

CHAPTER I: Preliminaries

Organization

- Contact:
 - Email: `sagnol -at- math.tu-berlin.de`
 - Office: MA 518
 - Assistant: Daniel Schmidt g.W., MA 524
- Timing:
 - Monday 10-12, H1029
 - Thursday 14-16, H3006
 - Lecture every Monday
 - Lecture or Exercises every second week on Thursday
 - Occasionally: Tutorials with python notebooks
- Resources:
 - There will be a handout, put online as and when the corresponding chapters are finished
 - The course is mainly based on the book “Convex Optimization”, S. Boyd & L. Vandenberghe, 2004, freely available online at <http://stanford.edu/boyd/cvxbook/>.
 - Selected chapters are also based on the following references:
 - * “Lecture on Modern Convex Optimization”, A. Ben-Tal & A. Nemirovski, 2001.
 - * “Topics in Convex Optimisation”, Lecture notes of H. Fawzi at Cambridge.
 - * “Approximation Algorithms and Semidefinite Programming”, Lecture notes of B. Gärtner & J. Matoušek at ETH Zurich.
 - * “Semidefinite Optimization”, Lecture notes of M. Laurent & F. Vallentin at Utrecht.
- Evaluation:
 - Exercises will be given on week in advance. At the beginning of exercise sessions, check the exercises you’ve prepared. One student will be asked to explain his solution. You need 50% of all exercises *checked* to pass the exam.
 - Oral examination. You should not learn everything *by heart*, but rather know which result exists and be able to find it in your handout if needed. You will not be asked to prove results from the course, but rather to solve some new exercises & show your modelling skills.
 - Mathematical rigor: the handout contains precise mathematical statements. On the blackboard however, I will handwavy sometimes. IMHO it’s OK to handwavy as long as:
 - * You are aware that you’re not handling all small details of the problem;
 - * You have an idea of how to solve the problem in a rigorous manner.

- Objectives:
 - Understand the concept of duality in convex optimization, for both the “standard” nonlinear programming approach and the “modern” conic programming approach.
 - Learn to recognize convexity in a mathematical optimization problem, and modelling tricks to represent or reformulate a problem as a linear program (LP), second-order cone program (SOCP), semidefinite program (SDP), or geometric program (GP).
 - Review many applications of convex optimization, in particular in the fields of engineering, data analysis, combinatorial optimization, and robust optimization.
 - Learn to use modern interface such as PICOS or CVXPY, to formulate conic optimization problems in python and solve them with state-of-the-art solvers.
 - Understand the algorithms (interior point methods).

Notation for scalars, vectors, and matrices

In this script, we use the following standard notation:

- $[n] := \{1, \dots, n\}$
- Scalar numbers are denoted by plain lower case letters, e.g. $c \in \mathbb{R}$.
- Vectors are denoted by boldface lower case letters, e.g. $\mathbf{v} \in \mathbb{R}^n$, with elements v_1, \dots, v_n . Unless stated otherwise, the symbol \mathbf{v} always denote a *column vector*. The associated row vector is \mathbf{v}^T . On the blackboard, we'll make no distinction between plain and boldface letters, but it should be clear from the context whether the symbol x denotes a scalar or a vector.
- Matrices are denoted by upper case letters, e.g. $A \in \mathbb{R}^{m \times n}$, with elements A_{ij} ($i \in [m]$, $j \in [n]$).
- Random variables are denoted by boldface upper case letters, e.g. \mathbf{X} . (Note that we do not distinguish between random scalars and random vectors, this should be clear from the context. The coordinates of a random vector $\mathbf{X} \in \mathbb{R}^n$ are the \mathbf{X}_i 's, also written in boldface since they are random variables).
- Sometimes, we will write $A = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{m \times n}$, which means that A is a matrix with columns $\mathbf{a}_j \in \mathbb{R}^m$ ($\forall j \in [n]$).
- Similarly, $A = [\mathbf{a}_1, \dots, \mathbf{a}_m]^T \in \mathbb{R}^{m \times n}$, means that A is a matrix with rows $\mathbf{a}_i \in \mathbb{R}^n$ ($\forall i \in [m]$).
- $\mathbb{R}_+ = \{x \in \mathbb{R} : x \geq 0\}$.
- $\mathbb{R}_{++} = \{x \in \mathbb{R} : x > 0\}$.
- $\mathbb{S}^n = \{X \in \mathbb{R}^{n \times n} : X = X^T\}$.
- $\mathbb{S}_+^n = \{X \in \mathbb{S}^n : X \text{ is positive semidefinite}\}$.
- $\mathbb{S}_{++}^n = \{X \in \mathbb{S}^n : X \text{ is positive definite}\}$.

The inequalities between vectors are *elementwise*, that is, if $A = [\mathbf{a}_1, \dots, \mathbf{a}_m]^T \in \mathbb{R}^{m \times n}$ is a matrix with rows \mathbf{a}_i^T 's, then $A\mathbf{x} \leq \mathbf{b}$ means

$$\mathbf{a}_i^T \mathbf{x} \leq b_i \quad (\forall i \in [m]).$$

The i th vector of the (canonical) basis of \mathbb{R}^n is $\mathbf{e}_i = [0, \dots, 0, 1, 0, \dots, 0]^T$, with the 1 in i th position.

Example:

For a vector $\mathbf{v} \in \mathbb{R}^n$, we have $\mathbf{e}_i^T \mathbf{v} = v_i$.

For a vector $A \in \mathbb{R}^{m \times n}$, we have $\mathbf{e}_i^T A \mathbf{e}_j = A_{ij}$. (here, it is implicit that the vectors \mathbf{e}_i and \mathbf{e}_j are of appropriate size, that is, $\mathbf{e}_i \in \mathbb{R}^m$ and $\mathbf{e}_j \in \mathbb{R}^n$.)

The matrix $E_{ij} = \mathbf{e}_i \mathbf{e}_j^T$ is a matrix (whose size depend on the dimensions of \mathbf{e}_i and \mathbf{e}_j , and should be clear from the context) with zeroes everywhere, except for a one in position (i, j) .

#1

The vector of all ones is $\mathbf{1} = \sum_i \mathbf{e}_i$. We sometimes write $\mathbf{1}_n$ to make it clear that $\mathbf{1}_n \in \mathbb{R}^n$.

The matrix of all ones is $J = \mathbf{1}\mathbf{1}^T$. When the dimension is not clear from the context, we can write $J_n = \mathbf{1}_n \mathbf{1}_n^T$ $J_{m,n} = \mathbf{1}_m \mathbf{1}_n^T$.

The identity matrix is I (or $I_n \in \mathbb{S}^n$). The zero vector is $\mathbf{0}$ (or $\mathbf{0}_n \in \mathbb{R}^n$), and the zero matrix is O (or $O_n \in \mathbb{S}^n$).

Let $\mathbf{u} \in \mathbb{R}^n$. The diagonal matrix with elements u_1, \dots, u_n is denoted by $\text{Diag}(\mathbf{u})$.

The vector of diagonal elements of a matrix $X \in \mathbb{R}^{n \times n}$ is denoted by $\text{diag } X$.

The image space of $A \in \mathbb{R}^{m \times n}$, i.e., the space spanned by the columns of A , is denoted by

$$\mathbf{Im } A := \{A\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\} \subseteq \mathbb{R}^m.$$

The nullspace of A is denoted by

$$\mathbf{Ker } A := \{\mathbf{y} \in \mathbb{R}^n : A\mathbf{y} = \mathbf{0}\}.$$

We also recall that the rank of $A \in \mathbb{R}^{m \times n}$ is the dimension of $\mathbf{Im } A$, and we say that A has *full column rank* whenever $\text{rank } A = n$, which means that the n columns of A are linearly independent. Similarly, we say that A has *full row rank* if A^T has full column rank, i.e., if the m rows of A are linearly independent.

Scalar products and norms

The scalar product of two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ is

$$\langle \mathbf{u}, \mathbf{v} \rangle := \mathbf{u}^T \mathbf{v} = \sum_{i=1}^n u_i v_i$$

The scalar product of two matrices of the same size is

$$\langle A, B \rangle := \text{trace } A^T B = \sum_j (A^T B)_{jj} = \sum_{i,j} A_{ij} B_{ij}.$$

In particular, note that when A and B are symmetric, we simply have $\langle A, B \rangle = \text{trace } AB$.

Example:

The sum of all entries of a vector $\mathbf{v} \in \mathbb{R}^n$ is $\sum_{i=1}^n v_i = \mathbf{1}^T \mathbf{v}$.

The sum of all entries of a matrix $A \in \mathbb{R}^{m \times n}$ is $\sum_{i=1}^m \sum_{j=1}^n A_{ij} = \langle J, A \rangle$.

#2

Unless stated otherwise, the symbol $\|\mathbf{v}\|$ denotes the Euclidean norm of the vector $\mathbf{v} \in \mathbb{R}^n$:

$$\|\mathbf{v}\| := \sqrt{\mathbf{v}^T \mathbf{v}} = \left(\sum_{i=1}^n v_i^2 \right)^{1/2}$$

When there might be an ambiguity, we use $\|\mathbf{v}\|_2$ to denote the Euclidean norm of \mathbf{v} . For all $p \geq 1$, the L_p -norm of \mathbf{v} is $\|\mathbf{v}\|_p := \left(\sum_i v_i^p \right)^{1/p}$.

The *Frobenius norm* of a matrix is

$$\|A\|_F := \sqrt{\langle A, A \rangle} = \left(\sum_{i,j} A_{ij}^2 \right)^{1/2}.$$

The *vectorization* of a matrix $A = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{m \times n}$ with columns \mathbf{a}'_j s is

$$\text{vec}(A) := \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{bmatrix} \in \mathbb{R}^{mn}.$$

In particular, note that $\langle A, B \rangle = \text{vec}(A)^T \text{vec}(B)$ and $\|A\|_F = \|\text{vec}(A)\|$.

Linear and Quadratic functions

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be *affine* if it is of the form $\mathbf{x} \mapsto \sum_{i=1}^n a_i x_i + b$. In vector notation, an affine function can always be written as

$$\mathbf{x} \mapsto \mathbf{a}^T \mathbf{x} + b.$$

where $\mathbf{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$. More generally, an affine function mapping \mathbb{R}^m to \mathbb{R}^n has the form

$$f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$$

for some matrix $A \in \mathbb{R}^{n \times m}$ and a vector $\mathbf{b} \in \mathbb{R}^n$. We sometimes abuse the language and say that f is *linear*, although this term is normally reserved for functions of the form $\mathbf{x} \mapsto A\mathbf{x}$ (without the constant \mathbf{b}). To emphasize that a linear function has no constant part ($\mathbf{b} = \mathbf{0}$), we will speak of a *linear form*.

Quadratic functions of \mathbb{R}^n to \mathbb{R} are of the form

$$\mathbf{x} \mapsto \sum_{ij} Q_{ij} x_i x_j + \sum_i a_i x_i + b = \mathbf{x}^T Q \mathbf{x} + \mathbf{a}^T \mathbf{x} + b.$$

Note that we can always assume without loss of generality (w.l.o.g.) that $Q \in \mathbb{S}^n$, because $\mathbf{x}^T Q \mathbf{x} = (\mathbf{x}^T Q \mathbf{x})^T = \mathbf{x}^T Q^T \mathbf{x}$, so:

$$\mathbf{x}^T Q \mathbf{x} = \frac{1}{2} \mathbf{x}^T (Q + Q^T) \mathbf{x}.$$

When a quadratic function has no affine part, i.e., $f(\mathbf{x}) = \mathbf{x}^T Q \mathbf{x}$, we speak of a *quadratic form*. An alternative formulation shows that any quadratic function can be assimilated with a quadratic form over $\mathbb{R}^n \times \{1\}$:

$$\mathbf{x}^T Q \mathbf{x} + \mathbf{a}^T \mathbf{x} + b = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}^T \begin{bmatrix} Q & \frac{1}{2} \mathbf{a} \\ \frac{1}{2} \mathbf{a}^T & b \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}.$$

The gradient of a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is

$$\nabla f : \mathbf{x} \mapsto \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^n,$$

and if f is twice differentiable, its Hessian is the function

$$\nabla^2 f : \mathbf{x} \mapsto \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{pmatrix} \in \mathbb{S}^n,$$

In particular, the gradient and Hessian of linear and quadratic forms read

$$\begin{aligned} \nabla(\mathbf{x} \mapsto \mathbf{a}^T \mathbf{x}) &= \mathbf{a} \\ \nabla^2(\mathbf{x} \mapsto \mathbf{a}^T \mathbf{x}) &= \mathbf{O} \\ \nabla(\mathbf{x} \mapsto \frac{1}{2} \mathbf{x}^T Q \mathbf{x}) &= Q \mathbf{x} \\ \nabla^2(\mathbf{x} \mapsto \frac{1}{2} \mathbf{x}^T Q \mathbf{x}) &= Q. \end{aligned}$$

Example:

Let $\mathbf{u} \in \mathbb{R}^m$ and $\mathbf{v} \in \mathbb{R}^n$. The function $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $X \mapsto \mathbf{u}^T X \mathbf{v}$ is a linear function of X . This becomes obvious if we use the following formula:

$$\mathbf{u}^T X \mathbf{v} = \langle X, \mathbf{u} \mathbf{v}^T \rangle. \quad (1)$$

The above is a simple consequence from the fact that the trace is *invariant to cyclic permutations*: $\text{trace } AB = \text{trace } BA$. Indeed,

$$\begin{aligned} \mathbf{u}^T X \mathbf{v} &= \text{trace } \mathbf{u}^T X \mathbf{v} && \text{(seen as a } 1 \times 1\text{-matrix)} \\ &= \text{trace } X \mathbf{v} \mathbf{u}^T && \text{(from the invariance to cyclic permutations)} \\ &= \langle X, \mathbf{u} \mathbf{v}^T \rangle && \text{(note that } \mathbf{u} \mathbf{v}^T \text{ is an } m \times n\text{-matrix)} \end{aligned}$$

#3

Random vectors

Let \mathbf{X} be a real-valued random variable, associated with a probability measure \mathbb{P} . The *expected value* (or the *expectation*) of \mathbf{X} is

$$\mathbb{E}[\mathbf{X}] = \int_{x \in \mathbb{R}} x d\mathbb{P}(x).$$

and its *variance* is

$$\mathbb{V}[\mathbf{X}] = \mathbb{E}[\mathbf{X}^2] - \mathbb{E}[\mathbf{X}]^2 = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2] \geq 0.$$

Similarly, when \mathbf{X} is a random vector in \mathbb{R}^n , its expected value is the vector

$$\mathbb{E}[\mathbf{X}] = \int_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x} d\mathbb{P}(\mathbf{x}) = \begin{bmatrix} \mathbb{E}[\mathbf{X}_1] \\ \vdots \\ \mathbb{E}[\mathbf{X}_n] \end{bmatrix}.$$

The *variance-covariance matrix* of \mathbf{X} is the matrix with the *covariance* of \mathbf{X}_i and \mathbf{X}_j in position (i, j) : $\text{cov}[\mathbf{X}_i, \mathbf{X}_j] = \mathbb{E}[\mathbf{X}_i \mathbf{X}_j] - \mathbb{E}[\mathbf{X}_i] \mathbb{E}[\mathbf{X}_j] = \mathbb{E}[(\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i])(\mathbf{X}_j - \mathbb{E}[\mathbf{X}_j])]$. In matrix notation,

$$\mathbb{V}[\mathbf{X}] := \mathbb{E}[\mathbf{X} \mathbf{X}^T] - \mathbb{E}[\mathbf{X}] \mathbb{E}[\mathbf{X}]^T = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T].$$

Note that the diagonal of i th diagonal element of $\mathbb{V}[\mathbf{X}]$ is nothing but $\mathbb{V}[\mathbf{X}_i]$. The last expression also shows

that a variance-covariance matrix is always a positive semidefinite matrix (this will become clear during the lecture): $\mathbb{V}[\mathbf{X}] \in \mathbb{S}_+^n$.

The expectation is linear, so in particular, for a scalar random variable \mathbf{X} and constants $a, b \in \mathbb{R}$, we have $\mathbb{E}[a\mathbf{X} + b] = a\mathbb{E}[\mathbf{X}] + b$. Similarly, for a random vector $\mathbf{X} \in \mathbb{R}^n$, a matrix $A \in \mathbb{R}^{m \times n}$ and a (constant) vector $\mathbf{b} \in \mathbb{R}^m$, it holds

$$\mathbb{E}[A\mathbf{X} + \mathbf{b}] = A\mathbb{E}[\mathbf{X}] + \mathbf{b}.$$

A consequence is that when \mathbf{X} is multiplied by a constant $a \in \mathbb{R}$, its variance (or variance-covariance matrix) gets multiplied by the square of this constant: $\mathbb{V}[a\mathbf{X}] = a^2\mathbb{V}[\mathbf{X}]$. More generally, we can show that if \mathbf{X} takes values in \mathbb{R}^n and A is an $m \times n$ -matrix,

$$\mathbb{V}[A\mathbf{X}] = A\mathbb{V}[\mathbf{X}]A^T.$$

Example:

Let \mathbf{X} and \mathbf{Y} be two real-valued random variables. We can use the above formula to derive the standard formula

$$\mathbb{V}[\mathbf{X} + \mathbf{Y}] = \mathbb{V}[\mathbf{X}] + \mathbb{V}[\mathbf{Y}] + 2\mathbf{cov}[\mathbf{X}, \mathbf{Y}],$$

which shows, in particular, that the variance of the sum equals the sum of the variances whenever the two variables are uncorrelated. Indeed, let $\mathbf{Z} = [\mathbf{X}, \mathbf{Y}]^T$ be the random vector of \mathbb{R}^2 obtained by concatenating \mathbf{X} and \mathbf{Y} . Then,

$$\mathbb{V}[\mathbf{X} + \mathbf{Y}] = \mathbb{V}[\mathbf{1}^T \mathbf{Z}] = \mathbf{1}^T \mathbb{V}[\mathbf{Z}] \mathbf{1} = \mathbf{1}^T \begin{bmatrix} \mathbb{V}[\mathbf{X}] & \mathbf{cov}[\mathbf{X}, \mathbf{Y}] \\ \mathbf{cov}[\mathbf{X}, \mathbf{Y}] & \mathbb{V}[\mathbf{Y}] \end{bmatrix} \mathbf{1},$$

which is nothing but the sum of the 4 entries of the last 2×2 -matrix.

#4

Example:

Let \mathbf{X} be a random vector which takes values in \mathbb{R}^n , with expected value $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu} \in \mathbb{R}^n$ and variance-covariance matrix $\mathbb{V}[\mathbf{X}] = \Sigma \in \mathbb{S}_+^n$.

The function $f : \mathbb{S}^n \rightarrow \mathbb{R}$, which associates Q to the expected value of $\mathbf{X}^T Q \mathbf{X}$ is a linear form. We will show that $f(Q)$ has a simple expression, which depends only on $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$ and $\Sigma = \mathbb{V}[\mathbf{X}]$. To see this, we can write

$$\begin{aligned} f(Q) &= \mathbb{E}[\mathbf{X}^T Q \mathbf{X}] \\ &= \mathbb{E}[\langle Q, \mathbf{X} \mathbf{X}^T \rangle] && \text{(cf. Eq (1) in Example #3)} \\ &= \langle Q, \mathbb{E}[\mathbf{X} \mathbf{X}^T] \rangle && \text{(linearity of expectation)} \\ &= \langle Q, \Sigma + \boldsymbol{\mu} \boldsymbol{\mu}^T \rangle && \text{(because } \Sigma = \mathbb{E}[\mathbf{X} \mathbf{X}^T] - \boldsymbol{\mu} \boldsymbol{\mu}^T \text{)}. \end{aligned}$$

#5