

# Optimization Methods

Fabian Bosshard

July 14, 2025

## Contents

<b>Preface</b>	<b>ii</b>
<b>References</b>	<b>ii</b>
<b>1 General formulation of an optimization problem</b>	<b>1</b>
<b>2 Useful definitions and properties</b>	<b>2</b>
<b>3 Classification of Problems</b>	<b>3</b>
3.1 Linear vs Nonlinear Programming . . . . .	3
3.2 Continuous vs Discrete Programming . . . . .	3
3.3 Convex vs Non-convex Optimization . . . . .	3
3.4 Constrained vs Unconstrained Optimization . . . . .	4
3.5 Smooth vs Non-smooth Optimization . . . . .	5
3.6 Global vs Local Optimization . . . . .	5
3.7 Stochastic vs Deterministic Optimization . . . . .	6
<b>4 Unconstrained Optimization</b>	<b>7</b>
4.1 Univariate Objective Functions . . . . .	7
4.1.1 Taylor Expansions for Univariate Functions . . . . .	7
4.1.2 First and Second Order Optimality Conditions for Univariate Functions . . . . .	8
4.1.3 Minimizers for Convex Univariate Functions . . . . .	9
4.2 Multivariate Objective Functions . . . . .	9
4.2.1 First Order Derivatives for Multivariate Functions . . . . .	9
4.2.2 First Order Taylor's Expansion for Multivariate Functions . . . . .	10
4.2.3 First Order Optimality Conditions for Multivariate Functions . . . . .	11
4.3 Gradient Descent . . . . .	11
4.3.1 Gradient Descent Algorithm . . . . .	11
4.3.2 Convergence Analysis . . . . .	12
4.3.3 Convergence Rate for Strongly Convex Functions . . . . .	13
4.4 Newton's Method . . . . .	14
4.4.1 Newton's Method for Univariate Functions . . . . .	14
4.4.2 Second Order Taylor's Expansion for Multivariate Functions . . . . .	15
4.4.3 Second Order Optimality Conditions for Multivariate Functions . . . . .	16
4.4.4 Minimization of Quadratic Functions . . . . .	16
4.4.5 Newton's Method for Multivariate Functions . . . . .	17
4.4.6 Convergence of Newton's Method . . . . .	18
4.5 Line Search Algorithms . . . . .	19
4.5.1 Wolfe Conditions . . . . .	20
4.5.2 Backtracking Line Search . . . . .	20
4.5.3 Global Convergence Theorem . . . . .	21
4.6 Variations of Newton's Method . . . . .	22
4.6.1 Newton with Hessian Correction . . . . .	22
4.6.2 Quasi-Newton Methods . . . . .	23
<b>5 Constrained Optimization</b>	<b>25</b>
5.1 Introduction and Terminology . . . . .	25
5.1.1 Classes of Constrained Optimization Problems . . . . .	25
5.1.2 Active Constraints . . . . .	25
5.1.3 Feasible Directions . . . . .	25
5.1.4 Constraint Qualification Condition . . . . .	26

5.2	Optimality Conditions . . . . .	26
5.2.1	Necessary First-Order Conditions . . . . .	26
5.2.2	Second-Order Conditions . . . . .	28
<b>6</b>	<b>Linear Programming</b>	<b>29</b>
6.1	General and Standard Form . . . . .	29
6.1.1	“Generic” Form . . . . .	29
6.1.2	General Form . . . . .	29
6.1.3	Standard Form . . . . .	30
6.2	Geometric Properties . . . . .	30
6.2.1	Corner Points . . . . .	30
6.2.2	Basic Solutions for Polyhedra in Standard Form . . . . .	31
6.2.3	Degeneracy . . . . .	32
6.2.4	Existence of Basic Feasible Solutions . . . . .	32
6.3	Optimality of Corner Points . . . . .	32
6.4	Simplex Algorithm . . . . .	33
6.4.1	Feasible and Basic Directions . . . . .	33
6.4.2	Reduced Costs . . . . .	34
6.4.3	Optimality Conditions . . . . .	34
6.4.4	Non-degenerate Case . . . . .	34
6.4.5	Degenerate Case . . . . .	35
6.4.6	Computational Complexity . . . . .	35

## Preface

This document is an unofficial student-made summary of the course Optimization Methods taught by Deborah Sulem in Spring 2025 at the Università della Svizzera italiana. It is based on the slides and lecture notes, as well as [1]. The summary is not exhaustive and may contain errors. If you find any, please report them to [fabianlucasbosshard@gmail.com](mailto:fabianlucasbosshard@gmail.com) or open an issue at [https://github.com/fabianbosshard/USI\\_summaries](https://github.com/fabianbosshard/USI_summaries).

This work is licensed under a Creative Commons “Attribution 4.0 International” license.



## References

- [1] Jorge Nocedal and Stephen J. Wright. Numerical Optimization. Second. Springer, 2006. ISBN: 9780387303031. DOI: 10.1007/978-0-387-40065-5. URL: <https://link.springer.com/book/10.1007/978-0-387-40065-5>.

# 1 General formulation of an optimization problem

main components:

- **objective function**  $f(\vec{x})$ : function to be minimized or maximized
- **variables**  $\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ : quantities of the problem that can be adjusted and need to be chosen **optimally**
- **constraint functions**  $c_i(\vec{x})$ : conditions that limit the possible values of  $\vec{x}$ . they can be **equality** or **inequality** constraints
- **feasible region**: set of all variables  $\vec{x}$  satisfying all constraints

general form of an optimization problem:

$$\min_{\vec{x} \in \mathbb{R}^n} f(\vec{x}) \quad \text{subject to} \quad \begin{cases} c_i(\vec{x}) = 0, & i \in \mathcal{E} \\ c_i(\vec{x}) \geq 0, & i \in \mathcal{I} \end{cases} \quad (1)$$

with

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  the objective function
- $\mathcal{E}$  the index set of equality constraints
- $\mathcal{I}$  the index set of inequality constraints

**Remark 1.**

- We do not lose generality by restricting our analysis to minimization problems, since maximizing  $f$  is **equivalent** to minimizing  $-f$ .
- We do not lose generality by considering constraints of the form in (1) since:
  - a constraint of the form  $c_i(\vec{x}) = b$  can be transformed into  $\bar{c}_i(\vec{x}) = 0$  with  $\bar{c}_i(\vec{x}) = c_i(\vec{x}) - b$ .
  - a constraint of the form  $c_i(\vec{x}) \leq 0$  can be transformed into  $\bar{c}_i(\vec{x}) \geq 0$  with  $\bar{c}_i(\vec{x}) = -c_i(\vec{x})$ .
- If  $n = 1$ , the function  $f$  is **univariate** and we also say that the optimization problem is univariate. If  $n = 2$ , it is a **bivariate** problem, and if  $n \geq 2$ , it is a **multivariate** problem. ◀

A point  $\vec{x} \in \mathbb{R}^n$  satisfying all the constraints is called a **feasible point** and the set of all feasible points is called the **feasible set** (or, feasible region) of the optimization problem. It is formally defined as

$$\Omega = \{\vec{x} \in \mathbb{R}^n : c_i(\vec{x}) = 0, i \in \mathcal{E}, c_i(\vec{x}) \geq 0, i \in \mathcal{I}\} \quad (2)$$

This also implies that the problem (1) can be re-written as

$$\min_{\vec{x} \in \Omega} f(\vec{x})$$

The general formulation of optimization problems includes a large variety of problems and each category of problem requires a specific methodology. There is no universal optimization algorithm that is applicable to all problems. Thus, one needs to understand the specificity of each problem and algorithm to identify the most suitable methodology.

## 2 Useful definitions and properties

**Theorem 1** (Cauchy–Schwarz Inequality). The **Cauchy–Schwarz inequality** states that for any vectors  $\vec{x}, \vec{y}$  of an inner product space,  $|\langle \vec{x}, \vec{y} \rangle|^2 \leq \langle \vec{x}, \vec{x} \rangle \langle \vec{y}, \vec{y} \rangle$  where  $\langle \cdot, \cdot \rangle$  is the inner product. Or written in terms of the induced norm  $\|\cdot\| := \sqrt{\langle \cdot, \cdot \rangle}$ ,  $|\langle \vec{x}, \vec{y} \rangle| \leq \|\vec{x}\| \|\vec{y}\|$ . For  $\mathbb{R}^n$  and the dot product, this means

$$|\vec{x} \cdot \vec{y}| \leq \|\vec{x}\| \|\vec{y}\|$$

$$\left| \sum_{i=1}^n x_i y_i \right| \leq \sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}$$

◁

The dot product can be expressed in terms of the angle  $\theta$  between the two vectors:

$$\vec{x} \cdot \vec{y} = \vec{x}^T \vec{y} = \|\vec{x}\| \|\vec{y}\| \cos(\theta)$$

**Definition 1** (Spectral Norm). The **spectral norm** of a matrix  $\underline{A} \in \mathbb{R}^{m \times n}$  is defined as

$$\|\underline{A}\|_2 := \max_{\vec{x} \neq \vec{0}} \frac{\|\underline{A}\vec{x}\|_2}{\|\vec{x}\|_2} = \max_{\|\vec{x}\|_2=1} \|\underline{A}\vec{x}\|_2 = \sqrt{\lambda_{\max}(\underline{A}^T \underline{A})}$$

where  $\lambda_{\max}(\underline{A}^T \underline{A})$  is the largest eigenvalue of  $\underline{A}^T \underline{A} \in \mathbb{R}^{n \times n}$ . If  $\underline{A}$  is symmetric, then the spectral norm is equal to the largest eigenvalue (in absolute value) of  $\underline{A}$ ,

$$\|\underline{A}\|_2 = \max_{i \in \{1, \dots, n\}} |\lambda_i(\underline{A})| \quad \text{if } \underline{A}^T = \underline{A}$$

where  $\lambda_i(\underline{A})$  is the  $i$ -th eigenvalue of  $\underline{A}$ . ◁

**Definition 2** (Condition Number). The **condition number** of a matrix  $\underline{A} \in \mathbb{R}^{m \times n}$  is defined as

$$\kappa(\underline{A}) := \|\underline{A}\|_2 \cdot \|\underline{A}^{-1}\|_2 \geq 1$$

If  $\underline{A}$  is symmetric positive definite,

$$\kappa(\underline{A}) = \frac{\lambda_{\max}(\underline{A})}{\lambda_{\min}(\underline{A})} \quad \text{if } \underline{A} \succ 0$$

◁

### 3 Classification of Problems

Note that “programming” is used synonymously with “optimization”.

#### 3.1 Linear vs Nonlinear Programming

An optimization problem is said to be **linear** if  $f$  is linear (or, affine) and the constraints are linear; otherwise, the problem is **nonlinear**.

**Definition 3** (Linear Function). A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **linear** if:

- for any  $\vec{x}, \vec{y} \in \mathbb{R}^n$ ,  $f(\vec{x} + \vec{y}) = f(\vec{x}) + f(\vec{y})$ ,
- for any  $\vec{x} \in \mathbb{R}^n$  and  $\lambda \in \mathbb{R}$ ,  $f(\lambda\vec{x}) = \lambda f(\vec{x})$ .

Equivalently,  $f$  is linear if it can be written as

$$f(\vec{x}) = \vec{c}^T \vec{x} = \sum_{i=1}^n c_i x_i$$

with  $\vec{c} \in \mathbb{R}^n$ . ◀

**Definition 4** (Affine Function). A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **affine** if

$$f(\alpha\vec{x} + (1 - \alpha)\vec{y}) = \alpha f(\vec{x}) + (1 - \alpha)f(\vec{y}), \quad \forall \vec{x}, \vec{y} \in \mathbb{R}^n, \forall \alpha \in [0, 1]$$

or equivalently, if there exist  $\vec{c} \in \mathbb{R}^n$  and  $b \in \mathbb{R}$  such that

$$f(\vec{x}) = \vec{c}^T \vec{x} + b$$
◀

**Example 1.** The following minimization problem is linear:

$$\min_{\vec{x} \in \mathbb{R}^n} \sum_{i=1}^n c_i x_i \quad \text{subject to} \quad \sum_{i=1}^n x_i = 1$$

with  $c_1, \dots, c_n \in \mathbb{R}$ . ◀

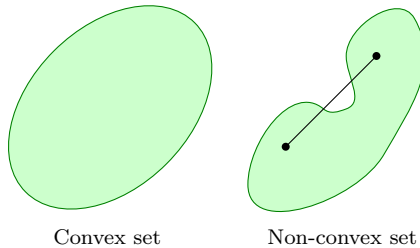
#### 3.2 Continuous vs Discrete Programming

If the components of  $\vec{x}$  are real numbers and in an uncountable set, it is a **continuous** optimization problem. Otherwise, it is a **discrete** programming problem (if the variables can only be integers, it is an **integer** programming problem).

**Example 2.** of a continuous as well as a discrete problem:

- If  $x_1, \dots, x_n$  can take any value in  $\mathbb{R}$  or in  $[a, b]$  with  $a < b$ , then the problem is continuous.
- If  $x_1, \dots, x_n$  can take values in  $\mathbb{N}$  or  $\mathbb{Z}$ , then the problem is discrete. ◀

#### 3.3 Convex vs Non-convex Optimization



An optimization problem is said to be **convex** if the objective function  $f$  is convex, the equality constraints are linear, and the inequality constraints are concave. Equivalently, the problem is convex if the feasible set  $\Omega$  is convex.

**Definition 5** (Convex Set). A set  $S \subset \mathbb{R}^n$  is **convex** if  $\forall \vec{x}, \vec{y} \in S, \forall t \in [0, 1]$ ,

$$t\vec{x} + (1-t)\vec{y} \in S \quad \text{◀}$$

**Example 3.** A ball defined as

$$\mathcal{B}_r(\vec{x}_0) = \{\vec{x} \in \mathbb{R}^n : \|\vec{x} - \vec{x}_0\| \leq r\}$$

with center  $\vec{x}_0 \in \mathbb{R}^n$  and radius  $r$ , is convex. ▶

Unless specified otherwise, the norm  $\|\cdot\|$  is the Euclidean norm:

$$\|\vec{x}\| = \|\vec{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

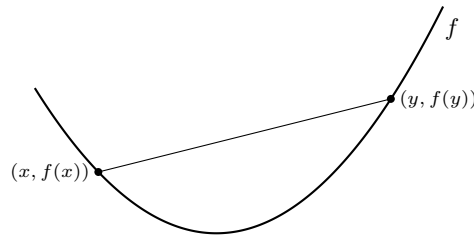
**Definition 6** (Convex Function). A function  $f : S \rightarrow \mathbb{R}$  is **convex** if its domain  $S$  is convex and  $\forall \vec{x}, \vec{y} \in S \forall t \in [0, 1]$ ,

$$f(t\vec{x} + (1-t)\vec{y}) \leq tf(\vec{x}) + (1-t)f(\vec{y}) \quad \text{◀}$$

Concavity is the opposite of convexity:  $f$  is **concave** if  $-f$  is convex.

**Definition 7** (Strictly Convex Function). A function  $f : S \rightarrow \mathbb{R}$  is **strictly convex** if its domain  $S$  is convex and  $\forall \vec{x}, \vec{y} \in S \forall t \in [0, 1]$ ,

$$f(t\vec{x} + (1-t)\vec{y}) < tf(\vec{x}) + (1-t)f(\vec{y}) \quad \text{◀}$$



Example of a convex function.

**Definition 8** ( $\mu$ -Strongly Convex Function). A function  $f : S \rightarrow \mathbb{R}$  is  **$\mu$ -strongly convex** (with  $\mu > 0$ ) if its domain  $S$  is convex and  $\forall \vec{x}, \vec{y} \in S \forall t \in [0, 1]$ ,

$$f(t\vec{x} + (1-t)\vec{y}) \leq tf(\vec{x}) + (1-t)f(\vec{y}) - \frac{\mu}{2}t(1-t)\|\vec{x} - \vec{y}\|^2 \quad \text{◀}$$

Observe that “strong convexity”  $\implies$  “strict convexity”  $\implies$  “convexity”.

**Example 4.** The following examples and properties of convex functions are important:

- The exponential function is convex, while the logarithm function is concave.
- Powers of positive numbers  $x^\alpha$  with  $x > 0$  are convex if  $\alpha \geq 1$  or  $\alpha \leq 0$ , and concave if  $0 \leq \alpha \leq 1$ .
- Affine functions are both convex and concave.
- Absolute value and norms, e.g.,  $\ell_p$ -norms are convex:
  - $\|\vec{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$
  - $\|\vec{x}\|_\infty = \max\{|x_1|, \dots, |x_n|\}$
- Positive part (Rectified Linear Unit)  $\max\{x, 0\}$  is convex, negative part  $\min\{0, x\}$  is concave.
- Positive linear combinations of convex functions are convex.
- If  $h$  is convex and  $g$  is convex and non-decreasing, then  $g(h(\vec{x}))$  is convex. ▶

### 3.4 Constrained vs Unconstrained Optimization

An optimization problem is **constrained** if  $\mathcal{E}$  OR  $\mathcal{I}$  are NOT empty. It is **unconstrained** if both  $\mathcal{E}$  AND  $\mathcal{I}$  are empty, i.e.,  $\mathcal{E} = \mathcal{I} = \emptyset$ .

### 3.5 Smooth vs Non-smooth Optimization

An optimization problem is said to be **smooth** if the objective function  $f$  is smooth, i.e.,  $f$  is continuously differentiable and its derivative is Lipschitz continuous; otherwise, the problem is **non-smooth**.

**Definition 9** (Continuity). A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is **continuous** on  $\mathbb{R}$  if for all  $x_0 \in \mathbb{R}$ ,

$$\lim_{x \rightarrow x_0} f(x) = f(x_0) \quad \text{✓}$$

**Definition 10** (Lipschitz Continuity). A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is  **$L$ -Lipschitz continuous** with  $L > 0$  if for all  $x, y \in \mathbb{R}$ ,

$$|f(x) - f(y)| \leq L|x - y| \quad \text{✓}$$

Observe that “Lipschitz continuity”  $\implies$  “continuity”.

**Definition 11** (Derivative). The derivative of  $f : S \rightarrow \mathbb{R}$  at  $x \in S \subset \mathbb{R}$  is defined as

$$f'(x) = \lim_{t \rightarrow 0} \frac{f(x+t) - f(x)}{t} = \frac{df}{dx}(x) \quad \text{✓}$$

**Definition 12** (Differentiability). A univariate function  $f : S \rightarrow \mathbb{R}$  with domain  $S \subseteq \mathbb{R}$  is **differentiable** if  $S$  is an open set and the derivative  $f'(x)$  exists for every  $x \in S$ . It is **continuously differentiable** if  $f'$  is continuous. ✓

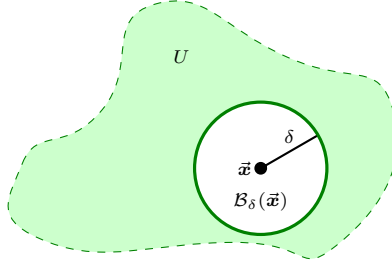


Illustration of an open set.

**Definition 13** (Open Set). An open set  $U \subseteq \mathbb{R}^n$  is a set that contains an open ball around each point, i.e.,

$$\forall \vec{x} \in U, \quad \exists \delta > 0, \quad B_\delta(\vec{x}) = \{\vec{y} \in \mathbb{R}^n : \|\vec{y} - \vec{x}\| < \delta\} \subseteq U \quad \text{✓}$$

### 3.6 Global vs Local Optimization

We say that an optimization problem is solved globally if a **global minimizer** is found; otherwise, if only a **local minimizer** is found, the problem is solved locally. Recall that the feasible set is denoted by  $\Omega$ .

**Definition 14** (Global Minimizer). A vector  $\vec{x}^* \in \Omega$  is a **global minimizer** of  $f$  if

$$f(\vec{x}^*) \leq f(\vec{x}), \quad \forall \vec{x} \in \Omega \quad \text{✓}$$

**Definition 15** (Local Minimizer). A vector  $\vec{x}^* \in \Omega$  is a **local minimizer** of  $f$  if

$$f(\vec{x}^*) \leq f(\vec{x})$$

for any  $\vec{x}$  in a neighborhood of  $\vec{x}^*$ . ✓

**Definition 16** (Strict Local Minimizer). A vector  $\vec{x}^* \in \Omega$  is a **strict local minimizer** of  $f$  if

$$f(\vec{x}^*) < f(\vec{x})$$

for any  $\vec{x}$  in a neighborhood of  $\vec{x}^*$ . ✓

**Definition 17** (Isolated Local Minimizer). A vector  $\vec{x}^* \in \Omega$  is an **isolated local minimizer** of  $f$  if it is the only local minimizer in a neighborhood of  $\vec{x}^*$ . ✓

**Remark 2.** An example for a neighborhood  $\mathcal{N}(\vec{x}^*)$  is the open ball

$$\mathcal{B}_r(\vec{x}^*) := \{\vec{x} \in \mathbb{R}^n : \|\vec{x} - \vec{x}^*\| < r\}$$

with some radius  $r > 0$ . ◀

**Remark 3.**

- An isolated local minimizer is always a strict minimizer but the converse is not true: there are strict minimizers that are not isolated.
- A global minimizer can also be strict (by replacing  $\leq$  with  $<$  in Definition 14). It can also be isolated (if there exists a neighborhood with no local minimizer). ◀

### 3.7 Stochastic vs Deterministic Optimization

The distinction here refers to the type of algorithms used to solve the optimization problem. A stochastic optimization algorithm incorporates some randomness in its iterations. Such algorithm may find different local minimizers even if it is initialised at the same point.



## 4 Unconstrained Optimization

### 4.1 Univariate Objective Functions

The general form of a univariate unconstrained optimization problem is

$$\min_{x \in \mathbb{R}} f(x)$$

with  $f : \mathbb{R} \rightarrow \mathbb{R}$  the objective function. Since there are no equality or inequality constraints,  $x$  may take any real value. In this course we restrict ourselves to smooth functions (see 3.5).

In this chapter (and later ones) we will see that the concept of derivatives is fundamental for solving optimization problems; in particular, derivatives allow us:

1. to construct local approximating models (via Taylor's theorem),
2. to characterise the solutions of the problem (i.e., by deriving optimality conditions).

#### 4.1.1 Taylor Expansions for Univariate Functions

We recall the definitions of the first and second derivatives for a univariate function.

**Definition 18** (First Derivative). For a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  and a point  $x \in \mathbb{R}$ , the first (or first-order) derivative is defined as

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

□

**Definition 19** (Second Derivative). The second (or second-order) derivative of  $f$  at  $x$  is defined as

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

□

Recall that  $f'(x)$  measures the rate of change (the slope of the tangent line) of  $f$  at  $x$ . Thus for small  $d$  one has

$$f(x+d) \approx f(x) + f'(x)d$$

In fact, the tangent line at  $x$  is given by

$$t(y; x) = f(x) + f'(x)(y - x), \quad y \in \mathbb{R}$$

or, equivalently, writing  $y = x + d$ ,

$$t(x+d; x) = f(x) + f'(x)d$$

This linear approximation is known as the **first-order Taylor approximation** of  $f$ .

By incorporating the second derivative, one can construct a local quadratic model (the second-order Taylor approximation):

$$f(x+d) \approx q(x+d; x) = f(x) + f'(x)d + \frac{1}{2}f''(x)d^2, \quad d \in \mathbb{R} \quad (3)$$

or equivalently,

$$f(y) \approx q(y; x) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(x)(y - x)^2, \quad y \in \mathbb{R} \quad (4)$$

Here,  $f''(x)$  is sometimes called the **curvature** of  $f$  at  $x$ . Note that these approximations are valid in a neighbourhood of  $x$  (i.e. for small  $d = y - x$ ).

Taylor's theorem provides an exact expansion.

**Theorem 2** (First and Second Order Taylor's Expansion for Univariate Functions). If  $f : \mathbb{R} \rightarrow \mathbb{R}$  is continuously differentiable, then for any  $x \in \mathbb{R}$  and  $d \in \mathbb{R}$  there exists  $t \in (0, 1)$  such that

$$f(x+d) = f(x) + f'(x+td)d$$

Moreover, if  $f$  is twice continuously differentiable, then for any  $x \in \mathbb{R}$  and  $d \in \mathbb{R}$  there exists  $t \in (0, 1)$  such that

$$f(x+d) = f(x) + f'(x)d + \frac{1}{2}f''(x+td)d^2$$

◁

A useful consequence of Taylor's theorem is a bound on the approximation errors of the linear and quadratic models.

**Corollary 3.** If  $f$  is twice continuously differentiable then for any  $x \in \mathbb{R}$  there exist constants  $L > 0$  and  $\epsilon > 0$  such that

$$|f(y) - t(y; x)| \leq L(y - x)^2$$

for all  $y \in \mathbb{R}$  with  $|y - x| \leq \epsilon$ . If  $f$  is three times continuously differentiable, there exist  $L > 0$  and  $\epsilon > 0$  such that

$$|f(y) - q(y; x)| \leq L|y - x|^3$$

for all  $y \in \mathbb{R}$  with  $|y - x| \leq \epsilon$ . ◁

That is, for  $|y - x|$  sufficiently small, the error of the linear model is of order  $O((y - x)^2)$  and the quadratic model approximates  $f$  with error of order  $O(|y - x|^3)$ .

Actually Taylor's expansion does not stop at the second derivatives. It can also include higher-order derivatives of  $f$  (if they exist) in order to construct increasingly better approximation of  $f$ :

$$f(x + d) \approx f(x) + \frac{1}{1!}f'(x)d + \frac{1}{2!}f''(x)d^2 + \frac{1}{3!}f^{(3)}(x)d^3 + \cdots + \frac{1}{k!}f^{(k)}(x)d^k$$

Nonetheless in optimization we will essentially need the first and second order approximations.

#### 4.1.2 First and Second Order Optimality Conditions for Univariate Functions

In unconstrained optimization we can derive “recipes” to characterise solutions of our problem. These “recipes” are optimality conditions.

**Theorem 4** (Necessary First-Order Optimality Condition for Univariate Functions). If  $f$  is continuously differentiable and  $x^*$  is a local minimizer of  $f$ , then  $f'(x^*) = 0$ , i.e.  $x^*$  is a **stationary point**. ◁

**Proof.** (Contradiction) Assume  $f'(x^*) \neq 0$ . We will use the fact that around  $x^*$  the tangent line approximates  $f$  so that if its slope is not null, we can decrease  $f$  by either decreasing or increasing  $x^*$ . Let  $d := -f'(x^*)$ . Thus  $df'(x^*) < 0$  and since  $f'$  is continuous there exist  $t_0$  such that  $\forall t \in [0, t_0]$

$$df'(x^* + td) < 0$$

Let  $t \in [0, t_0]$ . By the first part of Taylor's theorem, there exists  $\bar{t} \in (0, t)$  such that

$$f(x^* + td) = f(x^*) + tdf'(x^* + \bar{t}d)$$

Since  $\bar{t} \leq t \leq t_0$  we have  $df'(x^* + \bar{t}d) < 0$  which implies  $f(x^* + td) < f(x^*)$ . Since one can choose  $t$  arbitrarily close to 0, this proves that  $x^*$  is not a local minimizer. ◻

It is important to note that stationary points may correspond to local minimizers, local maximizers, or saddle points. Therefore, further analysis (using second derivatives) is often required.

**Theorem 5** (Necessary Second-Order Optimality Condition). If  $f$  is twice continuously differentiable and  $x^*$  is a local minimizer of  $f$ , then

$$f''(x^*) \geq 0$$

◁

**Proof.** (Contradiction) Assume  $f''(x^*) < 0$ . Note that if  $f'(x^*) \neq 0$  we already proved that  $x^*$  is NOT a local minimizer so we can assume that  $f'(x^*) = 0$ . Let  $d \neq 0$ . Since  $f''(x^*) < 0$ , we have  $f''(x^*)d^2 < 0$ . By continuity of  $f''$ , there exists  $t_0 > 0$  such that

$$f''(x^* + td)d^2 < 0, \quad \forall t \in (0, t_0)$$

Let  $t \in (0, t_0)$ . By the second part of Taylor's theorem there exists  $\bar{t} \in (0, t)$  such that

$$f(x^* + td) = f(x^*) + \frac{1}{2}f''(x^* + \bar{t}d)t^2d^2$$

Since  $\bar{t} \leq t \leq t_0$ ,  $f''(x^* + \bar{t}d) < 0$ , which implies that  $f(x^* + td) < f(x^*)$ . As before we can conclude that  $x^*$  is NOT a local minimizer. ◻

However, the above conditions are necessary but not sufficient. For example, if  $f'(x) = 0$  and  $f''(x) = 0$  the test is inconclusive.

A sufficient condition is obtained by strengthening the second-order condition.

**Theorem 6** (Sufficient Optimality Conditions). If  $f$  is twice continuously differentiable,  $f'(x^*) = 0$  and  $f''(x^*) > 0$ , then  $x^*$  is a strict local minimizer of  $f$ . ◀

**Proof.** Since  $f''$  is continuous, there exists  $r > 0$  such that for all  $d$  with  $|d| < r$  we have  $f''(x^* + d) > 0$ . Then by the second-order Taylor expansion, for all such  $d$  (with  $d \neq 0$ ) there exists  $t \in (0, 1)$  such that

$$f(x^* + d) = f(x^*) + \frac{1}{2} f''(x^* + td) d^2 > f(x^*)$$

Thus  $x^*$  is a strict local minimizer. ◻

**Remark 4.** Note that the sufficient condition is not necessary; for instance,  $f(x) = x^4$  has a strict local minimizer at  $x = 0$  despite  $f''(0) = 0$ . ◀

#### 4.1.3 Minimizers for Convex Univariate Functions

Identifying minimizers is considerably simpler when  $f$  is convex (Definition 6). If  $f$  is twice differentiable, the following characterisations are equivalent:

$$\begin{aligned} f \text{ is convex} &\iff f''(x) \geq 0, \quad \forall x \in \mathbb{R} \\ &\iff f(y) \geq f(x) + f'(x)(y - x), \quad \forall x, y \in \mathbb{R} \end{aligned}$$

Two important properties follow immediately:

**Proposition 7.** If  $f$  is convex, any local minimizer is a global minimizer. ◀

**Proposition 8.** If  $f$  is convex and differentiable, then any stationary point is a global minimizer. ◀

## 4.2 Multivariate Objective Functions

The general form of a multivariate unconstrained optimization problem is

$$\min_{\vec{x} \in \mathbb{R}^n} f(\vec{x})$$

with  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  the objective function, i.e. its argument is a vector with  $n$  coordinates:

$$f(\vec{x}) = f(x_1, \dots, x_n)$$

We will again assume that  $f$  is smooth (see 3.5): A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is said to be ***L-smooth*** if it is continuously differentiable and its gradient  $\nabla f$  is Lipschitz continuous with Lipschitz constant  $L > 0$ :

$$\|\nabla f(\vec{x}) - \nabla f(\vec{y})\| \leq L \|\vec{x} - \vec{y}\|, \quad \forall \vec{x}, \vec{y} \in \mathbb{R}^n \quad (5)$$

If  $f$  is twice continuously differentiable, this is equivalent to  $\|\underline{H}_f(\vec{x})\|_2 \leq L$  for all  $\vec{x} \in \mathbb{R}^n$ , where  $\|\cdot\|_2$  denotes the spectral norm (Definition 1). Additionally, we assume that  $f$  is bounded from below, i.e., there exists  $M \in \mathbb{R}$  such that  $f(\vec{x}) \geq M$  for all  $\vec{x} \in \mathbb{R}^n$ .

### 4.2.1 First Order Derivatives for Multivariate Functions

The first-order derivative of a multivariate function is given by its **gradient**. To define the gradient we first recall partial derivatives and introduce the canonical basis.

We denote by  $\{\vec{e}_i\}_{i=1, \dots, n}$  the canonical basis of  $\mathbb{R}^n$  where

$$\vec{e}_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \quad (\text{with the } i\text{-th coordinate equal to } 1)$$

**Definition 20** (Partial Derivative). Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . The partial derivative of  $f$  with respect to the  $i$ -th coordinate at a point  $\vec{x} \in \mathbb{R}^n$  is defined as

$$\frac{\partial f}{\partial x_i}(\vec{x}) = \lim_{h \rightarrow 0} \frac{f(\vec{x} + h \vec{e}_i) - f(\vec{x})}{h} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{h}$$

◀

**Definition 21** (Gradient). Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . The **gradient** of  $f$  at  $\vec{x}$  is the vector

$$\nabla f(\vec{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(\vec{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\vec{x}) \end{pmatrix} \in \mathbb{R}^n$$

◀

The gradient gives the direction of steepest increase of  $f$  at  $\vec{x}$  and each coordinate quantifies the rate of change of  $f$  in the corresponding direction. We say that  $f$  is differentiable if  $\nabla f(\vec{x})$  exists for all  $\vec{x}$ , and continuously differentiable if  $\nabla f$  is continuous.

The gradient can also be used to compute the change in any direction. If we define the univariate function

$$g(t) = f(\vec{x} + t \vec{d}), \quad t \in \mathbb{R}$$

for a direction  $\vec{d} \in \mathbb{R}^n$ , then by the chain rule we have

$$g'(t) = \frac{d}{dt} g(t) = \frac{d}{dt} (f(\vec{x} + t \vec{d})) = \nabla f(\vec{x} + t \vec{d})^T \vec{d}$$

Thus, the **directional derivative** of  $f$  at  $\vec{x}$  in the direction  $\vec{d}$  is given by

$$g'(0) = \nabla f(\vec{x})^T \vec{d}$$

Among all unit vectors  $\vec{d}$  (i.e. with  $\|\vec{d}\| = 1$ ), the maximum value of  $g'(0)$  is attained when  $\vec{d}$  is aligned with  $\nabla f(\vec{x})$ .

#### 4.2.2 First Order Taylor's Expansion for Multivariate Functions

A first-order Taylor expansion (or linear local approximation) of  $f$  at  $\vec{x}$  is

$$t(\vec{y}; \vec{x}) = f(\vec{x}) + \nabla f(\vec{x})^T (\vec{y} - \vec{x}), \quad \vec{y} \in \mathbb{R}^n$$

That is,

$$f(\vec{x} + \vec{d}) \approx f(\vec{x}) + \nabla f(\vec{x})^T \vec{d}$$

**Theorem 9** (First Order Taylor's Expansion for Multivariate Functions). Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable. Then for any  $\vec{x} \in \mathbb{R}^n$  and  $\vec{d} \in \mathbb{R}^n$  there exists  $t \in (0, 1)$  such that

$$f(\vec{x} + \vec{d}) = f(\vec{x}) + \nabla f(\vec{x} + t \vec{d})^T \vec{d}$$

If we set  $n = 1$  we recover the univariate case (Theorem 2).

◀

**Corollary 10.** If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable, then for any  $\vec{x} \in \mathbb{R}^n$  there exist constants  $L > 0$  and  $\epsilon > 0$  such that for all  $\vec{y}$  in

$$\mathcal{B}_\epsilon(\vec{x}) = \{\vec{x} \in \mathbb{R}^n : \|\vec{x} - \vec{x}\| \leq \epsilon\}$$

the error in the linear approximation satisfies

$$|f(\vec{y}) - t(\vec{y}; \vec{x})| \leq L \|\vec{y} - \vec{x}\|^2$$

If we set  $n = 1$  we recover the univariate case (Corollary 3).

◀

### 4.2.3 First Order Optimality Conditions for Multivariate Functions

Recall the definitions for local and global minimizers (Definitions 15 and 14).

**Theorem 11** (Necessary First-Order Optimality Condition). If  $f$  is continuously differentiable and  $\vec{x}^*$  is a local minimizer of  $f$  then  $\nabla f(\vec{x}^*) = \vec{0}$ , i.e.  $\vec{x}^*$  is a **stationary point**. If we set  $n = 1$  we recover the univariate case (Theorem 4).  $\triangleleft$

**Proof.** (Contradiction) Assume that  $\nabla f(\vec{x}^*) \neq \vec{0}$ . Define the direction  $\vec{d} = -\nabla f(\vec{x}^*)$ . Consider the univariate function

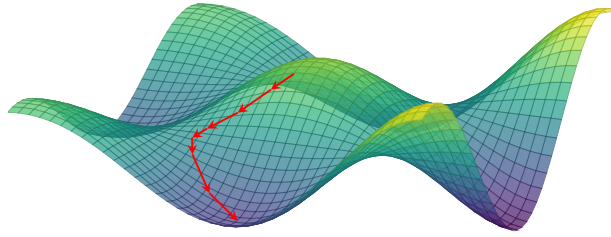
$$g(t) = f(\vec{x}^* + t\vec{d}), \quad t \in \mathbb{R}$$

Then,  $g'(0) = \nabla f(\vec{x}^*)^T \vec{d} = -\|\nabla f(\vec{x}^*)\|^2 < 0$ . Thus, for sufficiently small  $t > 0$  we have  $g(t) < g(0)$ , meaning

$$f(\vec{x}^* + t\vec{d}) < f(\vec{x}^*)$$

This contradicts the local minimality of  $\vec{x}^*$ . Hence, it must be that  $\nabla f(\vec{x}^*) = \vec{0}$ .  $\square$

### 4.3 Gradient Descent



Since the gradient  $\nabla f(\vec{x})$  points in the direction of steepest increase, the negative gradient  $-\nabla f(\vec{x})$  gives the direction of steepest descent. Recall that for any direction  $\vec{d} \in \mathbb{R}^n$  the directional derivative of  $f$  at  $\vec{x}$  in the direction  $\vec{d}$  is  $g'(0) = \nabla f(\vec{x})^T \vec{d}$ . We can write

$$\nabla f(\vec{x})^T \vec{d} = \cos(\theta) \|\nabla f(\vec{x})\| \|\vec{d}\|$$

where  $\theta$  is the angle between  $\vec{d}$  and  $\nabla f(\vec{x})$ . Among all unit directions (i.e.  $\|\vec{d}\| = 1$ ), the decrease is maximized when  $\cos(\theta) = -1$ , i.e. when  $\vec{d}$  is collinear with  $-\nabla f(\vec{x})$ , but points in the opposite direction. Since it is of unit norm, we have

$$\vec{d} = -\frac{\nabla f(\vec{x})}{\|\nabla f(\vec{x})\|}$$

which is called the **direction of steepest descent**.

A gradient descent algorithm then generates a sequence  $\{\vec{x}^{(k)}\}$  via

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \alpha \nabla f(\vec{x}^{(k)})$$

where  $\alpha > 0$  is the step size. A common stopping rule is to terminate when  $\|\nabla f(\vec{x}^{(k)})\|$  falls below a prescribed tolerance level  $\epsilon$  or after a maximum number of iterations  $k_{\max}$ .

#### 4.3.1 Gradient Descent Algorithm

---

##### Algorithm 1 Gradient Descent

---

- 1: **Input:** Starting point  $\vec{x}^{(0)}$ , step size  $\alpha > 0$ , tolerance  $\epsilon > 0$ , and maximum iterations  $k_{\max}$
  - 2: Set  $k \leftarrow 0$
  - 3: Compute  $\vec{g}^{(0)} = \nabla f(\vec{x}^{(0)})$
  - 4: **while**  $\|\vec{g}^{(k)}\| > \epsilon$  and  $k < k_{\max}$  **do**
  - 5:     Update:  $\vec{x}^{(k+1)} = \vec{x}^{(k)} - \alpha \vec{g}^{(k)}$
  - 6:     Set  $k \leftarrow k + 1$
  - 7:     Compute  $\vec{g}^{(k)} = \nabla f(\vec{x}^{(k)})$
  - 8: **Output:**  $\vec{x}^{(k)}$
-

### 4.3.2 Convergence Analysis

**Theorem 12** (Global Convergence of Gradient Descent). Let  $f$  be twice continuously differentiable and  $L$ -smooth (i.e. its gradient is Lipschitz continuous with constant  $L > 0$ ). Assume that  $f$  is bounded from below by  $m$  and let  $0 < \alpha < \frac{2}{L}$ . Then for any starting point  $\vec{x}^{(0)}$ , the iterates  $\{\vec{x}^{(k)}\}$  produced by Algorithm 1 satisfy

$$\lim_{k \rightarrow \infty} \nabla f(\vec{x}^{(k)}) = \vec{0} \quad (6)$$

◁

**Proof.** Since  $f$  is  $L$ -smooth, for any  $\vec{x}, \vec{x}' \in \mathbb{R}^n$  one has

$$f(\vec{x}') \leq f(\vec{x}) + \nabla f(\vec{x})^T (\vec{x}' - \vec{x}) + \frac{L}{2} \|\vec{x}' - \vec{x}\|^2 \quad (7)$$

Apply this inequality with  $\vec{x} = \vec{x}^{(k)}$  and  $\vec{x}' = \vec{x}^{(k+1)} = \vec{x}^{(k)} - \alpha \nabla f(\vec{x}^{(k)})$ . Then

$$\begin{aligned} f(\vec{x}^{(k+1)}) &\leq f(\vec{x}^{(k)}) - \alpha \|\nabla f(\vec{x}^{(k)})\|^2 + \frac{\alpha^2 L}{2} \|\nabla f(\vec{x}^{(k)})\|^2 \\ &= f(\vec{x}^{(k)}) - \underbrace{\left( \alpha - \frac{\alpha^2 L}{2} \right)}_{=: \beta} \|\nabla f(\vec{x}^{(k)})\|^2 \end{aligned} \quad (8)$$

Since  $\alpha < \frac{2}{L}$ , we have  $\beta = \alpha - \frac{\alpha^2 L}{2} = \alpha \left( 1 - \frac{\alpha L}{2} \right) > 0$ . Then,

$$f(\vec{x}^{(k+1)}) \leq f(\vec{x}^{(k)}) - \beta \|\nabla f(\vec{x}^{(k)})\|^2 \quad (9)$$

We can sum these inequalities until iteration  $k$  (accumulation of progress):

$$f(\vec{x}^{(k)}) \leq f(\vec{x}^{(0)}) - \beta \sum_{i=0}^{k-1} \|\nabla f(\vec{x}^{(i)})\|^2$$

Since  $f$  is bounded,  $m \leq f(\vec{x}^{(k)}) \leq f(\vec{x}^{(0)}) - \beta \sum_{i=0}^{k-1} \|\nabla f(\vec{x}^{(i)})\|^2$ . Therefore  $\sum_{i=0}^{k-1} \|\nabla f(\vec{x}^{(i)})\|^2 < \infty$ . But this implies that  $\lim_{k \rightarrow \infty} \|\nabla f(\vec{x}^{(k)})\| = 0$ , which in turn implies (6).  $\square$

**Remark 5.**

- Using the update inequality (9), the best step size  $\alpha^* \in (0, \frac{2}{L})$  is obtained by maximizing  $\beta(\alpha) = \alpha - \frac{\alpha^2 L}{2}$

$$\beta'(\alpha) = 1 - \alpha L \stackrel{!}{=} 0 \Rightarrow \alpha^* = \frac{1}{L}$$

- It is a **global** convergence result, since it is valid for any starting point  $\vec{x}^{(0)}$ .
- It does NOT necessarily imply that the iterates  $\{\vec{x}^{(k)}\}$  converge.
- Often,  $L$  is unknown and we need to carefully tune the step size  $\alpha$ .

◀

**Definition 22** (Global Convergence). An optimization algorithm is globally convergent if, from any starting point  $\vec{x}^{(0)}$ , its iterates  $\vec{x}^{(k)}$  satisfy

$$\nabla f(\vec{x}^{(k)}) \xrightarrow[k \rightarrow \infty]{} \vec{0}$$

◀

**Definition 23** (Convergence of Iterates). We say that the sequence of iterates  $\{\vec{x}^{(k)}\}$  converges to  $\vec{x}^*$  if

$$\|\vec{x}^{(k)} - \vec{x}^*\| \xrightarrow[k \rightarrow \infty]{} 0$$

◀

### 4.3.3 Convergence Rate for Strongly Convex Functions

**Lemma 13.** If  $f$  is differentiable, being  $\mu$ -strongly convex (Definition 8) is equivalent to

$$f(\vec{y}) \geq f(\vec{x}) + \nabla f(\vec{x})^T (\vec{y} - \vec{x}) + \frac{\mu}{2} \|\vec{y} - \vec{x}\|^2 \quad (10)$$

and we have the inequality

$$2\mu(f(\vec{x}) - f(\vec{x}^*)) \leq \|\nabla f(\vec{x})\|^2, \quad \forall \vec{x} \in \mathbb{R}^n \quad (11)$$

where  $\vec{x}^*$  is the unique global minimizer of  $f$ .  $\triangleleft$

Equation (11) says that the gradient norm  $\|\nabla f(\vec{x})\|$  grows with the **optimality**  $f(\vec{x}) - f(\vec{x}^*)$ . This allows us to quantify how close we are from the global minimum by looking at the gradient norm.

**Proof.** Rearranging equation (10) with  $\vec{y} = \vec{x}^*$  yields

$$\begin{aligned} f(\vec{x}) - f(\vec{x}^*) &\leq -\nabla f(\vec{x})^T (\vec{x}^* - \vec{x}) - \frac{\mu}{2} \|\vec{x}^* - \vec{x}\|^2 \\ 2\mu(f(\vec{x}) - f(\vec{x}^*)) &\leq -2\mu\nabla f(\vec{x})^T (\vec{x}^* - \vec{x}) - \mu^2 \|\vec{x}^* - \vec{x}\|^2. \end{aligned}$$

Adding  $0 = \|\nabla f(\vec{x})\|^2 - \|\nabla f(\vec{x})\|^2$  to the right-hand side, allows us to rewrite the last line as

$$\begin{aligned} 2\mu(f(\vec{x}) - f(\vec{x}^*)) &\leq \|\nabla f(\vec{x})\|^2 - \|\nabla f(\vec{x})\|^2 - 2\mu\nabla f(\vec{x})^T (\vec{x}^* - \vec{x}) - \mu^2 \|\vec{x}^* - \vec{x}\|^2 \\ &= \|\nabla f(\vec{x})\|^2 - \left( \|\nabla f(\vec{x})\|^2 + 2\mu\nabla f(\vec{x})^T (\vec{x}^* - \vec{x}) + \mu^2 \|\vec{x}^* - \vec{x}\|^2 \right) \\ &= \|\nabla f(\vec{x})\|^2 - \|\nabla f(\vec{x}) + \mu(\vec{x}^* - \vec{x})\|^2 \\ &\leq \|\nabla f(\vec{x})\|^2. \end{aligned}$$

**Theorem 14.** Let  $f$  be  $L$ -smooth and  $\mu$ -strongly convex with unique minimizer  $\vec{x}^*$ . If the step size satisfies  $\alpha \in (0, \frac{1}{L}]$ , then for every iteration  $k \geq 0$  of gradient descent,

$$f(\vec{x}^{(k+1)}) - f(\vec{x}^*) \leq (1 - \alpha\mu) (f(\vec{x}^{(k)}) - f(\vec{x}^*))$$

**Proof.** From the smoothness of  $f$  we have already derived (Equation 8) that

$$f(\vec{x}^{(k+1)}) \leq f(\vec{x}^{(k)}) - \alpha \left(1 - \frac{\alpha L}{2}\right) \|\nabla f(\vec{x}^{(k)})\|^2$$

Since  $\alpha \leq \frac{1}{L}$ , we have  $(1 - \frac{\alpha L}{2}) \geq (1 - \frac{1}{2}) = \frac{1}{2}$  and therefore

$$f(\vec{x}^{(k+1)}) - f(\vec{x}^*) \leq f(\vec{x}^{(k)}) - f(\vec{x}^*) - \frac{\alpha}{2} \|\nabla f(\vec{x}^{(k)})\|^2$$

Using Equation (11) from Lemma 13,  $2\mu(f(\vec{x}^{(k)}) - f(\vec{x}^*)) \leq \|\nabla f(\vec{x}^{(k)})\|^2$ , we obtain

$$f(\vec{x}^{(k+1)}) - f(\vec{x}^*) \leq f(\vec{x}^{(k)}) - f(\vec{x}^*) - \alpha\mu(f(\vec{x}^{(k)}) - f(\vec{x}^*))$$

which simplifies to the stated result.  $\square$

Theorem 14 shows that the error in function value decreases by a factor of  $1 - \alpha\mu$  at each iteration. In other words, gradient descent converges **linearly** with rate  $C = 1 - \alpha\mu$ .

**Definition 24** (Linear Convergence). Let  $\{z^{(k)}\}$  be a sequence with  $z^{(k)} \rightarrow z^*$  and there exists  $C \in [0, 1)$  and  $\hat{k}$  such that

$$\|z^{(k+1)} - z^*\| \leq C \|z^{(k)} - z^*\|$$

for all  $k \geq \hat{k}$ . Then the sequence converges **linearly** to  $z^*$  with rate  $C$ .  $\triangleleft$

**Remark 6.**

- The smaller the linear rate  $C$ , the faster is the convergence.
- The rate  $1 - \alpha\mu$  is the “worst-case” rate, but it can be faster for certain starting points.
- The best theoretical rate is when  $1 - \alpha\mu$  is smallest, i.e.,  $\alpha^* = \frac{1}{L}$  (we do not always know  $L$ !)

- Linear rate is generally considered slow, but unfortunately this is the best “worst-case” rate a first-order method can generally achieve!
- In fact if  $f$  is only convex, the worst-case rate is sub-linear, i.e., slower than any linear rate.

◀

For strongly convex and smooth objective functions, we can also determine how many iterations we need until  $f(\vec{x}^{(k)}) - f(\vec{x}^*) \leq \epsilon$  for a given  $\epsilon$ .

Since the error decreases

$$f(\vec{x}^{(k)}) - f(\vec{x}^*) \leq (1 - \alpha\mu)^k (f(\vec{x}^{(0)}) - f(\vec{x}^*))$$

then to guarantee that

$$f(\vec{x}^{(k)}) - f(\vec{x}^*) \leq \epsilon$$

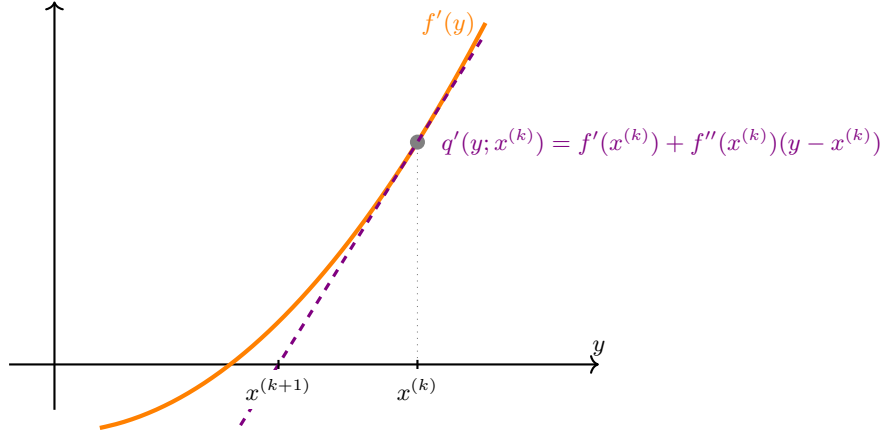
it suffices to have

$$k \geq \frac{1}{-\log(1 - \alpha\mu)} \log \left( \frac{f(\vec{x}^{(0)}) - f(\vec{x}^*)}{\epsilon} \right)$$

i.e. the number of iterations is  $O(\log \frac{1}{\epsilon})$ .

## 4.4 Newton's Method

### 4.4.1 Newton's Method for Univariate Functions



Recall (Equation 4) the second-order Taylor approximation for a univariate  $f$ . Since  $f \approx q(y; x^{(k)})$ , it makes sense to choose  $x^{(k+1)}$  as the minimizer of  $q(y; x^{(k)})$ , i.e.

$$\begin{aligned} x^{(k+1)} &= \arg \min_y q(y; x^{(k)}) \Rightarrow q'(y; x^{(k)}) = f'(x^{(k)}) + f''(x^{(k)})(y - x^{(k)}) \stackrel{!}{=} 0 \\ \implies x^{(k+1)} &= x^{(k)} - \frac{f'(x^{(k)})}{f''(x^{(k)})} \end{aligned} \quad (12)$$

and if  $q''(x^{(k+1)}; x^{(k)}) = f''(x^{(k)}) > 0$ , then  $x^{(k+1)}$  is the global minimizer of  $q(\cdot; x^{(k)})$ , since the latter is convex. The direction

$$d_N = -\frac{f'(x^{(k)})}{f''(x^{(k)})}$$

is called the Newton's direction and the step (12) called the Newton's step (see Algorithm 2).

---

#### Algorithm 2 Newton's Algorithm

---

- 1: **Input:** Starting point  $x^{(0)}$ , tolerance  $\epsilon > 0$ , and maximum iterations  $k_{\max}$
  - 2: Set  $k \leftarrow 0$
  - 3: **while**  $|f'(x^{(k)})| > \epsilon$  and  $k \leq k_{\max}$  **do**
  - 4:     Compute  $f'(x^{(k)})$  and  $f''(x^{(k)})$
  - 5:     Update:  $x^{(k+1)} = x^{(k)} - \frac{f'(x^{(k)})}{f''(x^{(k)})}$
  - 6:     Set  $k \leftarrow k + 1$
  - 7: **Output:**  $x^{(k)}$
-



**Remark 7.** Newton's algorithm does **not** work when  $f''(x^{(k)}) \leq 0$ , but there are possible remedies. If  $f$  is convex, then  $f''(x) \geq 0$  for any  $x$ , and if  $f''(x^{(k)}) > 0$ , the Newton's step is well defined. ◀

**Remark 8.** Originally, Newton's algorithm is a method for finding the root of a differentiable function, i.e., finding  $x$  such that  $g(x) = 0$ . The form of the iterates comes from the first-order Taylor approximation:

$$g(y) \approx t(y; x^{(k)}) = g(x^{(k)}) + g'(x^{(k)})(y - x^{(k)})$$

The next iterate  $x^{(k+1)}$  is chosen as the root of  $t(y; x^{(k)})$ , i.e.,  $t(x^{(k+1)}; x^{(k)}) = 0$ , which gives

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(k)})}$$

In optimization, we use Newton's method to find a stationary point (which **may** be a local minimizer and **is** a global minimizer if  $f$  is convex). Therefore, we use Newton's method to find a root of  $g = f'$ , which corresponds to the update

$$x^{(k+1)} = x^{(k)} - \frac{f'(x^{(k)})}{f''(x^{(k)})}$$

which is the same as (12), obtained from minimizing the second-order Taylor approximation. If the objective is not convex, we need to check that the stationary point found by Newton's method is a local minimizer. ◀

#### 4.4.2 Second Order Taylor's Expansion for Multivariate Functions

Assume  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable. The second-order Taylor approximation is

$$f(\vec{y}) \approx q(\vec{y}; \vec{x}) = f(\vec{x}) + \nabla f(\vec{x})^T (\vec{y} - \vec{x}) + \frac{1}{2} (\vec{y} - \vec{x})^T \underline{H}_f(\vec{x}) (\vec{y} - \vec{x}) \quad (13)$$

where  $\underline{H}_f(\vec{x})$  is the Hessian matrix defined as

$$\underline{H}_f(\vec{x}) = \left( \frac{\partial^2 f}{\partial x_i \partial x_j} \right)_{i,j \in \{1, \dots, n\}} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix} \in \mathbb{R}^{n \times n} \quad (14)$$

Since  $f$  is twice continuously differentiable, the Hessian is symmetric, i.e.,  $\underline{H}_f(\vec{x}) = \underline{H}_f(\vec{x})^T$ .

We have already seen the first order Taylor's expansion for multivariate functions (Theorem 9, Corollary 10). Now let's extend this to the second order:

**Theorem 15** (Second Order Taylor's Expansion for Multivariate Functions). If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable, then for any  $\vec{x} \in \mathbb{R}^n$  and  $\vec{d} \in \mathbb{R}^n$  there exists  $t \in (0, 1)$  such that

$$f(\vec{x} + \vec{d}) = f(\vec{x}) + \nabla f(\vec{x} + t\vec{d})^T \vec{d}$$

Moreover, if  $f$  is twice continuously differentiable, then for any  $\vec{x} \in \mathbb{R}^n$  and  $\vec{d} \in \mathbb{R}^n$  there exists  $t \in (0, 1)$  such that

$$f(\vec{x} + \vec{d}) = f(\vec{x}) + \nabla f(\vec{x})^T \vec{d} + \frac{1}{2} \vec{d}^T \underline{H}_f(\vec{x} + t\vec{d}) \vec{d}$$

For  $n = 1$ , this reduces to the univariate case (Theorem 2). ◀

**Corollary 16.** If  $f$  is three times continuously differentiable, then for any  $\vec{x} \in \mathbb{R}^n$  there exists  $L > 0$  and  $\epsilon > 0$  such that

$$|f(\vec{y}) - q(\vec{y}; \vec{x})| \leq L \|\vec{y} - \vec{x}\|^3$$

for any  $\vec{y} \in \mathbb{R}^n$  such that  $\|\vec{x} - \vec{y}\| \leq \epsilon$ . For  $n = 1$ , this reduces to the univariate case (Corollary 3). ◀

#### 4.4.3 Second Order Optimality Conditions for Multivariate Functions

**Theorem 17** (Second Order Necessary Optimality Condition). If  $f$  is twice continuously differentiable and  $\vec{x}^*$  is a local minimizer of  $f$ , then  $\underline{H}_f(\vec{x}^*) \succeq 0$ , i.e.,  $\vec{p}^T \underline{H}_f(\vec{x}^*) \vec{p} \geq 0$  for all  $\vec{p} \in \mathbb{R}^n$ .  $\triangleleft$

**Proof.** (Contradiction) Assume  $\vec{x}^*$  is a local minimizer of  $f$  and  $\underline{H}_f(\vec{x}^*) \not\succeq 0$ . Then there exists  $\vec{p} \in \mathbb{R}^n$  such that  $\vec{p}^T \underline{H}_f(\vec{x}^*) \vec{p} < 0$ . Since  $\underline{H}_f(\cdot)$  is continuous there exists  $t_0 > 0$  such that for all  $t \in (0, t_0)$ ,

$$\vec{p}^T \underline{H}_f(\vec{x}^* + t\vec{p}) \vec{p} < 0$$

Let  $t \in (0, t_0)$ . If  $\nabla f(\vec{x}^*) \neq 0$  we already know that  $\vec{x}^*$  cannot be a local minimizer (Theorem 11), so we assume  $\nabla f(\vec{x}^*) = 0$ . Theorem 15 says there exists  $\bar{t} \in (0, t)$  such that

$$f(\vec{x}^* + t\vec{p}) = f(\vec{x}^*) + \frac{1}{2} t^2 \vec{p}^T \underline{H}_f(\vec{x}^* + \bar{t}\vec{p}) \vec{p}$$

and since  $0 < \bar{t} < t < t_0$ , we have  $\vec{p}^T \underline{H}_f(\vec{x}^* + \bar{t}\vec{p}) \vec{p} < 0$ , which implies  $f(\vec{x}^* + t\vec{p}) < f(\vec{x}^*)$ , contradicting the assumption that  $\vec{x}^*$  is a local minimizer.  $\square$

**Theorem 18** (Sufficient Optimality Conditions). If  $f$  is twice continuously differentiable and  $\vec{x}^*$  is such that  $\nabla f(\vec{x}^*) = \vec{0}$  and  $\underline{H}_f(\vec{x}^*) \succ 0$ , then  $\vec{x}^*$  is a strict local minimizer of  $f$ .  $\triangleleft$

**Proof.** Since  $\underline{H}_f(\cdot)$  is continuous, there exists  $r > 0$  such that for any  $\vec{p}$  with  $\|\vec{p} - \vec{x}^*\| < r$ , we have  $\underline{H}_f(\vec{x}^* + \vec{p}) \succ 0$ . If  $\vec{p} \neq \vec{0}$ , then  $\vec{p}^T \underline{H}_f(\vec{x}^* + \vec{p}) \vec{p} > 0$ , and by Theorem 15 there exists  $t \in (0, 1)$  such that

$$f(\vec{x}^* + \vec{p}) = f(\vec{x}^*) + \frac{1}{2} \vec{p}^T \underline{H}_f(\vec{x}^* + t\vec{p}) \vec{p} > f(\vec{x}^*)$$

which shows that  $\vec{x}^*$  is a strict local minimizer.  $\square$

#### 4.4.4 Minimization of Quadratic Functions

For a multivariate function  $f$ , at iteration  $k$  of Newton's method (Algorithm 3), we need to find the minimizer  $\vec{x}^*$  of the second-order Taylor approximation (Equation 13) at  $\vec{x}^{(k)}$  and define the new iterate as  $\vec{x}^{(k+1)} = \vec{x}^*$ . Equation 13 is a multivariate quadratic function of the form

$$q(\vec{x}) = b + \vec{g}^T \vec{x} + \vec{x}^T \underline{Q} \vec{x} \quad (15)$$

$$= b + \sum_{i=1}^n g_i x_i + \sum_{i=1}^n \sum_{j=1}^n q_{ij} x_i x_j \quad (16)$$

where  $\vec{x} = \vec{y} - \vec{x}^{(k)}$ ,  $b = f(\vec{x}^{(k)})$ ,  $\vec{g} = \nabla f(\vec{x}^{(k)})$ , and  $\underline{Q} = \frac{1}{2} \underline{H}_f(\vec{x}^{(k)})$ . We assume that  $\underline{Q}$  is symmetric, which is not restrictive, since any quadratic model with a non-symmetric matrix can be re-written with a symmetric matrix  $\tilde{\underline{Q}} = \frac{1}{2}(\underline{Q} + \underline{Q}^T)$ .

Using Equation (16), it is easier to compute the derivatives:

$$\frac{\partial}{\partial x_k} q(\vec{x}) = \sum_{i=1}^n \frac{\partial}{\partial x_k} g_i x_i + \sum_{i=1}^n \sum_{j=1}^n \frac{\partial}{\partial x_k} q_{ij} x_i x_j = g_k + 2 \sum_{j=1}^n q_{kj} x_j$$

$$\frac{\partial^2}{\partial x_k \partial x_l} q(\vec{x}) = 2q_{kl}$$

or in vectorized form:

$$\nabla q(\vec{x}) = \vec{g} + 2\underline{Q}\vec{x}$$

$$\underline{H}_q(\vec{x}) = 2\underline{Q}$$

A stationary point of  $q$  thus satisfies

$$\nabla q(\vec{x}^*) = \vec{0} \iff 2\underline{Q}\vec{x}^* = -\vec{g}$$

so  $\vec{x}^*$  is the solution of a linear system of equations. If  $\underline{Q}$  is invertible, then  $\vec{x}^* = -\frac{1}{2}\underline{Q}^{-1}\vec{g}$ .

**Remark 9.** It is always less computationally expensive to solve a linear system  $\underline{A}\vec{x} = \vec{b}$  than inverting  $\underline{A}$ , i.e., computing  $\underline{A}^{-1}\vec{b}$ .  $\blacktriangleleft$

Observe that for quadratic functions to be a local minimizer the

- necessary second-order optimality condition is  $\underline{Q} \succeq 0$ ,
- sufficient second-order optimality condition is  $\underline{Q} \succ 0 \wedge 2\underline{Q}\vec{x}^* = -\vec{g}$ .

**Caution 10.** Some important matrix properties. Let  $\underline{A} \in \mathbb{R}^{n \times n}$ .

- The determinant  $\det(\underline{A})$  is equal to the product of the eigenvalues of  $\underline{A}$ .
- The trace  $\text{tr}(\underline{A})$  is equal to the sum of the eigenvalues of  $\underline{A}$ .
- $\underline{A}$  is positive definite ( $\underline{A} \succ 0$ ) if and only if it is symmetric and all its eigenvalues are positive.
- $\underline{A}$  is positive semi-definite ( $\underline{A} \succeq 0$ ) if and only if it is symmetric and all its eigenvalues are non-negative.
- $\underline{A}$  is negative definite ( $\underline{A} \prec 0$ ) if and only if it is symmetric and all its eigenvalues are negative.
- $\underline{A}$  is negative semi-definite ( $\underline{A} \preceq 0$ ) if and only if it is symmetric and all its eigenvalues are non-positive.
- If  $\underline{A}$  is symmetric and has both positive and negative eigenvalues, then it is indefinite.

Now consider  $\vec{x}^*$  a stationary point of  $q(\vec{x}) = b + \vec{g}^T \vec{x} + \vec{x}^T \underline{Q} \vec{x}$ , i.e. a point satisfying  $2\underline{Q}\vec{x}^* = -\vec{g}$ .

- If  $\underline{Q} \succ 0$ , then  $q$  is strictly convex and  $\vec{x}^*$  is the global minimizer of  $q$ .
- If  $\underline{Q} \prec 0$ , then  $q$  is strictly concave and  $\vec{x}^*$  is the global maximizer of  $q$ .
- If  $\underline{Q}$  has at least one negative eigenvalue then  $q$  is unbounded below and has no minimizer.
- If  $\underline{Q}$  is indefinite,  $\vec{x}^*$  is a saddle point of  $q$ .
- If  $\underline{Q} \succeq 0$ , then  $q$  is convex. If  $\underline{Q}$  is not positive definite, then it is singular and in this case  $q$  is either unbounded or has an infinite number of global minimizers.

#### 4.4.5 Newton's Method for Multivariate Functions

---

**Algorithm 3** Newton's Algorithm (Multivariate)

---

```

1: Input: Starting point  $\vec{x}^{(0)}$ , tolerance  $\epsilon > 0$ , and maximum iterations  $k_{\max}$ 
2: Set  $k \leftarrow 0$ 
3: while  $|\nabla f(\vec{x}^{(k)})| > \epsilon$  and  $k \leq k_{\max}$  do
4:   Compute  $\nabla f(\vec{x}^{(k)})$  and  $\underline{H}_f(\vec{x}^{(k)})$ 
5:   Compute  $\vec{u}_{\min}$  as the minimizer of  $q(\vec{u}; \vec{x}^{(k)})$ . ▷ If it exists
6:   Update:  $\vec{x}^{(k+1)} = \vec{u}_{\min}$  ▷ If it exists:  $\vec{x}^{(k+1)} = \vec{x}^{(k)} - (\underline{H}_f(\vec{x}^{(k)})^{-1} \nabla f(\vec{x}^{(k)}))$ 
7:   Set  $k \leftarrow k + 1$ 
8: Output:  $\vec{x}^{(k)}$ 

```

---

In Line 5 of Algorithm 3, we need to find the minimizer of the quadratic local model at  $\vec{x}^{(k)}$  given by the second-order Taylor approximation (Section 4.4.2, Equation 13). If we set  $\vec{y} = \vec{x}^{(k)} + \vec{d}$  in (13), we obtain

$$q(\vec{d}; \vec{x}^{(k)}) = f(\vec{x}^{(k)}) + \nabla f(\vec{x}^{(k)})^T \vec{d} + \frac{1}{2} \vec{d}^T \underline{H}_f(\vec{x}^{(k)}) \vec{d}$$

Denoting  $b = f(\vec{x}^{(k)})$ ,  $\vec{g} = \nabla f(\vec{x}^{(k)})$ , and  $\underline{Q} = \frac{1}{2} \underline{H}_f(\vec{x}^{(k)})$ , we have a quadratic function of the same form as in (15):

$$q(\vec{d}; \vec{x}^{(k)}) = b + \vec{g}^T \vec{d} + \vec{d}^T \underline{Q} \vec{d}$$

Thus a stationary point of  $q$  satisfies

$$\vec{d}^* = -\frac{1}{2} \underline{Q}^{-1} \vec{g} = -(\underline{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)})$$

If  $\underline{Q} \succ 0$ , then  $\vec{d}^*$  is the global minimizer of  $q(\vec{d}; \vec{x}^{(k)})$ , and is the **Newton's direction**, denoted  $\vec{d}_N$ . Hence in the original problem with  $\vec{y} = \vec{x}^{(k)} + \vec{d}$ , the minimizer  $\vec{x}^{(k+1)}$  is

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \vec{d}_N = \vec{x}^{(k)} - (\underline{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)})$$

This update rule is called the **Newton's step**.

**Caution 11.** Two important remarks:

- If  $\underline{H}_f(\vec{x}^{(k)})$  is not positive definite, then we cannot compute the Newton's direction. We need to detect this case by checking if all eigenvalues of  $\underline{H}_f(\vec{x}^{(k)})$  are positive. If they are not, an easy solution is to take a gradient step for that iteration.
- To compute  $\vec{d}_N$ , we should not invert  $\underline{H}_f(\vec{x}^{(k)})$ , but solve the linear system

$$2Q\vec{d}_N = -\vec{g}, \text{ i.e., } \underline{H}_f(\vec{x}^{(k)})\vec{d}_N = -\nabla f(\vec{x}^{(k)})$$

for  $\vec{d}_N$ . ◀

#### 4.4.6 Convergence of Newton's Method

**Theorem 19** (Convergence of Newton's Method). If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable,  $\vec{x}^*$  is a point such that  $\nabla f(\vec{x}^*) = \vec{0} \wedge \underline{H}_f(\vec{x}^*) \succ 0$ , and  $\underline{H}_f(\cdot)$  is  $L$ -Lipschitz continuous with  $L > 0$  in a neighborhood of  $\vec{x}^*$ , then there exists  $\epsilon > 0$  such that for any  $\vec{x}^{(0)}$  with  $\|\vec{x}^{(0)} - \vec{x}^*\| \leq \epsilon$ , at each iteration  $k \geq 0$  of Newton's method (Algorithm 3),

$$\|\vec{x}^{(k+1)} - \vec{x}^*\| \leq L \left\| (\underline{H}_f(\vec{x}^*))^{-1} \right\| \|\vec{x}^{(k)} - \vec{x}^*\|^2 \quad (17)$$

i.e., the iterates  $\{\vec{x}^{(k)}\}$  converge quadratically to  $\vec{x}^*$ . ◁

**Remark 12.** Theorem 19 shows that if Algorithm 3 is initialized sufficiently close to a local minimizer  $\vec{x}^*$  at which the second-order sufficient optimality conditions (Theorem 18) hold, then the iterates converge quadratically to  $\vec{x}^*$ . It is thus a **local convergence** result and not a global convergence result. ◀

**Definition 25** (Convergence Rates). A sequence of vectors  $\{\vec{z}^{(k)}\}$  with  $\vec{z}^{(k)} \rightarrow \vec{z}^*$  converges

- **linearly** if there exists  $c \in (0, 1)$  and  $\hat{k} \geq 0$  such that

$$\|\vec{z}^{(k+1)} - \vec{z}^*\| \leq c \|\vec{z}^{(k)} - \vec{z}^*\| \text{ for all } k \geq \hat{k}.$$

- **quadratically** if there exists  $c \in (0, \infty)$  and  $\hat{k} \geq 0$  such that

$$\|\vec{z}^{(k+1)} - \vec{z}^*\| \leq c \|\vec{z}^{(k)} - \vec{z}^*\|^2 \text{ for all } k \geq \hat{k}.$$

- **sublinearly** if there exists  $\{c^{(k)}\} \subset (0, 1)$  with  $c^{(k)} \rightarrow 1$  and  $\hat{k} \geq 0$  such that

$$\|\vec{z}^{(k+1)} - \vec{z}^*\| \leq c^{(k)} \|\vec{z}^{(k)} - \vec{z}^*\| \text{ for all } k \geq \hat{k}.$$

- **superlinearly** if there exists  $\{c^{(k)}\} \subset (0, \infty)$  with  $c^{(k)} \rightarrow 0$  and  $\hat{k} \geq 0$  such that

$$\|\vec{z}^{(k+1)} - \vec{z}^*\| \leq c^{(k)} \|\vec{z}^{(k)} - \vec{z}^*\| \text{ for all } k \geq \hat{k}. \quad \not\Rightarrow$$

**Remark 13.** Observe that: quadratic  $\implies$  superlinear  $\implies$  linear  $\implies$  sublinear. ◀

Theorem 19 says that when Algorithm 3 is initialized sufficiently close to the local minimizer  $x^*$ , we have  $C = \|(\underline{H}_f(\vec{x}^*))^{-1}\| > 0$  such that

$$\|x^{(k+1)} - x^*\| \leq C \|x^{(k)} - x^*\|^2 \quad (18)$$

Denoting the error at iteration  $k$  by

$$\text{err}_k = \|x^{(k)} - x^*\|$$

we then can rewrite (18) as

$$\text{err}_{k+1} \leq C (\text{err}_k)^2 \quad (19)$$

Recursively applying this estimate leads to

$$\text{err}_K \leq C \left( \underbrace{C \left( \underbrace{C \left( \underbrace{C \text{err}_0^2}_{\geq \text{err}_1} \right)^2}_{\geq \text{err}_2} \right)^2}_{\geq \text{err}_{K-1}} \right)^2 = C^{2^K - 1} (\text{err}_0)^{2^K}$$

For simplicity, assume  $\mathbf{err}_0 = 1$ . To achieve an error  $\mathbf{err}_K \leq \epsilon$  with  $0 < \epsilon \ll 1$ , it suffices that

$$C^{2^K - 1} \leq \epsilon$$

Taking logarithms gives

$$(2^K - 1) \log C \leq \log \epsilon$$

Since  $\log C < 0$  (because  $C < 1$  in the quadratic regime), we can rearrange to obtain

$$2^K \geq 1 + \frac{\log \epsilon}{\log C}$$

and taking logarithms once more yields

$$K \geq \log \left( 1 + \frac{\log \epsilon^{-1}}{\log C^{-1}} \right) = O(\log \log \epsilon^{-1})$$

Thus, to reduce the optimality error to order  $\epsilon$  we need only  $O(\log \log \epsilon^{-1})$  iterations. This is significantly fewer than the  $O(\log \epsilon^{-1})$  iterations typically required by gradient descent.

In its standard form, Algorithm 3 is not globally convergent, but it can be made so with a slight variation of the update rule using [line search](#). Also, computing second-order derivatives is computationally expensive, but also this can be avoided by using [quasi-Newton methods](#) while keeping some of the benefits of Newton's method.

## 4.5 Line Search Algorithms

To make Newton's algorithm globally convergent from any starting point we need more flexibility in the update rule. One possibility is to take a smaller step in the same direction, i.e.,

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha \vec{d}_N = \vec{x}^{(k)} + \underbrace{\alpha \left( -(\mathbf{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)}) \right)}_{\vec{d}_N}$$

with  $\alpha \in (0, 1)$  a step size. This is similar to Gradient Descent where instead of  $\vec{d}_N$  we have  $\vec{d}_S = -\nabla f(\vec{x}^{(k)})$ . In fact, both Gradient Descent and Newton's algorithms (and their variants) fall in the class of [line search](#) algorithms where finding a "good" step size is a sub-routine. A line search method has an update rule of the form

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)}$$

where:

- $\vec{d}^{(k)} \in \mathbb{R}^n$  is a search direction;
- $\alpha^{(k)} > 0$  is a chosen step size.

Performing a line search specifically refers to the sub-routine of choosing the best possible (or at least, good enough)  $\alpha^{(k)}$ , given  $\vec{d}^{(k)}$ . At a minimum we require that we achieve a decrease of  $f$ , i.e.,

$$f(\vec{x}^{(k+1)}) = f(\vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)}) < f(\vec{x}^{(k)})$$

This is always satisfied for sufficiently small  $\alpha^{(k)}$  if  $\vec{d}^{(k)}$  is a descent direction. Recall that if we define

$$g(\alpha) = f(\vec{x}^{(k)} + \alpha \vec{d}^{(k)})$$

then its derivative is  $g'(\alpha) = \nabla f(\vec{x}^{(k)} + \alpha \vec{d}^{(k)})^T \vec{d}^{(k)}$  with  $g'(0) = \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}$ . Since

$$g(\alpha) \approx f(\vec{x}^{(k)}) + g'(0) \alpha$$

we can guarantee  $g(\alpha) < f(\vec{x}^{(k)})$  for sufficiently small  $\alpha > 0$  if

$$g'(0) = \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)} < 0$$

i.e., if  $\vec{d}^{(k)}$  is indeed a descent direction. Unfortunately, this minimal requirement is not enough to choose a good step size in practice.

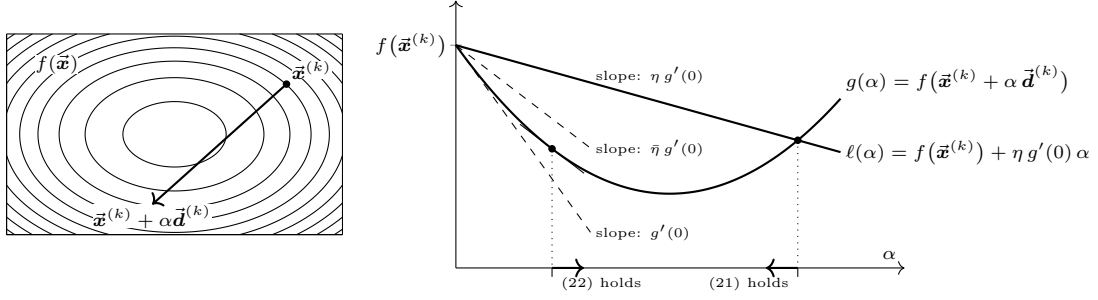
One possibility is to choose the step size that minimizes  $g(\alpha)$ , i.e.,

$$\min_{\alpha > 0} g(\alpha) \quad \text{with} \quad g(\alpha) = f(\vec{x}^{(k)} + \alpha \vec{d}^{(k)}) \quad (20)$$

which is called [exact line search](#); however, solving this subproblem can be expensive. Therefore, in practice,  $\alpha^{(k)}$  is chosen to satisfy less strong requirements, for instance the Wolfe conditions.

**Remark 14.** Whenever  $\alpha^{(k)}$  is [not](#) chosen by solving the minimization problem (20) exactly, we say it is an [inexact line search](#) method. ◀

#### 4.5.1 Wolfe Conditions



The Wolfe conditions are two conditions that guide the choice of a “good enough” step size  $\alpha^{(k)}$  at each iteration of a line search method to achieve global convergence. Consider the affine function

$$\ell(\alpha) = f(\vec{x}^{(k)}) + \eta \alpha \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}$$

with  $\eta \in (0, 1)$  a relaxation parameter (sometimes called the **relaxed tangent**). The first Wolfe condition (also called the **Armijo** or **sufficient decrease condition**) stipulates that  $g(\alpha^{(k)})$  should be no larger than  $\ell(\alpha^{(k)})$ . Hence the Armijo condition requires that the decrease in  $f$  is at least proportional to both  $\alpha^{(k)}$  and  $\nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}$ . In practice one typically chooses  $\eta$  to be rather small (e.g.,  $\eta = 10^{-4}$  or  $10^{-2}$ ).

The second Wolfe condition (also called the **curvature condition**) discards step sizes that are too small, thereby avoiding very slow progress.

**Definition 26** (First Wolfe Condition: Armijo or Sufficient Decrease Condition). Given a point  $\vec{x}^{(k)}$ , a direction  $\vec{d}^{(k)}$  and a parameter  $\eta \in (0, 1)$ , the step size  $\alpha^{(k)} > 0$  in a line search method should verify

$$g(\alpha^{(k)}) = f(\vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)}) \leq f(\vec{x}^{(k)}) + \eta \alpha^{(k)} \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)} = \ell(\alpha^{(k)}) \quad (21)$$

□

**Definition 27** (Second Wolfe Condition: Curvature Condition). Given a point  $\vec{x}^{(k)}$ , a descent direction  $\vec{d}^{(k)}$ , and a parameter  $\bar{\eta} \in (\eta, 1)$ , the step size  $\alpha^{(k)}$  should verify

$$g'(\alpha^{(k)}) = \nabla f(\vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)})^T \vec{d}^{(k)} \geq \bar{\eta} \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)} = \bar{\eta} g'(0) \quad (22)$$

□

Note that the left-hand side of the above inequality is  $g'(\alpha^{(k)})$  while the right-hand side is  $\bar{\eta} g'(0)$  (and  $g'(0) < 0$  since  $\vec{d}^{(k)}$  is a descent direction). In other words, the curvature condition ensures that the derivative of  $g$  at  $\alpha^{(k)}$  is not too small in magnitude; otherwise, one could (and should) take a longer step.

The Armijo and curvature conditions together are known as the **Wolfe conditions**. Although they give criteria for selecting a “good”  $\alpha^{(k)}$ , they do not, by themselves, indicate how to find such an  $\alpha^{(k)}$ . One popular approach is the backtracking line search algorithm.

#### 4.5.2 Backtracking Line Search

The main idea behind backtracking line search (Algorithm 4) is to start from a “large” initial value  $\bar{\alpha}$  (e.g.  $\bar{\alpha} = 10$ ) and then decrease it until the Armijo condition is met.

---

##### Algorithm 4 Backtracking Line Search

---

- 1: **Given:** Current point  $\vec{x}^{(k)}$ , descent direction  $\vec{d}^{(k)}$
  - 2: **Parameters:** Initial trial step size  $\bar{\alpha} > 0$ , relaxation parameter  $\eta \in (0, 1)$
  - 3: **Initialize:** Set  $\alpha^{(k)} \leftarrow \bar{\alpha}$
  - 4: **while**  $g(\alpha^{(k)}) > \ell(\alpha^{(k)})$  **do**
  - 5:      $\alpha^{(k)} \leftarrow \rho \alpha^{(k)}$  ▷ with  $\rho \in (0, 1)$ ; often  $\rho = \frac{1}{2}$  is used
  - 6: **Output:**  $\alpha^{(k)}$
-

The backtracking line search algorithm can be incorporated into both Newton's and Gradient Descent algorithms to choose the step size at each iteration. A general line search procedure is given in Algorithm 5. In this algorithm,  $\vec{d}^{(k)} = -\nabla f(\vec{x}^{(k)})$  (Gradient Descent; Algorithm 1) or  $\vec{d}^{(k)} = -(\underline{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)})$  (Newton's Method; Algorithm 3) or any other descent direction.

---

**Algorithm 5** General Line Search Algorithm

---

```

1: Given: Starting point  $\vec{x}^{(0)}$ 
2: Initialize: Set  $k \leftarrow 0$ 
3: while  $\|\nabla f(\vec{x}^{(k)})\| > \epsilon$  and  $k \leq k_{\max}$  do
4:   Compute search direction  $\vec{d}^{(k)}$  ▷ e.g. using  $-\nabla f(\vec{x}^{(k)})$  or  $-(\underline{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)})$ 
5:   Find  $\alpha^{(k)}$  using the backtracking line search (Algorithm 4)
6:   Update:  $\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)}$ 
7:   Increment:  $k \leftarrow k + 1$ 
8: Output:  $\vec{x}^{(k)}$ 

```

---

#### 4.5.3 Global Convergence Theorem

With the addition of the line search sub-routine, it is possible to prove that Newton's algorithm is globally convergent. In fact, a more general result holds for any line search method.

**Theorem 20** (Zoutendijk's Theorem). Assume that:

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable,  $L$ -smooth with  $L > 0$ , and bounded from below;
- For each iteration  $k$ , the search direction  $\vec{d}^{(k)}$  is a descent direction and the step size  $\alpha^{(k)}$  satisfies the Wolfe conditions (21) and (22).

Then

$$\sum_{k \geq 0} \cos^2 \theta^{(k)} \|\nabla f(\vec{x}^{(k)})\|^2 < \infty$$

where  $\theta^{(k)}$  is the angle between  $\vec{d}^{(k)}$  and  $-\nabla f(\vec{x}^{(k)})$ , i.e.,

$$\cos \theta^{(k)} = \frac{\nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}}{\|\nabla f(\vec{x}^{(k)})\| \|\vec{d}^{(k)}\|}$$

◁

**Proof.** Let  $k \geq 0$ . Since  $\alpha^{(k)}$  satisfies the curvature condition for some  $0 < \bar{\eta} < 1$ , we have

$$\nabla f(\vec{x}^{(k+1)})^T \vec{d}^{(k)} = \nabla f(\vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)})^T \vec{d}^{(k)} \geq \bar{\eta} \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}$$

which can be rewritten as

$$\left( \nabla f(\vec{x}^{(k+1)}) - \nabla f(\vec{x}^{(k)}) \right)^T \vec{d}^{(k)} \geq (\bar{\eta} - 1) \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}$$

and, since  $\nabla f$  is  $L$ -Lipschitz and by the Cauchy-Schwarz inequality, we have

$$\left( \nabla f(\vec{x}^{(k+1)}) - \nabla f(\vec{x}^{(k)}) \right)^T \vec{d}^{(k)} \leq L \alpha^{(k)} \|\vec{d}^{(k)}\|^2$$

Combining the two inequalities, we obtain

$$\alpha^{(k)} \geq \frac{(\bar{\eta} - 1) \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}}{L \|\vec{d}^{(k)}\|^2}$$

Since  $\vec{d}^{(k)}$  is a descent direction and  $\alpha^{(k)}$  satisfies the Armijo condition for some  $\eta \in (0, \bar{\eta})$ , it follows that

$$\begin{aligned}
f(\vec{x}^{(k+1)}) &= f(\vec{x}^{(k)} + \alpha^{(k)} \vec{d}^{(k)}) \\
&\leq f(\vec{x}^{(k)}) + \eta \alpha^{(k)} \nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)} \\
&\leq f(\vec{x}^{(k)}) - \frac{\eta(\bar{\eta} - 1)}{L} \frac{|\nabla f(\vec{x}^{(k)})^T \vec{d}^{(k)}|^2}{\|\vec{d}^{(k)}\|^2} \\
&= f(\vec{x}^{(k)}) - \frac{\eta(\bar{\eta} - 1)}{L} \cos^2 \theta^{(k)} \|\nabla f(\vec{x}^{(k)})\|^2
\end{aligned}$$

If we denote  $c = \frac{\eta(\bar{\eta}-1)}{L}$  and sum these inequalities over  $k$ , we obtain

$$f(\vec{x}^{(K)}) \leq f(\vec{x}^{(0)}) - c \sum_{k=0}^K \cos^2 \theta^{(k)} \|\nabla f(\vec{x}^{(k)})\|^2$$

Since  $f$  is bounded from below and  $c > 0$ , it follows that

$$\sum_{k \geq 0} \cos^2 \theta^{(k)} \|\nabla f(\vec{x}^{(k)})\|^2 < \infty$$

□

An immediate consequence of Theorem 20 is that

$$\cos^2 \theta^{(k)} \|\nabla f(\vec{x}^{(k)})\|^2 \rightarrow 0$$

If the cosine terms are bounded away from zero (i.e., there exists  $\delta > 0$  such that  $\cos \theta^{(k)} \geq \delta$  for all  $k$ ), then this implies  $\|\nabla f(\vec{x}^{(k)})\| \rightarrow 0$ ; in other words, the iterates are attracted to a stationary point. For example:

- In Gradient Descent, where  $\vec{d}^{(k)} = -\nabla f(\vec{x}^{(k)})$ , we have  $\cos \theta^{(k)} = -1$ .
- In Newton's method, when

$$\vec{d}^{(k)} = -(\underline{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)}),$$

one may show that

$$\cos \theta^{(k)} \geq \frac{\lambda_{\min}(\underline{H}_f(\vec{x}^{(k)}))}{\lambda_{\max}(\underline{H}_f(\vec{x}^{(k)}))}$$

where  $\lambda_{\min}(\underline{H}_f(\vec{x}^{(k)}))$  and  $\lambda_{\max}(\underline{H}_f(\vec{x}^{(k)}))$  denote the smallest and largest eigenvalues of  $\underline{H}_f(\vec{x}^{(k)})$ . Hence if the condition number

$$\kappa(\underline{H}_f(\vec{x}^{(k)})) = \frac{\lambda_{\max}(\underline{H}_f(\vec{x}^{(k)}))}{\lambda_{\min}(\underline{H}_f(\vec{x}^{(k)}))}$$

is uniformly bounded, then  $\cos \theta^{(k)}$  is bounded away from zero.

Thus, under appropriate assumptions, the line search methods (whether steepest descent, Newton's method with line search, or quasi-Newton methods with line search) yield

$$\|\nabla f(\vec{x}^{(k)})\| \rightarrow 0$$

Note, however, that this global convergence result does not necessarily imply that the iterates converge to a local minimizer without additional assumptions on the Hessian of  $f$ .

## 4.6 Variations of Newton's Method

Even with the addition of a line search, there remain two main drawbacks of Newton's algorithm:

- The update is well defined only if the Hessian matrix  $\underline{H}_f(\vec{x}^{(k)})$  is positive definite,
- It requires the computation of second-order derivatives, which can be computationally expensive.

We now discuss possible solutions.

### 4.6.1 Newton with Hessian Correction

Recall that if  $\underline{H}_f(\vec{x}^{(k)})$  is not positive definite, then the Newton direction

$$\vec{d}^{(k)} = -(\underline{H}_f(\vec{x}^{(k)}))^{-1} \nabla f(\vec{x}^{(k)})$$

may not be well defined or may fail to be a descent direction. A common remedy is to **correct** the Hessian by adding a multiple of the identity matrix. That is, one seeks  $\lambda > 0$  such that

$$\underline{H}_f(\vec{x}^{(k)}) + \lambda \underline{I}$$

is positive definite.



A practical method for checking whether a symmetric matrix is positive definite is to attempt its Cholesky decomposition. If the decomposition fails, then the matrix is not positive definite. Recall that the Cholesky decomposition of a positive definite matrix  $\underline{H}$  is a factorization of the form

$$\underline{H} = \underline{L}\underline{L}^T$$

where  $\underline{L}$  is a lower triangular matrix (the Cholesky factor). This decomposition is widely used to solve linear systems  $\underline{H}\vec{x} = \vec{b}$  (by first solving  $\underline{L}\vec{v} = \vec{b}$  via forward substitution and then  $\underline{L}^T\vec{x} = \vec{v}$  via backward substitution) and to compute the inverse by solving  $\underline{A}\underline{X} = \underline{I}$ .

Algorithm 6 describes a procedure for finding  $\lambda$  (the Hessian correction parameter) using the Cholesky decomposition.

---

**Algorithm 6** Hessian Correction

---

```

1: Given: Hessian matrix  $\underline{H}_f(\vec{x})$ 
2: Parameters: Initial correction parameter  $\bar{\lambda} > 0$ , increase factor  $c \in (0, 1)$ 
3: Try to compute the Cholesky factor  $\underline{L}$  of  $\underline{H}_f(\vec{x})$ 
4: if successful then
5:   Return  $\underline{L}$ 
6: Set  $\lambda \leftarrow \bar{\lambda}$ 
7: Try to compute the Cholesky factor  $\underline{L}$  of  $\underline{H}_f(\vec{x}) + \lambda \underline{I}$ 
8: while not successful do
9:   Increase  $\lambda \leftarrow c\lambda$ 
10:  Try to compute the Cholesky factor  $\underline{L}$  of  $\underline{H}_f(\vec{x}) + \lambda \underline{I}$ 
11: Output:  $\underline{L}$ , the Cholesky factor of the corrected Hessian

```

---

Using the Hessian correction subroutine above, one can incorporate it into Newton's method. The overall algorithm with line search and Hessian correction is given in Algorithm 7.

---

**Algorithm 7** Newton's Algorithm with Line Search and Hessian Correction

---

```

1: Given: Starting point  $\vec{x}^{(0)}$ 
2: Initialize:  $k \leftarrow 0$ 
3: while  $\|\nabla f(\vec{x}^{(k)})\| > \epsilon$  and  $k \leq k_{\max}$  do
4:   Compute  $\vec{g}^{(k)} = \nabla f(\vec{x}^{(k)})$  and  $\underline{H}_f(\vec{x}^{(k)})$ 
5:   Compute the Cholesky factor  $\underline{L}$  of  $\underline{H}_f(\vec{x}^{(k)})$  using the Hessian Correction subroutine (Algorithm 6)
6:   Solve  $\underline{L}\vec{v} = -\vec{g}^{(k)}$  (forward substitution)
7:   Solve  $\underline{L}^T\vec{d} = \vec{v}$  (backward substitution)
8:   Find  $\alpha^{(k)}$  using a line search procedure (e.g., Algorithm 4)
9:   Update:  $\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha^{(k)}\vec{d}^{(k)}$ 
10:  Increment:  $k \leftarrow k + 1$ 
11: Output:  $\vec{x}^{(k)}$ 

```

---

#### 4.6.2 Quasi-Newton Methods

Quasi-Newton methods aim to capture some of the advantages of Newton's method while avoiding the expensive computation of second-order derivatives. The main idea is to replace the Hessian  $\underline{H}_f(\vec{x}^{(k)})$  by an approximation that is updated using gradient information only.

A natural starting point is the first-order Taylor expansion of  $\nabla f$  at  $\vec{x}^{(k)}$ :

$$\nabla f(\vec{x}^{(k+1)}) \approx \nabla f(\vec{x}^{(k)}) + \underline{H}_f(\vec{x}^{(k)}) (\vec{x}^{(k+1)} - \vec{x}^{(k)})$$

which can be rearranged to yield the **secant equation**

$$\underline{H}_f(\vec{x}^{(k)}) \underbrace{(\vec{x}^{(k+1)} - \vec{x}^{(k)})}_{\vec{s}^{(k)}} \approx \nabla f(\vec{x}^{(k+1)}) - \nabla f(\vec{x}^{(k)}) \triangleq \vec{y}^{(k)}$$

That is, we wish to have

$$\tilde{\underline{H}}^{(k+1)} \vec{s}^{(k)} = \vec{y}^{(k)}$$

where  $\tilde{\underline{H}}^{(k+1)}$  is a symmetric approximation of the Hessian. Since a symmetric  $n \times n$  matrix has  $\frac{n(n+1)}{2}$  free parameters, the system is underdetermined and there exist infinitely many solutions. One usually

chooses the update so that  $\tilde{\underline{H}}^{(k+1)}$  is positive definite (to guarantee that the resulting search direction is a descent direction).

One of the earliest quasi-Newton methods is the DFP (Davidon-Fletcher-Powell) method. Its update is given by

$$\tilde{\underline{H}}^{(k+1)} = \left( \underline{I} - \frac{\tilde{\underline{y}}^{(k)} \tilde{\underline{s}}^{(k)T}}{\rho^{(k)}} \right)^T \tilde{\underline{H}}^{(k)} \left( \underline{I} - \frac{\tilde{\underline{y}}^{(k)} \tilde{\underline{s}}^{(k)T}}{\rho^{(k)}} \right) + \frac{\tilde{\underline{y}}^{(k)} \tilde{\underline{y}}^{(k)T}}{\rho^{(k)}}$$

with  $\rho^{(k)} = \tilde{\underline{y}}^{(k)T} \tilde{\underline{s}}^{(k)}$

A more popular approach is the BFGS (Broyden-Fletcher-Goldfarb-Shanno) method, which instead updates an approximation  $\underline{B}^{(k)}$  of the inverse Hessian. By considering the first-order Taylor expansion of  $\nabla f$ , one may wish to satisfy

$$\underline{B}^{(k+1)} \tilde{\underline{y}}^{(k)} = \tilde{\underline{s}}^{(k)}$$

and the BFGS update is given by

$$\underline{B}^{(k+1)} = \underline{B}^{(k)} - \frac{\underline{B}^{(k)} \tilde{\underline{s}}^{(k)} \tilde{\underline{s}}^{(k)T} \underline{B}^{(k)}}{\tilde{\underline{s}}^{(k)T} \underline{B}^{(k)} \tilde{\underline{s}}^{(k)}} + \frac{\tilde{\underline{y}}^{(k)} \tilde{\underline{y}}^{(k)T}}{\rho^{(k)}}$$

with  $\rho^{(k)} = \tilde{\underline{y}}^{(k)T} \tilde{\underline{s}}^{(k)}$

## 5 Constrained Optimization

### 5.1 Introduction and Terminology

Recall Equation (1), the general formulation of a constrained optimization problem, the definition of the feasible set (2), the definition of an affine function (Definition 4) and the types of minimizers (Section 3.6, Definitions 14, 15, 16).

#### 5.1.1 Classes of Constrained Optimization Problems

Depending on the type and form of  $f$  and  $c_i$ 's, the type of constrained problem can be radically different and require different classes of optimization algorithms:

- $f$  and  $c_i$ 's are **affine** functions: Linear Programming problem (Section 6)
- $f$  is **quadratic** and all  $c_i$ 's are **affine**: Quadratic Programming problem
- $f$  is **nonlinear**: Nonlinear Programming problem
- $f$  is **convex**, all inequality constraints  $c_{i,i \in \mathcal{I}}$  are **convex**, and all equality constraints  $c_{i,i \in \mathcal{E}}$  are **affine**: Convex Problem

**Lemma 21** (Properties of Convex Problems). We consider a constrained optimization problem of the form (1). If the problem is convex, then

- its feasible region  $\Omega$  is a convex set (Definition 5)
- any local minimizer is also a global minimizer ◁

#### 5.1.2 Active Constraints

are an important concept and in some cases allow us to reduce the number of constraints.

**Definition 28** (Active Set). The active set  $\mathcal{A}$  at any feasible  $\vec{x}$  consists of the indices of the constraints that are satisfied with equality at  $\vec{x}$ :

$$\mathcal{A}(\vec{x}) = \mathcal{E} \cup \{i \in \mathcal{I} \mid c_i(\vec{x}) = 0\}$$

At a feasible point  $\vec{x}$ , the inequality constraint  $i$  is said to be **active** if  $c_i(\vec{x}) = 0$  and **inactive** if the strict inequality  $c_i(\vec{x}) > 0$  is satisfied. ◁

Active inequality constraints often play a similar role as equality constraints.

If  $\vec{x}^*$  is a local minimizer, then if the non-active constraints at  $\vec{x}^*$  were removed,  $\vec{x}^*$  would remain a local minimizer.

#### 5.1.3 Feasible Directions

We need to design algorithms that produce iterates that stay within the feasible set. That is why we need the concept of **feasible directions**, which tells us in which directions we are allowed to move from a feasible point.

**Definition 29** (Feasible direction). Let  $\vec{x} \in \Omega$ . A vector  $\vec{d} \in \mathbb{R}^n$  is a **feasible direction** at  $\vec{x}$  if there exists  $\varepsilon > 0$  such that  $\vec{x} + t\vec{d} \in \Omega$ ,  $\forall t \in (0, \varepsilon]$ . In other words, by doing a small step in the direction  $\vec{d}$  from  $\vec{x}$ , we stay in  $\Omega$ . ◁

**Definition 30** (Set of linearized feasible directions). Given a feasible point  $\vec{x} \in \Omega$ , the **set of linearized feasible directions** is

$$\mathcal{F}(\vec{x}) := \left\{ \vec{d} \in \mathbb{R}^n \mid \nabla c_i(\vec{x})^\top \vec{d} = 0 \forall i \in \mathcal{E} \wedge \nabla c_i(\vec{x})^\top \vec{d} \geq 0 \forall i \in \mathcal{I} \cap \mathcal{A}(\vec{x}) \right\} \quad (23)$$

◁

Note that the inactive constraints are not relevant for finding the feasible directions.

**Definition 31** (Feasible Descent Direction). A direction  $\vec{d} \in \mathcal{F}(\vec{x})$  is a feasible descent direction at  $\vec{x} \in \Omega$  if  $\nabla f(\vec{x})^\top \vec{d} < 0$ . ◁

### 5.1.4 Constraint Qualification Condition

**Definition 32** (Constraint Qualification Condition). The Constraint Qualification Condition holds when all constraints are affine or if the set of active constraint gradients  $\{\nabla c_i(\vec{x}) \mid i \in \mathcal{A}(\vec{x})\}$  is linearly independent and  $\vec{x} \in \Omega$ .  $\triangleleft$

## 5.2 Optimality Conditions

### 5.2.1 Necessary First-Order Conditions

In Unconstrained Optimization, a necessary condition is that if  $\vec{x}^*$  is a local minimizer then there is no descent direction at  $\vec{x}^*$ , i.e.,

$$\nabla f(\vec{x}^*)^T \vec{d} \geq 0$$

for **all** directions  $\vec{d} \in \mathbb{R}^n$ , which implies the Necessary First-Order Optimality Condition  $\nabla f(\vec{x}^*) = 0$ . In constrained problems, a similar necessary condition holds:

**Lemma 22** (Fundamental Necessary Condition). If  $\vec{x}^*$  is a local minimizer and the Constraint Qualification Condition holds at  $\vec{x}^*$ , then  $\vec{x}^* \in \Omega$  and there exist no Feasible Descent Direction at  $\vec{x}^*$ , i.e.,

$$\nabla f(\vec{x}^*)^T \vec{d} \geq 0$$

for **all feasible** directions  $\vec{d} \in \mathcal{F}(\vec{x}^*)$ . Since we would need to examine every  $\vec{d} \in \mathcal{F}(\vec{x}^*)$ , this condition is not very practical.  $\triangleleft$

**Theorem 23** (Karush–Kuhn–Tucker conditions (KKT)). Suppose  $\vec{x}^*$  is a local minimizer of (1) where  $f$  and the  $c_i$ 's are continuously differentiable and the Constraint Qualification Condition holds at  $\vec{x}^*$ . Then there exists a unique  $\vec{\lambda}^* \in \mathbb{R}^{|\mathcal{E}|+|\mathcal{I}|}$  such that

$$\nabla f(\vec{x}^*) - \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i^* \nabla c_i(\vec{x}^*) = \vec{0} \quad (\text{stationarity}) \quad (24a)$$

$$c_i(\vec{x}^*) = 0, \quad \forall i \in \mathcal{E} \quad (\text{primal feasibility}) \quad (24b)$$

$$c_i(\vec{x}^*) \geq 0, \quad \forall i \in \mathcal{I} \quad (\text{primal feasibility}) \quad (24c)$$

$$\lambda_i^* \geq 0, \quad \forall i \in \mathcal{I} \quad (\text{dual feasibility}) \quad (24d)$$

$$\lambda_i^* c_i(\vec{x}^*) = 0, \quad \forall i \in \mathcal{E} \cup \mathcal{I} \quad (\text{complementary slackness}) \quad (24e)$$

Complementary Slackness (24e):

- $i \in \mathcal{A}(\vec{x}^*) \Rightarrow c_i(\vec{x}^*) = 0$
- $i \notin \mathcal{A}(\vec{x}^*) \Rightarrow \lambda_i^* = 0$

where  $\vec{\lambda}^*$  are called the **Lagrange multipliers** at  $\vec{x}^*$  and the pair  $(\vec{x}^*, \vec{\lambda}^*)$  is called a KKT point.  $\triangleleft$

Defining the **Lagrangian function** as

$$\mathcal{L}(\vec{x}, \vec{\lambda}) := f(\vec{x}) - \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i c_i(\vec{x}), \quad (\vec{x}, \vec{\lambda}) \in \mathbb{R}^n \times \mathbb{R}^{|\mathcal{E}|+|\mathcal{I}|} \quad (25)$$

the first condition (24a) can be expressed concisely as

$$\nabla_{\vec{x}} \mathcal{L}(\vec{x}^*, \vec{\lambda}^*) = \vec{0} \quad (26)$$

The last condition (24e) implies that if  $c_i$  is inactive at  $\vec{x}^*$ , i.e.  $c_i(\vec{x}^*) > 0$ , then  $\lambda_i^* = 0$ . Thus, condition (24a) can be rewritten as

$$\nabla f(\vec{x}^*) - \sum_{i \in \mathcal{A}(\vec{x}^*)} \lambda_i^* \nabla c_i(\vec{x}^*) = \vec{0} \quad (27)$$

Observe that (24e) also implies that

$$\mathcal{L}(\vec{x}^*, \vec{\lambda}^*) = f(\vec{x}^*) \quad (28)$$

A special case of (24e) is **strict complementarity**, which is satisfied if exactly one of  $\lambda_i^*$  and  $c_i(\vec{x}^*)$  is zero for each  $i \in \mathcal{I}$ . In other words,  $\lambda_i^* > 0$  for each  $i \in \mathcal{I} \cap \mathcal{A}(\vec{x}^*)$  and  $\lambda_i^* = 0$  for each  $i \in \mathcal{I} \setminus \mathcal{A}(\vec{x}^*)$ . Satisfaction of the strict complementarity property usually makes it easier for algorithms to determine the active set  $\mathcal{A}(\vec{x}^*)$  and converge rapidly to the solution  $\vec{x}^*$ .

**Definition 33** (Cone).  $S \subseteq \mathbb{R}^n$  is a **cone** if for all  $\vec{x} \in S$  and all  $\alpha \in \mathbb{R}_{\geq 0}$ , we have  $\alpha \vec{x} \in S$ .  $\triangleleft$

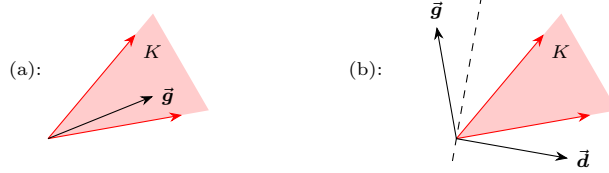
**Lemma 24** (Farkas). Let  $\underline{B} \in \mathbb{R}^{n \times m}$ ,  $\underline{C} \in \mathbb{R}^{n \times p}$  and  $\vec{g} \in \mathbb{R}^n$  and let  $K$  be the cone defined by

$$K := \{ \underline{B}\vec{y} + \underline{C}\vec{w} \mid \vec{y} \in \mathbb{R}_{\geq 0}^m, \vec{w} \in \mathbb{R}^p \} \quad (29)$$

where  $\mathbb{R}_{\geq 0}^m$  is the non-negative orthant<sup>1</sup> in  $\mathbb{R}^m$ , i.e., the set of vectors in  $\mathbb{R}^m$  whose components are all non-negative. Then for any  $\vec{g} \in \mathbb{R}^n$ , exactly one of the following conditions holds:

- (a)  $\vec{g} \in K$
- (b)  $\exists \vec{d} \in \mathbb{R}^n \mid \vec{g}^\top \vec{d} < 0, \underline{B}^\top \vec{d} \geq \vec{0}, \underline{C}^\top \vec{d} = \vec{0}$   $\triangleleft$

Geometrically Lemma 24 says that if  $\vec{g} \notin K$ , we can find a hyperplane with normal vector  $\vec{d}$  separating  $\vec{g}$  from the cone  $K$  such that  $\vec{g}$  lies on the opposite side of  $\vec{d}$ :



**Proof** (Theorem 23). Let  $\vec{x}^*$  be a local minimizer of (1) and assume that the Constraint Qualification Condition holds at  $\vec{x}^*$ . If we define as

$$N := \left\{ \sum_{i \in \mathcal{A}(\vec{x}^*)} \lambda_i \nabla c_i(\vec{x}^*) \mid \vec{\lambda} \in \mathbb{R}^{|\mathcal{A}(\vec{x}^*)|}, \lambda_i \geq 0, i \in \mathcal{A}(\vec{x}^*) \cap \mathcal{I} \right\} \quad (30)$$

it is of the form (29) with  $\underline{B} = [\nabla c_i(\vec{x}^*)]_{i \in \mathcal{I} \cap \mathcal{A}(\vec{x}^*)}$  and  $\underline{C} = [\nabla c_i(\vec{x}^*)]_{i \in \mathcal{E}}$ .

By Farkas, either  $\nabla f(\vec{x}^*) \in N$  or there exists  $\vec{d} \in \mathbb{R}^n$  such that  $\nabla f(\vec{x}^*)^\top \vec{d} < 0$ ,  $\nabla c_i(\vec{x}^*)^\top \vec{d} \geq 0$  for all  $i \in \mathcal{A}(\vec{x}^*) \cap \mathcal{I}$ ,  $\nabla c_i(\vec{x}^*)^\top \vec{d} = 0$  for all  $i \in \mathcal{E}$ . In the second case,  $\vec{d}$  would be both a descent direction and a feasible direction. By Lemma 22, if  $\vec{x}^*$  is a local minimizer, then there are no feasible descent directions at  $\vec{x}^*$  and therefore (b) cannot hold, implying that  $\nabla f(\vec{x}^*) \in N$ .

This means there exists  $\vec{\lambda} \in \mathbb{R}^{|\mathcal{A}(\vec{x}^*)|}$  such that  $\lambda_i \geq 0$  for  $i \in \mathcal{A}(\vec{x}^*) \cap \mathcal{I}$  and

$$\nabla f(\vec{x}^*) = \sum_{i \in \mathcal{A}(\vec{x}^*)} \lambda_i \nabla c_i(\vec{x}^*) \quad (31)$$

Defining the vector  $\vec{\lambda}^* \in \mathbb{R}^{|\mathcal{E}|+|\mathcal{I}|}$  as  $\lambda_i^* = \begin{cases} \lambda_i & i \in \mathcal{A}(\vec{x}^*) \\ 0 & i \in \mathcal{I} \setminus \mathcal{A}(\vec{x}^*) \end{cases}$ , we obtain

$$\nabla f(\vec{x}^*) - \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i^* \nabla c_i(\vec{x}^*) = \vec{0} \quad (32)$$

proving (24a).

(24b) and (24c) follow from the fact that  $\vec{x}^*$  is feasible, since it is a local minimizer of the constrained problem.

Note that  $\lambda_i^* \geq 0$  for all  $i \in \mathcal{I}$  since  $\lambda_i \geq 0$  for  $i \in \mathcal{A}(\vec{x}^*) \cap \mathcal{I}$  and by definition  $\lambda_i^* = 0$  for  $i \in \mathcal{I} \setminus \mathcal{A}(\vec{x}^*)$ , proving (24d).

(24e) follows from the definition of  $\vec{\lambda}^*$  since  $\lambda_i^* = 0$  whenever  $c_i(\vec{x}^*) > 0$ , i.e.  $i \notin \mathcal{A}(\vec{x}^*)$  and  $c_i(\vec{x}^*) = 0$  when  $i \in \mathcal{A}(\vec{x}^*)$ .  $\square$

When the KKT conditions are satisfied at a point  $\vec{x}^*$ , a (infinitesimal) move along any vector  $\vec{d} \in \mathcal{F}(\vec{x}^*)$  remains feasible and either increases the first-order approximation of  $f$ , if  $\vec{d}^\top \nabla f(\vec{x}^*) > 0$  (“decided”) or else keeps it constant, if  $\vec{d}^\top \nabla f(\vec{x}^*) = 0$  (“undecided”). A decrease is not possible, as we see if we express  $\nabla f(\vec{x}^*)$  as a linear combination of the active constraint gradients

$$\vec{d}^\top \nabla f(\vec{x}^*) = \sum_{i \in \mathcal{E}} \lambda_i^* \underbrace{\vec{d}^\top \nabla c_i(\vec{x}^*)}_{=0} + \sum_{i \in \mathcal{I} \cap \mathcal{A}(\vec{x}^*)} \overbrace{\lambda_i^*}^{\geq 0} \underbrace{\vec{d}^\top \nabla c_i(\vec{x}^*)}_{\geq 0} \geq 0 \quad (33)$$

where we used Definition 30 for the dot products and Condition (24d) for the multipliers.

<sup>1</sup>An orthant is the higher-dimensional analogue of a quadrant in the plane or an octant in three dimensions. The non-negative orthant is the generalization of the first quadrant in  $\mathbb{R}^2$ .

### 5.2.2 Second-Order Conditions

For undecided directions  $\vec{d}^\top \nabla f(\vec{x}^*) = 0$ , we canNOT determine from first derivative information whether a move along this direction will increase or decrease  $f$ . Second-order conditions examine the second derivative terms in the Taylor expansions of  $f$  and  $c_i$ , to see whether this extra information resolves the question.

**Definition 34** (Critical Cone). For a KKT point  $(\vec{x}, \vec{\lambda})$ , the **critical cone** at  $(\vec{x}, \vec{\lambda})$  is defined as

$$\mathcal{C}(\vec{x}, \vec{\lambda}) := \left\{ \vec{d} \in \mathcal{F}(\vec{x}) \mid \nabla c_i(\vec{x})^\top \vec{d} = 0 \text{ for all } i \in \mathcal{I} \cap \mathcal{A}(\vec{x}) \text{ with } \lambda_i > 0 \right\} \quad (34)$$

Equivalently,  $\vec{d} \in \mathcal{C}(\vec{x}, \vec{\lambda})$  if and only if all of the following hold:

- $\nabla c_i(\vec{x})^\top \vec{d} = 0$  for all  $i \in \mathcal{E}$
- $\nabla c_i(\vec{x})^\top \vec{d} = 0$  for all  $i \in \mathcal{I} \cap \mathcal{A}(\vec{x})$  with  $\lambda_i > 0$
- $\nabla c_i(\vec{x})^\top \vec{d} \geq 0$  for all  $i \in \mathcal{I} \cap \mathcal{A}(\vec{x})$  with  $\lambda_i = 0$   $\triangleleft$

From Definition 34, it follows that  $\mathcal{C}(\vec{x}, \vec{\lambda}) \subseteq \mathcal{F}(\vec{x})$  and, since  $\lambda_i = 0$  for inactive components  $i \in \mathcal{I} \setminus \mathcal{A}(\vec{x})$ , for any  $\vec{d} \in \mathcal{C}(\vec{x}, \vec{\lambda})$  we also have,

$$\lambda_i \nabla c_i(\vec{x})^\top \vec{d} = 0 \quad \forall i \in \mathcal{E} \cup \mathcal{I} \quad (35)$$

The critical cone  $\mathcal{C}(\vec{x}, \vec{\lambda})$  contains exactly those directions  $\vec{d} \in \mathcal{F}(\vec{x})$  for which it is not clear from first derivative information alone whether  $f$  will increase or decrease (“undecided” directions):

$$\vec{d}^\top \nabla f(\vec{x}) \stackrel{(24a)}{=} \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \vec{d}^\top \nabla c_i(\vec{x}) \stackrel{(35)}{=} 0 \quad \forall \vec{d} \in \mathcal{C}(\vec{x}, \vec{\lambda}) \quad (36)$$

To find out the change of  $f$  in directions  $\vec{d} \in \mathcal{C}(\vec{x}, \vec{\lambda})$ , we consider the second-order derivative (14) of  $\mathcal{L}(\vec{x}, \vec{\lambda})$  with respect to  $\vec{x}$ :

$$\nabla_{\vec{x}\vec{x}}^2 \mathcal{L}(\vec{x}, \vec{\lambda}) := \underline{H}_f(\vec{x}) - \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \underline{H}_{c_i}(\vec{x})$$

**Theorem 25** (Second-Order Necessary Optimality Condition). If  $\vec{x}^*$  is a local minimizer, the Constraint Qualification Condition holds at  $\vec{x}^*$  and  $\vec{\lambda}^*$  is the vector of Lagrange multipliers that satisfies the KKT, then:

$$\boxed{\vec{d}^\top \nabla_{\vec{x}\vec{x}}^2 \mathcal{L}(\vec{x}^*, \vec{\lambda}^*) \vec{d} \geq 0 \quad \forall \vec{d} \in \mathcal{C}(\vec{x}^*, \vec{\lambda}^*)} \quad (37) \quad \triangleleft$$

**Theorem 26** (Second-Order Sufficient Optimality Condition). Let  $(\vec{x}^*, \vec{\lambda}^*)$  be a KKT point. If

$$\boxed{\vec{d}^\top \nabla_{\vec{x}\vec{x}}^2 \mathcal{L}(\vec{x}^*, \vec{\lambda}^*) \vec{d} > 0 \quad \forall \vec{d} \in \mathcal{C}(\vec{x}^*, \vec{\lambda}^*) \setminus \{\vec{0}\}} \quad (38)$$

then  $\vec{x}^*$  is a strict local minimizer.  $\triangleleft$

If  $(\vec{x}^*, \vec{\lambda}^*)$  is a KKT point such that  $\nabla_{\vec{x}\vec{x}}^2 \mathcal{L}(\vec{x}^*, \vec{\lambda}^*) \succ 0$ , then  $\vec{x}^*$  is automatically a (strict) local minimizer!

Differences between Theorems 25 and 26:

- strict inequality in (37)
- absence of the Constraint Qualification Condition in (38)

## 6 Linear Programming

### 6.1 General and Standard Form

#### 6.1.1 “Generic” Form

In a Linear Programming (LP) problem,  $f$  and all  $c_i$  are affine functions. Thus (1) can be written as

$$\min_{\vec{x} \in \mathbb{R}^n} \vec{c}^\top \vec{x} \quad \text{subject to} \quad \begin{cases} \vec{a}_i^\top \vec{x} + b_i = 0, & i \in \mathcal{E} \\ \vec{a}_i^\top \vec{x} + b_i \geq 0, & i \in \mathcal{I} \end{cases} \quad (39)$$

with  $\vec{c} \in \mathbb{R}^n$ ,  $\vec{a}_i \in \mathbb{R}^n$ ,  $b_i \in \mathbb{R}$ ,  $i \in \mathcal{E} \cup \mathcal{I}$ . This is called the **generic form** of a LP problem.

Remarks:

- We do not lose generality by not including a constant term in (39) since the location is unaffected by an additive constant.
- Since affine functions are convex (Definition 6), any local minimizer of (39) is a global minimizer.
- $\vec{c}^\top \vec{x}$  is often called the **cost** and the value  $\vec{c}^\top \vec{x}^*$  at a minimizer  $\vec{x}^*$  the **optimal cost**.
- General Form and Standard Form are special cases of (39).

A LP problem is either **infeasible** ( $\Omega = \emptyset$ ), **unbounded** (the cost can be made arbitrarily small) or has **at least one minimizer**. Note that in last case, the set of minimizers could also be unbounded.

#### 6.1.2 General Form

**Transformation** from “Generic” Form to General Form:

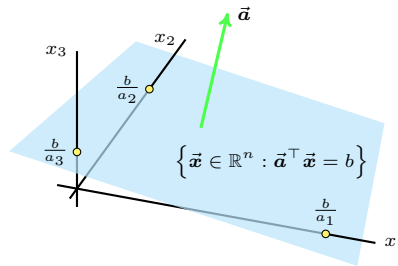
- $\vec{a}_i^\top \vec{x} + b_i = 0 \quad \longrightarrow \quad \vec{a}_i^\top \vec{x} + b_i \geq 0 \wedge -\vec{a}_i^\top \vec{x} - b_i \geq 0$
- $\vec{a}_i^\top \vec{x} + b_i \geq 0 \quad \longrightarrow \quad \vec{a}_i^\top \vec{x} \geq -b_i$

Applying these transformations to (1), we obtain  $m \geq |\mathcal{E} \cup \mathcal{I}|$  inequality constraints and (39) can be written as

$$\min_{\vec{x} \in \mathbb{R}^n} \vec{c}^\top \vec{x} \quad \text{subject to} \quad \underline{A}\vec{x} \geq \vec{b} \quad (40)$$

with  $\underline{A} \in \mathbb{R}^{m \times n}$ ,  $\vec{b} \in \mathbb{R}^m$  and  $\vec{c} \in \mathbb{R}^n$ . This is called the **general form** of a LP problem. This form is important to understand the geometry of the feasible set.

A set of the form  $\{\vec{x} \in \mathbb{R}^n : \vec{a}^\top \vec{x} = b\}$  is called a **hyperplane** in  $\mathbb{R}^n$  and it separates  $\mathbb{R}^n$  into the two sets  $\{\vec{x} \in \mathbb{R}^n : \vec{a}^\top \vec{x} \geq b\}$  and  $\{\vec{x} \in \mathbb{R}^n : \vec{a}^\top \vec{x} \leq b\}$ , which are called **halfspaces**.  $\vec{a}$  is orthogonal to the hyperplane:



The feasible set of a LP problem thus is an intersection of  $m$  halfspaces, i.e. a set of the form  $\{\vec{x} \in \mathbb{R}^n : \underline{A}\vec{x} \geq \vec{b}\}$ . It is a **polyhedron** (Definition 35).

**Definition 35** (Polyhedron). A set of the form

$$P = \left\{ \vec{x} \in \mathbb{R}^n : \vec{a}_i^\top \vec{x} = b_i, i \in \mathcal{E}, \vec{a}_i^\top \vec{x} \geq b_i, i \in \mathcal{I} \right\} \quad (41)$$

is called a polyhedron. A bounded polyhedron is called a **polytope**.  $\triangleleft$

### 6.1.3 Standard Form

The **standard form** of a LP problem is

$$\boxed{\min_{\vec{x} \in \mathbb{R}^n} \vec{c}^\top \vec{x} \quad \text{subject to} \quad \begin{cases} \underline{A}\vec{x} = \vec{b} \\ \vec{x} \geq \vec{0} \end{cases}} \quad (42)$$

with  $\underline{A} \in \mathbb{R}^{m \times n}$ ,  $\vec{b} \in \mathbb{R}^m$  and  $\vec{c} \in \mathbb{R}^n$ . This form is important for designing algorithms.

The standard form can be interpreted as follow:

- the columns  $\vec{A}_1, \dots, \vec{A}_n$  of  $\underline{A}$  are the **resource vectors**
- the non-negative variables  $x_j$  are the **amount/quantity** of resource  $j$
- the  $c_j$  are the **unit cost** of resource  $j$
- the  $\vec{b}$  is the **target** vector

We can therefore interpret a LP problem in standard form as a minimal cost “synthesis” problem where (non-negative) quantities of some resources (ingredients) are used to get a target product (food) containing  $b_i$  of nutrient  $i$ , where  $c_j$  is the unit cost of ingredient  $j$  and  $a_{ij}$  is the amount of nutrient  $i$  in ingredient  $j$ .

**Transformation** from Standard Form to General Form (easier):

$$\underline{A}' = \begin{pmatrix} \underline{A} \\ -\underline{A} \\ \underline{I}_n \end{pmatrix}, \quad \vec{b}' = \begin{pmatrix} \vec{b} \\ -\vec{b} \\ \vec{0}_n \end{pmatrix}$$

**Transformation** from General Form to equivalent<sup>2</sup> problem in Standard Form:

- non-positive variable:  $x_j \leq 0 \longrightarrow x'_j = -x_j \geq 0$
- free variable (neither  $x_j \leq 0$  nor  $x_j \geq 0$ ):  $x_j \in \mathbb{R} \longrightarrow x_j = x_j^+ - x_j^-, \quad x_j^+, x_j^- \geq 0$
- inequality constraint:  $\vec{a}_i^\top \vec{x} \geq b_i \longrightarrow \vec{a}_i^\top \vec{x} - s_i = b_i, \quad s_i \geq 0$  (slack/surplus variable)

## 6.2 Geometric Properties

### 6.2.1 Corner Points

A general property of LP problems is that minimizers can be found at “corner points”, which we will define in 3 equivalent ways (Definition 36, 37, 39).

**Definition 36** (Extreme Point).  $\vec{x} \in P$  is an **extreme point** of a polyhedron  $P$  if there exist no  $\vec{u}, \vec{v} \in P \setminus \{\vec{x}\}$  and  $\lambda \in [0, 1]$  such that  $\vec{x} = \lambda \vec{u} + (1 - \lambda) \vec{v}$ .  $\mathcal{L}_0$

Intuitively, this means that  $\vec{x}$  cannot be in-between two other points of  $P$ .

**Definition 37** (Vertex).  $\vec{x} \in P$  is a **vertex** of a polyhedron  $P$  if there exists  $\vec{c} \in \mathbb{R}^n$  such that  $\vec{c}^\top \vec{x} < \vec{c}^\top \vec{y}$  for all  $\vec{y} \in P \setminus \{\vec{x}\}$ .  $\mathcal{L}_0$

Intuitively, this means that  $\vec{x}$  is a unique solution (strict minimizer) of a LP problem with feasible set  $P$ .

Recall that  $\vec{c}^\top \vec{x} = b$  defines a family of parallel hyperplanes indexed by  $b$ , and that  $\vec{c}$  points in the direction where  $\vec{c}^\top \vec{x} = b$  increases. Thus, if  $\vec{x}^*$  is a vertex of  $P$ , the hyperplane  $\vec{c}^\top \vec{x} = b^*$  where  $b^* = \vec{c}^\top \vec{x}^*$  intersects  $P$  in only one point,  $\vec{x}^*$ , and  $b^*$  is the smallest possible value of parameter  $b$  such that the hyperplane intersects  $P$  in at least one point.

Definition 36 and 37 are purely geometric, i.e., they do not depend on a specific representation of the polyhedron  $P$ , but they are difficult to use for finding the corner points of  $P$ . Definition 39 uses the algebraic representation (41) of the polyhedron. It is less intuitive, but more practical, as it gives us a procedure for identifying all corner points.

Consider a polyhedron  $P$  defined in the “generic form” (Definition 35). Note that the general and standard form of  $P$  are special cases of this form.

<sup>2</sup>Strictly speaking, we obtain a different problem, with a different set of variables and thus different feasible set. However, the transformed problem is equivalent to the original in the sense that it has the same optimal cost (or is infeasible, in the case the original was also infeasible).



**Definition 38** (Basic Solution). The vector  $\vec{x} \in \mathbb{R}^n$  is a **basic solution** of  $P$  if all equality constraints are satisfied, i.e.,  $\vec{a}_i^\top \vec{x} = b_i$  for all  $i \in \mathcal{E}$  and there are  $n$  constraints in  $\mathcal{A}(\vec{x})$  such that their gradients are linearly independent.  $\triangleleft$

Remarks:

- If the number of constraints  $|\mathcal{E} \cup \mathcal{I}|$  defining  $P$  is less than  $n$ , there cannot be  $n$  linearly independent constraints, thus there exists no basic solution of  $P$ .
- If we choose  $n$  linearly independent constraints, there exists a unique point  $\vec{x} \in \mathbb{R}^n$  at which they are all active (the intersection of  $n$  non-parallel hyperplanes in  $\mathbb{R}^n$  is a unique point). Thus, there can only be a finite number of basic solutions given a finite number of constraints and this number is bounded by  $\binom{|\mathcal{E} \cup \mathcal{I}|}{n}$ .
- A basic solution of a polyhedron  $P$  does not necessarily belong to  $P$ , since we do not require that all inequality constraints are satisfied at  $\vec{x}$ .

**Definition 39** (Basic Feasible Solution). The vector  $\vec{x} \in \mathbb{R}^n$  is a **Basic Feasible Solution (BFS)** of  $P$  if it is a basic solution and  $\vec{x} \in P$ .  $\triangleleft$

Remarks:

- As mentioned, the 3 characterizations of corner points (Definition 36, 37, 39) are equivalent.
- Although the number of corner points is finite, it can be very large.

**Example 5** (Unit Cube). The unit cube  $\{\vec{x} \in \mathbb{R}^n : 0 \leq x_j \leq 1, j = 1, \dots, n\}$  has  $2^n$  BFSs.  $\triangleleft$

- Not all polyhedra have BFSs (e.g. one constraint in  $\mathbb{R}^n$  with  $n > 1$ )<sup>3</sup>.

### 6.2.2 Basic Solutions for Polyhedra in Standard Form

We now consider a polyhedron represented in standard form, i.e.,

$$P = \{\vec{x} \in \mathbb{R}^n : \underline{A}\vec{x} = \vec{b}, \vec{x} \geq \vec{0}\} \quad (43)$$

where  $\underline{A} \in \mathbb{R}^{m \times n}$  and  $\vec{b} \in \mathbb{R}^m$  (the total number of constraints is  $|\mathcal{E} \cup \mathcal{I}| = m + n$ ).

We assume that  $\text{rank}(\underline{A}) = m$ , i.e.,  $\underline{A}$  has “full row rank” (rows of  $\underline{A}$  are linearly independent). In particular, this requires  $m \leq n$ . This is not restrictive because, if  $P$  is non-empty (i.e. the constraints are compatible), linearly dependent rows correspond to redundant constraints that can be discarded without changing the polyhedron.

Consider a basic solution  $\vec{x}$  (Definition 38). If the  $m$  equality constraints are satisfied, they are active at  $\vec{x}$  and they are linearly independent by the assumption on the rows of  $\underline{A}$ . Thus,  $n - m$  inequality constraints must also be active at  $\vec{x}$ , i.e.,  $n - m$  coordinates of  $\vec{x}$  need to be zero. However, we need to choose those variables so that the  $n$  active constraints are linearly independent.

**Theorem 27** (Basic Solutions - Standard Form). Assume  $P$  is a polyhedron in standard form (43) with  $\text{rank}(\underline{A}) = m$ .  $\vec{x} \in \mathbb{R}^n$  is a basic solution of  $P$  if and only if  $\underline{A}\vec{x} = \vec{b}$  and there exist indices  $j_1, \dots, j_m \in \{1, \dots, n\}$  (**basic indices**) such that

- $\underline{B} = [\vec{A}_{j_1}, \dots, \vec{A}_{j_m}] \in \mathbb{R}^{m \times m}$  is non-singular (**basis matrix**)
- $x_i = 0$  for  $i \notin \{j_1, \dots, j_m\}$

where  $\vec{A}_j$  is the  $j$ -th column of  $\underline{A}$ .  $\triangleleft$

We can thus construct all basic solutions of a polyhedron in standard form by enumerating all sets of  $m$  linearly independent columns of  $\underline{A}$ , setting  $x_i = 0$  for  $i \notin \{j_1, \dots, j_m\}$ , and solving

$$\underline{A}\vec{x} = \sum_{j=1}^n x_j \vec{A}_j = \sum_{i=1}^m x_{j_i} \vec{A}_{j_i} = \underline{B}\vec{x}_B = \vec{b} \quad (44)$$

where  $\vec{x}_B = (x_{j_1}, \dots, x_{j_m})^\top \in \mathbb{R}^m$  (**basic variables**). Denoting the set of basic indices by  $\mathcal{B} = \{j_1, \dots, j_m\}$  and the set of non-basic indices by  $\mathcal{N} = \{1, \dots, n\} \setminus \mathcal{B}$ , we thus have

$$\vec{x}_B = \underline{B}^{-1}\vec{b}, \quad \vec{x}_N = \vec{0}_{n-m} \quad (45)$$

<sup>3</sup>see Section 6.2.4

With this procedure we can also find all BFSs: for each basic solution  $\vec{x}$  found, it is a BFS if  $\vec{x}_B \geq \vec{0}$  since  $x_i = 0$  for  $i \notin \{j_1, \dots, j_m\}$  (non-basic variables).

A set of  $m$  basic indices uniquely defines a basic solution (the corresponding basis matrix  $\underline{B}$  is invertible). However, several sets of basic indices may lead to the same basic solution if the latter is degenerate (see 6.2.3).

### 6.2.3 Degeneracy

is an important concept which can affect the behavior of algorithms.

**Definition 40.** A basic solution  $\vec{x}$  of a generic polyhedron  $P$  is said to be degenerate if more than  $n$  linearly independent constraints are active at  $\vec{x}$ . If there are exactly  $n$  linearly independent active constraints,  $\vec{x}$  is non-degenerate.  $\triangleleft$

For a standard form polyhedron, more than  $n$  active constraints (with the equality constraints satisfied) means that more than  $n - m$  coordinates of  $\vec{x}$  are zero. This means that some of the basic variables in  $\vec{x}_B = \underline{B}^{-1}\vec{b}$  are zero.

### 6.2.4 Existence of Basic Feasible Solutions

In “generic” polyhedra of the form (41), Corner Points (i.e., Extreme Points / Vertexes / Basic Feasible Solutions) do not necessarily exist. Obviously, if  $|\mathcal{E} \cup \mathcal{I}| < n$ , there cannot be  $n$  active constraints, thus there are no basic solutions and no BFSs.

A useful necessary and sufficient geometric condition is

**Theorem 28.** A non-empty polyhedron  $P$  has at least one BFS if and only if it does not contain a line, i.e., a set of the form  $\{\vec{x} + \lambda \vec{d} : \lambda \in \mathbb{R}\}$ .  $\triangleleft$

**Example 6.** The polyhedron  $P = \{\vec{x} \in \mathbb{R}^n : \vec{x} \geq \vec{0}\}$  does not contain a line.  $\triangleleft$

**Corollary 29.** A non-empty polyhedron in standard form (43) has at least one BFS. A non-empty polytope (bounded polyhedron) has at least one BFS.  $\triangleleft$

**Proof.** Since the positive orthant does not contain a line (Example 6), a set of the form (43), a fortiori, cannot contain a line either (since the latter is a subset of former). A bounded set cannot contain a line either, since a line is unbounded.  $\square$

## 6.3 Optimality of Corner Points

The following theorem formalizes the intuition that “an optimal solution can be found at a corner point”:

**Theorem 30.** Consider a LP problem  $\min_{\vec{x} \in P} \vec{c}^\top \vec{x}$ , where  $P \subseteq \mathbb{R}^n$  is a polyhedron. If there exists a minimizer and  $P$  has at least one BFS, then there exists a minimizer that is a BFS of  $P$ .  $\triangleleft$

**Remark:** If  $P$  is in Standard Form, it has at least one BFS (Corollary 29) and thus if there is a minimizer, there exists a minimizer that is a BFS of  $P$ .

**Proof.** Let  $\vec{x}^* \in P$  be a minimizer and  $c^* = \vec{c}^\top \vec{x}^*$  be the optimal cost. Then the set of minimizers

$$Q = \{\vec{x} \in P : \vec{c}^\top \vec{x} = c^*\} = \{\vec{x} \in \mathbb{R}^n : \underline{A}\vec{x} \geq \vec{b}, \vec{c}^\top \vec{x} = c^*\}$$

is a polyhedron and not empty, because  $\vec{x}^* \in Q$ .

If  $P$  has at least one BFS, it does not contain a line (Theorem 28) and thus, a fortiori, neither does  $Q$ , since  $Q \subseteq P$ . Therefore  $Q$  has at least one BFS, henceforth denoted  $\vec{y}^*$ .

We now show that  $\vec{y}^*$  is also a BFS of  $P$  employing the definition of a Vertex (Definition 37). If  $\vec{y}^*$  were not a BFS of  $P$ , there would exist  $\vec{u}, \vec{v} \in P \setminus \{\vec{y}^*\}$  and  $\lambda \in (0, 1)$  such that  $\vec{y}^* = \lambda \vec{u} + (1 - \lambda) \vec{v}$ . Then

$$c^* = \vec{c}^\top \vec{y}^* = \lambda \underbrace{\vec{c}^\top \vec{u}}_{\geq c^*} + (1 - \lambda) \underbrace{\vec{c}^\top \vec{v}}_{\geq c^*} \geq \lambda c^* + (1 - \lambda) c^* = c^*$$

which implies that  $\vec{c}^\top \vec{u} = \vec{c}^\top \vec{v} = c^*$  and thus  $\vec{u}, \vec{v} \in Q$ . This contradicts the fact that  $\vec{y}^*$  is a BFS of  $Q$ , and hence  $\vec{y}^*$  must be a BFS of  $P$ .  $\square$

A slightly stronger result is

**Theorem 31.** Consider a LP problem  $\min_{\vec{x} \in P} \vec{c}^\top \vec{x}$ , where  $P \subseteq \mathbb{R}^n$  is a polyhedron. If  $P$  has at least one BFS, then either the optimal cost is  $-\infty$  or there exists a minimizer that is a BFS of  $P$ .  $\triangleleft$

Remark: By transforming the problem into Standard Form, we often change the variables and the feasible set. Thus, if a point is a BFS of the transformed problem, it is not necessarily a BFS of the original problem. Nonetheless, if we find a minimizer of the transformed problem at a BFS, we also obtain a minimizer of the original problem.

## 6.4 Simplex Algorithm

is based on Corollary 29 and Theorem 30: it considers a LP problem in standard form and explores the BFSs of the feasible set, moving along its edges in directions that reduce the cost, until a minimizer is found.

We consider a LP problem in the form (42) with the feasible set

$$P = \{\vec{x} \in \mathbb{R}^n : \underline{A}\vec{x} = \vec{b}, \vec{x} \geq \vec{0}\} \quad (46)$$

where  $\underline{A} \in \mathbb{R}^{m \times n}$  is a matrix with linearly independent rows.

### 6.4.1 Feasible and Basic Directions

The set of feasible directions (Definition 29, 30) in the context of LP becomes:

**Theorem 32** (Feasible Direction). Let  $\vec{x} \in P$  be a feasible point and  $\vec{d} \in \mathbb{R}^n$  a direction. Then

$$\vec{d} \in \mathcal{F}(\vec{x}) \iff \begin{cases} \underline{A}\vec{d} = \vec{0} \\ d_i \geq 0, \quad i \in I_0 \end{cases} \quad (47)$$

where  $I_0$  are the indices in  $\{1, \dots, n\}$  for which  $x_i = 0$ .  $\triangleleft$

At a BFS  $\vec{x}^*$ , we can express  $\mathcal{F}(\vec{x}^*)$  using a finite number of directions:

**Definition 41** (Basic Direction). Given the basic indices  $\mathcal{B} = \{j_1, \dots, j_m\}$  and the corresponding basis matrix  $\underline{B} = [\vec{A}_{j_1}, \dots, \vec{A}_{j_m}]$  at  $\vec{x}^*$ , we construct a basic direction  $\vec{d}$  as follows:

1. select a non-basic index  $j \in \mathcal{N} = \{1, \dots, n\} \setminus \mathcal{B}$ .
2. set  $d_j = 1$  and  $d_i = 0$  for all other non-basic indices  $i \in \mathcal{N} \setminus \{j\}$ . i.e., moving in  $\vec{d}$  from  $\vec{x}^*$ , the non-basic variable  $x_j$  is moved by one unit while the other non-basic variables are kept unchanged.
3. set components at basic indices such that  $\underline{A}\vec{d} = \vec{0}$ :

$$\underline{A}\vec{d} = \sum_{i=1}^n d_i \vec{A}_i = \vec{A}_j + \sum_{i=1}^m \vec{A}_{j_i} d_{j_i} = \vec{A}_j + \underline{B}\vec{d}_{\mathcal{B}} = \vec{0} \implies \boxed{\vec{d}_{\mathcal{B}} = -\underline{B}^{-1} \vec{A}_j} \quad (48)$$

where  $\vec{d}_{\mathcal{B}} = (d_{j_1}, \dots, d_{j_m})$  are the basic components of  $\vec{d}$ .

This defines the  $j^{\text{th}}$  basic direction at  $\vec{x}^*$ .  $\triangleleft$

Remarks:

- There are as many basic directions as indices  $i$  such that  $x_i^* = 0$ .
- If  $\vec{x}^*$  is non-degenerate, there are  $n - m$  basic directions.
- If  $\vec{x}^*$  is non-degenerate, all basic directions are feasible, but if  $\vec{x}^*$  is degenerate, there may be a basic direction  $\vec{d}$  and a basic index  $j_i$  such that  $x_{j_i}^* = 0$  and  $d_{j_i} < 0$ , which means that  $\vec{d}$  is **not** feasible!<sup>4</sup>
- At a non-degenerate BFS, moving along a basic direction corresponds to moving along an edge of the feasible set.
- At a BFS (degenerate or not), any  $\vec{d} \in \mathcal{F}(\vec{x}^*)$  can be expressed as a linear combination of basic directions.

---

<sup>4</sup>important for 6.4.5

### 6.4.2 Reduced Costs

We now look at the change of the cost function  $f(\vec{x}) = \vec{c}^\top \vec{x}$  along the  $j^{\text{th}}$  basic direction  $\vec{d}$ .

Directional derivative of  $f(\vec{x})$  at BFS  $\vec{x}^*$  in direction  $\vec{d}$ :

$$\nabla f(\vec{x}^*)^\top \vec{d} = \vec{c}^\top \vec{d} = \sum_{i=1}^n c_i d_i = \vec{c}_B^\top \vec{d}_B + c_j \stackrel{(48)}{=} -\vec{c}_B^\top \underline{B}^{-1} \vec{A}_j + c_j \quad (49)$$

since  $d_i = 0$  for any non-basic index  $i \neq j$  and  $\vec{d}_B = -\underline{B}^{-1} \vec{A}_j$ .

**Definition 42** (Reduced Cost). For every index  $j \in \{1, \dots, n\}$  the **reduced cost** of variable  $x_j$  at  $\vec{x}^*$  is defined as

$$c_j^r = c_j - \vec{c}_B^\top \underline{B}^{-1} \vec{A}_j \quad (50)$$

and the **vector of reduced costs** can be written as  $\vec{c}^r = \vec{c} - \underline{A}^\top \underline{B}^{-1} \vec{c}_B$ .  $\triangleleft$

The reduced cost of a basic index  $j_i \in \mathcal{B}$  is zero:

$$c_{j_i}^r = c_{j_i} - \vec{c}_B^\top \underline{B}^{-1} \vec{A}_{j_i} = c_{j_i} - \vec{c}_B^\top \vec{e}_i = c_{j_i} - c_{j_i} = 0 \quad (51)$$

### 6.4.3 Optimality Conditions

The Fundamental Necessary Condition (Lemma 22) in the context of LP becomes

**Theorem 33** (Fundamental Necessary Condition). If  $\vec{x}^*$  is a local minimizer, then

$$\vec{c}^\top \vec{d} \geq 0, \quad \forall \vec{d} \in \mathcal{F}(\vec{x}^*)$$

in the context of LP.  $\triangleleft$

Using reduced costs we can restate this as

**Theorem 34** (Necessary Optimality Condition). Let  $\vec{x}^*$  be a BFS of  $P$  with a corresponding basis matrix  $\underline{B}$  and reduced cost vector  $\vec{c}^r$ . If  $\vec{x}^*$  is a minimizer and is **non-degenerate**, then  $\vec{c}^r \geq \vec{0}_n$ .  $\triangleleft$

**Theorem 35** (Sufficient Optimality Condition). Let  $\vec{x}^*$  be a BFS of  $P$  with a corresponding basis matrix  $\underline{B}$  and reduced cost vector  $\vec{c}^r$ . If  $\vec{c}^r \geq \vec{0}_n$ , then  $\vec{x}^*$  is a minimizer.  $\triangleleft$

A basis matrix  $\underline{B}$  is said to be optimal if

1.  $\vec{x}_B = \underline{B}^{-1} \vec{b} \geq \vec{0}_m$  (i.e., the corresponding basic solution is feasible)
2.  $\vec{c}^r = \vec{c} - \underline{A}^\top \underline{B}^{-1} \vec{c}_B \geq \vec{0}_n$  (i.e., the reduced costs are non-negative)

### 6.4.4 Non-degenerate Case

---

#### Algorithm 8 One Iteration of the Simplex Method

---

**Require:** BFS  $\vec{x}$ , corresponding basis  $\mathcal{B}$ , linear program  $(\underline{A}, \vec{b}, \vec{c})$

- 1: **for all**  $j \in \mathcal{N}$  **do**  $\triangleright$  compute vector of reduced costs  $\vec{c}^r$
  - 2:    $c_j^r \leftarrow c_j - \vec{c}_B^\top \underline{B}^{-1} \vec{A}_j$   $\triangleright$  reduced cost for non-basic index  $j$  (Definition 42)
  - 3: **if**  $\vec{c}_N^r \geq \vec{0}$  **then**  $\triangleright$  not necessary to check  $\vec{c}_B^r$ , as it is zero (51)
  - 4:   **return**  $\vec{x}$   $\triangleright$   $\vec{x}$  is optimal, since all reduced costs are non-negative (Theorem 35)
  - 5:  $\triangleright$  otherwise, we select one of the edges of the feasible set along which the cost decreases  $\triangleleft$
  - 6: choose  $j \in \mathcal{N}$  such that  $c_j^r < 0$   $\triangleright$  entering index: corresponding basic direction is a descent direction
  - 7:  $\vec{d}_B \leftarrow -\underline{B}^{-1} \vec{A}_j$   $\triangleright$  compute corresponding  $j^{\text{th}}$  basic direction (Definition 41)
  - 8:  $\triangleright$  see how far we can move along this direction before some basic variable hits zero (leaving index  $\ell$ )  $\triangleleft$
  - 9: **if**  $\vec{d}_B \geq \vec{0}$  **then**  $\triangleright$  if all basic variables increase, there is no limiting constraint
  - 10:   **PRINT**("unbounded problem")
  - 11:  $\theta \leftarrow \min_{i \in \mathcal{B}: d_i < 0} \frac{x_i}{-d_i}$   $\triangleright$  largest feasible step until a basic variable ( $x_\ell$ ) hits zero
  - 12: choose  $\ell \in \mathcal{B}$  that attains the minimum in the previous line  $\triangleright$  leaving index
  - 13:  $\vec{y} \leftarrow \vec{x} + \theta \vec{d}$   $\triangleright$  update BFS
  - 14:  $\mathcal{B} \leftarrow (\mathcal{B} \setminus \{\ell\}) \cup \{j\}$   $\triangleright$  update basic indices
-

### 6.4.5 Degenerate Case

When the current BFS  $\vec{x}$  is degenerate (see 6.2.3), at least one basic variable is already at 0. Consequently the step size chosen in Line 11 of Algorithm 8 can be

$$\theta = \min_{i \in \mathcal{B}: d_i < 0} \frac{x_i}{-d_i} = 0$$

so that the new point  $\vec{y} = \vec{x} + \theta \vec{d}$  coincides with the old one ( $\vec{y} = \vec{x}$ ) even though the basis has changed. If this situation happens repeatedly the **plain** simplex method can **cycle**, i.e. visit the same degenerate vertex with different bases forever.

To guarantee finite termination of Algorithm 8 one restricts the free choices made in Lines 6 and 12. A widely used rule is **Bland's (smallest-index) rule**:

- choose as entering index the smallest  $j \in \mathcal{N}$  with  $c_j^r < 0$  in Line 6
- choose as leaving index the smallest  $\ell \in \mathcal{B}$  that attains  $\theta = \frac{x_\ell}{-d_\ell}$  in Line 12

Bland's rule prevents cycling even in the presence of degeneracy and therefore restores the finite-termination guarantee of the simplex algorithm.

Even with anti-cycling rules, degeneracy may slow the practical progress of the method because several consecutive iterations can keep the objective value unchanged before a strictly improving move is found.

### 6.4.6 Computational Complexity

In a naïve implementation of Algorithm 8, the two linear solves cost  $O(m^3)$ , scanning all non-basic columns for reduced costs adds  $O(mn)$ :

$$O(m^3 + mn)$$

A **revised** variant that stores and updates  $\underline{\mathbf{B}}^{-1}$  improves the cost to:

$$O(m^2 + mn)$$

The number of possible bases is  $\binom{n}{m}$ , thus the theoretical worst case is exponential in  $m$ . In practice the revised variant with good pivot rules is fast for large, sparse problems.