# Introduction to Optimization

## Introduction to Continuous Optimization III / Gradient-Based Algorithms

November 20, 2015

École Centrale Paris, Châtenay-Malabry, France

Dimo Brockhoff

INRIA Lille – Nord Europe

INVENTORS FOR THE DIGITAL WORLD

# 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 *[salle Proto]* |
| Fri, 6.11.2015 | D | Approximation algorithms and heuristics *[S205/S207]* |
| Mon, 9.11.2015 | C | Introduction to Continuous Optimization I *[S118]* |
| Fri, 13.11.2015 | C | Introduction to Continuous Optimization II *[from here onwards always: S205/S207]* |
| **Fri, 20.11.2015** | **C** | **Gradient-based Algorithms [+ finishing the intro]** |
| Fri, 27.11.2015 | C | ~~End of Gradient-based Algorithms + Linear Programming~~ *Stochastic Optimization and Derivative Free Optimization I* |
| Fri, 4.12.2015 | C | Stochastic Optimization and Derivative Free Optimization II |
| Tue, 15.12.2015 | | Exam |

# Lecture Overview Continuous Optimization

**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*

**Question:** Is the Hessian matrix always symmetric?

**Answer:** No, but $f$ having continuous second order partial derivatives is a sufficient condition for the Hessian to be symmetric ("Schwarz' theorem").

# Remark on Last Lecture II

**Question:** How do we prove in general that the gradient is orthogonal to the level sets?
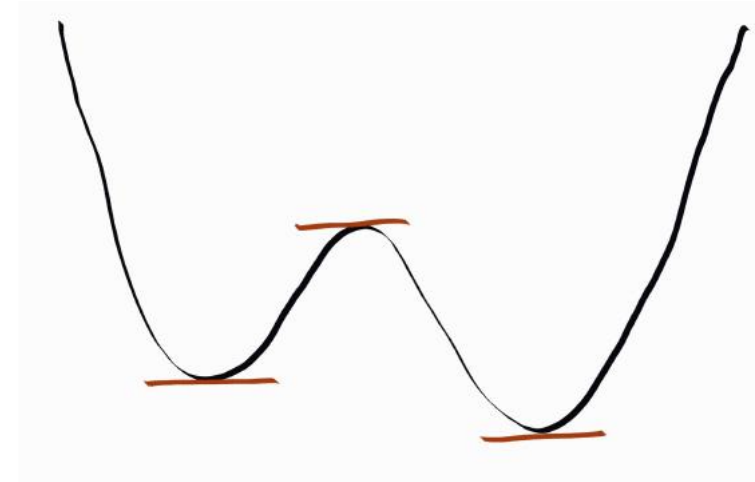
**Answer:**

- similar to what we did for two variables
- take any curve within the level set, parametrized by $t \longmapsto c(t)$
- clear: $f(c(t)) = c$ for all $t$
- derivative wrt to $t$: $\frac{d}{dt} f\big(c(t)\big) = 0$
- but also $\frac{d}{dt} f(c(t)) = \nabla(f(c(t)) \frac{d}{dt} c(t)$

  [via chain rule, $\frac{d}{dt} c(t)$ is a vector, tangent to the curve in $t$]

# Mathematical Tools to Characterize Optima

# Mathematical Characterization of Optima

**Objective:** Derive general characterization of optima

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



**Final Goal:**

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

# Optimality Conditions for Unconstrained Problems

**For 1-dimensional optimization problems $f\colon \mathbb{R} \to \mathbb{R}$**

Assume $f$ is differentiable

- $x^*$ is a local optimum $\implies f'(x^*) = 0$

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

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

- points $y$ such that $f'(y) = 0$ are called <span style="color:red">critical</span> or <span style="color:red">stationary</span> points

**Generalization to $n$-dimensional functions**

If $f\colon U \subset \mathbb{R}^n \longmapsto \mathbb{R}$ is differentiable

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

  *proof via Taylor formula*

If $f$ is twice continuously differentiable

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

  *proof via Taylor formula at order 2*

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

**Proof of Sufficient Condition:**

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

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

$$> \frac{\lambda}{2} ||h||^2 + o(||h||^2) = \left( \frac{\lambda}{2} + \frac{o(||h||^2)}{||h||^2} \right) ||h||^2$$

# Convex Functions

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

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

## Theorem

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

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

*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 $\nabla^2 f(\boldsymbol{x})$ is positive semi-definite for all $\boldsymbol{x}$.

**Examples of Convex Functions:**

- $f(\boldsymbol{x}) = a^T \boldsymbol{x} + b$

- $f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^T A \boldsymbol{x} + a^T \boldsymbol{x} + b$, $A$ symmetric positive definite

- the negative of the entropy function (i.e. $f(\boldsymbol{x}) = \sum_{i=1}^{n} \boldsymbol{x}_i \ln(\boldsymbol{x_i})$ for positive $\boldsymbol{x}$)

**Exercise:**

Let $f \colon U \to \mathbb{R}$ be a convex and differentiable function on a convex open $U$.
Show that if $\nabla f(\boldsymbol{x}^*) = 0$, then $\boldsymbol{x}^*$ is a global minimum of $f$

**Why convexity?** local minima are also global under convexity assumption.

# Constrained Optimization

**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 derivative is continuous*

**Theorem:**

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

Let $a \in E$ satisfy

$$\begin{cases} f(a) = \inf \{f(x) \mid x \in U, g(x) = 0\} \\ \qquad\qquad g(a) = 0 \end{cases}$$

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

$$\underbrace{\nabla f(a) + \lambda \nabla g(a) = 0}$$

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

Note: $a$ need not be a global minimum but a local one

**Exercise:**

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

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

1) Plot the level sets of $f$, plot $g = 0$
2) Compute $\nabla f$ and $\nabla g$
3) Find the solutions with $\nabla f + \lambda \nabla g = 0$
   
   *equation solving with 3 unknowns $(x, y, \lambda)$*
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.

**Theorem**

- Assume $f: U \to \mathbb{R}$ and $g_k: U \to \mathbb{R}$ $(1 \le k \le p)$ are $\mathcal{C}^1$.

- Let $a$ be such that

$$\begin{cases} f(a) = \inf \{ f(x) \mid x \in \mathbb{R}^n, \quad g_k(x) = 0, \quad 1 \le k \le p \} \\ \qquad\qquad g_k(a) = 0 \text{ for all } 1 \le k \le p \end{cases}$$

- If $\left( \nabla g_k(a) \right)_{1 \le k \le p}$ are linearly independent, then there exist $p$ real constants $(\lambda_k)_{1 \le k \le p}$ 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

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

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

# Inequality Constraints: 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 $g_k(a) = 0$.

**Theorem (Karush-Kuhn-Tucker, KKT):**

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

Furthermore, let $a \in U$ satisfy

$$\begin{cases} f(a) = \inf(f(x) \mid x \in U, g_k(x) = 0 \text{ (for } k \in E), g_k(x) \le 0 \text{ (for } k \in \mathrm{I}) \\ \qquad\qquad g_k(a) = 0 \text{ (for } k \in E) \\ \qquad\qquad g_k(a) \le 0 \text{ (for } k \in I) \end{cases}$$

also works again for $a$
being a local minimum

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

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

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

**Theorem (Karush-Kuhn-Tucker, KKT):**

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

Furthermore, let $a \in U$ satisfy

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

also works again for $a$
being a local minimum

Let $I_a^0$ be the set of constraints that are active in $a$ and assume that $\left(\nabla g_k(a)\right)_{k \, \in \, E \, \cup \, I_a^0}$ 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**

&#10102;  choose an initial point $\boldsymbol{x}_0$, set $t = 1$

&#10103;  while not happy

- choose a descent direction $\boldsymbol{d}_t \neq 0$

- line search:

  - choose a step size $\sigma_t > 0$
  - set $\boldsymbol{x}_{t+1} = \boldsymbol{x}_t + \sigma_t \boldsymbol{d}_t$

- set $t = t + 1$

**Remaining questions**

- how to choose $\boldsymbol{d}_t$?
- how to choose $\sigma_t$?

# Gradient Descent

**Rationale:** $\boldsymbol{d}_t = -\nabla f(\boldsymbol{x}_t)$ is a descent direction

indeed for $f$ differentiable

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

## Step-size

- optimal step-size: $\sigma_t = \underset{\sigma}{\mathrm{argmin}}\, f(\boldsymbol{x}_t - \sigma \nabla f(\boldsymbol{x}_t))$

- **Line Search:** total or partial optimization w.r.t. $\sigma$
  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 slide and exercise

## Stopping criteria:

norm of gradient smaller than $\epsilon$

**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 $\sigma$ 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 $\sigma_k$ is done in direction $\boldsymbol{d}$: $\sigma_k \nabla f(x_k)^T \boldsymbol{d}$
    - actual decrease: $f(x_k) - f(x_k + \sigma_k \boldsymbol{d})$
    - 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 $\sigma$

Initialize $\sigma$: $\sigma \leftarrow \sigma_0$
**while** $f(\mathbf{x} + \sigma\mathbf{d}) > f(\mathbf{x}) + \theta\sigma\nabla f(\mathbf{x})^T\mathbf{d}$ **do**
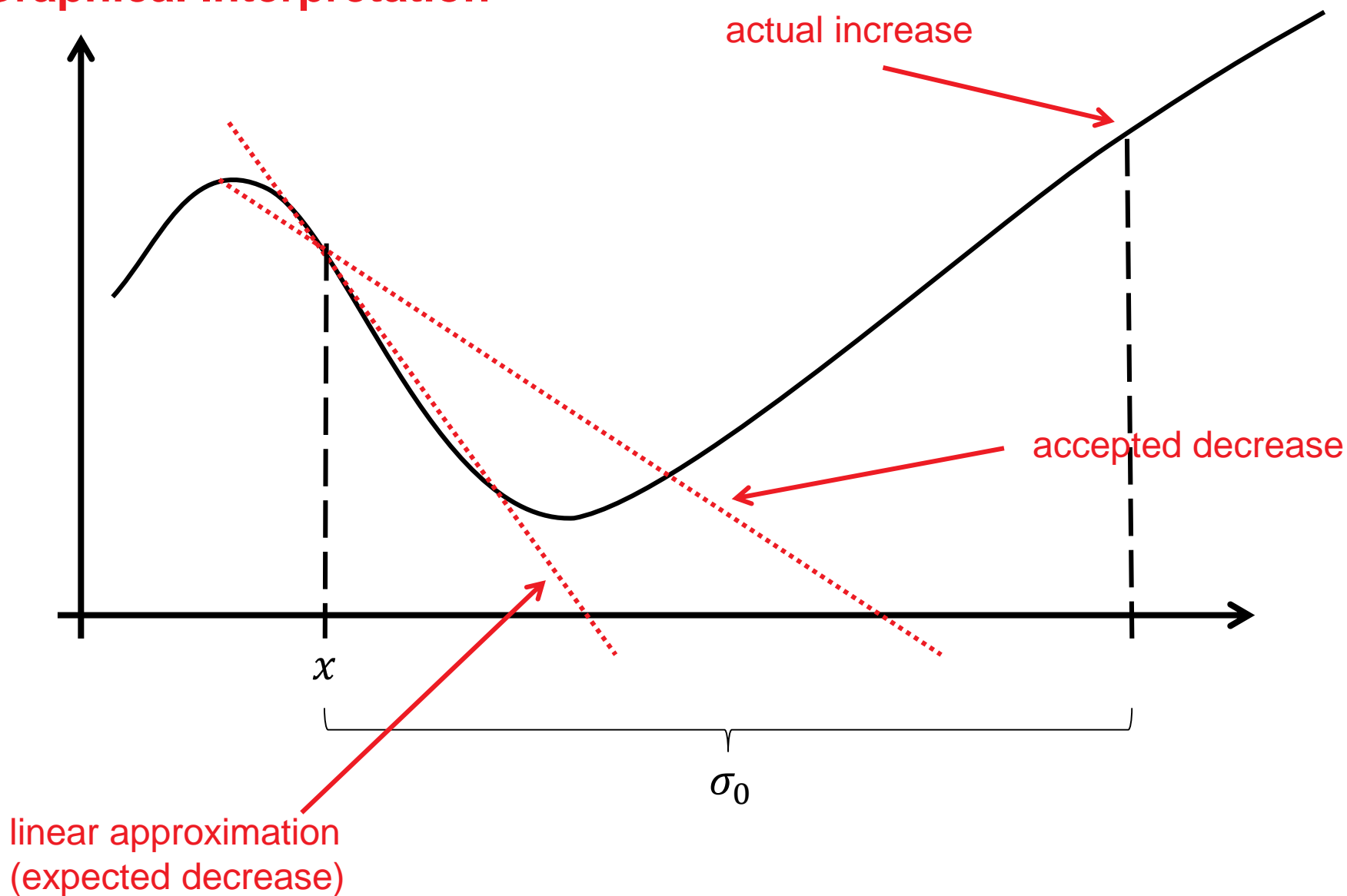    $\sigma \leftarrow \beta\sigma$
**end while**

---

Armijo, in his original publication chose $\beta = \theta = 0.5$.
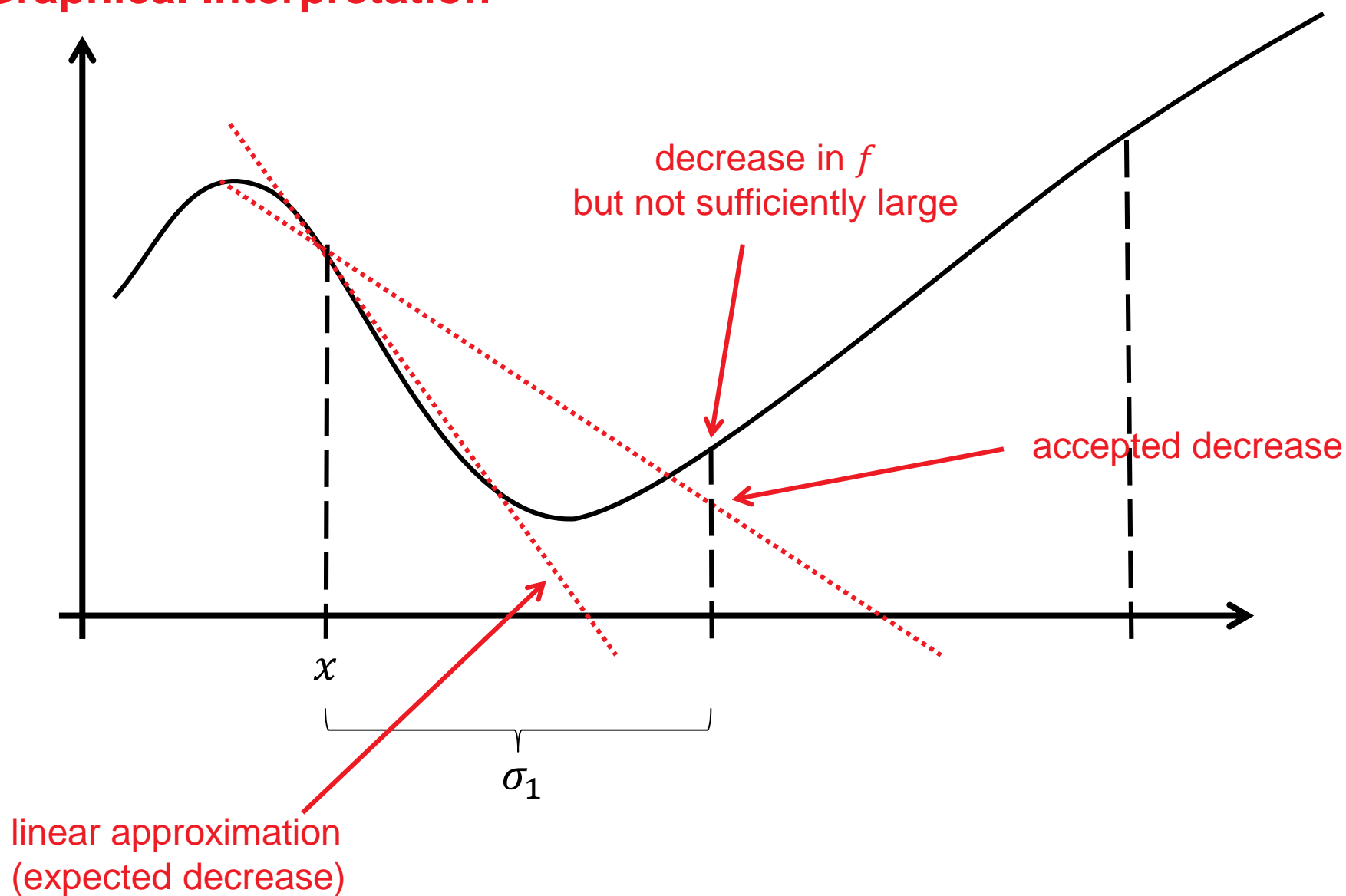
Choosing $\theta = 0$ means the algorithm accepts any decrease.
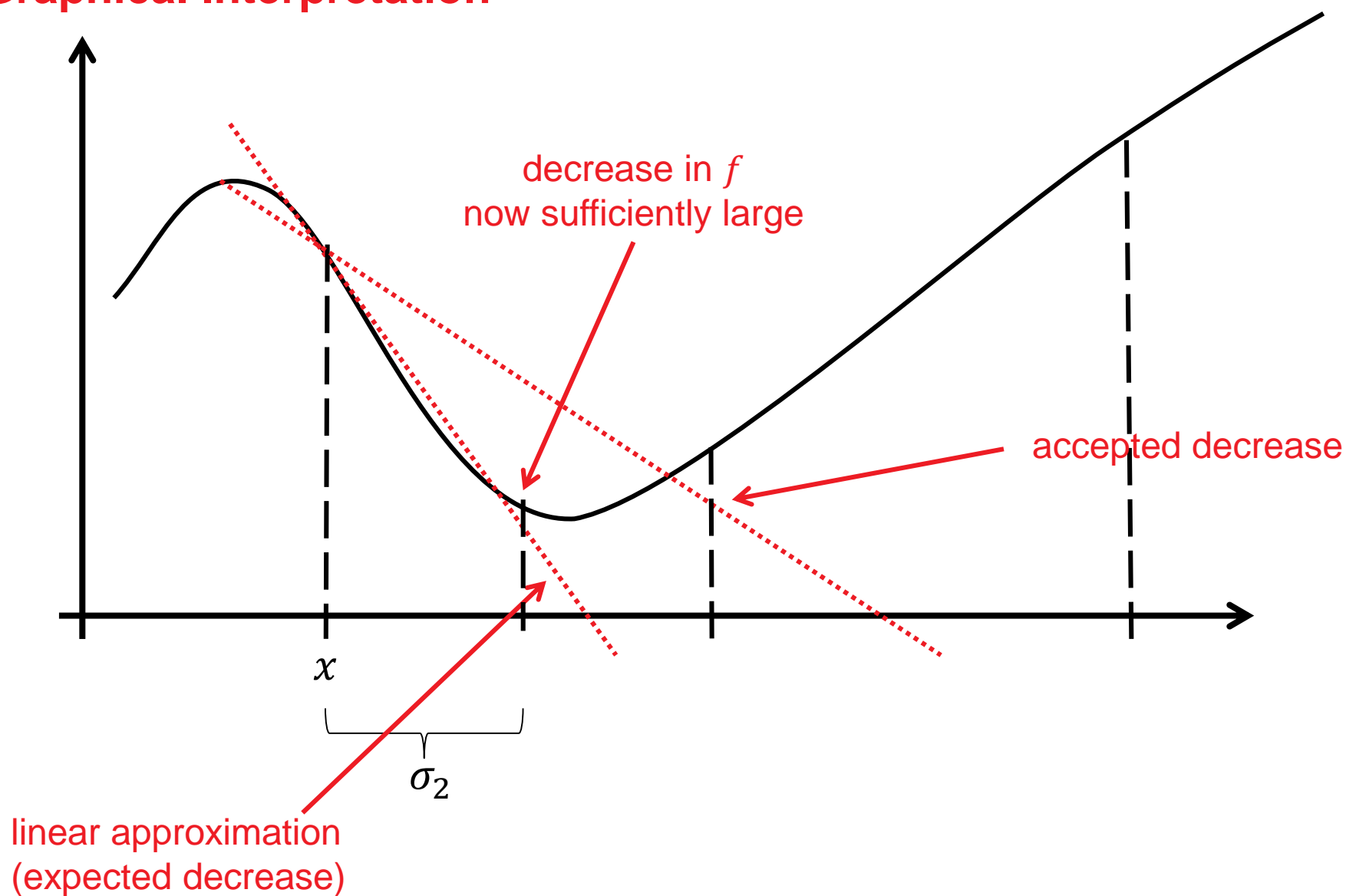
**Graphical Interpretation**



actual increase

accepted decrease

$x$

$\sigma_0$

linear approximation
(expected decrease)

## Graphical Interpretation



decrease in $f$
but not sufficiently large

accepted decrease

$x$

$\sigma_1$

linear approximation
(expected decrease)

## Graphical Interpretation



decrease in $f$
now sufficiently large

accepted decrease

$x$

$\sigma_2$

linear approximation
(expected decrease)

# Gradient Descent: Simple Theoretical Analysis

Assume $f$ is twice continuously differentiable, convex and that

$\mu I_d \preccurlyeq \nabla^2 f(x) \preccurlyeq L I_d$ with $\mu > 0$ holds, assume a fixed step-size $\sigma_t = \frac{1}{L}$

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

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

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

Hence $||x_{t+1} - x^*||^2 \leq |||I_d - \frac{1}{L} \nabla^2 f(y_t)|||^2 \ ||x_t - x^*||^2$

$$\leq \left( 1 - \frac{\mu}{L} \right)^2 ||x_t - x^*||^2$$

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: $-[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$ [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*

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

**Affine Invariance:** same behavior on $f(x)$ and $f(Ax + b)$ for $A \in \text{GLn}(\mathbb{R})$

- 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

$x_{t+1} = x_t - \sigma_t H_t \nabla f(x_t)$ 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`

I hope it became clear...

    ...what are <span style="color:red">gradient</span> and <span style="color:red">Hessian</span>

    ...what are sufficient and necessary conditions for optimality

    ...what is the difference between <span style="color:red">gradient</span> and <span style="color:red">Newton direction</span>

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