

Lecture Notes for Statistics and Estimation Theory

Paul Ferrand

2019

Prerequisites

Linear algebra

We will consider matrices and vectors, noted as uppercase bold \mathbf{A} and lowercase bold \mathbf{v} respectively. A matrix is a collection of elements arranged in m rows and n columns. Vectors are matrices with a single column. We have

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{1m} & \dots & a_{mn} \end{pmatrix} \quad (1)$$

The transpose of a matrix is

$$\mathbf{A}^T = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{1m} & \dots & a_{mn} \end{pmatrix} \quad (2)$$

A matrix is symmetric if $\mathbf{A}^T = \mathbf{A}$, and it is square if $m = n$. You can add matrices and vectors if they have the same size. Matrices can be multiplied as follows. Assume that \mathbf{A} is $m \times n$ and \mathbf{B} is $n \times p$. Note that the number of columns of \mathbf{A} needs to be equal to the number of columns of \mathbf{B} . Then

$$\mathbf{C} = \mathbf{A} \cdot \mathbf{B} \text{ and } c_{ij} = \sum_{k=1}^n a_{ik}b_{kj} \quad (3)$$

When $\mathbf{A} = \mathbf{a}^T$ is actually a *row* vector and \mathbf{b} is a *column* vector, this multiplication outputs a scalar, and the operation is called a scalar product. It is sometimes denoted as

$$\mathbf{a}^T \mathbf{b} = \langle \mathbf{a} | \mathbf{b} \rangle \quad (4)$$

A matrix multiplication can thus be seen as a collection of scalar products between the rows of a matrix and the columns of another matrix—this can come in handy at times.

Scalar products Scalar products are used to check *orthogonality* in arbitrary spaces. In particular, we will say that 2 vectors are orthogonal if their scalar product is equal to 0. For vectors in 2 dimensions, this directly corresponds to the natural orthogonality that we learn since

childhood, and extends in arbitrary dimensions. We can also define the *norm* of a vector through scalar products. We will say that

$$\|v\| = \sqrt{v^T v} \quad (5)$$

which in turns, allows to define an *angle* between arbitrary vectors through their cosine:

$$\cos[\angle(a, b)] = \frac{a^T b}{\|a\| \|b\|} \quad (6)$$

Orthogonal vectors have their cosine equal to 0—once again, as we learned before. In contrast, *colinear vectors* have their cosines equal to 1, and thus the scalar product of colinear vectors is equal to the product of their norms.

Gram-Schmidt There is more to orthogonality; consider 2 vectors a and b that are just arbitrary. We can split b into a part that is colinear with a and a part that is orthogonal to a as follows

$$b = \frac{a^T b}{\|a\|^2} a + \left(b - \frac{a^T b}{\|a\|^2} a \right) \quad (7)$$

The first term represent the colinear part and the second one between parentheses the orthogonal one. Through standard algebra you can verify that the above is obviously true. This can be used to create a set of orthogonal vectors through a set of arbitrary ones. Let such a set be denoted $\{u_k\}$. We proceed iteratively, forming the set $\{v_k\}$ as

$$v_0 = u_0 \quad (8)$$

$$v_1 = u_1 - \frac{(v_0^T u_1)}{\|v_0\|^2} v_0 \quad (9)$$

$$v_2 = u_2 - \frac{(v_0^T u_2)}{\|v_0\|^2} v_0 - \frac{(v_1^T u_2)}{\|v_1\|^2} v_1 \quad (10)$$

$$\vdots \quad (11)$$

This operation is sometimes called an orthogonalization or a *Gram-Schmidt* procedure, and we say that $\{v_k\}$ is a set of linearly independent vectors. One can indeed verify that $v_i^T v_j = 0$ by construction for any i , and j . If we normalize the vectors in the set, they will form a *basis*. The number of vectors in the basis is related to the dimension of the linear subspace described by the vectors, and is intuitively linked to the number of *degrees of freedom* of the linear space.

Verify that this indeed creates a set of orthogonal vectors.

Rank and identity matrix Back to matrices, we define the *rank* of a matrix as the number of linearly independent rows or columns, whichever is the smallest. There are alternative characterizations but

this one is practical. In particular, the number of linearly independent columns are linked to the descriptive power of the matrix since we can form a basis with the columns. Among specific matrices, we call the one that is neutral with respect to multiplication the identity matrix. It is a square matrix with ones on the diagonal, and it verifies that

$$AI = IA = A \quad (12)$$

where the sizes of I are chosen so that the multiplication makes sense. This neutral element allows one to define the inverse of a matrix, which is *any* matrix such that

$$AA^{-1} = A^{-1}A = I \quad (13)$$

This definition only really makes sense when A is square, although there are *pseudo-inverses* in some other cases. When A is square, the inverse of the matrix will exist and be unique *if and only if (iff)* the rank of A is equal to n . If not, A is singular or *rank-deficient* and basically has no inverse. A matrix is said to be *orthogonal* if its columns and rows are orthogonal pairwise, and if

$$A^{-1} = A^T \quad AA^T = A^T A = I \quad (14)$$

Determinant and trace Next we define the determinant of a square matrix as the scalar value constructed in a recursive manner as

$$\det A = \sum_{j=1}^n a_{ij} C_{ij} \quad (15)$$

The values C_{ij} are the *cofactors at ij* of the matrix A obtained as

$$C_{ij} = (-1)^{i+j} \det M_{ij} \quad (16)$$

The matrix M_{ij} is obtained by deleting the i -th row and j -th column of $\det A$ —hence the recursive definition. The cofactors are often considered as a matrix C , which is important because the inverse of a matrix is defined through this matrix, as

$$A^{-1} = \frac{C^T}{\det A} \quad (17)$$

It is the only general analytical formula to invert an arbitrary matrix; when A has some specific structure, there may be easier inverse formulas. We can also define the *trace* of a matrix as the sum of its diagonal elements

$$\text{tr} A = \sum a_{ii} \quad (18)$$

Both the determinant and the trace express profound things about a matrix; we will see shortly what that is.

Quadratic forms and definiteness A quadratic form is a function of a vector x defined through a matrix A as

$$Q(x) = \sum_i \sum_j A_{ij} x_i x_j = x^T A x \quad (19)$$

The matrix A is supposed to be symmetric, which entails no loss in generality since any quadratic form may be expressed in this way. If for all $x \neq \mathbf{0}$ we have

$$x^T A x \geq 0 \quad (20)$$

Exercise: show why.

the matrix A is said to be positive semi-definite. If the inequality is strict, A is said to be positive definite.

Properties of matrices, determinants and traces

- $(AB)^T = B^T A^T$
- $(A^{-1})^T = (A^T)^{-1} = A^{-T}$
- $(AB)^{-1} = B^{-1} A^{-1}$
- $\det A^T = \det A$
- $\det(cA) = c^n \det A$
- $\det(AB) = \det A \det B$
- $\det A^{-1} = (\det A)^{-1}$
- $\text{tr}(AB) = \text{tr}(BA)$
- $\text{tr}(A^T B) = \sum_i \sum_j A_{ij} B_{ij}$
- $\det I = 1$
- If D is a diagonal matrix, then $\det D = \prod_i d_{ii}$
- $\text{tr}(c) = c$ for all scalars c

Some theorems

1. A square matrix A is invertible *if and only if* its determinant is non-zero.
2. A square matrix A is positive definite iff it can be written as $A = CC^T$ for some matrix C that is also full-rank and invertible.
3. A square matrix A is positive definite iff all its principal minors are positive¹
4. If C is not full-rank or one of the principal minors is equal to 0, A is only positive semidefinite.

¹ The principal minors are the determinants of the top-leftmost square matrices of increasing size.

5. If A is positive definite, then for some full rank $m \times n$ matrix B with $m \leq n$, BAB^T is also positive definite.
6. If A is positive definite, then its diagonal elements are positive and the determinant of A is positive. If A is only positive semi-definite, the previous properties are relaxed to non-negative—i.e. they can be equal to 0.

Matrix inversion lemma A very common and useful property is stated as follows

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(DA^{-1}B + C^{-1})^{-1}DA^{-1} \quad (21)$$

Proof of this?

It relates the inverse of sum and product of matrices with the sum and products of the component inverse. Since the inverse is a complex operation, it is very useful numerically. As an example, it can be used to compute the so-called *rank-1* update of the inverse of the matrix A

$$(A + uu^T)^{-1} = A^{-1} - \frac{A^{-1}uu^T A^{-1}}{1 + u^T A^{-1}u} \quad (22)$$

where u is a column vector and uu^T is the *outer product* of u with itself². This matrix has rank 1 since its columns are all colinear with u . This operation happens usually to integrate a newly measured value into a matrix A for which we already computed the inverse. Here we never invert a matrix, but we use the already known value of A^{-1} to compute the new value after adding a rank-1 matrix. If furthermore A is symmetric, then $A^{-1} = A^{-T}$ and denoting $v = A^{-1}u$ we can rewrite (22) as

² uu^T is a matrix and its component at row i and column j is equal to $u_i u_j$

$$(A + uu^T)^{-1} = A^{-1} - \frac{vv^T}{1 + u^T v} \quad (23)$$

Eigendecompositions An eigenvector v of a matrix is the square satisfies

$$Av = \lambda v \quad (24)$$

for some scalar λ , which may be complex in general even if A only has real coefficients. We call λ the eigenvalue associated to the eigenvector v . By convention, $\|v\| = 1$. There are many properties on eigenvalues depending on the matrix./ In general, our interest will be on symmetric and possibly positive definite matrices, for which the eigenvalues are respectively real and positive, and whose eigenvectors are orthogonal if they correspond to different eigenvalues. In particular, this means that we can write

Exercise: what are the eigenvectors of the identity matrix?

$$A[v_1 \cdots v_n] = [\lambda_1 v_1 \cdots \lambda_n v_n] \quad (25)$$

or written differently

$$AV = V\Lambda \quad (26)$$

with

$$V = [v_1 \cdots v_n] \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

Since V is formed by orthogonal vectors you can rewrite (26) as

$$A = V\Lambda V^T = \sum \lambda_i v_i v_i^T \quad (27)$$

EXERCISE How does the eigen decomposition helps in computing the inverse of a matrix?

We have

$$A^{-1} = V^{-T} \Lambda^{-1} V^{-1} = V \Lambda^{-1} V^T = \sum \frac{1}{\lambda_i} v_i v_i^T$$

EXERCISE What is the determinant of an orthogonal matrix? What can you say about the determinant of a square matrix with respect to its eigenvalues? What about the trace?

If V is orthogonal, then $V^T = V^{-1}$ which is equivalent to

$$\det V = (\det V)^{-1} \Leftrightarrow \det V = \pm 1$$

In turn, this means that

$$\det A = \det V \det \Lambda (\det V^T) = \det \Lambda = \prod \lambda_i$$

For the trace, remark that

$$\text{tr}(V\Lambda V^T) = \text{tr}(V^T V\Lambda) = \text{tr}(\Lambda) = \sum \lambda_i$$

Probability

We will consider a simplistic view of probability, one that serves our interest without getting in the way or introducing hard concepts unless they provide a very strong intuition about a problem. A *random variable* can be considered as a mathematical object that can take any value on a specific domain. This domain may be purely discrete, purely continuous, or something in-between. This last possibility creates a host of practical issues that require something called *measure theory* to solve adequately. As such, we will try to avoid these random variables and only consider the pure ones, remembering that there is a unifying view behind these.

Random variables are usually noted as upper-case letters as X and Y , and the values they take on their domains as lowercase letters. A *realization* of a random variable takes a value, which in the discrete case is naturally noted as the *event* $X = x$. In a continuous domain, there is an uncountable infinity of possible values for the random

variable to take, and therefore there is not much interest in a specific values, but rather in ranges of values. In this case, we are more interested in events $a \leq X \leq b$ for some $a < b$. If the domain of the r.v. is $(-\infty, \infty)$ we will mostly consider the event $X \leq x$.

Cumulative and probability functions To each event we associate a probability, meaning that we are able to measure the relative weight of the event among all possible realizations of the random variable, assuming that the sum of event weights are equal to 1. For discrete variables, this measure of weight or probability *law* is given as a function of the $x \in \mathcal{D}$ where \mathcal{D} is the definition domain, as

$$\mathbb{P}(X = x) \quad (28)$$

This function can sometimes be parametrized, something that will be of great interest in the rest of the course. A common discrete probability law is the Poisson law, where for an integer $k > 0$ we have

$$\mathbb{P}(X = k) = \frac{\theta^k}{k!} e^{-\theta} \quad (29)$$

In the continuous case, the weight of the event is related to the range and defined through a continuous *probability density function* or *p.d.f.* $f(x)$, so that

$$\begin{aligned} F(b) - F(a) &= \mathbb{P}(a \leq x \leq b) = \int_a^b f(x) dx \\ &= \int_{-\infty}^b f(x) dx - \int_{-\infty}^a f(x) dx \end{aligned} \quad (30)$$

The cumulative distribution function $F(x)$ is thus defined as

$$F(x) = \int_{-\infty}^x f(y) dy \quad (31)$$

Once again, care must be taken but when this integral is proper, and $f(x)$ is well behaved, it usually poses no issue. Note that you should definitely not think about $f(x)$ as the probability of the event $X = x$; as said before, the probability of this event is virtually zero and uninteresting. In particular, although in the discrete case we have by construction that $P(X = x) \leq 1$ since the sum of the probability weights over all the possible events sums to 1, in the continuous case, assuming the domain is $(-\infty, \infty)$, we want

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (32)$$

which does not preclude $f(x)$ to be strictly larger than 1 for some x .

Show that this is true for the Poisson law; solution uses the fact that

$$e^{\theta} = \sum_{k=0}^{\infty} \frac{\theta^k}{k!}$$

Usual laws Among the laws used the most for continuous r.v., the most common one is definitely the Gaussian or so called *normal* law, denoted as $\mathcal{N}(\mu, \sigma^2)$ and whose p.d.f. is

$$f(x; \sigma, \mu) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (33)$$

We include the parameters in the function using the notation $f(x; \mu, \sigma)$. As we gravitate towards estimation, this notation will prove useful. Another common law is the continuous distribution between end-points a and b , usually written as $\mathcal{U}(a, b)$

$$f(x; a, b) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases} \quad (34)$$

The exponential law $\mathcal{E}(\theta)$ has density

$$f(x; \theta) = \frac{1}{\theta} e^{-\frac{x}{\theta}} \quad (35)$$

and is defined on $[0, \infty]$. The parameter θ is sometimes defined as $\mu = \theta^{-1}$. In this case, the density is

$$f(x; \mu) = \mu e^{-\mu x}. \quad (36)$$

Finally, the χ^2 distribution regularly appears throughout probability and statistics. In particular, if you have r.v.s. X_i distributed as standard Gaussian $\mathcal{N}(0, 1)$ and are independent³ then the sum of their squares

³ We'll see later what that means exactly

$$Y = \sum_{i=0}^n X_i^2 \quad (37)$$

follows a $\chi^2(n)$ law. The parameter K is the number of independent Gaussians in the sum and is dubbed *degrees of freedom* in most literature. The p.d.f. of a χ^2 random variable is

$$f(y; n) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} y^{\frac{n}{2}-1} \exp\left(-\frac{1}{2}y\right) & y \geq 0 \\ 0 & y < 0 \end{cases} \quad (38)$$

The function $\Gamma(\cdot)$ is the Gamma integral; it is a common special function that behaves like the continuous equivalent of the factorial function.

Expectation, moments and quantiles The expectation of a random variable is defined as

$$\mathbb{E}[X] = \sum_i x_i \mathbb{P}(X = x_i) \quad (39)$$

if X is a discrete r.v. and takes on values in a set $\{x_i\}$ and

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} xf(x)dx \quad (40)$$

if X is a continuous r.v. The integration is limited to the domain if it is less than the entire real line. This basic definition is usually called the mean, although expectation is actually defined for any function $g(X)$, for example in the continuous case⁴

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x)f(x)dx. \quad (41)$$

Since it is defined as an integral, the expectation is linear, in the sense that if $g(X) = aX + b$ for some scalars a and b , we have

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b. \quad (42)$$

Among possible functions $g(\cdot)$, monomials x^k are called k -th order moments, and

$$\mathbb{E}[(X - \mathbb{E}[X])^k] \quad (43)$$

are called *centered moments*. The centered moment of order 2 is called *variance* and holds a very special place, so it sometimes has its own notation

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] \quad (44)$$

There is an interesting interplay between the mean and variance. Consider a r.v. X such that $\mathbb{E}[X^2] < \infty$. Then for all c in \mathbb{R} we have that

$$\begin{aligned} \mathbb{E}[(X - c)^2] &= \mathbb{E}[X^2] - 2\mathbb{E}[X]c + c^2 \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]c + c^2 + \text{Var}(X) \\ &= (\mathbb{E}[X] - c)^2 + \text{Var}(X) \end{aligned}$$

This also says that there is an extremal for the mean with respect to the variance, in the sense that $\mu = \mathbb{E}[X]$ iff

$$\mathbb{E}[(X - \mu)^2] = \min_{c \in \mathbb{R}} \mathbb{E}[(X - c)^2] \quad (45)$$

Any other choice of c beyond the mean leads to a higher value for $\mathbb{E}[(X - c)^2]$.

Among other characterizations of an r.v. an interesting one is the *quantile*, defined as the solution⁵ to

$$F(q_p) = p \quad (46)$$

for some $0 \leq p \leq 1$. Within quantiles the value $M = q_{1/2}$ is called the *median* and has special properties. As per the definition, we have that

$$\mathbb{P}(X \geq M) \geq 1/2 \text{ and } \mathbb{P}(X \leq M) \geq 1/2 \quad (47)$$

Note that the expectation may not exist, in which case usually the integral definition will diverge or be improper. Sometimes you can salvage this with measure theory, but there exists even common cases where you can't—e.g. Cauchy r.v.s.

⁴ Unless specified, the discrete case mirrors all definitions.

$$\begin{aligned} &\mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \end{aligned}$$

⁵ This definition only really makes sense for well-behaved continuous r.v.

Like the mean, the median characterizes the *position* or *location* of the r.v., whereas values like the variance characterize its *scale*. Unlike the mean however the median is always defined. It also has an extremal property, where

$$\mathbb{E} [|X - M|] = \min_{c \in \mathbb{R}} \mathbb{E} [|X - c|] \quad (48)$$

However this point may not be unique, since in the very general case there may be many possible values of a verifying (47).

Inequalities There are interesting inequalities in probability theory that should be known. If anything, remember that they exist so that you can check whether they can help you derive bounds on quantities of interest involving random variables. They are also used in proofs a lot, for obvious reason: in mathematics, bounding a value is a very powerful statement.

Proposition 1 (Markov inequality). *Let $h(\cdot)$ be an increasing and positive function and X a non-negative r.v. such that $\mathbb{E}[h(X)] < \infty$. Then for any a such that $h(a) > 0$ we have*

$$\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[h(X)]}{h(a)} \quad (49)$$

Proof. $h(\cdot)$ is an increasing positive function, so

$$\mathbb{P}(X \geq a) \leq \mathbb{P}(h(X) \geq h(a)) = \int \mathbb{1}_{h(X) \geq h(a)} f(x) dx \quad (50)$$

Then

$$\int \mathbb{1}_{h(X) \geq h(a)} f(x) dx = \mathbb{E} [\mathbb{1}_{h(X) \geq h(a)}] \leq \mathbb{E} \left[\frac{h(X)}{h(a)} \mathbb{1}_{h(X) \geq h(a)} \right] \quad (51)$$

Then

$$\frac{\mathbb{E} [h(X) \mathbb{1}_{h(X) \geq h(a)}]}{h(a)} = \frac{\mathbb{E} [h(X)]}{h(a)} - \frac{\mathbb{E} [h(X) \mathbb{1}_{h(X) < h(a)}]}{h(a)} \quad (52)$$

The last term is positive since $h(\cdot)$ is positive; removing it thus leads to the expected result. \square

Proposition 2 (Chebyshev inequality). *Let X be an r.v. with $\mathbb{E}[X^2] < \infty$. Then assuming $a > 0$*

$$\mathbb{P}(|X| \geq a) \leq \frac{\mathbb{E}[X^2]}{a^2} \quad \mathbb{P}(|X - \mathbb{E}[X]| \geq a) \leq \frac{\text{Var} X}{a^2} \quad (53)$$

Proof. Use Markov's inequality. \square

TODO: understand the derivation

This notation uses indicators rather than bound; notice that

$$\mathbb{P}(h(X) \geq h(a)) = \mathbb{E} [\mathbb{1}_{h(X) \geq h(a)}]$$

which is the basis of a whole way to define probability theory through expectations.

Proposition 3 (Hölder inequality). Let $1 < r < \infty$ and $1/r + 1/s = 1$. Let X and Y be 2 r.v. such that $\mathbb{E}[|X|^r] < \infty$ and $\mathbb{E}[|Y|^s] < \infty$. Then $\mathbb{E}[|XY|] < \infty$ and

$$\mathbb{E}[|XY|] \leq \mathbb{E}[|X|^r]^{\frac{1}{r}} \mathbb{E}[|Y|^s]^{\frac{1}{s}} \quad (54)$$

Proof. Remark that

$$\frac{1}{r} \log a + \frac{1}{s} \log b \leq \log \left(\frac{a}{r} + \frac{b}{s} \right) \quad (55)$$

for $a > 0$ and $b > 0$, since \log is a concave function. This is equivalent to

$$a^{\frac{1}{r}} b^{\frac{1}{s}} \leq \frac{a}{r} + \frac{b}{s} \quad (56)$$

Then, substitute

$$a = \frac{|X|^r}{\mathbb{E}[|X|^r]} \quad b = \frac{|Y|^s}{\mathbb{E}[|Y|^s]} \quad (57)$$

which leads to

$$|XY| \leq \mathbb{E}[|X|^r]^{\frac{1}{r}} \mathbb{E}[|Y|^s]^{\frac{1}{s}} \left(\frac{|X|^r}{r\mathbb{E}[|X|^r]} + \frac{|Y|^s}{s\mathbb{E}[|Y|^s]} \right) \quad (58)$$

Taking the expectation on both sides gives the proof. \square

Proposition 4 (Lyapunov inequality). Let $0 < v < t$ and $1/r + 1/s = 1$. Let $X \mathbb{E}[|X|^t] < \infty$. Then $\mathbb{E}[|X|^v] < \infty$ and

$$\mathbb{E}[|X|^v]^{\frac{1}{v}} \leq \mathbb{E}[|X|^t]^{\frac{1}{t}} \quad (59)$$

Proof. Use Hölder's with $X = X^v$, $Y = 1$, $r = t/v$. \square

Proposition 5 (Jensen inequality). If $g(\cdot)$ is a convex function and $\mathbb{E}[|X|] < \infty$, then

$$g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)] \quad (60)$$

Proof. Since $g(\cdot)$ is convex, for any point a in \mathbb{R} there is an affine function $\ell_a(x)$ which is dominated by $g(x)$ and for which $\ell_a(a) = g(a)$. If $g(\cdot)$ is differentiable, this is the tangent of $g(\cdot)$ at a . Then for the r.v. X we have that

$$\ell_a(X) \leq g(X) \implies \mathbb{E}[\ell_a(X)] \leq \mathbb{E}[g(X)]$$

Since $\ell_a(\cdot)$ is linear,

$$\mathbb{E}[\ell_a(X)] = \ell_a(\mathbb{E}[X])$$

Then setting $a = \mathbb{E}[X]$ finishes the proof. \square

Proposition 6 (Cauchy-Schwarz inequality). Let X and Y be 2 r.v. with finite variance. Then

$$\mathbb{E}[XY]^2 \leq \mathbb{E}[|XY|]^2 \leq \mathbb{E}[X^2]\mathbb{E}[Y^2] \quad (61)$$

This inequality implies that

$$\mathbb{E}[|X|] \leq \mathbb{E}[|X|^2]^{\frac{1}{2}} \leq \dots \leq \mathbb{E}[|X|^k]^{\frac{1}{k}}$$

For real variables this means that $\mathbb{E}[|X|] < \infty$ is always true if $\mathbb{E}[X^2] < \infty$.

Proof. Note that the second inequality can be obtained through Hölder's. By Jensen's inequality, we have that

$$|\mathbb{E}[X]| \leq \mathbb{E}[|X|] \quad (62)$$

which can be used to prove the first inequality. \square

Independence Let us start by defining the joint repartition of a couple of continuous r.v. X and Y as

$$F(x, y) = \mathbb{P}(X \leq x, Y \leq y) \quad (63)$$

Notice once again that we are weighing a set of values defined by bounds on said values. The distribution function provides the weight to assign to the set, in a way. If everyone is well behave, we can define the joint p.d.f. as

$$f_{X,Y}(x, y) = \frac{\partial^2 F(x, y)}{\partial x \partial y} \quad (64)$$

The density should verify

$$\iint_{\mathbb{R}^2} f_{X,Y}(x, y) dx dy = 1 \quad (65)$$

as in the single r.v. case—we are still normalizing our probabilities to 1. The *marginal* p.d.f. of X and Y are obtained by integrating over a single variable

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy \quad f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx \quad (66)$$

Knowing the marginals does not give you the joint distribution, but the opposite is true. In the discrete case the joint p.d.f. is naturally defined through $\mathbb{P}(X = x, Y = y)$.

We can now define independence as follows: we say that 2 r.v. are independent iff their *joint* p.d.f. factors in separate parts depending only on their respective variables, i.e.

$$f_{X,Y}(x, y) = f_X(x) f_Y(y) \quad (67)$$

The separate terms are in fact the marginals; you can check this by definition. We can revisit our previous statement: knowing the marginals gives you the joint distribution iff the variables are independent. If we also look at the joint expectation of X and Y through the straightforward definition, we see that

$$\mathbb{E}[XY] = \iint xy f(x, y) dy dx \quad (68)$$

Whenever X and Y are independent, this can be turned into

$$\begin{aligned}\mathbb{E}[XY] &= \iint xyf(x)f(y)dydx \\ &= \int xf(x) \left(\int yf(y)dy \right) dx \\ &= \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

The separation of the integrals is known as Fubini's theorem or Tonelli's variant. It requires conditions on $f(x)$ and $f(y)$ which are a bit beyond the scope of this class. We will just say that for most non-pathological cases at hand they are indeed satisfied. Note here that while independence implies that $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$, having $\mathbb{E}[X]\mathbb{E}[Y]$ is not a sufficient condition for independence of X and Y .

Conditioning If the variable are not independent, we can develop on how they depend on one another. Here, the joint p.d.f. will always tell you everything you would need, but in practice it is not always the most practical to manipulate or observe. In this case, you can try to study the conditioned r.v. $Y|X$, with its p.d.f. defined as

$$f(y|x) = \begin{cases} \frac{f(x,y)}{f(x)} & \text{if } f(x) > 0 \\ f(y) & \text{if } f(x) = 0 \end{cases} \quad (69)$$

Check that $Y|X$ is indeed an r.v., and check what happens when X and Y are independent.

We also define the conditional expectation as a function of x through

$$\mathbb{E}(Y|X = x) = \int yf(y|x)dy \quad (70)$$

This last value is a scalar for any x . However, X is itself a r.v. so in full generality we may also see $\mathbb{E}[Y|X]$ as a r.v. As an expectation, it also has many properties that we expect, among which linearity, as well as

1. $\mathbb{E}[Y|X] = \mathbb{E}[Y]$ if X and Y are independent
2. $\mathbb{E}[h(X)|X] = h(X)$ for well-behaved functions $h(\cdot)$
3. $\mathbb{E}[h(X, Y)|X = x] = \mathbb{E}[h(x, Y)|X = x]$

TODO: can the last one be shown without measure getting in the way? It's more or less OK in the discrete case.

Covariance Dependence and conditioning can be a bit annoying to manipulate in practice, and in many cases they are very hard to characterize. Let us define the covariance between 2 r.v. with finite variance as

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \quad (71)$$

There are some specific and useful properties for the covariance

Blitz exercises

1. $\text{Cov}(X, X) = \text{Var}(X)$
2. $\text{Cov}(aX, bY) = ab\text{Cov}(X, Y)$
3. $\text{Cov}(X + a, Y) = \text{Cov}(X, Y)$
4. $\text{Cov}(Y, X) = \text{Cov}(X, Y)$
5. $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$
6. $\text{Cov}(X, Y) = 0$ if X and Y are independent

Random vectors Beyond couples of random variables, we will encounter a lot of random vectors within this course. A random vector is just a collection of r.v. indexed and stacked into a vector $\mathbf{x} = (X_1, \dots, X_n)^T$. We can define its mean as the vectors of means of the individual components, i.e.

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] = \begin{pmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_n] \end{pmatrix} \quad (72)$$

Similarly, we can look at the covariance between the components of the vector, e.g. $\mathbb{E}[X_i X_j]$ for $1 \leq i \leq n$ and $1 \leq j \leq n$. This also has a nice matrix through what is sometimes called the *dyadic* or *outer* product

$$\mathbf{C} = \mathbb{E} \left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right] \quad (73)$$

and we take the expectation on each component of the matrix.

Among the possible random vectors an astounding majority of practical—or should we say solvable—problems are related to the multivariate Gaussian distribution, which is the extension of the Gaussian distribution to larger vectors using the aforementioned mean and covariance. The p.d.f. of this distribution is

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{N}{2}} \sqrt{|\mathbf{C}|}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (74)$$

You can see that this indeed looks like a *p.d.f.*; the matrix operations within the exponential are a quadratic form between the *centered* vector \mathbf{x} and the *inverse* covariance matrix \mathbf{C} . This interpretation will come in handy for us later.

A very important property of the multivariate Gaussian is that all linear transformations of a multivariate Gaussian random vector are still multivariate Gaussian random vectors. In other words, if we define

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} \quad (75)$$

Exercise: show that the covariance matrix is positive semidefinite. Consider a r.v. \mathbf{c} that is centered w.l.o.g. and has covariance matrix \mathbf{C} . For any \mathbf{x} we have that $\mathbf{x}^T \mathbf{C} \mathbf{x} = \mathbf{x}^T \mathbb{E}[\mathbf{c}\mathbf{c}^T] \mathbf{x} = \mathbb{E}[(\mathbf{x}^T \mathbf{c})(\mathbf{c}^T \mathbf{x})]$. This is the variance of some random r.v. $z = \mathbf{x}^T \mathbf{c}$, and since the variance is always positive we deduce that \mathbf{C} is positive definite.

Exercise: what happens where the components of \mathbf{x} are independent?

and \mathbf{x} follows a multivariate Gaussian distribution, then so does \mathbf{y} . Obviously, the parameters of the distribution are different, but we can compute them easily using the linearity of the expectation:

$$\mathbb{E}[\mathbf{y}] = A\mathbb{E}[\mathbf{x}] + \mathbf{b} \quad (76)$$

and

$$\mathbb{E}[(\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^T] = A\mathbf{C}A^T \quad (77)$$

These expressions are actually not conditioned on the fact that \mathbf{y} and \mathbf{x} are multivariate Gaussian; they are properties of the expectation for any random vectors, and analogous to the scalar cases $\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$ and $\text{Var}(aX + b) = a^2\text{Var}(X)$.

EXERCISE Assuming X_1, \dots, X_n are independent variables, let

$$X_{\min} = \min(X_1, \dots, X_n) \quad X_{\max} = \max(X_1, \dots, X_n)$$

Show that

$$\mathbb{P}(X_{\min} > x) = \prod_i \mathbb{P}(X_i > x) \quad \mathbb{P}(X_{\max} < x) = \prod_i \mathbb{P}(X_i < x)$$

If X_1, \dots, X_n are uniform random variables, compute $\mathbb{E}[X_{\max}]$ and $\text{Var}(X_{\max})$.

EXERCISE Let X be a positive r.v. with finite expectation. Show that

$$\mathbb{E}[X] = \int_0^\infty (1 - F(x))dx = \int_0^\infty P(X > x)dx \quad (78)$$

Hint: the expectation of an indicator is the probability.

EXERCISE Let X, Y_1 and Y_2 be independent r.v. such that Y_1 and Y_2 are $\mathcal{N}(0, 1)$. Let

$$Z = \frac{Y_1 + XY_2}{\sqrt{1 + X^2}}$$

Show that Z is $\mathcal{N}(0, 1)$ using $\mathbb{P}(Z \leq z | X = x)$.

EXERCISE Let X and N be 2 r.v. with finite absolute expectation and N takes values in the strictly positive integers. Let X_1, X_2, \dots be the sequence of r.v. with the same law as X . Using conditioning, show Wald's inequality

$$\mathbb{E} \left[\sum_{i=1}^N X_i \right] = \mathbb{E}[N]\mathbb{E}[X]$$