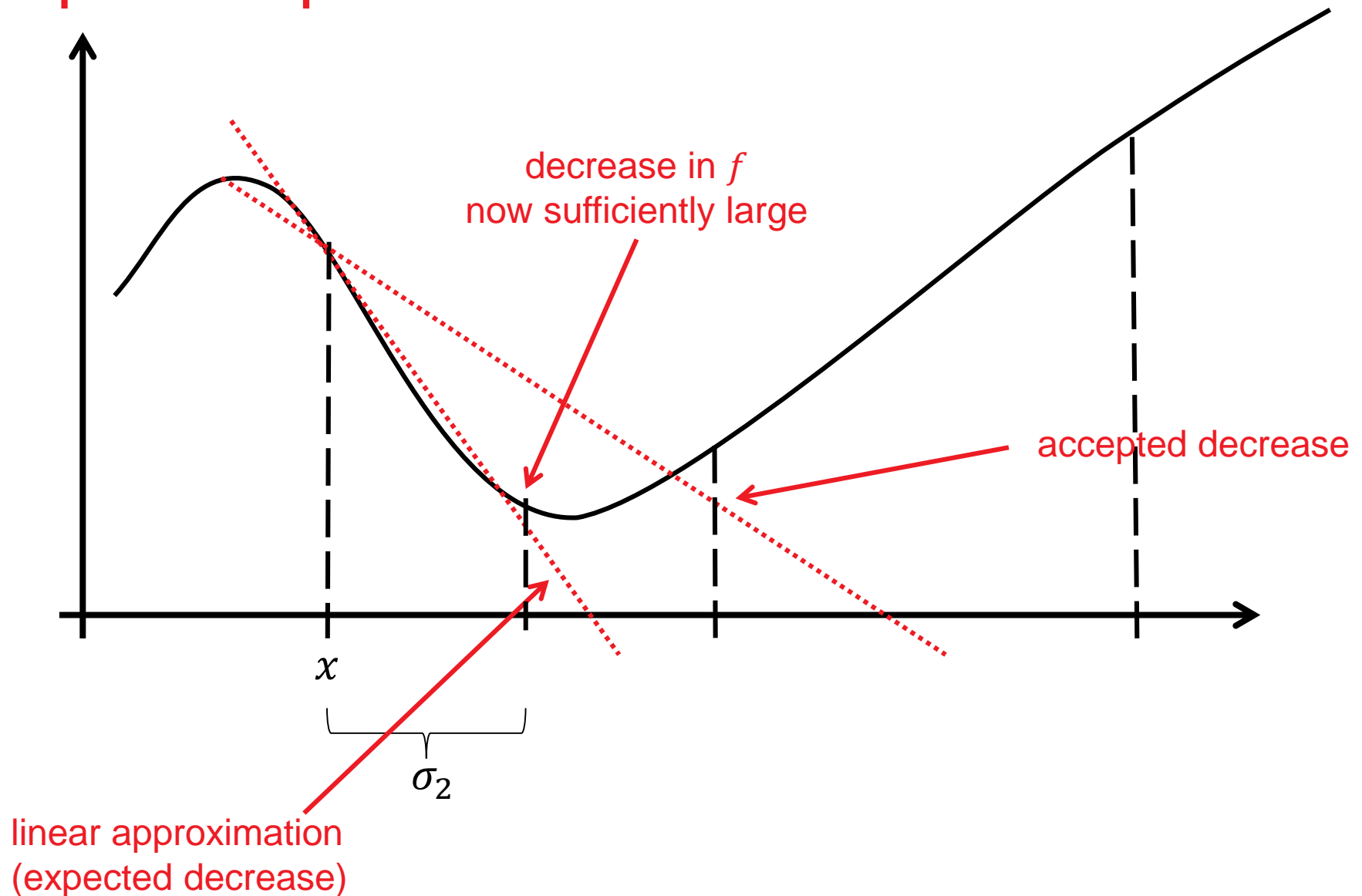
The Armijo-Goldstein Rule

Graphical Interpretation



The Armijo-Goldstein Rule

Graphical Interpretation



Gradient Descent: Simple Theoretical Analysis

Assume f is twice continuously differentiable, convex and that $\mu I_d \preceq \nabla^2 f(x) \preceq L I_d$ with $\mu > 0$ holds, assume a fixed step-size $\sigma_t = \frac{1}{L}$

Note: $A \preceq B$ means $x^T A x \leq x^T B x$ for all x

$$x_{t+1} - x^* = x_t - x^* - \sigma_t \nabla^2 f(y_t)(x_t - x^*) \text{ for some } y_t \in [x_t, x^*]$$

$$x_{t+1} - x^* = \left(I_d - \frac{1}{L} \nabla^2 f(y_t) \right) (x_t - x^*)$$

$$\begin{aligned} \text{Hence } \|x_{t+1} - x^*\|^2 &\leq \|I_d - \frac{1}{L} \nabla^2 f(y_t)\|^2 \|x_t - x^*\|^2 \\ &\leq \left(1 - \frac{\mu}{L}\right)^2 \|x_t - x^*\|^2 \end{aligned}$$

$$\text{Linear convergence: } \|x_{t+1} - x^*\| \leq \left(1 - \frac{\mu}{L}\right) \|x_t - x^*\|$$

algorithm slower and slower with increasing condition number

Non-convex setting: convergence towards stationary point

Newton Algorithm

Newton Method

- descent direction: $-\left[\nabla^2 f(x_k)\right]^{-1} \nabla f(x_k)$ [so-called **Newton direction**]
- The Newton direction:
 - minimizes the best (locally) quadratic approximation of f :
$$\tilde{f}(x + \Delta x) = f(x) + \nabla f(x)^T \Delta x + \frac{1}{2} (\Delta x)^T \nabla^2 f(x) \Delta x$$
 - points towards the optimum on $f(x) = (x - x^*)^T A (x - x^*)$
- however, Hessian matrix is expensive to compute in general and its inversion is also not easy

quadratic convergence

$$\left(\text{i.e. } \lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^2} = \mu > 0 \right)$$

Remark: Affine Invariance

Affine Invariance: same behavior on $f(x)$ and $f(Ax + b)$ for $A \in \text{GLn}(\mathbb{R}) = \text{set of all invertible } n \times n \text{ matrices over } \mathbb{R}$

- Newton method is affine invariant

see http://users.ece.utexas.edu/~cmcaram/EE381V_2012F/Lecture_6_Scribe_Notes.final.pdf

- same convergence rate on all convex-quadratic functions
- Gradient method not affine invariant

Quasi-Newton Method: BFGS

$x_{t+1} = x_t - \sigma_t H_t \nabla f(x_t)$ where H_t is an **approximation** of the inverse Hessian

Key idea of Quasi Newton:

successive iterates x_t, x_{t+1} and gradients $\nabla f(x_t), \nabla f(x_{t+1})$ yield second order information

$$q_t \approx \nabla^2 f(x_{t+1}) p_t$$

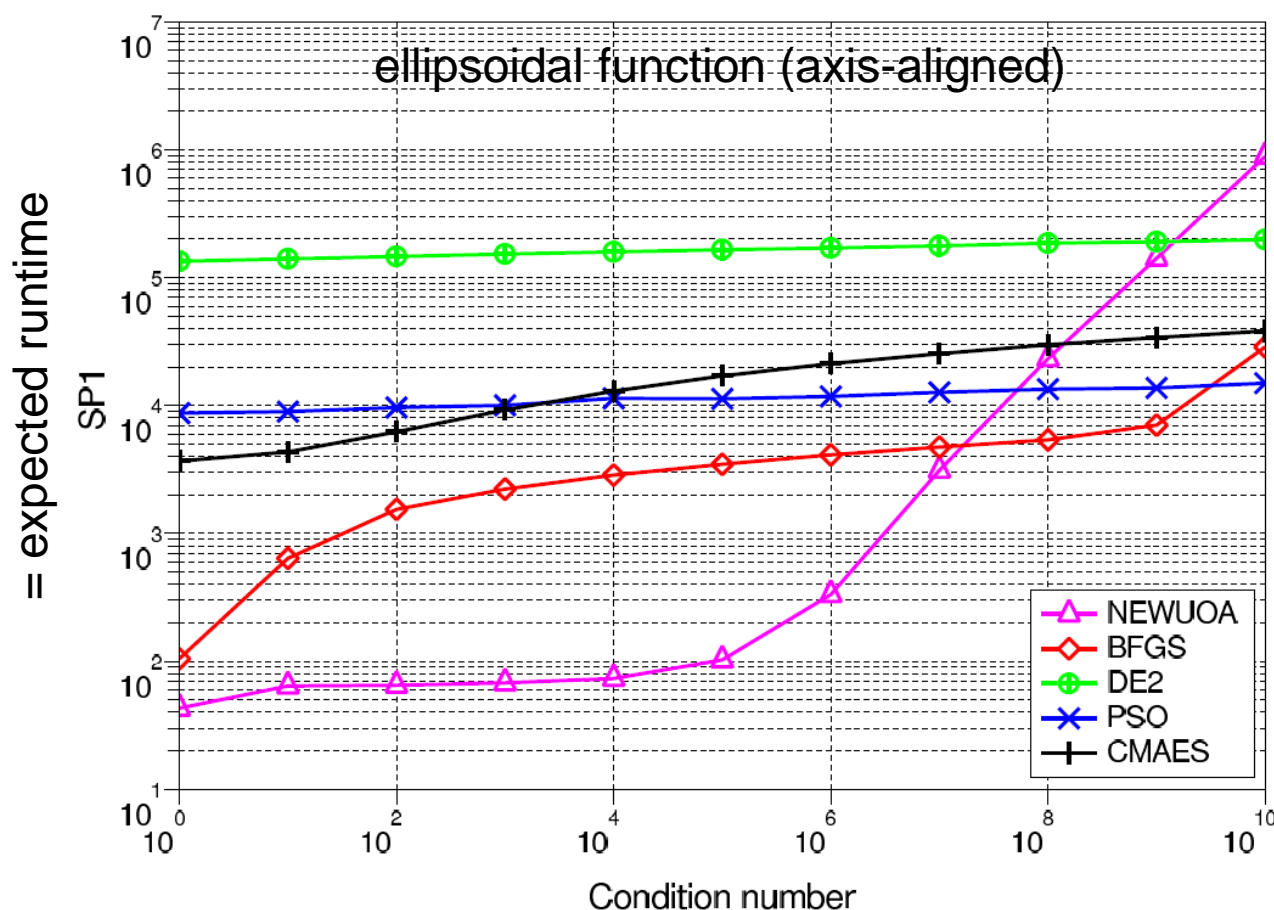
where $p_t = x_{t+1} - x_t$ and $q_t = \nabla f(x_{t+1}) - \nabla f(x_t)$

Most popular implementation of this idea: **Broyden-Fletcher-Goldfarb-Shanno (BFGS)**

- default in MATLAB's `fminunc` and python's `scipy.optimize.minimize`

Once Again: Invariances

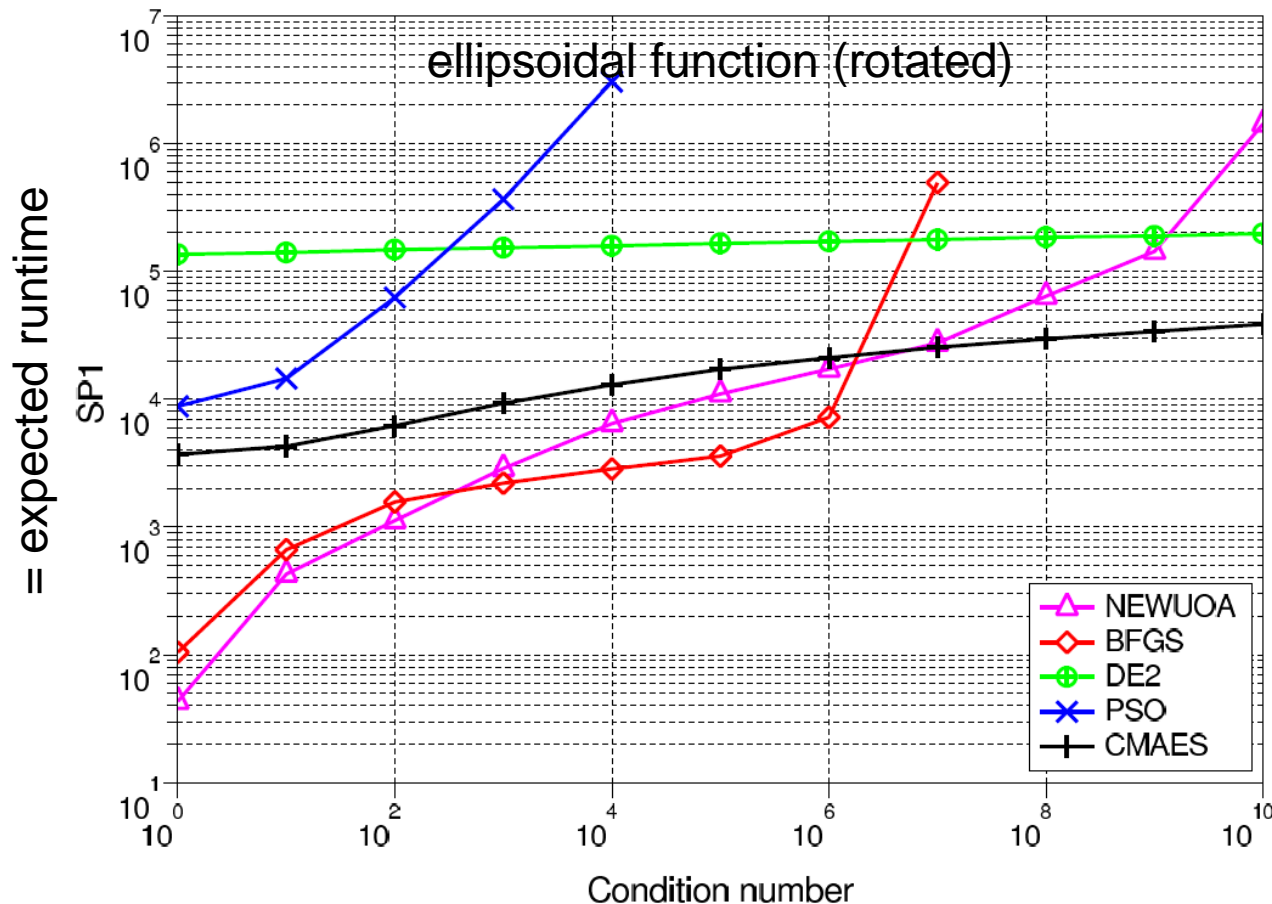
- Newton as well as BFGS are affine invariant in theory
- However, numerics might play an important role



Picture taken from “Experimental Comparisons of Derivative Free Optimization Algorithms by A. Auger, N. Hansen, J. M. Perez Zerpa, R. Ros, and M. Schoenauer. In: *International Symposium on Experimental Algorithms*. Springer, Berlin, Heidelberg, 2009.

Once Again: Invariances

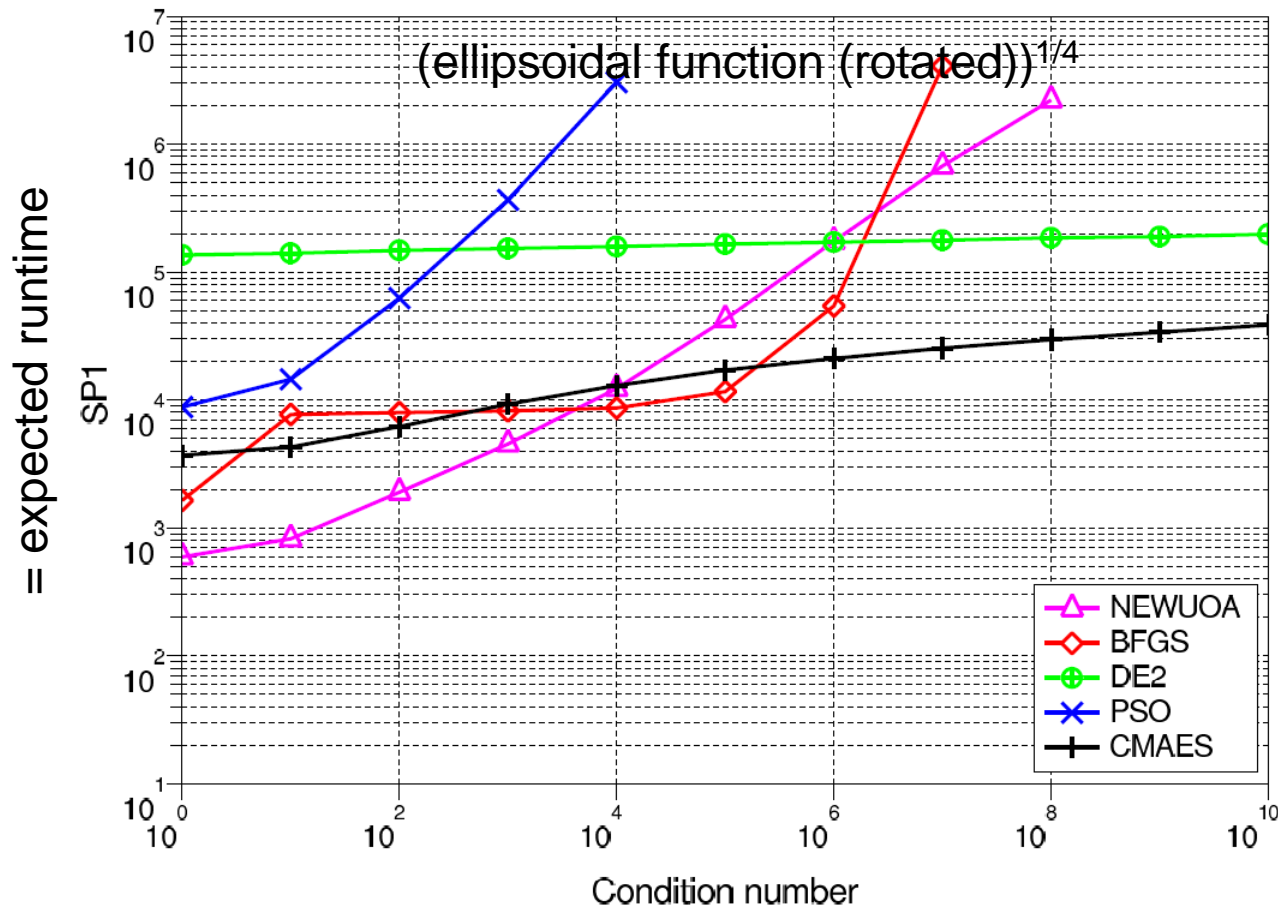
- Newton as well as BFGS are affine invariant in theory
- However, numerics might play an important role



Picture taken from “Experimental Comparisons of Derivative Free Optimization Algorithms by A. Auger, N. Hansen, J. M. Perez Zerpa, R. Ros, and M. Schoenauer. In: *International Symposium on Experimental Algorithms*. Springer, Berlin, Heidelberg, 2009.

Once Again: Invariances

Other invariances might be nice as well, for example in function-value-free algorithms...



Picture taken from “Experimental Comparisons of Derivative Free Optimization Algorithms by A. Auger, N. Hansen, J. M. Perez Zerpa, R. Ros, and M. Schoenauer. In: *International Symposium on Experimental Algorithms*. Springer, Berlin, Heidelberg, 2009.

Conclusions

I hope it became clear...

...what is the difference between **gradient** and **Newton direction**
and **gradient descent** and **(quasi-)Newton algorithms**

... that adapting the step size in descent algorithms is crucial and

... what invariances are and why they are important in
optimization (more details on that in the CMA-ES lecture)