# 1

PROBABILITY SPACES

## THE IDEA

We begin prob with the vague notion of a random experiment. A repeatable process that does not have a determinate outcome.

We begin to abstract this by formalising the:
- **Sample Space:** $\Omega$ set of all outcomes
- **Outcomes:** An element $\omega \in \Omega$
- **Events:** Subsets $A \subset \Omega$ such that we can define probability on it.

But note that in the general case we cannot use ARBITRARY subsets as events.

## SETS

We denote the union of disjoint sets $D \cup E = D + E$

The symmetric difference of a set is: $A \Delta B = A \backslash B + B \backslash A$
$$= (A \cup B) \backslash (A \cap B)$$

**INDICATORS:** For an arbitrary set $A$ the indicator function
$$\mathbb{1}_A(\omega) = \begin{cases} 1, & \omega \in A \\ 0, & \omega \notin A \end{cases}$$

For two sets $A$ & $B$: $\mathbb{1}_{A^c} = 1 - \mathbb{1}_A$, $\mathbb{1}_{A \cup B} = \max\{\mathbb{1}_A, \mathbb{1}_B\}$
$$\mathbb{1}_{A \cap B} = \mathbb{1}_A \mathbb{1}_B, \quad \mathbb{1}_{A \Delta B} = |\mathbb{1}_A - \mathbb{1}_B|$$

## EVENTS & ALGEBRAS

We need that our events are closed under certain operations so that when we manipulate events we still have events.

A family $\mathcal{F}$ of subsets of $\Omega$ is a $\sigma$-Algebra if it satisfies

A.1) $\Omega \in \mathcal{F}$
A.2) $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$
A.3) $A_1, A_2, \ldots \in \mathcal{F} \Rightarrow \bigcup_{n \geq 1} A_n \in \mathcal{F}$

Note that from De'Morgans Law this also implies closed under countable intersections.

Thus we will now call the elements of an appropriate $\sigma$-Algebra generated on subsets of $\Omega$ **Events**

**T:** For two $\sigma$-Algebras $\mathcal{F}_1$ & $\mathcal{F}_2$ on a common sample space $\Omega$, then $\mathcal{F}_1 \cap \mathcal{F}_2$ is also a $\sigma$-Algebra.

**T:** $\mathcal{F}_n$, $n \in \mathbb{N}$, $\sigma$-Algebras on a common sample space $\Rightarrow \bigcap_{n \geq 1} \mathcal{F}_n$ is a $\sigma$-Algebra.

Given some $\Omega$ how do we create a $\sigma$-Algebra?

**$\sigma$ ALGEBRA GENERATED:**
① For a single $A \subset \Omega$. $\sigma(A) = \{\emptyset, A, A^c, \Omega\}$
② For $\mathcal{G} = \{A_1, \ldots, A_n\}$ a finite partition of $\Omega$
$$\sigma(\mathcal{G}) = \left\{ \sum_{i \in I} A_i : I \subset \{1, \ldots, n\} \right\}$$
For an arbitrary collection of sets, consider that we can form a partition by taking all possible intersections.

**T:** For any family of subsets of $\Omega$, $\mathcal{G}$, $\exists$ a unique $\sigma$-Algebra, $\sigma(\mathcal{G})$ s.t. $\mathcal{G} \subset \sigma(\mathcal{G})$
& $[\mathcal{H}$ a $\sigma$-Algebra on $\Omega$ & $\mathcal{G} \subset \mathcal{H} \Rightarrow \sigma(\mathcal{G}) \subset \mathcal{H}]$

## BOREL SET:
This is the cannonical $\sigma$-Alg generated used on $\mathbb{R}$.
$$\mathcal{B}(\mathbb{R}) = \sigma\{(a, b] \mid a, b \in \mathbb{R}, a < b\}$$
$$\mathcal{B}(\mathbb{R}^m) = \sigma\left\{ \prod_{i=1}^{m} (a_i, b_i] \mid a_i, b_i \in \mathbb{R}, a_i < b_i \right\}$$

$A_n$ occur i.o. slide 11.

## PROBABILITY SPACE

A set, $\Omega$, paired with a $\sigma$-Alg generated on its subsets, $\mathcal{F}$, is called a **measurable space** $(\Omega, \mathcal{F})$

A **probability** on $(\Omega, \mathcal{F})$ is a function $P: \mathcal{F} \longrightarrow \mathbb{R}$ satisfying
- P.1) $P(A) \geq 0$, $A \in \mathcal{F}$
- P.2) $P(\Omega) = 1$
- P.3) For any pairwise disjoint $A_1, A_2, \ldots \in \mathcal{F}$ $P(\bigcup_{j \geq 1} A_j) = \sum_{j \geq 1} P(A_j)$

The tripple $(\Omega, \mathcal{F}, \mathbb{P})$ is called a **prob space**.

**EX.** Degenerate Dist at a fixed $\omega \in \Omega$
$$\varepsilon_\omega(A) = \mathbb{1}_A(\omega)$$

**E.** Counting measure on $\mathbb{N}$
$$M(B) \equiv \sum_{n \geq 1} \varepsilon_n(B), \quad B \in P(\mathbb{N})$$

From these axioms we can further deduce these properties of the probability measure:

**T:** $P(\emptyset) = 0$, **T:** $P(\bigcup_{n=1}^{m} A_j) = \sum_{n=1}^{m} P(A_j)$

**T:** $P(A^c) = 1 - P(A)$ For pairwise disjoint $A_1, \ldots, A_m$

**T:** $A \subset B \Rightarrow P(B \backslash A) = P(B) - P(A) \Rightarrow P(A) \leq P(B)$

**T:** $P(A \cup B) = P(A) + P(B) - P(B \cap A)$

**T: Booles Ineq:** $P(\bigcup_{j \geq 1} A_j) \leq \sum_{j \geq 1} P(A_j)$

**T: Borel-Cantelli:** $\sum_{n \geq 1} P(A_n) < \infty \Rightarrow P(A_n \text{ i.o.}) = 0$

CONTINUITY PROPERTIES:

The infinite case of the countable additivity property of our probability is responsible for important continuity properties of the $\mathbb{P}$.

$A_n \uparrow A$ $\iff$ $A_1 \subset A_2 \subset \cdots$ & $\bigcup_{n \geq 1} A_n = A$

$A_n \downarrow A$ $\iff$ $A_1 \supset A_2 \supset \cdots$ & $\bigcap_{n \geq 1} A_n = A$

**T:** A function $P: \mathcal{F} \longrightarrow \mathbb{R}$ satisfies P.1, P.2 & has finite additivity. THEN (the following are $\iff$) $\mathbb{P}$ has property P.3
$$\iff [A_n \uparrow A \Rightarrow P(A_n) \uparrow P(A)] \iff [A_n \downarrow \emptyset \Rightarrow P(A_n) \downarrow 0]$$
$$\iff [A_n \downarrow A \Rightarrow P(A_n) \downarrow P(A)]$$

# 2

**PROBABILITIES ON R.**

Probabilities on $\mathbb{R}$ are defined on the measurable space $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$.

BUT $\mathcal{B}(\mathbb{R})$ is HUGE! How do we specify the probability on all events in $\mathcal{B}(\mathbb{R})$ for a $P$?

## DISTRIBUTION FUNCTION:
It turns out we can entirely specify a $P$ by its distribution function.

**Distribution Function** (D.F) of a probability $P$ on $\mathbb{R}$ is the function $F_P : \mathbb{R} \to \mathbb{R}$, $t \mapsto P(-\infty, t]$

**T:** For any probability $P$ on $\mathbb{R}$ its D.F. $F$ satisfies all of...
- D.1) Non Decreasing: $s < t \Rightarrow F(s) \leq F(t)$
- D.2) Right continuous: $F(t) = F(t+)$
- D.3) $\lim_{t \to -\infty} F(t) = 0$ & $\lim_{t \to \infty} F(t) = 1$

**T:** $F_P = F_{P'} \iff P = P'$  (The probability & its distribution function entirely determine one another)

**T:** For any $F: \mathbb{R} \to \mathbb{R}$ satisfying D.1-3 there corresponds EXACTLY ONE Probability $P$ on $\mathcal{B}(\mathbb{R})$.

## CLASSIFYING P on R

$P$ is **Discrete** on $\mathbb{R}$
$\iff (\exists C \subset \mathbb{R})(C$ countable & $P(C) = 1)$

**T:** $\iff \exists \{t_i\}_{i \geq 1} \subset \mathbb{R}$ & $\exists \{p_i > 0\}_{i \geq 1}$ such that
$$\sum_i p_i = 1 \quad \& \quad P = \sum_i p_i \, \varepsilon_{t_i}$$
$\iff \exists \{t_i\}_{i \geq 1} \subset \mathbb{R}$ & $\exists \{p_i > 0\}_{i \geq 1}$ such that
$$\sum_i p_i = 1 \quad \& \quad F_P(t) = \sum_i p_i \, \mathbb{1}(t_i \leq t)$$

$P$ is **absolutely continuous** (A.C) on $\mathbb{R}$ if
$\exists f_P : \mathbb{R} \to \mathbb{R}$ such that $F_P(t) = \int_{-\infty}^{t} f_P(x) \, dx$.

So A.C. distributions are the ones with densities. Note that $f_P = F_P'$ almost everywhere. We must be careful about points of discontinuity.

**T:** Any function $f: \mathbb{R} \to \mathbb{R}$ that is
① $f \geq 0$ ② Integrable ③ $\int f \, dx = 1$
specifies a probability on $\mathbb{R}$.

A **mixed distribution** $P$ is one such that for some $p \in (0,1)$ $P = p P_d + (1-p) P_a$ where $P_d$ is discrete & $P_a$ is A.C.

A **singular distributions** are continuous but not A.C., they Dont have a density however no single point has a positive probability.

**T:** Any probability on $\mathbb{R}$ has a unique representation of the form
$$P = \alpha_d P_d + \alpha_a P_a + \alpha_s P_s \,, \quad \alpha_i \geq 0, \sum_i \alpha_i = 1$$
$\quad \quad \hookrightarrow$ Discrete $\quad \searrow$ A.C $\quad \searrow$ Singular.

# 3

**RANDOM VARIABLES**

We naively think of a R.V. as a function of an outcome of a Random experiment that captures some information about the experiment. We also however want to be able to calculate the probability of our R.V. mapping to a certain value (of set of values)

Thus we define a **Random Variable** (R.V) as a function $X : \Omega \longrightarrow \mathbb{R}$ such that $\forall B \in \mathcal{B}(\mathbb{R}) \quad X^{-1}(B) \in \mathcal{F}$.

$$X^{-1}(B) = \{ \omega \in \Omega \mid X(\omega) \in B \}$$

**T:** For an arbitrary family of subsets of $\mathbb{R}$ $\{B_\alpha \mid \alpha \in I\}$
- $B_\alpha \subset B_\beta \implies X^{-1}(B_\alpha) \subset X^{-1}(B_\beta)$
- $\bigcup_{\alpha \in I} X^{-1}(B_\alpha) = X^{-1}\left(\bigcup_{\alpha \in I} B_\alpha\right)$  (similarly for $\cap$)
- $B_\alpha \cap B_\beta = \emptyset \implies X^{-1}(B_\alpha) \cap X^{-1}(B_\beta)$
- $X^{-1}(B_\alpha^c) = \left[ X^{-1}(B_\alpha) \right]^c$

Ex. Random indicator: For any event $A$, $\mathbb{1}_A$ is a R.V.

Simple R.V.: $\sum_{i=1}^{n} a_i \mathbb{1}_{A_i}$ , $a_i \in \mathbb{R}$, $A_i \in \mathcal{F}$, $i \leq n < \infty$

**T:** For a R.V $X$ , $\sigma(X) = \{ X^{-1}(B) \mid B \in \mathcal{B}(\mathbb{R}) \}$ is a $\sigma$-Alg.

## DISTRIBUTIONS OF R.V.

The **distribution** of a R.V $X$ on $(\Omega, \mathcal{F}, \mathbb{P})$ is defined as $P_X(B) = P(X \in B)$, $P_X : \mathcal{B}(\mathbb{R}) \longrightarrow \mathbb{R}$. This $P_X$ is in fact a probability on $\mathbb{R}$.

Thus the **distribution function** of $X$ is
$$F_X(t) = P_X((-\infty, t]) = \mathbb{P}(X \leq t)$$

$X$ is **discrete/A.C/singular** if $F_X$ is D/A.C/S. (same for R.vec)
The **survival function** of $X$ is $S_X(t) = 1 - F_X(t)$
Two R.V. $X$ & $Y$ are **identically distributed** $\iff P_X = P_Y$

## FUNCTIONS OF R.V.

**T:** $X$ a R.V, $g$ increasing $(g'' \geq 0)$ & continuous on $\mathbb{R}$ $\implies Y = g(X)$ is a R.V with D.F $F_Y(t) = F_X(g^{-1}(t))$

**T:** $X$ an A.C R.V, $g$ continuously diff on open set $U$ such that $P(X \in U) = 1$. $\implies Y = g(X)$ is A.C R.V. $f_Y(t) = |\frac{d}{dt} g^{-1}(t)| f_X(g^{-1}(t))$

Note that we must also have that $g$ is invertable & that inverse is differentiable.

For a D.F $F$ the **quantile function** is
$$Q(\alpha) = \inf\{ t : F(t) \geq \alpha \}, \quad \alpha \in [0,1)$$

**T:** $U \sim U[0,1] \implies X = Q(U) \sim F$

## RANDOM VECTORS

A **Random Vector** (R.vec) $X = (x_1, ..., X_n) : \Omega \longrightarrow \mathbb{R}^n$ is a function such that $X^{-1}(B) \in \mathcal{F} \quad \forall B \in \mathcal{B}(\mathbb{R})$

**T:** $X = (X_1, ..., X_n)$ R.vec $\iff (\forall i) X_i$ is a R.V)

**T:** $X = (X_1, ..., X_n)$ a R.vec & $g$ measurable $g : \mathbb{R}^n \to \mathbb{R}^n$ $\implies g(X)$ is a R.vec
$g$ measurable $\iff \forall B \in \mathcal{B}(\mathbb{R})$ $g^{-1}(B) \in \mathcal{B}(\mathbb{R})$

We have the **D.F. in the multivariate** case as
$$F_X(t_1, ..., t_n) = \mathbb{P}(X_1 \leq t_1, ..., X_d \leq t_d), \quad (t_1, ..., t_n) \in \mathbb{R}^n.$$

**T:** An A.C dist has a density $f_X$ satisfying
$$F_X(t_1, ..., t_n) = \int_{-\infty}^{t_1} \cdots \int_{-\infty}^{t_n} f_X(x_1, ..., x_n) dx_1 \cdots dx_n$$

**T:** $(X_1, ..., X_n)$ is discrete $\iff (\forall i)(X_i$ is discrete)

**T:** $(X_1, ..., X_n)$ is A.C. $\implies (\forall i)(X_i$ is AC) & 
Integrate out all but the desired variable.
$$f_{X_i}(x) = \int \cdots \int f_X(s_1, ..., x, ..., s_n) ds_1 \cdots ds_n$$

**T:** $X$ a R.vec, $g : \mathbb{R}^n \to \mathbb{R}^n$ has a smooth inverse $X$ A.C $\implies Y = g(X)$ is A.C R.V such that
$$f_Y(t) = |\det [g^{-1}(t)]| f_X(g^{-1}(t))$$
$\det(h(t)) = J(h(t)) = \det\left[\frac{dh_i}{dt_i}\right]$ is the **Jacobian**.

## INDEPENDENCE

A collection of R.V. $X_1, ..., X_n$ are **independent** if
$$(\forall B_1, ..., B_n \in \mathcal{B}(\mathbb{R}))\left( P(X_1 \in B_1, ..., X_n \in B_n) = \prod_{i=1}^{n} P(Y_i \in B_i) \right)$$

This is a very general definition of independence however it is not convenient as a test. For showing some R.V.s are independent we will often use the following.

**T:** $X_1, ..., X_n$ R.V independent $\iff$ $\forall t_1, ..., t_n \in \mathbb{R}$ $F_{X_1, ..., X_n}(t_1, ..., t_n) = \prod_{i=1}^{n} F_{X_i}(t_i)$

**T:** Discrete $X_1, ..., X_n$ R.V independent $\iff$ $\forall t_1, ..., t_n \in \mathbb{R}$ $P(X_1 = t_1, ..., X_n = t_n) = \prod_{i=1}^{n} P(X_i = t_i)$

**T:** AC R.V. $X_1, ..., X_n$ independent $\iff$ $\forall t_1, ..., t_n \in \mathbb{R}$ $f_{X_1, ..., X_n}(t_1, ..., t_n) = \prod_{i=1}^{n} f_{X_i}(t_i)$

**T:** $g_i$ measurable functions, $X_1, ..., X_n$ independent $\implies Y_j = g_j(X_j)$ are also independent.

**Events** $A_1, ..., A_n$ are **independent**
$\iff \mathbb{1}_{A_1}, ..., \mathbb{1}_{A_n}$ are independent as R.V.

**T:** $\iff (\forall I \subseteq \{1, ..., n\})\left( P(\bigcap_{i \in I} A_i) = \prod_{i \in I} P(A_i) \right)$.

**T:** $\iff A_1^c, ..., A_n^c$ are independent

# 4

**EXPECTATIONS**

## DEFINING EXPECTATIONS

In second year we talked about $E(X)$ for A.C & discrete as simply $\int x f(x)\,dx$. This however is just a computational tool, not an informative definition. It is also not general enough, what about singular or mixed distributions? We also want our def to align with frequentist intuition.

For $\mathbb{1}_A$, $A \in \mathcal{F}$ we define $E(X) = P(A)$    *Indicator*

For $X = \sum_{i=1}^{n} a_i \mathbb{1}_{A_i}$ we define $E(X) = \sum_{i=1}^{n} a_i P(A_i)$    *Simple R.V.*

Now note that any non-negative R.V can be approximated by an increasing sequence of simple R.V. $\{X_n\}_{n \geq 1}$ in the following way

$$\forall_w \in \Omega \qquad X_n(w) \underset{n \to \infty}{\uparrow} X(w) . \not\exists$$

We will use our current definition & this approximation to define

$X \geq 0$ Arbitrary, $E(X) = \lim_{n \to \infty} E(X_n)$   Where $\{X_n\}_{n \geq 1}$ is a sequence of simple R.V as $m \nearrow$

**T:** This definition is consistent. Different sequences will give the same expectation.

Let $X^+ = \max\{X, 0\}$ & $X^- = -\min\{X, 0\}$ & Note that for an arbitrary $X = X^+ - X^-$

A R.V. is **integrable** if $E(|X|) < \infty \iff X \in L^1$
If $X$ is integrable $E(X) = E(X^+) - E(X^-)$
The expectation of a R.V. $X$ over an event $A$ is $E(X; A) = E(X\mathbb{1}_A)$.

## UNDERSTANDING EXPECTATION

**T:** Expectation as a function is
- Monotone : $X \leq Y$ & $E(Y) < \infty \Rightarrow E(X) \leq E(Y)$
- Linear: $\forall a, b \in \mathbb{R}$, $X, Y \in L^1 \Rightarrow E(aX + bY) = aE(X) + bE(Y)$

We can denote $E$ using Lebesgue integrals
$$E(X) = \int_\Omega X(w) \mathbb{P}(dw) = \int_\Omega X(w) \, dP(w) = \int_\Omega X \, dP$$

Moving probability spaces from $(\Omega, \mathcal{F}, \mathbb{P})$ (general) to $(\mathbb{R}^d, B(\mathbb{R}^d), P_X)$ allows us to shift

**T:** $E(X) = \int_\Omega X(w) \, dP(w) = \int_\mathbb{R} x \, dP_X(x)$   when $E(X)$ is defined   *distribution func of $P_X$.*

Notation: $\int g(x) \, dP_X(x)$ often denoted $\int g(x) \, dF_X(x)$

**T:** If $F$ is A.C with density $f = F'$ (a.e) and both $f$ & $g$ are piecewise continuous then
$$\underbrace{\int g(x) dF(x)}_{\text{Lebesgue}} = \underbrace{\int_{-\infty}^{\infty} g(x) f(x) \, dx}_{\text{Riemann}}$$

## USING EXPECTATIONS

**T:** $Y = g(x) \sim P_Y$, $X \sim P_X \Rightarrow E(Y) = \int g(x) \, dP_X(x)$

**T:** $X \geq 0$, $E(X) = \int_0^\infty (1 - F_X(x)) \, dx$

**T:** $Y = g(X)$ for nice $g \Rightarrow E(Y) = \int g(x) \, dF_X(x)$

commonly $E(g(X)) = \int g(x) f_x(x) \, dx$ for A.C.
$\phantom{commonly\ E(g(X))} = \sum_{t_i \in C_X} g(t_i) P(X = t_i)$.

**T:** $X_1$ & $X_2$ independent & $g_i(X_i) = L^1$
$$\Rightarrow E(g_1(X_1) g_2(X_2)) = E(g_1(X_1)) E(g_2(X_2))$$

**INEQUALITIES:**

**T:** Jensens: $X \in L^1$, $g$ convex $\Rightarrow g[EX] \leq E[g(X)]$

**T:** Lyapunov : For $0 < r \leq s$   $(E[|X|^r])^{\frac{1}{r}} \leq (E[|X|^s])^{\frac{1}{s}}$

**T:** Chebyshev : $g$ positive nondecreasing on $\mathbb{R} \to \mathbb{R}$
$$\Rightarrow \text{For any } X \quad P(X \geq a) \leq \frac{E(g(X))}{g(a)}$$

**T:** Cauchy: $E|XY| \leq \sqrt{E(X^2) E(Y^2)}$

**MOMENTS:**
The $k^{th}$ moment of $X$ is $E(X^k)$
The $k^{th}$ central moment of $X$ is $E[X - E(X)]^k$
The $2^{nd}$ central moment is $V(X) = E(X^2) - [E(X)]^2$

The mixed moments of $X$ & $Y$ are $E(X^n Y^m)$.
For $X, Y \in L^2$   $Cov(X,Y) = E[(X - E(X))(Y - E(Y))]$
$\phantom{For X, Y \in L^2 \quad Cov(X,Y)} = E[XY] - E(X)E(Y)$

$Corr(X,Y) = \dfrac{Cov(X,Y)}{\sqrt{V(X)V(Y)}}$

**T:** $V(X+Y) = V(X) + V(Y) + 2Cov(X,Y)$

**T:** $|Corr(X,Y)| = 1 \iff P(Y = aX + b) = 1$ for $a \neq 0, b \in \mathbb{R}$

**MULTI - DIMENSION:**
When considering R-Vectors $X = (X_1, \ldots, X_d)$   $d \geq 3$
Covariance becomes a matrix:-   *A row vector*
$$C_X^2 = [Cov(X_i, X_j)] = E[(X - E(X))^T (X - E(X))]$$
$C_X^2(i,i) = V(X_i)$, $C_X^2$ is symetric, $C_X^2$ is semi-positive definite.
$\quad C_X^2(i,j) = C_X^2(j,i) \qquad \forall x \in \mathbb{R}^d \quad x C_X^2 x^T \geq 0$

**T:** If $X = (X_1, \ldots, X_n)$ has iid components $X_i \sim N(0,1)$
& $Y = \mu + XA$, $\mu \in \mathbb{R}^m$, $A \in \mathbb{R}^{n \times m}$
$$\Rightarrow Y \sim MVN(\mu, A^T A), \text{ For } Y \sim MVN(\mu, C_Y^2) \text{ :f will have density...}$$

$$f_Y(y) = \frac{1}{\sqrt{(2\pi)^m \det(C_Y^2)}} \exp\left[-\frac{1}{2}(y - \mu)[C_Y^2]^{-1}(y - \mu)^T\right]$$

# 5

CONDITIONAL
EXPECTATIONS

## DEFINING CE

event A occurred. We define CE in this context as...
The **Conditional Expectation (CE)** of R.V. X given event A
is $E(X|A) = \dfrac{E(X \mathbb{1}_A)}{P(A)}$.

Next consider if we have a partition of $\{A_1, ... A_n\}$
of the sample space $\Omega$. If all we know is which of these
events occurred then we have simply the value of
a simple R.V. $Y = \sum_{i=1}^{n} y_i \mathbb{1}_{A_i}$. Say Y takes value $y_i$.

Our best guess for X is now $E(X|A_i)$.
Since $A_i = \{Y = y_i\}$ & $\hat{X} = E(X|A_i) = h(y)$
let $E(X|Y) = h(Y)$. $\leftarrow$ SIMPLE R.V.

**T:** Let $X \in L^1$ & Y be R.V on common prob space
$\implies \exists \hat{X}$ a R.V satisfying    ($\hat{X}$ is a RV on $\sigma(Y)$)
   CE.1) $\hat{X}$ flat on atoms of $\sigma(Y)$ $\leftarrow$
   CE.2) $E(\hat{X}:A) = E(X:A)$ $\forall A \in \sigma(Y)$
   That is unique up to values on sets of
   zero probability.

$E(X:A) = E(X \mathbb{1}_A)$.

We call this unique R.V $E(X|Y)$.

Note if $\mathcal{F} \subset \sigma(Y)$ & we replace CE.1) with the
same condition on $\mathcal{F}$ the theorem is still true,
we call this **CE of X given $\sigma$-Alg $\mathcal{F}$** $E(X|\mathcal{F})$.

## PROPERTIES OF CE

**T:** $\varphi$ is 1-1 function (injective) $\implies E(X|Y) = E(X|\varphi(Y))$

**T:** Linearity: $(\forall a, b \in \mathbb{R})(E(aX + bZ|Y) = aE(X|Y) + bE(Z|Y))$

**T:** Monotone: $X \leq Z$ a.s. $\implies E(X|Y) \leq E(Z|Y)$ a.s.

**T:** $Z = g(Y) \implies E(ZX|Y) = Z E(X|Y)$

**T:** X & Y independent $\implies E(X|Y) = E(X)$

**T:** Double $\mathbb{E}$: $E[E(X|Y_1, Y_2)|Y_1] = E(X|Y_1)$

   In particular $E[E(X|Y)] = E(X)$.

## OTHER CONDITIONALS

**Conditional probabilities** are defined for an event
A & R.V Y $P(A|Y) = E(\mathbb{1}_A|Y)$.

**Conditional distributions** are non-trivial
however it can be proved that conditional
distributions $P_{X|Y}(B|y) = P(X \in B|Y) = E(\mathbb{1}(X \in B)|Y)$
always exist.

When $(X, Y)$ is A.C we can use conditional
densities $f_{X|Y}(x|y) = \dfrac{f_{X,Y}(x,y)}{f_Y(y)}$

Where $f_Y(y) = \int f_{X,Y}(x,y) \, dx$.

Thus $E(X|Y) = \int x f_{X|Y}(x|y) \, dx$.

# 6

**SUFFICIENT STATISTICS**

## THE MODEL

For observed data we make the assumptions that the underlying RE is given by $(\Omega, \mathcal{F}, P_\theta)$. $P_\theta$ is a probability depending on parameter $\theta \in \Theta \subset \mathbb{R}^d$ whose value we dont know. We observe a Random vector $X(\omega) = X \in \mathbb{R}^n$. $P_\theta$ is the distribution on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ induced by $X$ & $P_\theta$.

## STATISTICS

A measurable function $S(X)$ is a **statistic**. Statistics are also then R.vec.

**Estimators** of a parameter $\theta$ are statistics with codomain $\Theta$.

## SUFFICIENCY

A statistic $S$ is **sufficient (SS) for a parameter** $\theta$ if the conditional distribution $P_\theta(X \in B | S)$, $B \in \mathcal{B}(\mathbb{R})$ doesn't depend on $\theta$.

**T:** For $\varphi$ 1-1 function (Injection) & $S$ a SS for $\theta$ $\Rightarrow \varphi(S)$ is also SS for $\theta$.

How can we find & show statistics are sufficient? It is most convenient to use densities. Importantly only A.C. distributions have densities, however A.C. is relative to some measure! So discrete distributions are A.C. relative to the counting measure.

**T:** Suppose all $P_\theta$ are A.C. with respect to some measure $\mu$, with densities $f_\theta(x) = \frac{dP_\theta}{d\mu}(x)$.

$S$ is a SS for $\theta$ $\iff \exists \psi(s,\theta) \; \exists h(x)$
$$f_\theta(x) = \psi(S(x), \theta) h(x)$$
*Neyman-Fisher*   Note that $S(x)$ may be a R.Vec

NOTE: For $X = (x_1, \dots, x_n)$ iid $\Rightarrow f_\theta(x) = \prod f_\theta(x_i)$

**T:** If $T$ is a statistic and $S = \varphi(T)$ for some $\varphi$ is a ss for $\theta \Rightarrow T$ is SS for $\theta$.

$\theta^*(X) = \hat{\theta}^* = \operatorname*{argmax}_{\theta \in \Theta} f_\theta(X)$ is the **maximum likelihood estimator (MLE)** of $\theta$ from $X$.

**T:** $S$ a SS for $\theta \Rightarrow \hat{\theta}^*$ is a function of $S$ only.

## BIAS

How can we compare estimators? We know that there is not a perfect estimator for anything other than the degenerate distribution ie. $\nexists \theta^*$ such that $E(\theta^* - \theta)^2$ is minimised $\forall \theta$.

We need to ask for less, so we compare estimators within certain classes of estimators.

$\theta_0^* \in K$, a class of estimators for $\theta$, is **efficient in** $K$ $\iff \forall \theta^* \in K$, $E_\theta(\theta_0^* - \theta)^2 \leq E_\theta(\theta^* - \theta)^2$ $\forall \theta \in \Theta$.

A common class is the class of estimators with **bias** $b(\theta)$. $K_b = \{\theta^* \mid E_\theta(\theta^*) = \theta + b(\theta), \forall \theta \in \Theta\}$.

$K_0$ is the class of **unbiased estimators**.

**T:** An estimator efficient in $K_b$ is unique up to values on a set of $\theta$ probability.

**T:** Rao-Blackwell: $\theta^* \in K_b$, $S$ a SS for $\theta$ $\Rightarrow \theta_S^* = E_\theta(\theta^* | S)$ has properties
- $\theta_S^*$ is a function of $S$ only
- $\theta_S^* \in K_b$
- $E_\theta(\theta_S^* - \theta)^2 \leq E_\theta(\theta^* - \theta)^2$, $\forall \theta \in \Theta$.

For $\theta \in \mathbb{R}^d$ we can measure the performance of an estimator using $E_\theta(\theta^* - \theta, a)$ for $a \in \mathbb{R}^d$ where $(\cdot, \cdot)$ is the scalar product. *Dispersion.* We prefer an estimator if its dispersion is lower $\forall a$.

**T:** MV R-B: Same except the last condition now...
- $E_\theta(\theta_S^* - \theta, a)^2 \leq E_\theta(\theta^* - \theta, a)^2$, $\forall \theta \in \Theta$

# 7

**CONVERGENCE OF RANDOM VAR:**

We know what it formally means for a sequence of numbers to converge to a point. What might it mean for a sequence of functions, specifically R.V.s to converge to a single function. We have several interesting & useful notions.

$$X_n \xrightarrow[n\to\infty]{a.s.} X \iff (\exists A \subset \Omega)(P(A)=1 \ \& \ \forall \omega \in A \ X_n(\omega) \xrightarrow{n\to\infty} X(\omega))$$
Pointwise convergence on set of prob 1.

$$X_n \xrightarrow[n\to\infty]{P} X \iff (\forall \varepsilon > 0)(P(|X_n - X| \geq \varepsilon) \xrightarrow{n\to\infty} 0)$$

$$X_n \xrightarrow[n\to\infty]{L^2} X \iff (X_n, X \in L^2)(E(X_n - X)^2 \xrightarrow{n\to\infty} 0)$$

$$X_n \xrightarrow[n\to\infty]{L^1} X \iff (X_n, X \in L^1)(E|X_n - X| \xrightarrow{n\to\infty} 0)$$

$$X_n \xrightarrow[n\to\infty]{d} X \iff \lim_{n\to\infty} F_{X_n}(t) = F_X(t) \text{ at all continuity points of } F_X.$$

T: $\iff \forall f$ continuous & bounded $E[f(X_n)] \to E[f(X)]$

## RELATIONS BETWEEN CONVERGENCE

T: $\xrightarrow{a.s.} \Rightarrow \xrightarrow{P} \Rightarrow \xrightarrow{d}$

$\xrightarrow{L^2} \Rightarrow \xrightarrow{P}$ ; $(X_n, X \in L^2)(\xrightarrow{P} \Rightarrow \xrightarrow{L^2})$

## TRANSFORMATIONS

T: $g: \mathbb{R} \to \mathbb{R}$ continuous
- $X_n \xrightarrow{a.s.} X \implies g(X_n) \xrightarrow{a.s.} g(X)$
- $X_n \xrightarrow{P} X \implies g(X_n) \xrightarrow{P} g(X)$
- $X_n \xrightarrow{d} X \implies g(X_n) \xrightarrow{d} g(X)$

# CONVERGENCE THMS:

T: Monotone Convergence: $X_n \geq 0$ R.Vs on a common probability space, $X_n \uparrow X \implies E(X_n) \uparrow E(X)$

T: Fatous Lemma: $X_n \geq 0 \implies E(\liminf_{n\to\infty} X_n) \leq \liminf_{n\to\infty} E(X_n)$

T: Dominated Convergence: $(\forall n)(|X_n| \leq Y \text{ a.s. } \& E(Y) < \infty)$
$X_n \xrightarrow[n\to\infty]{a.s.} X \implies \lim_{n\to\infty} E(X_n) = E(X)$

$\{X_n\}$ is iid sequence of $B(p)$ R.Vs
$S_n = \sum_{i=1}^{n} X_i$

T: Weak LLN: $\dfrac{S_n}{n} \xrightarrow[n\to\infty]{P} p$

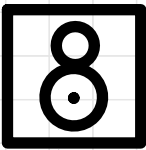T: Strong LLN: $\dfrac{S_n}{n} \xrightarrow[n\to\infty]{a.s.} p$.

# 8

**Characteristic Functions**

For any R.V. $X$ its characteristic function (ChF) is $\varphi_x : \mathbb{R} \to \mathbb{C}$, $\varphi_x(t) = \mathbb{E}(e^{itX})$

Because of this definition the chF always exists, and is finite.

**T:** $|\varphi_x(t)| \le 1$    **T:** $\varphi_x(0) = 1$

**T:** $Y = aX + b$, $a, b \in \mathbb{R} \implies \varphi_y(t) = e^{itb} \varphi_x(at)$.

**T:** $\overline{\varphi_x(t)} = \varphi_x(-t) = \varphi_{-x}(t)$

**T:** ChF is real valued $\iff X$ is symmetric $\iff X \overset{d}{=} -X$.

**T:** Any chF is uniformly continuous.

**T:** $X$ & $Y$ independent $\implies \varphi_{x+y}(t) = \varphi_x(t) \varphi_y(t)$

**T:** $k \in \mathbb{N}$, $E|X|^k < \infty \implies \varphi_x(t)$ is $k$ times cont' differentiable & $\mathbb{E}(X^k) = i^{-k} \frac{d^k}{dt^k} \varphi_x(t)\big|_{t=0}$

**T:** Inversion: $\int |\varphi_x(t)| \, dt < \infty \implies X$ has continuous density
$$f_x(x) = \frac{1}{2\pi} \int e^{-itx} \varphi_x(t) \, dt.$$

**T:** ChF uniquely specify the distribution.

**T:** $\int |t^k \varphi_x(t)| \, dt < \infty \implies X$ has $k$ times diff' continuous density

**T:** $X_n \overset{d}{\to} X \iff (\forall t \in \mathbb{R})(\varphi_{x_n}(t) \to \varphi_x(t))$

**T:** $(\forall t \in \mathbb{R})(\varphi_{x_n}(t) \to \varphi(t))$ where $\varphi_{x_n}$ are ChF & $\varphi(t)$ is continuous at $0 \implies \varphi_x$ is chF of some R.V. $X$ & $X_n \overset{d}{\to} X$.

Clearly then the chF contains a lot of information about the distribution. This is why we use them because they are compact & have plentiful info.

# APPLICATIONS TO STATS

Because of the property that the chF of a sum is the product of chF, they are convenient for proofs in statistics about sums.

**T:** WLLN: $X_1, X_2, \ldots$ iid $\implies \dfrac{\sum_{i=1}^{n} X_i}{n} \xrightarrow[n\to\infty]{P} \mathbb{E}(X_1)$

$E|X_1| < \infty$

**T:** CLT: Further $E(X_1^2) < \infty$ & $V(X_1) = \sigma^2 > 0$ $\implies \dfrac{\sum_{i=1}^{n} X_i - n E(X_1)}{\sigma \sqrt{n}} \xrightarrow[n\to\infty]{d} N(0,1)$

We can relax the condition of iid in the above to a more general condition. One example is for $E(X_i) = 0$ (rescale if necessary).
Lyapunov Condition: $B_n^2 = V(\sum_{i=1}^n X_i)$, need $B_n^{-3} \sum_{i=1}^n E|X_i|^3 \xrightarrow[n\to\infty]{} 0$

**T:** Poisson LT: $X_{n,1}, \ldots, X_{n,n}$ independent R.V.

$P(X_{n,j} = 1) = 1 - P(X_{n,j} = 0) = P_n$, $j = 1, \ldots, n$ & $np_n \to \lambda \in (0, \infty)$
$\implies \sum_{i=1}^{n} X_{n,i} \xrightarrow{d} P(\lambda)$.

# FOR R·VECTORS

$X = (X_1, \ldots, X_d)$, $t = (t_1, \ldots, t_d) \in \mathbb{R}^d$ then
$\varphi_x : \mathbb{R}^d \to \mathbb{C}$, $\varphi_x(t) = \mathbb{E}(e^{i\langle t, X\rangle}) = \mathbb{E}(\exp[i \sum_{j}^{d} t_j X_j])$ ← Dot product

All key results carry over.

**T:** $Y = XA + b$, $A$ a $d \times m$ matrix $b \in \mathbb{R}^m$ $\implies \varphi_y(s) = e^{i\langle s, b\rangle} \varphi_x(sA^T)$

**T:** $\dfrac{\partial^{k_1 + k_2}}{\partial t_1^{k_1} \partial t_2^{k_2}} \varphi_x(t) = i^{k_1 + k_2} \mathbb{E}[X_1^{k_1} X_2^{k_2} e^{i\langle t, X\rangle}]$   Ensuring we have appropriate finite moments

**T:** $(\forall b \in \mathbb{R}^d)(\varphi_{\langle b, x\rangle}(t) = \varphi_x(tb))$   projection of $X$ onto $b$.

**T:** WLLN & SLLN

**T:** CLT : $X_1, X_2, \ldots$ iid R-vect, $E\|X\|^2 < \infty$ & $C_x \leftarrow$ covariance matrix. exists. $\implies \dfrac{\sum_{i=1}^{n} X_i - \mu n}{\sqrt{n}} \to N(0, C_x^2)$

$x^2$ testing...?

**T:** $Z \sim N(0, I_d)$, $b_1, \ldots, b_d$ orthonormal system
$\implies Y = (\langle b_1, Z\rangle, \ldots, \langle b_d, Z\rangle) \sim N(0, I_d)$

# 9

Distribution Free
tests & MLE's

## EMPIRICAL DF

For $X_1, \ldots, X_n$ an iid sample we know that $S = (X_{(1)}, \ldots, X_{(n)})$ is a ss for $F$, the DF of $X_i$. The same info is captured in the

empirical distribution function

$$F_n^*(t) = \frac{1}{n} \sum_{j=1}^{n} \mathbb{1}(X_j \leq t) = \frac{1}{n} \sum_{j=1}^{n} \mathbb{1}(X_{(j)} \leq t)$$

The order stats are all the points of discontinuity of $F_n^*$ but we can also find other statistics

$$\overline{X} = \int t \, dF_n^*(t), \quad s^2 = \frac{1}{n} \sum_{j=1}^{n} (X_j - \overline{X})^2 = \int t^2 dF_n^*(t) - \left( \int t \, dF_n^*(t) \right)^2$$

If there is a parameter $\theta = G(F)$ then we can have a good estimator given by $\theta^* = G(F_n^*)$

**T** Glivenko – Cantelli: $X_1, X_2, \ldots$ iid, DF $F$.

$$\implies D_n = \sup_t |F_n^*(t) - F(t)| \xrightarrow{\text{a.s.}} 0$$

**T** For $X_1, X_2, \ldots$ iid DF $F$ & $U_1, U_2, \ldots \sim U[0,1]$

& $R_n^*(u) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(U_j \leq u)$ ← uniform EDF.

$$\implies D_n = \sup_t |F_n^*(t) - F(t)| = \sup_{u \in [0,1]} |R_n^*(u) - u| \quad \text{← Independent of } F.$$

Further $\sqrt{n} \left( R_n^*(u_1) - u_1, \ldots, R_n^*(u_d) - u_d \right) \longrightarrow N(0, C^2(u))$

$$C^2(u) = \left[ \min\{u_j, u_k\} (1 - \max\{u_j, u_k\}) \right]_{j,k = 1, \ldots, d}$$

We can use this for

Kolmogrov Test: $\lim\limits_{n \to \infty} P(\sqrt{n} D_n \leq x) = 1 + 2 \sum\limits_{k=1}^{\infty} (-1)^k e^{-2k^2 x^2}$

Mises-Smirnov $w^2$-Test: $w_n^2 \xrightarrow{d} \int_0^1 \left[ \sqrt{n} (R_n^*(u) - u) \right]^2 du$

$$\lim_{n \to \infty} P(w_n^2 \leq x) = P\left( \int_0^1 V^2(u) \, du \leq x \right)$$

with $V(u) \sim N(0, u(1-u))$

## MLE's

$X = (X_1, \ldots, X_n)$ $X_j$ have density $f_\theta(x)$ Then the MLE of $\theta$ is

$$\hat{\theta} = \arg\max_\theta f_\theta(X) = \arg\max_\theta \log(f_\theta(X))$$

**T** Gibbs Inequality: $f, g$ densities with respect to $\mu$, on common space $\int f(x) \log(f(x)) \mu(dx) \geq \int f(x) \log(g(x)) \mu(dx)$ when both integrals are finite.

**T** $\hat{\theta}_n \xrightarrow{P} \vartheta$ ← True value & $\sqrt{n} (\hat{\theta}_n - \vartheta) \xrightarrow{d} N\left(0, \frac{1}{I(\vartheta)}\right)$

Where $I(\vartheta) = \int \frac{[f_\vartheta'(x)]^2}{f_\vartheta(x)} \mu(dx)$

**T** $E_\vartheta(\theta_n^* - \vartheta)^2 \geq \frac{1}{n I(\vartheta)}$