

Probability

σ -algebra: Let Ω be a set. A family $\mathcal{F} \subset 2^\Omega$ is a σ -algebra if

- $\emptyset, \Omega \in \mathcal{F}$
- $\forall A \in \mathcal{F}, A^c \in \mathcal{F}$
- $\forall (A_n) \subset \mathcal{F}$ sequence in \mathcal{F} , $\bigcup_{i \in \mathbb{N}} A_i \in \mathcal{F}$

Then, the set (Ω, \mathcal{F}) is called a measurable space

Given a topological space $(\Omega, \mathcal{T}_\Omega)$, the Borel σ -algebra is the smaller σ -algebra that contains \mathcal{T}_Ω

Measurable function: Let $(\Omega_1, \mathcal{F}_1)$ and $(\Omega_2, \mathcal{F}_2)$ be measurable spaces. A function $f: \Omega_1 \rightarrow \Omega_2$ is called measurable if $f^{-1}(A) \in \mathcal{F}_1, \forall A \in \mathcal{F}_2$

Probability measure: A probability measure on (Ω, \mathcal{F}) measurable space is a function $P: \mathcal{F} \rightarrow [0, +\infty]$ such that

- $P[\emptyset] = 0$
- $P[\Omega] = 1$ (without it, P is just a measure on (Ω, \mathcal{F}))
- Given a sequence $(A_n) \subset \mathcal{F}$ two-by-two distinct, then

$$P\left[\bigcup_{n \in \mathbb{N}} A_n\right] = \sum_{n \in \mathbb{N}} P[A_n] \quad \begin{matrix} \leftarrow & \text{the infinite } +\infty \text{ for countable union arises from} \\ & \text{the divergence of this series} \end{matrix}$$

(Ω, \mathcal{F}, P) is called probability space

$$\text{Note that } P(A_1) = P(A_2) + P(A_1 - A_2) \quad \forall A_1, A_2 \in \mathcal{F}, P(A_1) < \infty$$

Lemma: Let $(A_n) \subset \mathcal{F}$ a crescent sequence, i.e., $A_1 \subseteq A_2 \subseteq \dots$. Then

$$P\left[\bigcup_{n \in \mathbb{N}} A_n\right] = \lim_{n \rightarrow \infty} P[A_n]$$

The integral of f with respect to P is

$$\int_{\Omega} f dP := \int_{\Omega} f^+ dP - \int_{\Omega} f^- dP$$

The set of all the functions with finite integral with respect to P is called $L(P)$

Lemma: Every $f \in L(P)$ induces a probability measure given by

$$\tilde{P}(A) = \left| \frac{\int_A f dP}{\int_{\Omega} f dP} \right| \quad \forall A \in \mathcal{F}$$

Theorem: A measurable function $f \in L(P) \iff |f| \in L(P)$

In special, if $g \in L(P)$ and $|f| \leq |g| \Rightarrow f \in L(P)$ and

$$\left| \int_{\Omega} f dP \right| \leq \int_{\Omega} |f| dP \leq \int_{\Omega} |g| dP$$

Monotone Convergence Theorem: If (f_n) is a monotone increasing sequence of functions in $M^+(\Omega, \mathcal{F})$ which converges to f , then

$$\int_{\Omega} f dP = \lim_{n \rightarrow \infty} \int_{\Omega} f_n dP$$

Fatou's Lemma: If $(f_n) \subset M^+(\Omega, \mathcal{F})$, then

$$\int_{\Omega} (\liminf f_n) dP \leq \liminf \left(\int_{\Omega} f_n dP \right)$$

Dominated Convergence Theorem: Let $(f_n) \in L(p)$ which converges almost everywhere to a function f . If $\exists g \in L(p)$ such that $|f_n| \leq g \quad \forall n \in \mathbb{N}$, then $f \in L(p)$ and

$$\int_{\Omega} f \, dP = \lim_{n \rightarrow \infty} \int_{\Omega} f_n \, dP$$

Useful bounds

- $|e^t - 1| \leq |t| e^{|t|} \quad \forall t \in \mathbb{R}$

- $|t| \leq e^{|t|} \quad \forall t \in \mathbb{R}$

Uniform convergence: The sequence (f_n) converges uniformly to f if, for every $\epsilon > 0$, $\exists N(\epsilon) > 0$ such that

$$N(\epsilon, x) = N(\epsilon)$$

$$n > N(\epsilon) \text{ and for } x \in \Omega \Rightarrow |f_n(x) - f(x)| < \epsilon$$

Radon-Nikodym Theorem: Let λ and μ be σ -finite measures defined on \mathcal{F} and suppose that $\lambda \ll \mu$. Then, $\exists \frac{d\lambda}{d\mu} \in M^+(\Omega, \mathcal{F})$ (not necessarily integrable) function such that

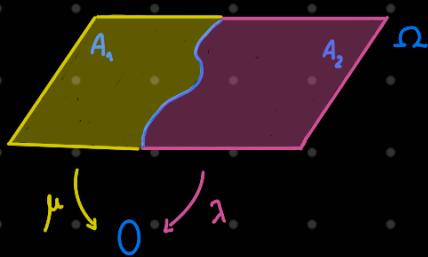
$$\lambda(A) = \int_A \frac{d\lambda}{d\mu} \, d\mu, \quad \forall A \in \mathcal{F},$$

$\lambda \ll \mu \Leftrightarrow \forall A \in \mathcal{F} \text{ such that } \mu(A) = 0, \text{ it implies that } \lambda(A) = 0$

(λ is absolutely continuous with resp. to μ)

and $\frac{d\lambda}{d\mu}$ is unique up to sets of measure μ null

Lebesgue Decomposition Theorem: Let λ and μ be σ -finite measures on (Ω, \mathcal{F}) . Then, \exists unique measures $\lambda_1 \ll \mu$, $\lambda_2 \perp \mu$ such that $\lambda = \lambda_1 + \lambda_2$.



$\lambda \perp \mu \Leftrightarrow \exists A_1, A_2 \in \mathcal{F}, A_1 \cap A_2 = \emptyset, \text{ such that}$

- $A_1 \cup A_2 = \Omega$

- $\mu(A_1) = 0 = \lambda(A_2)$
(λ and μ are mutually singular)

μ -equivalence class: Let $f \in L(\mu)$. Then,

$$[f] := \{g \in L(\mu) : f(x) = g(x) \quad \forall x \in A \text{ such that } A \in \mathcal{F} \text{ and } \mu(A) \neq 0\}$$

L^p -spaces: For $1 \leq p < \infty$,

$$L^p(\mu) = \left\{ [f] : f \text{ is } \mu\text{-measurable and } \left[\int_{\Omega} |f|^p d\mu \right]^{\frac{1}{p}} \right\}$$

are the complete normed vector spaces with the norm

$$\|f\|_p = \left[\int_{\Omega} |f|^p d\mu \right]^{\frac{1}{p}} \quad \begin{matrix} \leftarrow \text{the representant in } [f] \\ \text{doesn't matter} \end{matrix}$$

L^∞ -spaces: $L^\infty(\mu)$ is the set of all μ -equivalent classes which are bounded up to sets of measure zero, i.e,

$$L^\infty(\mu) := \left\{ [f] : f \text{ is measurable and } \forall A \in \mathcal{G} \text{ such that } \mu(A) \neq 0, \sup \{|f(x)| : x \in A\} < \infty \right\}$$

This space is a normed vector space. In fact, for $A \in \mathcal{G}$ such that $\mu(A) = 0$, for $f \in L^\infty(\mu)$ define

$$S(A) := \sup \{|f(x)| : x \notin A\}$$

$$\|f\|_\infty := \inf \{S(A) : A \in \mathcal{G}, \mu(A) \neq 0\} \quad \begin{matrix} \text{lower number} \\ \text{that limits all} \\ \text{sets with measure} \\ \text{non-null} \end{matrix}$$

Bounded linear function: A linear functional on $L^p(\mu)$, for $1 \leq p < \infty$, is a function $G : L^p(\mu) \rightarrow \mathbb{R}$ such that

$$G(af + bg) = aG(f) + bG(g) \quad \forall a, b \in \mathbb{R}, \forall f, g \in L^p(\mu).$$

G is bounded if $\exists M$ constant such that $|G(f)| \leq M\|f\|_p$. In this case, the norm or bound of G is

$$\|G\| := \sup \{|G(f)| : f \in L^p(\mu) \text{ such that } \|f\|_p \leq 1\}$$

Riesz Representation Theorem: If G is a bounded linear function on $L^p(\mu)$ for $1 < p < \infty$, $\exists g \in L^q(\mu)$, where $\frac{1}{p} + \frac{1}{q} = 1$, such that

$$G(f) = \int_{\Omega} fg d\mu \quad \forall f \in L^p(\mu) \quad \leftarrow \text{We represent } G \text{ as the integral of a function}$$

and $\|G\| = \|g\|_q$.

For $p \in \{1, \infty\}$, if μ is σ -finite measure and G is a bounded linear functional on $L^p(\mu)$, then $\exists g \in L^\infty(\mu)$ such that

$$G(f) = \int_{\Omega} fg d\mu \quad \forall f \in L^p(\mu)$$

Moreover, $\|G\| = \|g\|_\infty$ and $g \geq 0$ if G is positive, i.e., $\text{Im}(G) \subseteq \mathbb{R}_{\geq 0}$.

Algebra: A family $\tilde{\mathcal{F}} \subset 2^\Omega$ is an algebra of subsets of Ω if

- $\emptyset, \Omega \in \tilde{\mathcal{F}}$
- $A \in \tilde{\mathcal{F}} \Rightarrow A^c \in \tilde{\mathcal{F}}$
- $A_1, \dots, A_n \in \tilde{\mathcal{F}} \Rightarrow \bigcup_{i=1}^n A_i \in \tilde{\mathcal{F}}$ It distinguishes an algebra of a σ -algebra,
because here the union is finite

Then, a measure on $\tilde{\mathcal{F}}$ is a function $\mu: \tilde{\mathcal{F}} \rightarrow [0, \infty]$ such that

- $\mu(\emptyset) = 0$
- If $(A_n)_{n \in \mathbb{N}} \subset \tilde{\mathcal{F}}$ is a sequence of two-by-two disjoint sets such that $\bigcup_{i \in \mathbb{N}} A_i \in \tilde{\mathcal{F}}$, then

$$\mu\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \sum_{i \in \mathbb{N}} \mu(A_i)$$

Outer measure: Let $H \in 2^\Omega$, we define the outer measure

$$\mu^*(H) = \inf \left\{ \sum_{i \in \mathbb{N}} \mu(A_i) : (A_i)_{i \in \mathbb{N}} \subset \tilde{\mathcal{F}} \text{ such that } H \subset \bigcup_{i \in \mathbb{N}} A_i \right\}$$

\hookrightarrow It's subadditive: $(H_n)_{n \in \mathbb{N}} \subset 2^\Omega$, $\mu^*\left(\bigcup_{i \in \mathbb{N}} H_i\right) \leq \sum_{i \in \mathbb{N}} \mu^*(H_i)$

μ^* -measurable set: A set $H \subset \Omega$ is called μ^* -measurable if

$$\mu^*(A) = \mu^*(A \cap H) + \mu^*(A \setminus H) \quad \forall A \subset \Omega$$

The set of μ^* -measurable sets is called \mathcal{F}^*

Carathéodory Extension Theorem: A collection \mathcal{F}^* of all μ^* measurable sets is a σ -algebra containing \mathcal{F} . Moreover, if $(A_n)_{n \in \mathbb{N}} \subset \mathcal{F}^*$ is a sequence of two-by-two disjoint sets, then

$$\mu^*\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \sum_{i \in \mathbb{N}} \mu^*(A_i)$$

Hahn Extension Theorem: Suppose that μ is a σ -finite measure on an algebra \mathcal{F} . Then, there exists a unique extension of μ to a measure on \mathcal{F}^* .

Set-up for probability theory:

(Ω, \mathcal{F}, P)

↳ σ -algebra (space of events)

abstract space (sample space)

$$P: \mathcal{F} \longrightarrow \mathbb{R}_+ \cup \{\infty\}$$

$$\cdot P(\emptyset) = 0$$

$$\cdot \{A_j : j \geq 1\}, A_j \cap A_k = \emptyset$$

$$\cdot P\left[\bigcup_j A_j\right] = \sum_j A_j$$

$$\cdot P(\Omega) = 1 \quad (\text{without it, } P \text{ is a measure})$$

Random variable: $X: \Omega \longrightarrow \mathbb{R}$

With $(\mathbb{R}, \mathcal{B})$, X is measurable: $X^{-1}(A) \in \mathcal{F} \quad \forall A \in \mathcal{B}$

Probability distribution measure: X has a probability distribution measure μ (in the Borel σ -algebra of \mathbb{R}) if

$$\boxed{\mu = X_* P} \quad \text{i.e.} \quad \mu(A) = P[X^{-1}(A)] \stackrel{\text{notation}}{=} P[X \in A] \quad \forall A \in \mathcal{B}$$

Distribution function: The distribution function associated to X is

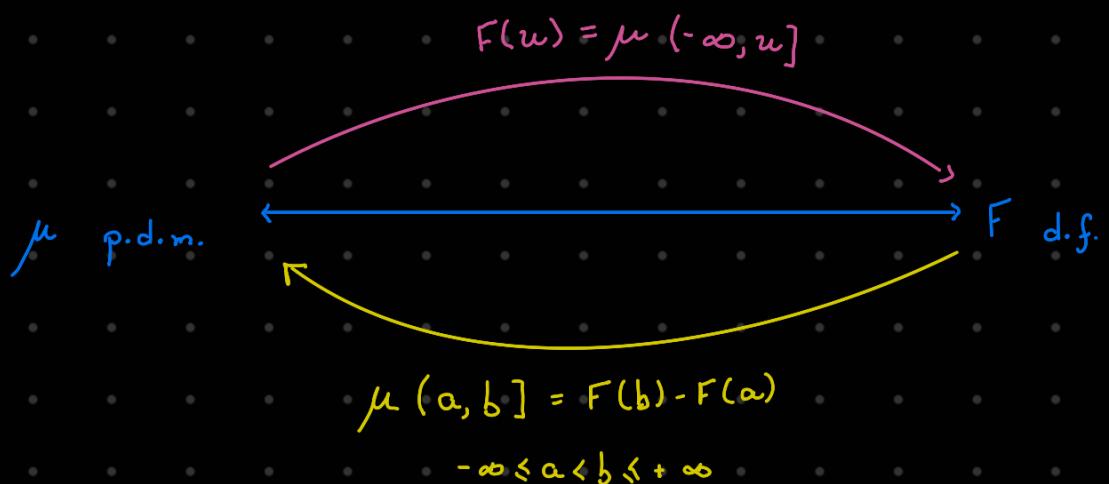
$$F_x : \mathbb{R} \longrightarrow [0, 1]$$

$$x \longmapsto P[X^{-1}((-\infty, x])] = P[X \leq x]$$



- F_x is monotone: $x < y \Rightarrow F_x(x) \leq F_x(y)$
- F_x is right-continuous, left-limit exists
- $\lim_{x \rightarrow \infty} F_x(x) = 1$, $\lim_{x \rightarrow -\infty} F_x(x) = 0$

There is an injection between the probability distribution measures and the distribution functions



$$S = \{(a, b] : -\infty < a < b < +\infty\}$$

is semi-algebra over \mathbb{R}

- $\emptyset, \mathbb{R} \in S$
- $\forall A, B \in S \Rightarrow A \cap B \in S$
- $\forall A, B \in S \Rightarrow \exists \{C_i\}_{i=1}^n$ pairwise-disjoint such that $A \setminus B = \bigcup_{i=1}^n C_i$

μ is σ -finite:

- $\exists \{C_i\}_{i \in \mathbb{N}} \subset \mathcal{B}$ such that
 - $\mu(C_i) < \infty$
 - $\sum_i C_i = \mathbb{R}$

Discrete random variables:

- X is a discrete random variable $\iff \exists A \subset \mathbb{R}$ countable such that $\sum_{x \in A} P(X=x) = 1$

Discrete distribution function:

F is a discrete d.f. $\iff \exists A = \{x_i\}_{i \geq 1} \subset \mathbb{R}$ and $\exists \{p_i\} \subset [0,1]$ such that

$$\sum_{i \geq 1} p_i = 1, \quad F(x) = \sum_{\substack{i \geq 1 \\ x_i \leq x}} p_i \quad \begin{matrix} \text{of } j \text{ such that} \\ x_j \leq x \end{matrix} \quad \text{summing}$$

Absolute continuous density function: F is absolutely continuous if $\exists f: \mathbb{R} \rightarrow \mathbb{R}_+$ measurable such that $F(b) - F(a) = \int_a^b f(t) dt \quad \forall a < b \in \mathbb{R}$

Remark: F is differentiable almost everywhere and $F' = f$

Singular continuous density function: F is singular $\iff \begin{cases} F \text{ is continuous} \\ F' = 0 \text{ almost everywhere} \end{cases}$

Theorem: Let F be a distribution function. Then, F can be written as a convex combination of a discrete, an absolute continuous and a singular density function.

$\exists \alpha, \beta \geq 0, \alpha + \beta \leq 1$, and $\exists F_d, F_{ac}, F_s$ such that

$$F = \alpha F_d + \beta F_{ac} + (1 - \alpha - \beta) F_s$$

Expectation: Let $X: \Omega \rightarrow \mathbb{R}$ a random variable such that

$$\int_{\Omega} |X| dP < \infty, \quad \leftarrow \text{ie, } X \text{ is integrable}$$

The same than suppose $X \in L(\Omega)$

then, we define

$$\mathbb{E}[X] = \int_{\Omega} X \, dP$$

Proposition: If $X \geq 0$ (ie, $X(\Omega) \subset \mathbb{R}_{\geq 0}$), then

$$\sum_{n \in \mathbb{N}} P[X \geq n] \leq \mathbb{E}[X] \leq 1 + \sum_{n \in \mathbb{N}} P[X \geq n]$$

$$\sum_{n \in \mathbb{N}} P[X^{-1}[n, +\infty)] \leq \mathbb{E}[X] \leq 1 + \sum_{n \in \mathbb{N}} P[X^{-1}[n, +\infty)]$$

Jensen inequality: Let $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ be a convex function, ie,

$$\varphi(\theta x + (1-\theta)y) \leq \theta \varphi(x) + (1-\theta)\varphi(y) \quad \forall \theta \in [0,1], \forall x, y \in \mathbb{R}.$$

Let X be a measurable random variable, ie, $\mathbb{E}[|X|] < \infty$. Assume that $\mathbb{E}[|\varphi(X)|] < \infty$. Then,

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$$

Chebyshew inequality: Let $X \geq 0$ be a random variable. Let $f: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ a monotonic non-decreasing function ($x \leq y \Rightarrow f(x) \leq f(y)$). Then,

$$P[X \geq a] \leq \frac{1}{f(a)} \mathbb{E}[f(X)] \quad P[X^{-1}[a, +\infty)] \leq \frac{1}{f(a)} \mathbb{E}[f(X)]$$

\uparrow
can be ∞

Independent: Let $A \in \mathcal{F}$. A is called an event.

events . . . Let $A_1, A_2, \dots, A_N \in \mathcal{F}$. They are called independent if, for all $\{n_1, n_2, \dots, n_k\} \subseteq \{1, \dots, N\}$, $n_j \neq n_k$ for $j \neq k$, holds

$$P[\bigcap_{k=1}^K A_{n_k}] = \prod_{k=1}^K P[A_{n_k}]$$

Independent random variables: Let X_1, \dots, X_n be random variables. They are called independent if, $\forall A_1, \dots, A_n \in \mathcal{B}$, holds

$$P\left[\bigcap_{j=1}^n \{X_j \in A_j\}\right] = \prod_{j=1}^n P[X_j \in A_j]$$

$$P\left[\bigcap_{j=1}^n X_j^{-1}(A_j)\right] = \prod_{j=1}^n P[X_j^{-1}(A_j)]$$

Independent family of random variables: Let $\{X_\alpha : \alpha \in I\}$ a family of random variables. We say this family is independent if, $\forall N \in \mathbb{N}$ and $\forall \{\alpha_1, \dots, \alpha_N\} \subset I$ with $\alpha_j \neq \alpha_k$ for $j \neq k$ (ie, for all finite subfamily of $\{X_\alpha\}$), $X_{\alpha_1}, \dots, X_{\alpha_N}$ are independent.

Remark: Any subfamily of an independent family is independent

Random vector: Consider \mathbb{R}^n with its σ -algebra \mathcal{B}^n . A n -dimensional random vector is a measurable function $X: \Omega \rightarrow \mathbb{R}^n$, ie, $\forall A \in \mathcal{B}^n$, $X^{-1}(A) \in \mathcal{I}$

Basically, it's a vector of random variables: $x = (x_1, \dots, x_n)$

Its distribution function is $F_x: \mathbb{R}^n \rightarrow [0, 1]$ \uparrow this entrance takes the first coordinate of A

$$F_x: \mathbb{R}^n \longrightarrow [0, 1] \\ (x_1, \dots, x_n) \longmapsto P[X_1 \leq x_1, \dots, X_n \leq x_n] = P\left[\bigcap_{j=1}^n \{X_j \leq x_j\}\right]$$

And its probability distribution measure is

$$\mu_x = X_* P, \text{ ie, } \mu(A) = P[X^{-1}(A)] = P[X \in A] \quad \forall A \in \mathcal{B}^n$$

Lemma: Let $X = (X_1, \dots, X_n)$ a random vector. Each X_j is associated to a F_{X_j}

$$X_1, \dots, X_n \text{ are independent} \iff F(x_1, \dots, x_n) = \prod_{j=1}^n F_{X_j}(x_j) \quad \forall x = (x_1, \dots, x_n) \in \mathbb{R}^n$$

By the identification between F and μ

Lemma: Let $X = (X_1, \dots, X_n)$ a random vector. Each X_j is associated to a μ_{X_j}

$$X_1, \dots, X_n \text{ are independent} \iff \mu_X(A_1 \times \dots \times A_n) = \prod_{j=1}^n \mu_{X_j}(A_j) \quad \forall A_1 \times \dots \times A_n \in \mathcal{B}^n$$

Theorem: Independence is preserved by measurable functions. Given X_1, \dots, X_n random variables and $f_1, \dots, f_n: \mathbb{R} \rightarrow \mathbb{R}$ measurable functions, then $f_1(X_1), \dots, f_n(X_n)$ are independent



we can take functions from more than one random variable and it still holds

Basically, independence is preserved by measurable functions

Theorem: Let X, Y independent random variables such that

$$\cdot \mathbb{E}[|X|], \mathbb{E}[|Y|] < \infty$$

Then,

$$\boxed{\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]}$$

Identically distributed random variables: Let X_1, \dots, X_n be random variables. They are identically distributed if

$$\mu_{X_1} = \mu_{X_2} = \dots = \mu_{X_n} = (X_1)_* P$$

Weak law of large numbers: Let X_1, \dots, X_n iid random variables. Suppose that

- $\mathbb{E}[X_1] < \infty$ ← as they are identically distributed, it's true for the others

Then,

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\left(\frac{X_1 + \dots + X_n}{n} - \mathbb{E}[X_1] \right)^2 \right] = 0 \quad \begin{matrix} \text{The empirical average of } \{X_i\}_{i=1}^n \\ \text{converges in } L^2 \text{ to } \mathbb{E}[X_1] \text{ for } n \text{ large enough} \end{matrix}$$

More general, for a sequence $\{X_i\}_{i=1}^n$ of iid random variables, then

$$\lim_{n \rightarrow \infty} X_n \xrightarrow{\text{function zero}} 0 \implies \lim_{n \rightarrow \infty} \mathbb{E}[|X_n|^\alpha] = 0 \quad \forall \alpha > 0$$

Consequently, $\forall \varepsilon > 0$, $\lim_{n \rightarrow \infty} P[|X_n| > \varepsilon] = 0$

and, specially,

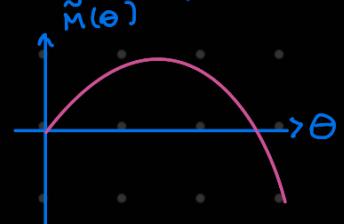
$$\lim_{n \rightarrow \infty} P \left[\left| \frac{X_1 + \dots + X_n}{n} - \mathbb{E}[X_1] \right| > \varepsilon \right] = 0$$

Upper bound for large deviations: Let X_1, \dots, X_n be iid random variables and $\theta \in \mathbb{R}$ such that $\mathbb{E}[e^{\theta X_1}] < \infty$. Define $M(\theta) = \log \mathbb{E}[e^{\theta X_1}]$. Then, $\forall a > 0$,

$$P \left[\frac{X_1 + \dots + X_n}{n} - m > a \right] \leq e^{n[-(a + \mathbb{E}[X_1])\theta + M(\theta)]} \leq e^{-nI(a)},$$

where $I(a) = \sup_{\theta > 0} \{a\theta - \tilde{M}(\theta)\}$ is the Legendre transformation for

the convex function $\tilde{M}(\theta) = M(\theta) - \mathbb{E}[X_1]\theta$



$\mathbb{E}_\theta[X]$: When someone writes \mathbb{E}_θ , he is parametrizing \mathbb{E} changing the measure P by a positive and unnormalized factor e^θ , $\theta \in \mathbb{R}$:

$$\mathbb{E}_\theta[f(x)] = \frac{\mathbb{E}[f(x)e^{\theta x}]}{\mathbb{E}[e^{\theta x}]} = \int_{\Omega} \frac{f(x)e^{\theta x} dP}{\int_{\Omega} e^{\theta x} dP},$$

If we have a family of densities $\{p_\theta : \theta \in \mathbb{R}\}$, essa notação pode significar $\mathbb{E}_\theta[x] = \int_{\Omega} p_\theta(x) \log p_0(x) dx$

for f measurable, continuous and bounded

$$\mathbb{E}[e^x] < \infty \Rightarrow \mathbb{E}[x^2] < \infty \Rightarrow \mathbb{E}[|x|] < \infty$$

Lower bound for large deviations: Given X_1, \dots, X_n iid random variables, for all open $G \subset \mathbb{R}$,

$$\liminf_{n \rightarrow \infty} \left(\frac{1}{n} \log P \left[\frac{X_1 + \dots + X_n}{n} \in G \right] \right) \geq -\inf_{\alpha \in G} I(\alpha),$$

$$\text{where } I(\alpha) = \sup_{\lambda \in \mathbb{R}} \{ \lambda \alpha - M(\lambda) \}$$

Almost Sure convergence: A sequence of random variables (X_n) converges almost surely to $X \Leftrightarrow \exists A \in \mathcal{F}, P[A] = 1$, such that, $\forall \omega \in A, X_n(\omega) \rightarrow X(\omega)$, ie, $\forall \omega \in A,$

$$X_n(\omega) \rightarrow X(\omega) \Leftrightarrow \forall \kappa > 1, \exists n \geq 1 \text{ such that } \forall m \geq n, |X_m(\omega) - X(\omega)| \leq 1/\kappa$$

It's weaker than $\forall \omega \in \Omega, X_n(\omega) \rightarrow X(\omega)$

Convergence in probability: A sequence of random variables (X_n) converges in probability to X , ie,

$$X_n \xrightarrow{P} X \Leftrightarrow \forall \varepsilon > 0, \lim_{n \rightarrow \infty} P[|X_n - X| > \varepsilon] = 0$$

Lemma: $X_n \xrightarrow{\text{almost sure}} X \Rightarrow X_n \xrightarrow{P} X$

Convergence almost surely
is stronger than convergence
in probability

Lemma: $X_n \xrightarrow{P} X \Rightarrow \exists \text{ subsequence } (X_{n_k}) \subset (X_n) \text{ such that } X_{n_k} \xrightarrow{\text{almost sure}} X$

Convergence in L^p : For $0 < p < \infty$, given a sequence (X_n) of random variables,

$$X_n \xrightarrow{L^p} X \iff \lim_{n \rightarrow \infty} \mathbb{E}[|X_n - X|^p] = 0$$

Lemma: $X_n \xrightarrow{L^p} X$ for $0 < p < \infty \Rightarrow X_n \xrightarrow{P} X$

Theorem: Let (X_n) be a sequence of random variables. Given $0 < p < \infty$, suppose that

- $X_n \xrightarrow{P} X$
- $\exists Y$ random variable bounding X_n : $|X_n| \leq Y \forall n \in \mathbb{N}$
- $\mathbb{E}[Y^p] < \infty$

Then, $X_n \xrightarrow{L^p} X$

Weak convergence: Given a sequence (μ_n) of probability distribution measures, it converges weakly to μ , i.e.,

$$\mu_n \xrightarrow{d} \mu \iff \forall (a, b] \text{ such that } \mu_n((a, b]) = 0 = \mu((a, b]),$$

$$\mu_n((a, b]) \rightarrow \mu((a, b])$$

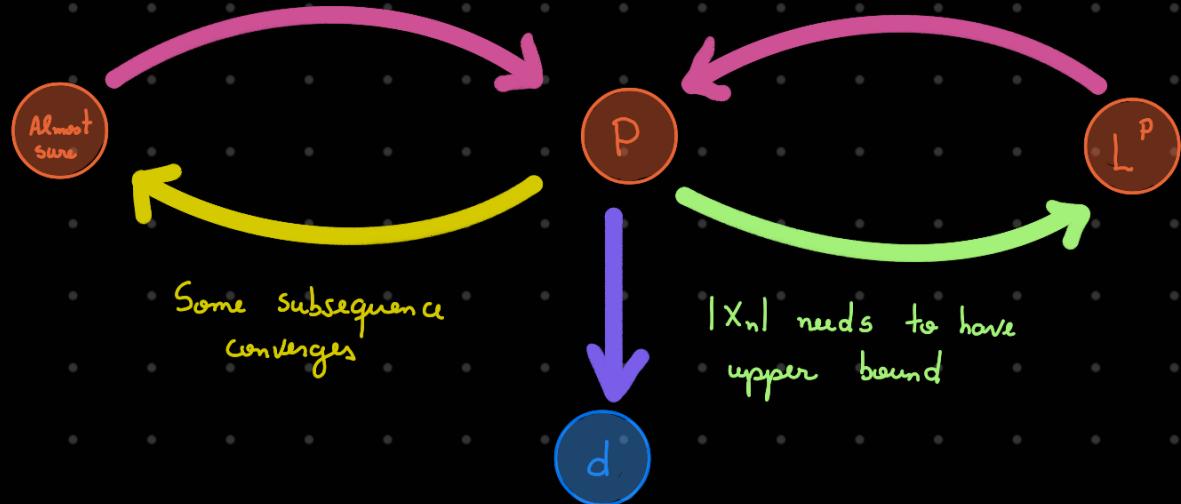
$\iff F_n(x) \rightarrow F(x) \quad \forall x \in \mathbb{R}$, where F_n, F are the distribution functions for μ_n and μ , resp.

the definitions are equivalent by the identification $\mu \leftrightarrow F$

this hypothesis is to deal with the δ distributions, because without this, each δ distribution would converge to each others

Convergence in distribution: (X_n) sequence of random variables converges in distribution to X if its sequence (F_{X_n}) of distribution functions converges to F_X

$$X_n \xrightarrow{d} X \Leftrightarrow F_{X_n} \longrightarrow F_X$$



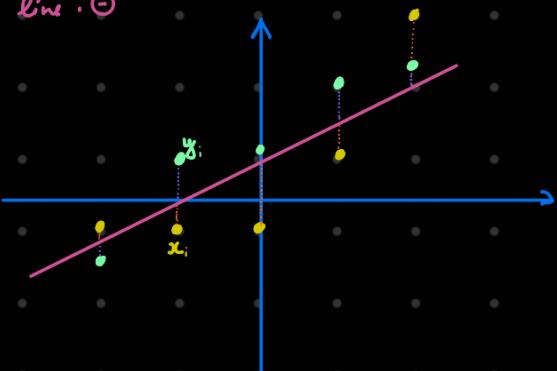
Least squared method: Given two random variables $X, Y : \Omega \rightarrow \mathbb{R}$, we estimate the intercept β_0 and the slope β_1 parameters as

$$(\beta_0, \beta_1) := \arg \min \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$$

and the error term of this estimation is defined as

$$\epsilon_i = y_i - \beta_0 - \beta_1 x_i \leftarrow \text{with sign } \begin{cases} \text{upper the line: } + \\ \text{under the line: } - \end{cases}$$

This method minimizes the sum of the areas of the squares which side lengths is given by the errors ϵ_i . The error of each point is given by the length of the vertical line connecting it and the line



It's not a good method when we have different degrees of numerical roundings. For example, for 4 observations, $y_1 \in [4, 6)$, $y_2 \in [3, 4)$, $y_3 \in [5, \infty)$, $y_4 \in [-\infty, 3)$ obtained from a normal process which we want to estimate. Y has distribution measure $N(\mu, \sigma)$.

1^o problem: Interval data: With intervals we can not compute ϵ_i because y_i is not known.

2^o problem: There is no independent variable X in this problem to estimate μ .

Joint probability density: Let $X: \Omega_1 \rightarrow \mathbb{R}$, $Y: \Omega_2 \rightarrow \mathbb{R}$ be random variables having said that, we obtain a new measure space $(\Omega_1 \times \Omega_2, \Sigma^*, P^*)$, where Σ^* is the product σ -algebra and $P^* := P_1 \otimes P_2$, i.e., $P^*(x, y) = P_1(x)P_2(y)$.

The joint probability density is the function $\rho_{x,y}: \Omega_1 \times \Omega_2 \rightarrow \mathbb{R}^2$

$$P^*((x, y) \in A) = \int_A \rho_{x,y}(x, y) dx dy \quad \forall A \in \mathcal{B}_{\mathbb{R}^2},$$

i.e., \hookrightarrow It implies $\rho_{x,y} \geq 0$ and $\int_{\mathbb{R}^2} \rho_{x,y} dx dy$

$$(x \otimes y)^*(P)(A) = P((x \otimes y)^{-1}(A)) = \int_A \rho_{x,y}(x, y) dx dy \quad \forall A \in \mathcal{B}_{\mathbb{R}^2}$$

\hookrightarrow If X, Y are independent, $\rho_{x,y} = \rho_x \cdot \rho_y$, where the probability density ρ_x and ρ_y are defined analogously.

Note it's different from probability distribution measures. In that case, we define a measure on $(\mathbb{R}, \mathcal{B})$ in terms of a probability measure P . Here, we have two measures (P on Σ and $\mu := dx dy$ on $\mathcal{B}_{\mathbb{R}^2}$) and we are defining a density function

\hookrightarrow It's just the Radon-Nikodin derivative

Marginal density: Given $\rho_{x,y}$ joint distribution density, we define the marginal density with respect to X and Y , respectively, as

$$\rho_x(x) = \int_{\mathbb{R}} \rho_{x,y}(x, y) dy, \quad \rho_y(y) = \int_{\mathbb{R}} \rho_{x,y}(x, y) dx.$$

Likelihood function: Let X_1, \dots, X_n be random variables. Let Θ a set called the parameter space. Fixed $\theta \in \Theta$, let $p_{x_1, \dots, x_n; \theta} = p(x_1, \dots, x_n | \theta)$ be the joint probability density of X_1, \dots, X_n . Given $X_1 = x_1, \dots, X_n = x_n$, the likelihood function is defined as

(we start with some model which gives the dependence in θ)

$$L : \Theta \longrightarrow \mathbb{R}$$

$$\theta \longmapsto L(\theta | x_1, \dots, x_n) := p(x_1, \dots, x_n | \theta)$$

If, for a sample x_1, \dots, x_n , $L(\theta_1 | x_1, \dots, x_n) > L(\theta_2 | x_1, \dots, x_n)$, then x_1, \dots, x_n is more likely if $\theta = \theta_1$, than if $\theta = \theta_2$

Likelihood Principle: If x_1, \dots, x_n and y_1, \dots, y_m are two sample points such that $L(\theta | x_1, \dots, x_n) \propto L(\theta | y_1, \dots, y_m) \forall \theta \in \Theta$, then $\{x_1, \dots, x_n\} = \{y_1, \dots, y_m\}$

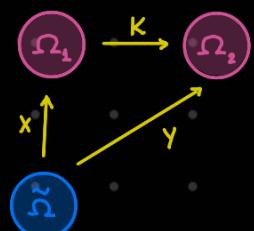
Infinitely divisible random variable: A random variable X is infinitely divisible if, $\forall n \in \mathbb{N}$, $\exists X_1, \dots, X_n$ iid random variables such that $X = \sum_{i=1}^n X_i$

Sufficient statistics: Let Ω_1, Ω_2 be finite non-empty sets and let $K: \Omega_1 \rightarrow \Omega_2$ a statistic which transforms the random variable X to $Y = K(X)$. Then, given a density $p(x|\theta)$ for X , it induces a density for Y : $q(y|\theta) = \sum_{x \in K^{-1}(y)} p(x|\theta)$

If the quantity below doesn't depend on θ , we say K is a sufficient statistics for X

$$r(x|\theta) = \frac{p(x|\theta)}{q(K(x)|\theta)}$$

In this case: $p(x|\theta) = r(x) q(K(x)|\theta)$ and K contains all the dependence of p in θ and it preserves all the information about θ .



Minimal sufficient statistics: A sufficient statistical K for X is minimal if any other sufficient statistics for X is a function of K , i.e., if \bar{K} is a sufficient statistics for X , then $\exists f: \Omega_2 \rightarrow \Omega_2$ measurable function such that $\bar{K} = f \circ K$ a.e. with respect to the underlined measure

$$K \text{ is minimal} \iff \left(\frac{p(x, \theta)}{p(y, \theta)} \text{ is independent of } \theta \iff K(x) = K(y) \quad \forall x, y \in \Omega_s \right)$$

Neyman-Fisher factorization criterion: Let $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ be a family of probability measures on a measurable space (Ω, \mathcal{F}) absolutely continuous with respect to a σ -finite measure μ . Let T a statistic between $(\Omega, \mathcal{F}, P_\theta)$ and another measurable space $(\tilde{\Omega}, \tilde{\mathcal{F}})$

$$\begin{aligned} T \text{ is sufficient with respect to } \mathcal{P} &\iff \forall \theta \in \Theta, \forall x \in \Omega, \exists g_\theta : \tilde{\Omega} \rightarrow \mathbb{R}_{>0} \\ &\quad \forall x \in \Omega, \exists h : \Omega \rightarrow \mathbb{R}_{>0}, h(x) \neq 0 \text{ almost everywhere} \\ &\quad (\text{both measurable function and in } L(\Omega, \mu)) \\ &\quad \text{such that} \\ &\quad \frac{dP_\theta}{d\mu}(x) = p(x | \theta) = g_\theta(T(x)) \cdot h(x) \text{ almost everywhere} \end{aligned}$$

Random matrix: Let $\{X_i : \Omega \rightarrow \mathbb{R}^n : 1 \leq i \leq p\}$ be a family of random vectors. A random matrix of dimension $n \times p$ is the matrix

$$M := [X_1 \ X_2 \ \dots \ X_p] = \begin{bmatrix} X_1^1 & X_2^1 & \dots & X_p^1 \\ X_1^2 & X_2^2 & \dots & X_p^2 \\ \vdots & \vdots & \ddots & \vdots \\ X_1^n & X_2^n & \dots & X_p^n \end{bmatrix}$$

Multivariate normal distribution: Let X be a n -dimensional random vector $X = (X_1, \dots, X_n)$. The n -variate normal distribution is the probability measure $\mathcal{N}_n(\mu, \Sigma)$ such that the probability density is given by

$$\rho_X(\underbrace{x_1, \dots, x_n}_{\bar{x}}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left(-\frac{1}{2} (\bar{x} - \mu)^\top \Sigma^{-1} (\bar{x} - \mu) \right),$$

omitted the parameters. Precisely, it should be $\rho(x_1, \dots, x_n | \mu, \Sigma)$

Notice it's well defined because $\mathcal{N}_n(\mu, \Sigma)(X^{-1}(A)) = \int_A \rho_x(x) dx \quad \forall A \in \mathcal{B}_{\mathbb{R}^n}$

If X follows $\mathcal{N}_n(\mu, \Sigma)$, so holds

$$\mu = \mathbb{E}[X]$$

$$\Sigma_{i,j} = \mathbb{E}[(x_i - \mu_i)(x_j - \mu_j)] = \text{Cov}(x_i, x_j)$$

The inverse of Σ is called precision matrix : $Q = \Sigma^{-1}$

Theorem (Characterization of normal random vector): Let $X = (x_1, \dots, x_n)$ a random vector. Then,

$$X \sim \mathcal{N}_n(\mu, \Sigma)$$



$\exists Z = (Z_1, \dots, Z_K)$ iid random vector

$\exists \mu = (\mu_1, \dots, \mu_n) \in \mathbb{R}^n$

$\exists A \in M_{n \times K}(\mathbb{R})$

such that

$$X = AZ + \mu, \quad Z_i \sim \mathcal{N}(0, I) \quad \forall i \in K, \quad AA^T = \Sigma$$

$\hookrightarrow X$ has the standard probability measure
up to an affine map

χ^2 -distribution: Let X_1, \dots, X_n be random variables following the standard normal distribution, ie,

$$\mathcal{N}(0, I)(X_i^{-1}(A)) = \mu(A) \quad \forall A \in \mathcal{B}_{\mathbb{R}}, \mu \text{ Lebesgue measure in } \mathbb{R}$$

Then, consider $Q := \sum_{i=1}^n X_i^2$. This new random variable $Q: \Omega \rightarrow \mathbb{R}$ will follow a new distribution, called χ^2 with n degrees of freedom, ie,

$$\chi_n^2(Q^{-1}(A)) = \mu(A) \quad \forall A \in \mathcal{B}_{\mathbb{R}}$$

Its probability density is

$$p_{\Omega}(x|n) = \begin{cases} \frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}, & x > 0 \\ 0, & \text{otherwise} \end{cases}$$

$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$
 $\Gamma(z) = (z-1)!$
 $\Gamma(z+1) = z \Gamma(z)$
 $\Gamma(\frac{1}{2}) = \sqrt{\pi}$

Sample following a distribution: A sample (x_1, x_2, \dots, x_n) which follows a distribution \mathcal{D} is just a point in the image set of some random vector X which follows \mathcal{D} . In other, $\exists X: \Omega \rightarrow \mathbb{R}^n$ random vector and $\omega \in \Omega$ such that $X(\omega) = (x_1, \dots, x_n)$ and such that $\mu = X_* \mathcal{D}$, where μ is the Lebesgue measure in \mathbb{R}^n :

$$\mathcal{D}(X^{-1}(A)) = \mu(A)$$

Generating according to one distribution: Generate a random variable following a probability distribution \mathcal{D} is the process of seek a random variable $X: \Omega \rightarrow \mathbb{R}$ such that $\mu = X_* \mathcal{D}$. It's the same for random vectors. Then, generate a sample of points following \mathcal{D} , formally, is the process of take a point in the image of a random variable X such that $\mu = X_* \mathcal{D}$

Left/Right invariant measures: Let (G, \cdot) be a locally compact Hausdorff topological group ($\forall g \in G$, $\exists K$ compact set and U open set such that $x \in U \subset K$). Let $g \in G$ and $S \subset G$. Then, define the translations

- Left: $gS = \{gs : s \in S\}$
 - Right: $Sg = \{sg : s \in S\}$
- ↗ It preserves Borel sets. As the group operation is continuous, it preserves the topology and, consequently, the σ -algebra \mathcal{B}



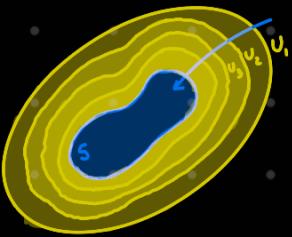
A measure μ on (G, \mathcal{B}) is

- Left-translation-invariant if: $\mu(gS) = \mu(S) \quad \forall g \in G, \forall S \in \mathcal{B}$ (Notation: LTI)
- Right-translation-invariant if: $\mu(Sg) = \mu(S) \quad \forall g \in G, \forall S \in \mathcal{B}$ (Notation: RTI)

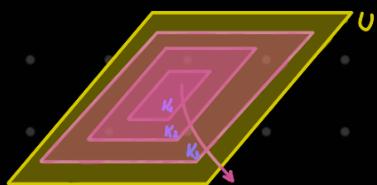
Haar's theorem: Up to a multiplicative constant, \exists ! nontrivial measure on \mathcal{B} of G such that

- μ is LTI: $\mu(gS) = \mu(S) \quad \forall g \in G, \forall S \in \mathcal{B}$
- μ is finite on compacts: $\mu(K) < \infty \quad \forall K \subset G$ compact
- μ is outer regular on Borel sets: $\mu(S) = \inf\{\mu(U): S \subset U, U \in \mathcal{T}_G\} \quad \forall S \in \mathcal{B}$
- μ is inner regular on open sets: $\mu(U) = \sup\{\mu(K): K \subset U, K \text{ compact}\} \quad \forall U \in \mathcal{T}_G$

Such a measure is called the left Haar measure on G

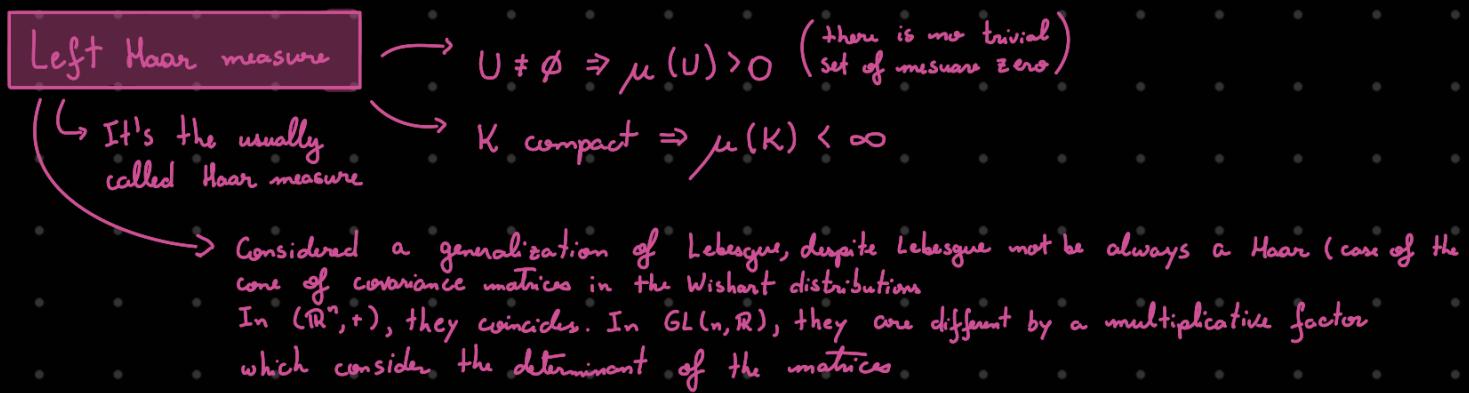


Outer regular



inner regular

This measure is unique with the normalization condition: $\mu(G) = 1$



Similarly, $\exists!$ right Haar measure ν up to a positive cte

Right Haar measure: $\nu(S) = a \mu(S^{-1})$, $a > 0, S^{-1} = \{g \in G : g^{-1}g = id \quad \forall g \in S\} \quad \forall S \in \mathcal{B}$

Modular function: $\exists \Delta: G \longrightarrow (0, +\infty)$, $g \mapsto \Delta(g)$ such that, for each $S \in \mathcal{B}$ fixed,

$$\nu(g^{-1}S) = \Delta(g) \nu(S)$$

↳ Modular function

Then, if $\Delta(g) = 1$, ν is right and left invariant and, consequently, $\nu = \mu$

The groups which $\Delta(g)$ are called unimodular groups

Radon measure: Let (Ω, \mathcal{B}) a measurable space with the Borel σ -algebra, where Ω is Hausdorff. A Radon measure in (Ω, \mathcal{B}) is inner regular and locally finite measure.

Every Haar measure is a Radon measure

Dirac measures are Radon

$\text{In } (\mathbb{R}^n, +) : \text{Lebesgue} \equiv \text{Haar} \equiv \text{Radon}$

	Lebesgue	Haar	Radon
Domínio de definição	\mathbb{R}^n	(G, \cdot) Locally compact Hausdorff topological groups	(Ω, τ) Locally compact Hausdorff topological spaces
Invariance	Translation	Group operation by left	

The cone $M_+(\Omega)$ of finite measures over a (finite) set Ω is a set of Radon measures, specially for the subset $P_+(\Omega)$ of probability measures

Polish space: A Polish space is a topological space (Ω, τ) which is

- Separable: It contains a countable dense subset
- Completely metrizable: It's homeomorphic to a complete metric space

Probability measures on the Borel σ -algebra of a Polish space is a Radon measure

Raw moment: Let $X: \Omega \rightarrow \mathbb{R}$ be a random variable with density function (with respect to the underlying σ -algebras and measures) ρ_x . Then, the n -th raw moment is given by

$$\langle X^n \rangle := \int_{\mathbb{R}} x^n \rho_x(x) d\mu(x)$$

Central moment: Let $X: \Omega \rightarrow \mathbb{R}$ be a random variable with density function (with respect to the underlying σ -algebras and measures) ρ_x . Given $c \in \mathbb{R}$, consider the integral

$$\langle X_c^n \rangle = \int_{\mathbb{R}} (x - c)^n \rho_x(x) d\mu(x)$$

When $c = \langle x \rangle =: \mathbb{E}_x[x]$, $\langle x_c^n \rangle$ is called the n -th central moment of X
 ↳ We introduce this notation to distinguish when we are considering the density p_x of X vs when we are just computing $\int_{\Omega} x \times dx$

Standardized moment: Let $X: \Omega \rightarrow \mathbb{R}$ be a random variable with the underlying density function p_x .
 The n -th standardized moment of X is the n -th central moment divided by the n -th standard deviation

$$\langle x_s^n \rangle = \frac{\mathbb{E}_x[(x - \mathbb{E}_x[x])^n]}{\sqrt{\mathbb{E}_x[(x - \mathbb{E}_x[x])^2]}}$$

#	Raw	Central	Standardized
1	Mean	0	0
2		Variance	1
3			Skewness
4			Kurtosis
5			Hyperskewness
6			Hypertailedness



Mixed moments: Let X_1, \dots, X_n be random variables with the underlying joint density function p_x .

For any $K_i \in \mathbb{N}$, we call

- $\mathbb{E}_x[X_1^{K_1} \dots X_n^{K_n}]$ the mixed moment of order $K = \sum_{i=1}^n K_i$
- $\mathbb{E}_x[(X_1 - \mathbb{E}_x[X_1])^{K_1} \dots (X_n - \mathbb{E}_x[X_n])^{K_n}]$ the central mixed moment of order $K = \sum_{i=1}^n K_i$

If $K=2$, there is only one central mixed moment of each 2 random variables, say X_1, X_2 ,

$$\text{Cov}(X_1, X_2) = \mathbb{E}_x[(X_1 - \mathbb{E}_x[X_1])(X_2 - \mathbb{E}_x[X_2])]$$

So, the covariance of X_1 and X_2 is the second central mixed moment of them (or, equivalently, the second central moment of the random vector $\vec{X} = (X_1, X_2)$)

↳ Obviously, $\text{Cov}(X_1, X_2)$ is unique and depends on the measures in the domain and codomain of \vec{X}

If X_1 and X_2 are random vectors of same size, $\text{Cov}(\vec{X}_1, \vec{X}_2)$ is the matrix of covariance

If $K=3$, we have the freedom to choose

- 2 random variables X_1 and X_2

and, consequently, we have

$$\mathbb{E}_{\vec{x}}[(X_1 - \mathbb{E}_{\vec{x}}[X_1])(X_2 - \mathbb{E}_{\vec{x}}[X_2])^2]$$

$$\mathbb{E}_{\vec{x}}[(X_1 - \mathbb{E}_{\vec{x}}[X_1])^2(X_2 - \mathbb{E}_{\vec{x}}[X_2])]$$

- 3 random variables X_1, X_2 and X_3

and, consequently, we have

$$\mathbb{E}_{\vec{x}}[(X_1 - \mathbb{E}_{\vec{x}}[X_1])(X_2 - \mathbb{E}_{\vec{x}}[X_2])(X_3 - \mathbb{E}_{\vec{x}}[X_3])]$$

Consequently, we can define the coskewness and cokurtosis in the same way as the covariance, but it's not unique for $K > 1$.

Cumulant: Let $X: \Omega \rightarrow \mathbb{R}$ be a random variable. We define the moment generating function as

$$M_x: \mathbb{R} \rightarrow \mathbb{R}, \quad t \mapsto \mathbb{E}_{\vec{x}}[e^{tx}],$$

when this expectation is finite. This name is justified because

$$M_x(t) = 1 + t \mathbb{E}_{\vec{x}}[X] + \frac{t^2}{2!} \mathbb{E}_{\vec{x}}[X^2] + \frac{t^3}{3!} \mathbb{E}_{\vec{x}}[X^3] + \dots$$

$\int x \rho_{\vec{x}}(dx) d\mu(x)$

If \vec{X} is a random vector, then $M_{\vec{X}}$ is defined as $M_{\vec{X}}(\vec{t}) = \mathbb{E}[e^{\langle \vec{t}, \vec{X} \rangle}]$, $t \in \mathbb{R}^n$

When M_x is defined, we call the cumulant-generating the logarithm of that one

$$K: \mathbb{R} \rightarrow \mathbb{R}, \quad t \mapsto \ln \mathbb{E}[e^{tx}],$$

The n -th cumulant of X is the n -th coefficient of the MacLaurin expansion of K

$$K(t) = \sum_{n \in \mathbb{N}} \frac{t^n}{n!} K_n, \quad K_n = \frac{d^{(n)} K(0)}{dt^{(n)}}$$

$$K_n(x+c) = \begin{cases} K_n(x) + c, & \text{if } n=1 \\ K_n(x), & \text{if } n > 1 \end{cases} \quad (\text{invariant by translation})$$

$$K_n(cx) = c^n K_n(x) \quad (n\text{-homogeneous})$$

$$K_n(x_1 + \dots + x_k) = \sum_{i=1}^k K_n(x_i) \quad (\text{linear on the sum})$$

$$\left. \begin{array}{l} \bullet K_1(x) = \langle x \rangle \\ \bullet K_2(x) = \text{Var}(x) \end{array} \right\}$$

↑
when $\{p_\theta\}$ is parametric, we write $\text{Var}_\theta(x)$

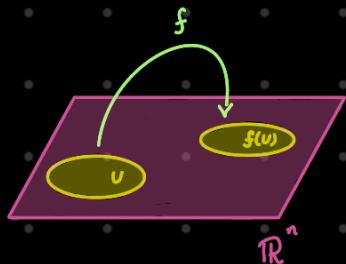
Important: The only probability distribution such that all the cumulants of order > 3 of a random variable which follows this distribution is the normal distribution (trivial δ -distribution if $K_2=0$)

In general, it's not unique. Then, usually, people write K_{r_1, \dots, r_k}

For random vector, we will have the cumulants obtained by the coefficients of the MacLaurin series for multiple variables. Then, if $M_{\vec{x}}(\vec{t})$ is the moment-generating function and it's finite, then a r -th cumulant of \vec{x} is

$$K_r(\vec{x}) = \left. \frac{\partial^{r_1}}{\partial t_1^{r_1}} \dots \frac{\partial^{r_n}}{\partial t_n^{r_n}} K(t) \right|_{t=0}; \quad r = \sum_{i=1}^n r_i$$

Invariance of domain: If U is an open set of \mathbb{R}^n and $f: U \rightarrow \mathbb{R}^m$ is a continuous injective map, then $f(U)$ is open in \mathbb{R}^m and f is homeomorphism between U and $f(U)$



↳ Continuous injective functions which preserve the dimension of the space are open

Complete statistics: A statistic T is complete for a family $\mathcal{P} = \{P_\theta : \theta \in \Xi\}$ of probability measures if, \forall measurable function f ,

$$\mathbb{E}_\theta [f \circ T] = c \quad \forall \theta \in \Xi \Rightarrow f \circ T = c \text{ a.e. with respect to all } P_\theta \in \mathcal{P}$$

If the expectation of $f \circ T$ is constant for all θ , then $f \circ T$ need to be equal to this constant (a.e. with respect to the distributions of \mathcal{P})

Note that, for each f , we can have another constant c ; c is constant in Ξ .

There is non-constant invisible functions of T for \mathcal{P} , i.e., with \mathbb{E}_θ null for all the parameters

Theorem: If T is complete and sufficient, then T is minimal sufficient

	Sufficient	Minimal	Complete
Concept	$T(X)$ contains all the needed information about the parameter present in the "data of X "	It's the "lowest" statistics which contains all the needed information about parameter present in the "data of X "	Guarantees there is no functions non-trivial (cte) of $T(X)$ with constant expectation for all θ
Role	Reduction of data without loss of information about the parameter	Best reduction of data without loss of information about the parameter	Provides the uniqueness of unbiased estimator of θ based on $T(X)$ with minimum variance

we can define it on a manifold, but we need to care about the support of the functions

Estimator: Let (Ω, \mathcal{F}) and $(\Xi \subset \mathbb{R}^n, \mathcal{G})$ be measurable spaces. An estimator is a measurable function $\hat{\theta}: \Omega \rightarrow \Xi$. Ω is called sample space and Ξ is called sample estimates space or sometimes just parameter space.

Given a random variable $X: \Omega_1 \rightarrow \Omega_2$, an estimator based on X is an estimator $\hat{\theta}: X(\Omega_1) \rightarrow \Xi$

We define the bias of $\hat{\theta}$ as the function

$$B: \Xi \rightarrow \mathbb{R}, \quad \theta \mapsto B(\theta) = E_\theta[\hat{\theta}(X)] - \theta$$

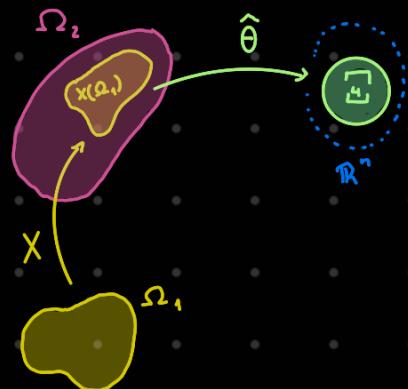
An estimator $\hat{\theta}$ is called unbiased for some $\theta \in \Xi$ if $B(\theta) = 0$, i.e., if $E_\theta[\hat{\theta}] = \theta$. In this case, usually is called "the true value of θ ".

The error of $\hat{\theta}$ for some $\theta \in \Xi$ is defined as the function

$$\epsilon_\theta: \Omega_1 \rightarrow \mathbb{R}, \quad x \mapsto \epsilon_\theta(x) = \hat{\theta}(x) - \theta$$

Consequently, the mean squared error for some $\theta \in \Xi$ is defined as

$$\bar{E}_\theta^2(\hat{\theta}) = E_\theta[\epsilon_\theta(x)^2] = E_\theta[(\hat{\theta}(x) - \theta)^2]$$



Lehmann-Scheffé' Theorem: Let $X: \Omega_1 \rightarrow \Omega_2$ be a random variable and let $T: \Omega_2 \rightarrow \Omega_3$ be a complete and sufficient statistic for a family $\mathcal{P} = \{P_\theta : \theta \in \Xi\}$ of probability measures in $(\Omega_2, \mathcal{F}_2)$. Let $\hat{\theta}(X)$ be an unbiased estimator of $\theta \in \Xi$ such that $\hat{\theta} = f \circ T(X)$, for some measurable function $f: \Omega_3 \rightarrow \Xi$. Then, $\hat{\theta}$ is the unique (up to sets of measure zero) minimum variance unbiased estimator of θ , ie, $\text{Var}_\theta(\hat{\theta}(X)) \leq \text{Var}_\theta(\hat{\theta}'(X))$, for all unbiased estimator $\hat{\theta}'$ of θ .

Rao-Blackwell Theorem: Let $X: \Omega_1 \rightarrow \Omega_2$ be a random variable and let $T: \Omega_2 \rightarrow \Omega_3$ be a sufficient statistic for a family $\mathcal{P} = \{P_\theta : \theta \in \Xi\}$ of probability measures in $(\Omega_2, \mathcal{F}_2)$. Let $\hat{\theta}$ be an unbiased estimator of $\theta \in \Xi$. Then, the unbiased estimator for θ given by

$$\hat{\theta}^*(x): \Omega_2 \rightarrow \Xi, \quad \hat{\theta}^*(x) = \mathbb{E}[\hat{\theta}(x) | T(x)]$$

is such that

$$\text{Var}_\theta(\hat{\theta}(x)) \geq \text{Var}_\theta(\hat{\theta}^*(x))$$

Statistical model: Let (Ω, \mathcal{F}) be a measurable space. A statistical model for Ω is a subset $\mathcal{P}' \subset \mathcal{P}$ of the set \mathcal{P} of probability measures in (Ω, \mathcal{F}) . Moreover, let Ξ a non-empty set. A parametric statistical model is a set $\mathcal{P}' = \{P_\theta : \theta \in \Xi\} \subset \mathcal{P}$ indexed by Ξ . In this case, Ξ is called a space of parameters for \mathcal{P}' . When another measure μ in (Ω, \mathcal{F}) is underlying and the Radon-Nikodym derivative of $dP/d\mu$ of P and μ exists for all $P \in \mathcal{P}'$, the set of this functions may be called a statistical model for Ω in some literatures.

Regular statistical model: Let $\mathcal{P}' = \{P_\theta : \theta \in \Xi\}$ be a parametric statistical model on (Ω, \mathcal{F}) .

It's called regular if all the conditions below hold:

- Global differentiation: $\Xi \subset \mathbb{R}^n$, $\Xi \in \mathcal{T}_{\mathbb{R}^n}$

- Differentiation of the Likelihood function: The function $\theta \mapsto L(\theta; \cdot)$ is $C^1(\Xi)$ (or sufficiently differentiable) μ -almost everywhere in Ω

