# Probability measure

- Given space $\Omega$ ("abstract space").

- Need a $\sigma$-algebra $\mathcal{A}$ on $\Omega$   ("measurable events")

    - $A \in \mathcal{A} \implies A^c \in \mathcal{A}$

    - $(A_i)_{i \in \mathbb{N}} \subset \mathcal{A} \implies \bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$   ("countable unions")

    - $\emptyset, \Omega \in \mathcal{A}$

    - countable intersections

- A measure $\mu$ on $(\Omega, \mathcal{A})$ is a function

$$\mu : \mathcal{A} \to [0, \infty]$$

that is countably additive: If $(A_i)_{i \in \mathbb{N}}$ is a sequence of pairwise disjoint sets, then

$$\mu \left( \bigcup_{i=1}^{\infty} A_i \right) = \sum_{i=1}^{\infty} \mu(A_i) .$$

A measure $P$ on a measurable space $(\Omega, \mathcal{A})$ is called a probability measure if $P(\Omega) = 1$.

The elements in $\mathcal{A}$ are called events.

Then $(\Omega, \mathcal{A}, P)$ is called a probability space.

# Example (1):  Throw one die

$\Omega = \{1, 2, \ldots 6\}$,  $\mathcal{A} = \mathcal{P}(\Omega)$  ($\sigma$-algebra generated by the "elementary events" $\{1\}, \{2\}, \ldots, \{6\}$).

$P$ can also be defined uniquely by assigning

$$P(\{1\}) = P(\{2\}) = \ldots = P(\{6\}) = \frac{1}{6}$$

For example  $P(\{1, 5\}) = P(\{1\}) + P(\{5\}) = \frac{1}{3}$

## Throw two dice:

$$\Omega = \{1, 2, \ldots 6\} \times \{1, 2, \ldots, 6\}$$

first die ↓ , second → $= \{(1,1), (1,2), (1,3), \ldots\}$

$$\mathcal{A} = \mathcal{P}(\Omega)$$
$$P(\{(i,j)\}) = \frac{1}{36}$$

# Example (2)   Normal distribution
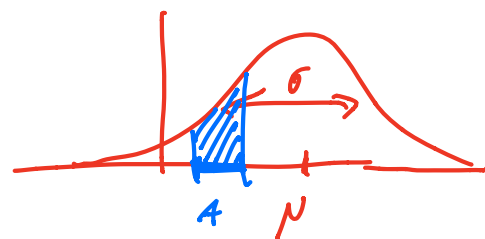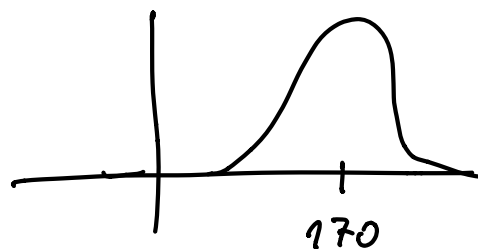
$\Omega = \mathbb{R}$

$\mathcal{A} =$ Borel-$\sigma$-algebra

$f_{\mu, \sigma} : \mathbb{R} \to \mathbb{R}$,

$$x \longmapsto \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(- \frac{(x-\mu)^2}{2\sigma^2}\right)$$



170



$\sigma$

$A$   $\mu$

$P : \mathcal{A} \to [0, 1]$

$$P(A) := \int_A f_{\mu, \sigma}(x)\, dx$$

# Different types of probability measures

## Discrete measure:

$\Omega = \{x_1, x_2, \ldots\}$ finite or at most countable.

$\mathcal{A} = \mathcal{P}(\Omega)$

We define a probability measure $P: \mathcal{A} \to [0,1]$ by assigning probabilities to the "elementary events":

$$P(\{x_i\}) =: p_i$$

with $0 \leq p_i \leq 1$, $\sum_i p_i = 1$.

For $A \in \mathcal{A}$ we assign

$$P(A) = \sum_{\{i \mid x_i \in A\}} p_i.$$

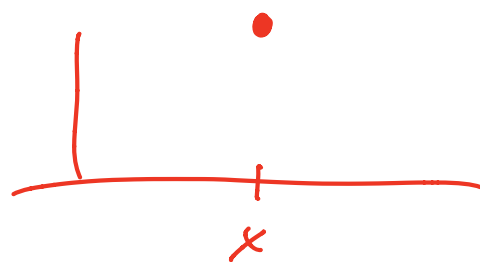Examples: toss a coin; distribution on $\mathbb{Q}$

## Dirac measures:

For $x \in \mathbb{R}$, we define the <u>Dirac measure</u> $\delta_x$ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ by setting

$$\delta_x(A) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}$$

Sometimes this is called a point mass at a point $x$.

A discrete measure on $\mathbb{R}$ can be written as a sum of Dirac measures. For example, throwing a die can be described as $\frac{1}{6}\left(\delta_1 + \delta_2 + \dots + \delta_6\right)$

## Measures with a density

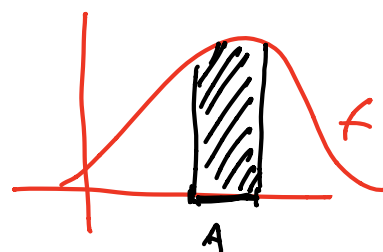Consider $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ and the Lebesgue measure $\lambda$.

Consider a function $f: \mathbb{R}^n \to \mathbb{R}_{\geq 0}$ that is measurable and satisfies $\int f \, d\lambda = 1.$ $\quad \left(= \int f(x) \, dx \right)$

Then we define a measure $\nu$ on $\mathbb{R}^n$ by setting, for all $A \in \mathcal{A}$,

$$\nu(A) := \int_A f(x) \, dx.$$



$\nu$ is the probability measure on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ with density $f$.
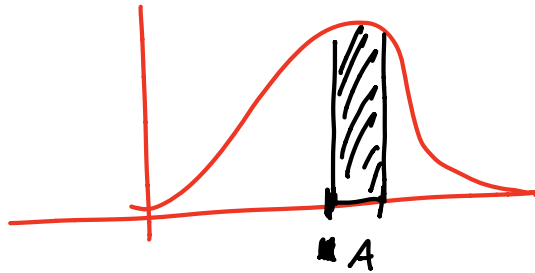
Notation: $\nu = f \cdot \lambda$

Question?. Can we describe every prob measure on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ in terms of a density?   Answer: no!

Counterexample: $\delta_0$ Dirac measure

**Def.** A prob. measure $\nu$ on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ is called
**absolutely continuous** with respect to another measure $\mu$ on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$
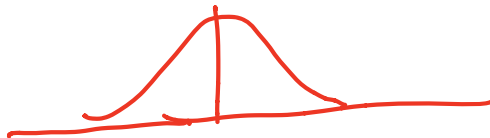if every $\mu$-null set is also a $\nu$-null set:

$$\forall B \in \mathcal{B}(\mathbb{R}^n): \quad \mu(B) = 0 \implies \nu(B) = 0.$$

Notation: $\nu \ll \mu$



$$\mu(A) = 0 \implies \underbrace{\int_A f \, d\mu}_{\nu(A)} = 0$$

Example: $N(0,1) \ll \lambda$



Example: $\delta_0 \not\ll \lambda$ because

$$\lambda(\{0\}) = 0 \quad \text{but} \quad \delta_0(\{0\}) = 1.$$

**Theorem (Radon–Nikodym)**

Consider two prob. measures $\nu, \mu$ on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$. Then the
following two statements are equivalent:

(see next page)

(1) $\nu$ has a density wrt $\mu$.

(2) $\nu$ is absolutely continuous wrt $\mu$.

## Proof idea

(1) $\Rightarrow$ (2)  easy

(2) $\Rightarrow$ (1)  We need to construct a density!

Consider the set $\mathcal{G}$ of all functions $g$ with the following properties:

$(*)$ $\begin{cases} \bullet \ g \text{ is measurable}, \ g \geq 0 \\ \\ \bullet \ g \cdot \mu \leq \nu, \text{ that is} \\ \qquad \forall A \in \mathcal{B}(\mathbb{R}^n): \ \int_A g \, d\mu \leq \nu(A). \end{cases}$

- Observe: $g \equiv 0$ satisfies $(*)$, so $\mathcal{G}$ is not empty.

- If $g, h$ both satisfy $(*)$, then $\sup(g, h)$ satisfies $(*)$.

- Define $\gamma := \sup_{g \in \mathcal{G}} \int g \, d\mu$ and construct a

  sequence $(g_n)_{n \in \mathbb{N}}$ such that $\lim \int g_n \, d\mu = \gamma$.

- Define "density" $f := \sup g_n$

- Now prove: $f$ does the job.

<u>Def</u> $\mu, \nu$ measures on $(\Omega, \mathcal{A})$. $\nu$ is called <u>singular</u>

wrt $\mu$ if there exists $A \in \mathcal{A}$ such that

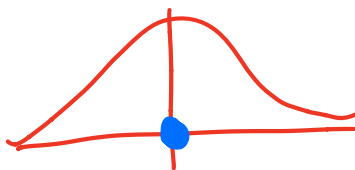$\mu(A) = 0$ but $\nu(A^c) = 0$. Notation: $\mu \perp \nu$.

$$\overset{\nu}{\text{|||||||||}} \qquad \bullet \qquad \mathbb{R}$$
$$\mu = \delta_0$$

Example: $\lambda \perp \delta_0$

<u>Theorem</u> ( Decomposition by Lebesgue)

$\mu, \nu$ prob. measures on $(\Omega, \mathcal{A})$. Then there exists a unique

decomposition $\nu = \nu_1 + \nu_2$ such that

$\nu_1 \ll \mu$ and $\nu_2 \perp \mu$.

Example: $\nu = \frac{1}{2}\left(N(0,1) + \delta_0\right)$

$\nu = \nu_1 + \nu_2$ where $\nu_1 = \frac{1}{2} N(0,1)$ , $\nu_2 = \frac{1}{2} \delta_0$.

<u>Proof</u> Let $\mathcal{N}_\mu$ be the set of all null-sets wrt $\mu$. $\subset \mathcal{A}$.

$\alpha := \sup \{\nu(A) \mid A \in \mathcal{N}_\mu\}$

Can construct a countable sequence $(A_n)_{n \in \mathbb{N}}$ , $A_n \in \mathcal{N}_\mu$,

such that $\nu(A_n) \nearrow \alpha$. By countable additivity

we get $\nu\left(\underbrace{\bigcup_{n \in \mathbb{N}} A_n}_{=:N}\right) = \alpha$.

Define $\nu_1 : A \longmapsto \nu(A \cap N^c)$

$\nu_2 : A \longmapsto \nu(A \cap N)$

Don the job.

Cantor - distribution: non-trivial distribution that is singular w.r.t. $\lambda$

Construct the Cantor set:

- Start with $C_0 := [0, 1]$

  "Remove middle part"



- $C_1 := [0, \frac{1}{3}] \cup [\frac{2}{3}, 1]$.

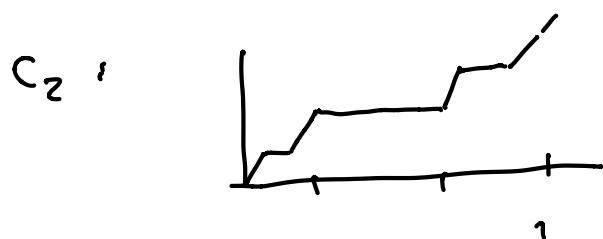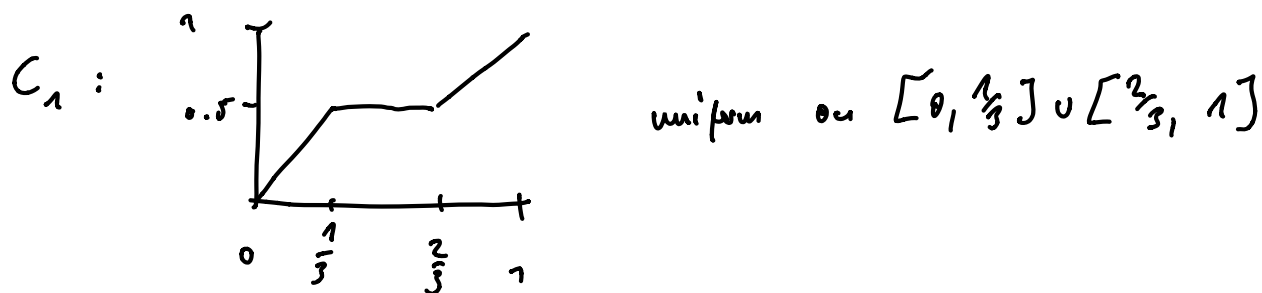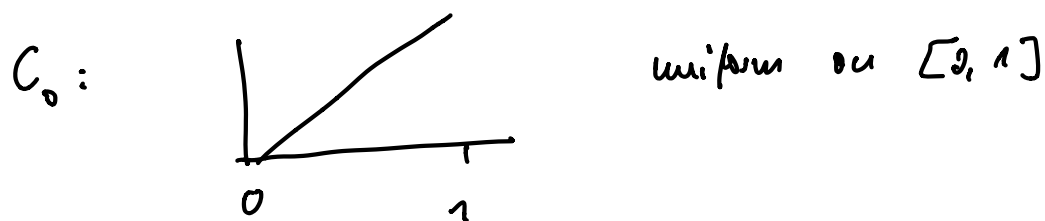  "Remove middle parts from all intervals"



- $C_2 = $

  $\vdots$



The Cantor set is the limit in this process.

Now construct a probability distribution:

Consider the cdfs of the sets $C_0, C_1, C_2 \ldots$

$C_0$ :

uniform on $[0, 1]$

$C_1$ :

uniform on $[0, \frac{1}{3}] \cup [\frac{2}{3}, 1]$

$C_2$ :

$\vdots$

Take limit. Can prove many strange properties:

- Cantor-set is compact, non-empty, empty interior.

- The cdf of "$\mu$" is continuous. $\mu$ is a prob. measure.

- But: $\lambda(c) = 0$.

$\Rightarrow \quad \lambda \perp \mu$

# Cumulative distribution function

Let $P$ is a prob. measure on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$.

Define the function $F: \mathbb{R} \to \mathbb{R}$, $x \mapsto P(]-\infty, x])$.

We say that $F$ is a cumulative distribution function (cdf), that is it satisfies the following ~~pob~~ properties:

(i) $F$ is monotonically increasing,
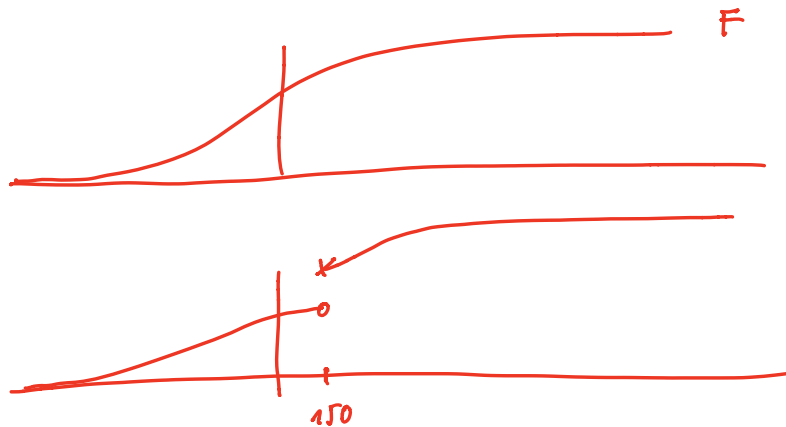
$$\lim_{x \to -\infty} F(x) = 0 \quad , \quad \lim_{x \to +\infty} F(x) = 1.$$

(ii) $F$ is continuous from the right:

$(x_n)_n$ sequence with $x_n \searrow x$

(i.e. $x_n \geq x_{n+1}$ and $x_n \to x$) then also

$F(x_n) \to F(x)$.

*pdf* = prob. density fct
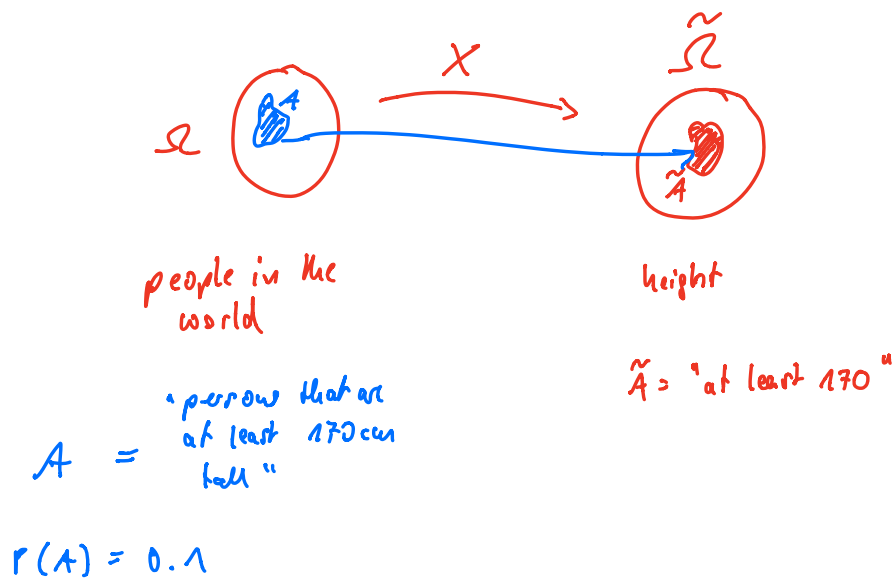
$f$ density of normal distr.

"volume"

$x_0$

*cdf*

$x_-$

Let $F: \mathbb{R} \to \mathbb{R}$ be a function with properties (i) and (ii).
Then there exist a unique prob. measure $P$ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$
such that $P(]-\infty, x]) := F(x)$.

# Random variable

**Def** Let $(\Omega, \mathcal{A}, P)$ be a probability space, $(\tilde{\Omega}, \tilde{\mathcal{A}})$ be another measurable space. A mapping: $X: \Omega \to \tilde{\Omega}$ is called a <u>random variable</u> if $X$ is measurable, i.e.

$$\forall \tilde{A} \in \tilde{\mathcal{A}}: \quad X^{-1}(\tilde{A}) := \{\omega \in \Omega \mid X(\omega) \in \tilde{A}\} \in \mathcal{A}.$$



people in the world

height

$A = $ "person that are at least 170cm tall"

$\tilde{A} = $ "at least 170"

$P(A) = 0.1$

Example: sum of two dice

$$\Omega = \{(i,j) \mid i,j \in \{1, \dots 6\}\}$$

$$\mathcal{A} = P(\Omega)$$

$$P(\{(i,j)\}) = \frac{1}{36}$$

$$\tilde{\Omega} = \{2, \dots, 12\}$$
$$\tilde{\mathcal{A}} = P(\tilde{\Omega})$$

$X$ "sum of the two values"

$X: \Omega \to \{2, \dots, 12\}$, $(i,j) \mapsto i+j$

Is measurable.

**Def** A random variable $X : \Omega \to \tilde{\Omega}$ induces a measure on the target space:

For $\tilde{A} \in \tilde{\mathcal{A}}$ we define

$$\underline{P_X}(\tilde{A}) := P\left(X^{-1}(\tilde{A})\right)$$

This is a probability measure on $(\tilde{\Omega}, \tilde{\mathcal{A}})$ and it is called the **distribution of** X.

**Def** $X : (\Omega, \mathcal{A}, P) \to (\tilde{\Omega}, \tilde{\mathcal{A}})$. Then the family

$$\sigma(X) := \left\{ X^{-1}(\tilde{A}) \mid \tilde{A} \in \tilde{\mathcal{A}} \right\}$$

is a $\sigma$-algebra on $\Omega$ and it is called the $\underline{\sigma\text{-algebra}}$ $\underline{\text{induced by } X}$

(it is the smallest $\sigma$-algebra on $\Omega$ that makes X measurable).

# Conditional probabilities

Notation: $P(A \cap B) = P(\text{"A and B"})$



$P(A \cup B) = P(\text{"A or B"})$



Def. $(\Omega, \mathcal{A}, P)$ probability space, $A, B \in \mathcal{A}$, $P(B) > 0$. Then

$$P(A \mid B) := \frac{P(A \cap B)}{P(B)} \quad \text{is called the}$$

conditional probability of $A$ given $B$.

Theorem: The mapping $P_B : \mathcal{A} \to [0, 1]$, $A \mapsto P(A \mid B)$ is a probability measure on $(\Omega, \mathcal{A})$, it is called the conditional distribution of $P$ with respect to $B$.

Example :    $\Omega$ = all persons on earth

$\mathcal{A}$ = $\mathcal{P}(\Omega)$

$P$ = "uniform"

Event  $A := $ " person has been vaccinated"

$B := $ " person has disease"

$P \left( \text{disease} \mid \text{vaccinated} \right)$



all persons

$P \left( \text{vaccinated} \mid \text{disease} \right)$

Example: two dice

$P \left( \text{"sum is 10"} \mid \text{"first die was 5"} \right)$

# Bayes formula

<u>Law of total probability</u> : Let $B_1, B_2, ..., B_k$ be a disjoint
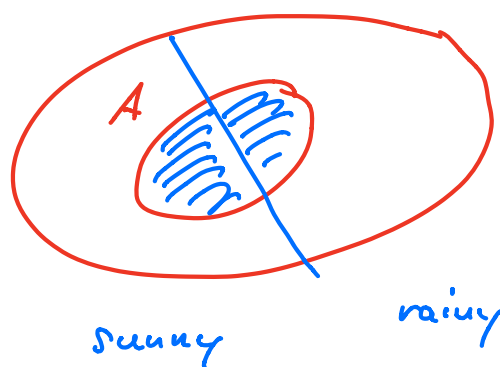
partition of $\Omega$ with $B_i \in \mathcal{A}$ for all $i$, and $A \in \mathcal{A}$. Then

$$P(A) = \sum_{i=1}^{n} P(A \mid B_i) \cdot P(B_i) \qquad = \sum_{i=1}^{n} P(A \cap B_i)$$



A

sunny        rainy

<u>Bayes formula</u> :

$$P(B_i \mid A) = \frac{P(A \mid B_i) \cdot P(B_i)}{\sum_i P(A \mid B_i) \cdot P(B_i)} \qquad = \frac{P(A \cap B_i)}{P(A)}$$

<u>Example</u> : breast cancer screening

Assume 1% of all women above 40 have breast cancer.

90% of women with breast cancer ba will be test positive. ("true positives")

8% of women without breast cancer will receive a positive
result as well   ("false positives")

Given that a woman receives a positive test result, what is
the likelihood that she has breast cancer?

$$P(\text{cancer} \mid \text{positive}) = \frac{P(\text{positive} \mid \text{cancer}) \cdot P(\text{cancer})}{P(\text{pos.} \mid \text{cancer}) P(\text{cancer}) + P(\text{pos} \mid \text{not}_{\text{cancer}}) \cdot P(\text{not}_{\text{cancer}})}$$

$$= \frac{0.9 \cdot 0.01}{0.9 \cdot 0.01 + 0.08 \cdot 0.99} \approx 10\%$$

# Independence

Consider a probability space $(\Omega, \mathcal{A}, P)$. Two events $A, B \in \mathcal{A}$ are called independent if

$$P(A \cap B) = P(A) \cdot P(B)$$

Observation: $A$ is independent of $B \iff P(A \mid B) = P(A)$

A family of events $(A_i)_{i \in I}$ is called independent if for all finite subsets $J \subset I$ we have

$$P\left( \bigcap_{i \in J} A_i \right) = \prod_{i \in J} P(A_i).$$

( Family is called pairwise independent if $\forall i, j \in I$:
$$P(A_i \cap A_j) = P(A_i) \cdot P(A_j).$$ This does not
imply independence! )

Two random variables $X : \Omega \to \Omega_1$, $Y : \Omega \to \Omega_2$ are called independent if their induced $\sigma$-algebras $\sigma(X)$, $\sigma(Y)$ are independent:
$$\forall A \in \sigma(X), \; B \in \sigma(Y): \quad P(A \cap B) = P(A) \cdot P(B).$$

Notation for independence:

$$A \perp\!\!\!\perp B$$

$$X \perp\!\!\!\perp Y$$

# Expectation (discrete case)

Consider a discrete random variable $X : \Omega \to \mathbb{R}$
(that is, $X(\Omega)$ is at most countable).

**Definition** $(\Omega, \mathcal{A}, P)$ prob. space, $S \subset \mathbb{R}$ at most countable, $X : \Omega \to S$ random variable.

If $\sum_{r \in S} |r| \cdot P(X = r) < \infty$, then

$$E(X) := \sum_{r \in S} r \cdot P(X = r) \quad \text{is called the expectation of } X.$$

(sometimes people write $EX$, $\mathbb{E}X$ or $\mathbb{E}(X)$).

## Examples

- Toss a coin. $\Omega = \{ head, tail \}$, $\mathcal{A} = \mathcal{P}(\Omega)$, $P(head) = p$
  $P(tail) = 1-p$.
  $0 < p < 1$.
  $X : \Omega \to \{0, 1\}$, $head \mapsto 1$, $tail \mapsto 0$.
  $E(X) = 0 \cdot \underbrace{P(X=0)}_{1-p} + 1 \cdot \underbrace{P(X=1)}_{p} = p$.

- Test error of a classifier.

**Def** A rv is called "centred" if $E(X) = 0$.

## Important properties:

- Linear:  $E(a \cdot \overset{rv}{X} + b \cdot \overset{rv}{Y}) = a \cdot E(X) + b \cdot E(Y).$

  $a \in \mathbb{R} \qquad b \in \mathbb{R}$

- $X, Y$ independent $\Rightarrow E(X \cdot Y) = E(X) \cdot E(Y)$

$$\sum_{ij} |x_i y_j| \, P(X = x_i, Y = y_j) \overset{ind.}{=}$$

$$= \sum_{ij} \underbrace{|x_i y_j|}_{|x_i| |y_j|} \, P(X = x_i) \cdot P(Y = y_j)$$

$$= \left( \underbrace{\sum_i |x_i| P(X = x_i)}_{< \infty} \right) \left( \underbrace{\sum_j |y_j| P(Y = y_j)}_{< \infty} \right)$$
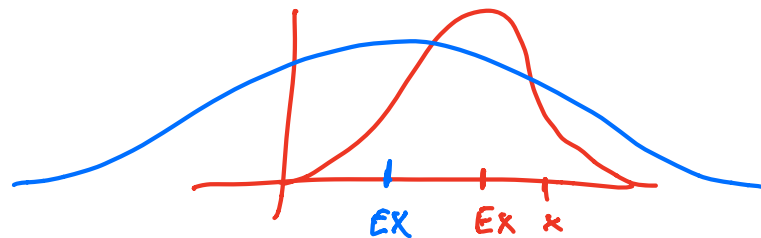
# Variance, covariance, correlation
## (discrete case)

**Def** $X, Y: (\Omega, \mathcal{A}, P) \to \mathbb{R}$ discrete rvs with
$E(X^2) < \infty$, $E(Y^2) < \infty$.

Then
$$\text{Var}(X) := E\left((X - E(X))^2\right)$$

is called the variance of $X$

and $\sqrt{\text{Var}(X)} =: \sigma_X$

is called the standard deviation.

EX  Ex  x

high variance

moderate variance

$$\text{Cov}(X,Y) := E\left((X - E(X)) \cdot (Y - E(Y))\right) \quad \text{is called}$$

the covariance of $X$ and $Y$.
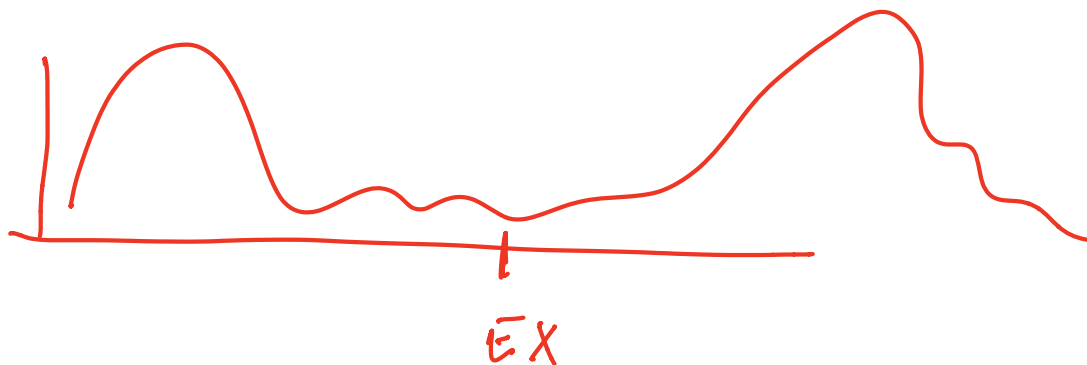
$$\rho_{XY} := \frac{\text{Cov}(X,Y)}{\sigma_X \cdot \sigma_Y} \quad \in [-1, 1] \quad \text{is called the}$$
correlation coefficient.

If $\text{Cov}(X,Y) = 0$, then $X$ and $Y$ are called uncorrelated.

More generally, for $k \in \mathbb{N}$ we define the terms
$$E(X^k) \quad (\text{"}k\text{-th moment"}),$$
$$E\left((X - E(X))^k\right) \quad (\text{"}k\text{th centered moment"})$$

$\overline{E}X$

## Intuition about covariance

$$Cov(X, Y) = E\left( (X - E(X)) \cdot (Y - E(Y)) \right)$$



$y$ body weight

$X$ = shoe size

positive, large covariance

$\rho \approx 0.9$



demand

price

negative cov, larger in absolute values

$\rho \approx -0.9$



$y$

$X$

$Cov \approx 0$
(uncorrelated).

$\triangle$ ! Uncorrelated $\not\Rightarrow$ independence!

$\Leftarrow$



independence

## Properties

- $Var(X) = E(x^2) - \left(E(x)\right)^2$

- $Cov(X,Y) = E(X \cdot Y) - E(X) \cdot E(Y)$

- $E(aX + b) = a \cdot E(X) + b$

- $Var(a \cdot X + b) = a^2 Var(X)$

- $Cov(X,Y) = Cov(Y, X)$

- $Var(X + Y) = Var(X) + Var(Y) + Cov(X,Y)$

- $X, Y$ independent $\Rightarrow Cov(X,Y) = 0$

$\not\Leftarrow$

- $X, Y$ independent $\Rightarrow Var(X + Y) = Var(X) + Var(Y)$.

# Expectation and variance in the general setting

$$L^k(\Omega, \mathcal{A}, P) := \left\{ X: \Omega \to \mathbb{R} \mid X \text{ measurable and } \int_\Omega |X^k| \, dP < \infty \right\}$$

$(\Omega, \mathcal{A}, P)$ prob. space, $X: \Omega \to \mathbb{R}$ with distribution $P_X = X(P)$, $X \in L^1(\Omega, \mathcal{A}, P)$. The ==expectation== of $X$ is then defined as

$$\boxed{E(X)} := \int_\Omega X \, dP = \int_\mathbb{R} x \, dP_X(x)$$

<span style="color:red">(case of density $f$:</span>
$$\int_\mathbb{R} x \, f(x) \, dx \quad )$$

If $X^k \in L^1(\Omega, \mathcal{A}, P)$ then

$$E(X^k) = \int x^k \, dP \text{ is called the } k\text{-th moment of } X.$$

If $X \in L^2(\Omega, \mathcal{A}, P)$ we define
$$\text{Var}(X) = E((X - E(X))^2)$$
$$\text{Cov}(X,Y) = E((X - E(X)) \cdot (Y - E(Y))).$$

# Markov and Chebyshev inequalities

## Cauchy-Schwartz-inequality

$X, Y \in L^2(\Omega, \mathcal{A}, P)$.    Then:

$$E(X \cdot Y)^2 \leq E(X^2) \cdot E(Y^2)$$

## Markov inequality:   $\varepsilon > 0$, $f: [0, \infty[ \rightarrow [0, \infty[$,

$f$ monotonically increasing. Then

$$P(|Y| > \varepsilon) \leq \frac{E(f(|Y|))}{f(\varepsilon)}$$

In particular,

$$P(|Y| > \varepsilon) \leq \frac{E(|Y|)}{\varepsilon}$$

## Chebyshev inequality:   $\varepsilon > 0$, $X \in L^2(\Omega, \mathcal{A}, P)$. Then:

$$P(|X - E(X)| > \varepsilon) \leq \frac{Var(X)}{\varepsilon^2}$$

key quantity in learning theory

# Examples of probability distributions

## Discrete distributions

- Uniform distr. on $\{1,...,n\}$: $P(\{i\}) = \frac{1}{n}$

- Binomial distribution on $\{0,...,n\}$

   Toss a coin $n$ times, independently, each time with probability $p$ of observing head. Denote head = 1, tail 0,

   $X :=$ # heads

   $$P(X=k) := \binom{n}{k} p^k (1-p)^{n-k}$$

- Poisson distribution on $\mathbb{N}$

   Parameter $\lambda > 0$

   $$P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

   Intuition: number of incoming calls at a hotline.

## Continuous distributions

Uniform distribution on $[a,b]$: constant density

# Normal distribution on $\mathbb{R}$

Density: parameter $\mu$ (mean), $\sigma$ (std. deviation)

$$f_{\mu, \sigma}(x) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



Notation: $N(\mu, \sigma^2)$

Some first properties:

- $X \sim N(\mu_1, \sigma_1^2)$ , $Y \sim N(\mu_2, \sigma_2^2)$ , $X, Y$ independent.

  Then $X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$

## Normal distribution in higher dimension

$$X: \Omega \to \mathbb{R}^n, \quad X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad \mu_i \in E(X_i), \quad \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$$

$\Sigma \in \mathbb{R}^{n \times n}$ with $\Sigma_{ij} = Cov(X_i, X_j)$, called covariance matrix.

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{n/2} (\det \Sigma)^{1/2}} \exp\left(- \frac{1}{2}(x-\mu)^t \Sigma^{-1} (x-\mu)\right)$$

Notation: $\mathcal{N}(\mu, \Sigma)$

**Prop** $\Sigma$ is psd and symmetric.

Consequence: $\Sigma$ has real-valued, non-negative eig.



Contour lines of $f_{\mu, \Sigma}$

• $X_1, \ldots, X_n$ are independent $\iff$ $\Sigma = \begin{pmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \\ & & \ddots & \\ 0 & & & \sigma_n^2 \end{pmatrix}$

- $X \sim N(\mu_1, \Sigma_1)$, $Y \sim N(\mu_2, \Sigma_2)$, independent, then

$$X + Y \sim N\left(\mu_1 + \mu_2, \ \Sigma_1 + \Sigma_2\right)$$

## Mixture of Gaussians

Consider $\pi_1, \pi_2, \ldots, \pi_k$ with $0 \le \pi_i \le 1$ and $\sum \pi_i = 1$

Consider the following density:

$$f(x) = \sum_{i=1}^{k} \pi_i \cdot f_{\mu_i, \Sigma_i}(x)$$

# Convergence of random variables

Consider rv $X_i : \Omega \to \mathbb{R}$, $i \in \mathbb{N}$, $X : \Omega \to \mathbb{R}$, $(\Omega, \mathcal{A}, P)$ a probability space.

(1) $(X_i)_{i \in \mathbb{N}}$ converges to $X$ **almost surely** : $\iff$

$$P\left( \left\{ \omega \in \Omega \mid \lim_{i \to \infty} X_i(\omega) = X(\omega) \right\} \right) = 1$$

Notation: $X_i \to X$ a.s.



(2) $(X_i)_{i \in \mathbb{N}}$ converges to $X$ **in probability** : $\iff$

$$\forall \varepsilon > 0 \quad P\left( \left\{ \omega \in \Omega \mid |X_i(\omega) - X(\omega)| > \varepsilon \right\} \right) \longrightarrow 0$$

Let us check that these definitions make sense. We need to prove that the events in (1) and (2) are in fact in the $\mathcal{A}$.

Case (1):

$$\lim_i X_i(\omega) = X(\omega)$$

$$\iff \forall k \in \mathbb{N} \; \exists N \in \mathbb{N} \; \forall n > N : \; |X_n(\omega) - X(\omega)| < \frac{1}{k}$$

So we get:

$$\{\omega \mid X_i(\omega) \to X(\omega)\} =$$

$$= \bigcap_{k \in \mathbb{N}} \bigcup_{N \in \mathbb{N}} \bigcap_{n \geq N} \left\{ \omega \mid |X_n(\omega) - X(\omega)| < \frac{1}{k} \right\} \quad \textcolor{red}{\in \mathcal{A}}$$

<span style="color:red">
$\underbrace{\phantom{countable}}$ countable unions and intersections
</span>

<span style="color:red">
$\underbrace{\phantom{Xn, X}}$ $X_n, X$ are measurable $\Rightarrow$ $|X_n - X|$ is measurable

so $\{\dots\} \in \mathcal{A}$
</span>

(3) $X_n \to X$ in $\underline{L^p}$ ("in the $p$-th mean") $\quad :\Leftrightarrow$

$X_n, X \in L^p$ and $\|X_i - X\|_p \to 0$.

(4) Let $M^1(\mathbb{R}^n)$ be the set of all probability measures on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$. Assume $(\mu_n)_n \subset M^1(\mathbb{R}^n)$, $\mu \in M^1(\mathbb{R}^n)$.
$C_b(\mathbb{R}^n) :=$ space of bounded continuous functions.

$$\mu_n \to \mu \quad \underline{\text{weakly}} \quad :\Leftrightarrow$$

$$\forall f \in C_b(\mathbb{R}^n) : \int f \, d\mu_n \longrightarrow \int f \, d\mu$$

⌜ Excursion:

In functional analysis, a sequence $(x_n)_n$ in a Banach space $B$ converges weakly if for all bounded lin. functionals $f$, we have that $f(x_n) \to f(x)$. (i.e. for all $f \in B'$).

Space $M^1(\mathbb{R}^n)$ itself is not a Banach space, but $\subset M(\mathbb{R}^n)$, space of all bounded measures. The dual space of $M(\mathbb{R}^n)$ is $C_b(\mathbb{R}^n)$.

⌞

(5) $X_i, X : (\Omega, \mathcal{A}, P) \to \mathbb{R}^n$. The sequence $X_n$ converges in __distribution__ to $X$ :$\iff$ the distributions $P_{X_n}$ converge to $P_X$ weakly.

We have the following implications.( but none of the missing implications is true in general):

almost surely                              in $L^1$ $\Leftarrow$ in $L^p$
                                                            $(p > 1)$

⟱                    ⤲

in probability

⟱

in distribution

## Example (converge a.s., in prob., but not in $L^1$)

$$X_n : \mathbb{R} \to \mathbb{R}, \qquad X_n(\omega) = \begin{cases} n & \text{for} \quad 0 \le x \le \frac{1}{n} \\ 0 & \text{otherwise} \end{cases}$$

area 1

$$\forall x > 0: \quad X_n(x) \to 0.$$

$\frac{1}{n} \qquad 1$

Can formally see: a.s., in prob.

But: no converge in $L^1$.

## Example (converge in prob., in $L^1$, but not a.s.)

"sliding blocks"

$$f_1 = \mathbb{1}_{[0,1]}$$

$$f_2 = \mathbb{1}_{[0, \frac{1}{2}]}, \qquad f_3 = \mathbb{1}_{[\frac{1}{2}, 1]}$$

$\frac{1}{2} \qquad 1 \qquad\qquad \frac{1}{2} \qquad 1$

$$f_4 = \mathbb{1}_{[0, \frac{1}{3}]}, \qquad f_5 = \mathbb{1}_{[\frac{1}{3}, \frac{2}{3}]}, \qquad f_6 = \mathbb{1}_{[\frac{2}{3}, 1)}$$

**Example** (Conv. in distribution, but not in prob.)

- $X_n : [0, 1] \to \mathbb{R}$, $X_1 = X_2 = \dots = \mathbb{1}_{[0, \frac{1}{2}[}$

- $X = \mathbb{1}_{[\frac{1}{2}, 1]}$

Obviously $X_n \not\to X$ in prob., but:

$$P_{X_1} = \frac{1}{2}\left(\delta_0 + \delta_1\right) = P_{X_2} = P_{X_3} = \dots = P_X$$

so $X_n \to X$ in distribution.

# Theorem of Borel - Cantelli

$(\Omega, \mathcal{A}, P)$ prob. space, $(A_n)_n$ sequence of events in $\mathcal{A}$.

$$P(A_n \text{ infinitely often }) := P(A_n \text{ i.o.})$$

$$= P(\{\omega \in \Omega \mid \omega \in A_n \text{ for infinitely many } n\})$$

**Proposition** : $X_n, X$ r.v. on $(\Omega, \mathcal{A}, P)$.

$$X_n \to X \text{ a.s.} \quad \Longleftrightarrow$$

$$\forall \varepsilon > 0 : \quad P(\{|X_n - X| > \varepsilon\} \text{ inf. often}) = 0$$

**Proof intuition** :

$$\{\lim X_n = X\}$$

$$= \{\forall k : |X_n - X| > \tfrac{1}{k} \text{ at most finitely often}\}$$

$$= \bigcap_{k \in \mathbb{N}} \{|X_n - X| > \tfrac{1}{k} \text{ at most fin. often}\}$$

$$\left(\bigcup_{k \in \mathbb{N}} \{|X_n - X| > \tfrac{1}{k} \text{ inf. often}\}\right)^{\text{complement}}$$

**Theorem :** Consider a sequence of events $(A_n)_n \subset \mathcal{A}$.

(1) If $\sum_{n=1}^{\infty} P(A_n) < \infty$, then $P(A_n \text{ i.o.}) = 0$.

(2) If $\sum_{n=1}^{\infty} P(A_n) = \infty$, and if $(A_n)_n$ are independent,

then $P(A_n \text{ i.o.}) = 1$.

**Application in learning theory:**

Assume that $P\left(|X_n - X| > \frac{1}{n}\right) < \delta_n$, and

assume that $\sum_{n=1}^{\infty} \delta_n < \infty$. Then you can use

Borel - Cantelli to prove that

$$P\left(|X_n - X| > \frac{1}{n} \text{ i.o.}\right) = 0,$$

thus $X_n \to X$ a.s.

# Limit Theorems: LLN and CLT

## Strong law of large numbers

$X_n : (\Omega, \mathcal{A}, P) \to \mathbb{R}$ iid (identically distributed and independent). Assume the mean $\mu := E(X_n) < \infty$, and $\text{Var}(X_1) =: \sigma^2 < \infty$. Then:

$$\lim \frac{1}{n} \sum_{i=1}^{n} X_i = \mu \quad \text{a.s. and in } L^2.$$

## Remarks:

- Many versions of this theorem exist. (slightly relaxing iid)

- "Strong law" $\Leftrightarrow$ converge a.s.

   "Weak law" $\Leftrightarrow$ convergence in probability

## Central Limit theorem

$(X_i)_{i \in \mathbb{N}}$ iid rv with mean $\mu$, variance $\sigma^2 < \infty$.
Consider the rv $S_n := \sum_{i=1}^{n} X_i$. We normalize it to

$$Y_n := \frac{S_n - n \cdot \mu}{\sqrt{n} \, \sigma}$$

(which has mean $0$ and standard dev. $1$).

Then $Y_n \to Y$ in distribution where $Y \sim N(0, 1)$.

Illustration: $X_i$ coin, head $\hat{=} 1$, tail $\hat{=} 0$

$$S_n = \sum X_i \in [0, n]$$

histogram



0   1   2   $\cdots$   n

$N(0,1)$

# Concentration inequalities

Motivation: random projections

$\mathbb{R}^d$, $d$ large

want to project in $\mathbb{R}^\ell$, $\ell$ "small"



## Theorem of Johnson - Lindenstrauss:

Can guarantee (for certain parameters $\varepsilon, k$)

$$(1-\varepsilon) \| x_i - x_j \|_{\mathbb{R}^d} \leq \| \pi(x_i) - \pi(x_j) \|_{\mathbb{R}^\ell}$$

$$\leq (1+\varepsilon) \| x_i - x_j \|_{\mathbb{R}^d}$$

Construction / proof steps:

(1) Assume you know $\| x_i - x_j \|_{\mathbb{R}^d} = 1$.

Compute $E\left( \| \pi(x_i) - \pi(x_j) \|_{\mathbb{R}^\ell} \right)$, "easy".

(2) $P\left( \left| \| \pi(x_i) - \pi(x_j) \| - E(\cdots) \right| > t \right)$ ?

# Hoeffding inequality

**Theorem (Hoeffding):** $X_1, \ldots, X_n$ rv $: (\Omega, \mathcal{A}, P) \to (\mathbb{R}, \mathcal{B})$, independent,

assume that $X_i \in [a_i, b_i]$ a.s. for $i = 1, \ldots, n$.

Let $S_n := \sum_{i=1}^{n} (X_i - E(X_i))$. Then for any $t > 0$,

$$P(S_n \geq t) \leq \exp\left( - \frac{2t^2}{\sum_{i=1}^{n} (b_i - a_i)^2} \right).$$



$-t$    $0$     mean of $S_n$    $t$

# Application of Hoeffding: SLLN

**Prop** $(X_i)_{i \in \mathbb{N}}$ iid rv, $a \leq X_i \leq b$, let $X$ have the same distribution as the $X_i$.

Then: $\frac{1}{n} \sum_{i=1}^{n} X_i \to E(X)$ a.s.

**Proof** Hoeffding $\Rightarrow$

- $P\left( \frac{1}{n} \sum X_i - E(X) > t \right) \leq \exp\left( - \frac{2nt^2}{(b-a)^2} \right)$

- $P\left( \frac{1}{n} \sum X_i - E(X) < -t \right)$

$$= P\left(\frac{1}{n}\sum(-x_i) - E(-x) > t\right) \le \exp\left(-\frac{2nt^2}{(b-a)^2}\right)$$

Combined we get

$$P\left(\left|\frac{1}{n}\sum x_i - E(x)\right| > t\right) \le 2\exp\left(-\frac{2nt^2}{(b-a)^2}\right).$$

Now want to apply Borel-Cantelli to get a.s. convergence:

$$z_n := \frac{1}{n}\sum_{i=1}^{n} x_i$$

$$\sum_{n=0}^{\infty} P(z_n - E(x) > t) \le 2\cdot\underbrace{\sum_{n=0}^{\infty} \exp\left(-\frac{2nt^2}{(b-a)^2}\right)}_{=:\text{ sum}} \overset{!}{\underset{(\ast)}{\le}} \infty$$

(∗) Substitute: $r := \exp\left(-\frac{2t^2}{(b-a)^2}\right) \in [0,1[$

Observe: $\exp\left(-\frac{2nt^2}{(b-a)^2}\right) = r^n$

$$\text{sum} = 2\sum_{n=0}^{\infty} r^n = 2\cdot\frac{1}{1-r} < \infty.$$

Now Borel-Cantelli gives almost sure convergence. □

Remark: Hoeffding is tight (cannot be improved without further assumptions). For fair coin tosses it is tight.

But: not tight if coin is biased ⤳ need other inequalities

# Bernstein inequality

**Theorem ( Bernstein):** $x_1, \ldots, x_n$ independent with $0$ mean, $|x_i| < 1$ a.s. Let $\sigma^2 := \frac{1}{n} \sum_{i=1}^{n} \text{Var}(x_i)$. Then for all $t > 0$,

$$P\left(\frac{1}{n} \sum_{i=1}^{n} x_i > t\right) \leq \exp\left(-\frac{n t^2}{2(\sigma^2 + t/3)}\right)$$

# Concentration inequality for functions with bounded differences

Consider a function $f: \mathbb{R}^n \to \mathbb{R}$ (or more generally, $f: \mathcal{X}^n \to \mathbb{R}$ for some "arbitrary" space $\mathcal{X}$).

We say that $f$ has the bounded differences property if there exist constants $c_1, \ldots, c_n$ such that

(⊛) $\quad \sup_{\substack{x_1 \ldots x_n \in \mathcal{X} \\ \tilde{x}_i \in \mathcal{X}}} \left| f(x_1, \ldots, x_{i-1}, x_i, x_{i+1}, \ldots, x_n) - f(x_1, \ldots, x_{i-1}, \tilde{x}_i, x_{i+1}, \ldots, x_n) \right| \leq c_i$

Example: $f(x_1, \ldots, x_n) = \sum_{i=1}^{n} x_i$, and $a \leq x_i \leq b \ \forall i$, then $f$ satisfies ⊛ with $c_i = b - a$.

**Theorem (Mc Diarmid)** $X_1, \ldots, X_n$ independent rv, $X_i \in \mathcal{X}_i$, $f: \mathcal{X}_1 \times \ldots \times \mathcal{X}_n \to \mathbb{R}$ function with bounded difference property.

Then, for any $t > 0$,

$$P\left( f(X_1, \ldots, X_n) - E\left( f(X_1, \ldots, X_n) \right) > t \right)$$

$$\leq \exp\left( - \frac{2 t^2}{\sum_{i=1}^{n} c_i^2} \right)$$

Applications:

- stability in ML

- standard theoretical CS, randomized algorithms

- largest eigenvalue of a random symmetric matrix



$A = \begin{pmatrix} \phantom{xxxxx} \end{pmatrix}$ $\to X_{27}$, $X_{39}$ $\sim$ draw iid

# Glivenko - Cantelli Theorem

$F$ cdf : $\quad F(a) = P(X \le a)$

$X_1, \dots, X_n \sim F$, iid

$F_n : \mathbb{R} \to [0,1]$

$F_n(a) := \frac{1}{n} \sum_{i=1}^{n} 1\{X_i \le a\}$



Now fix our particular $a_0 \in \mathbb{R}$.

$F_n(a_0) \longrightarrow F(a_0)$ by the law of large numbers.

$\qquad$ Because $1\{X_i \le a_0\}$ is a Binomial rv with

$\qquad p = P(X_i \le a_0).$

So it is clear that $F_n \xrightarrow{a.s.} F$ pointwise (i.e. $\forall a_0$)

Now let's look at uniform convergence.

**Theorem** $\quad X_1, \dots, X_n$ iid random variables with cdf $F$.
Let $F_n$ be the empirical cdf induced by the sample. Then:

$$P\left( \sup_{a \in \mathbb{R}} | F_n(a) - F(a) | > \varepsilon \right) \le$$



$$\le 8 \cdot (n+1) \cdot \exp\left( - \frac{n \varepsilon^2}{32} \right) .$$

In particular, $\sup | F_n - F | \longrightarrow 0$ a.s.,
i.e. $F_n \to F$ uniformly a.s.

**Proof**

Observe: LLN $\Rightarrow$ $P\left( |F_n(a_0) - F(a_0)| > \varepsilon \right) \to 0$

for any fixed $a_0$.

Problem: need to look at

$$P\left( \sup_{a \in \mathbb{R}} |F_n(a) - F(a)| > \varepsilon \right)$$

difficult because $\mathbb{R}$ is uncountable

If we take a supremum over a finite set, it is easier:

$$P\left( \max_{i=1\ldots n} |U_i| > \varepsilon \right) =$$

$$= P\left( |U_1| > \varepsilon \text{ or } |U_2| > \varepsilon \text{ or} \ldots \text{ or } |U_n| > \varepsilon \right)$$

$$\leq \sum_{i=1}^{n} P\left( |U_i| > \varepsilon \right)$$

Trick of the proof: convert $\sup_{a \in \mathbb{R}}$ to something "finite".

How could we achieve this?



true fact

$F$

$F_n$ induced by given sample

$F_n'$ induced by a ghost sample

$|red - green|$

$|red - blue| \leq 2 |green - blue|$

**Step 1** : Symmetrization by ghost sample

Assume $X_1', \ldots, X_n' \sim F$ independently ("ghost sample"),

Denote by $F_n'$ the empirical cdf induced by ghost sample

Now it is easy to prove:

$$P\left( \sup_a \, | F_n(a) - F(a) | > \varepsilon \right)$$

$$\leq 2 \, P\left( \sup_a \, | F_n(a) - F_n'(a) | > \frac{\varepsilon}{2} \right)$$

**Step 2** : Want to split this in two terms

$$| F_n(a) - F_n'(a) | = \left| \frac{1}{n} \sum_{i=1}^n \left( \mathbb{1}_{\{X_i \leq a\}} - \mathbb{1}_{\{X_i' < a\}} \right) \right| \quad \overbrace{\qquad\qquad\qquad\qquad\qquad}^{\textcircled{*}}$$

Introduce Rademacher random variables $\sigma_1, \ldots, \sigma_n$ :

$$\sigma_i \left( \{-1\} \right) = \sigma_i \left( \{1\} \right) = 1/2 .$$

Distribution of $\textcircled{*}$ is the same as the distr. of the following:

$$\left| \frac{1}{n} \sum_{i=1}^n \sigma_i \left( \mathbb{1}_{\{X_i \leq a\}} - \mathbb{1}_{\{X_i' \leq a\}} \right) \right| \qquad = \textcircled{**}$$

Now we have:

$$2P\left(\sup_a |F_n(a) - f_n'(a)| > \frac{\varepsilon}{2}\right)$$

$$= 2P\left(\sup_a \left|\frac{1}{n}\sum \sigma_i \left(\mathbb{1}_{X_i \leq a} - \mathbb{1}_{X_i' \leq a}\right)\right| > \frac{\varepsilon}{2}\right)$$

$$\leq 2P\left(\sup_a \left|\frac{1}{n}\sum \sigma_i \underbrace{\mathbb{1}_{X_i \leq a}}\right| > \frac{\varepsilon}{4}\right) + 2P\left(\sup_a \left|\frac{1}{n}\sum \sigma_i \underbrace{\mathbb{1}_{X_i' \leq a}}\right| > \frac{\varepsilon}{4}\right)$$

$$\phantom{xxxxxxxxxx} u \phantom{xxxxxxxxxxxxxxxxxxxxxxxx} v$$

<span style="color:red">Observe:</span>

<span style="color:red">$$P\left(|u-v| > \frac{\varepsilon}{2}\right) \leq P\left(|u| > \frac{\varepsilon}{4} \quad \text{or} \quad |v| > \frac{\varepsilon}{4}\right)$$</span>

<span style="color:red">$\uparrow$ right side is necessary for left side</span>

$$= 4 \cdot P\left(\sup_a \left|\frac{1}{n}\sum \sigma_i \mathbb{1}_{\{X_i \leq a\}}\right| > \frac{\varepsilon}{4}\right)$$

**Step 3** Exploit "finite structure":

Fix $X_1, \ldots, X_n$ ($\equiv$ condition on $X_1, \ldots, X_n$)

We look at $\mathbb{1}_{X_i \leq a}$ <span style="color:red">~~~~⊢—┼———┼——┼——┼——┼——┤</span>

The rvs $\mathbb{1}_{X_1 \leq a}, \ldots, \mathbb{1}_{\{X_n \leq a\}}$ for fixed $a$ can only

have $n+1$ realizables

$$P\left(\sup_a \frac{1}{n}\left|\sum \sigma_i \mathbb{1}_{X_i \leq a}\right| > \frac{\varepsilon}{4} \,\Big|\, X_1 \cdots X_n\right) \leq$$

$$\leq (n+1)\sup_a \underbrace{P\left(\frac{1}{n}\left|\sum \sigma_i \mathbb{1}_{\{X_i \leq a\}}\right| > \frac{\varepsilon}{4} \,\Big|\, X_1 \cdots X_n\right)}$$

<span style="color:blue">use Hoeffding ($\#$)</span>

Step 4   apply Hoeffding to (H)

Via:

$$P\left( \frac{1}{n}\left| \sum_{i} \sigma_i \, \mathbb{1}_{x_i \leq a}\right| > \frac{\varepsilon}{4} \;\middle|\; X_1 \ldots X_n \right)$$

$$\leq 2 \exp\left(- \frac{n \, \varepsilon^2}{32}\right)$$

Combining everything gives the theorem.

# Product space, joint distributions

Consider two measurable spaces $(\Omega_1, \mathcal{A}_1)$, $(\Omega_2, \mathcal{A}_2)$.

Define the product space $(\Omega_1 \times \Omega_2, \mathcal{A}_1 \otimes \mathcal{A}_2)$ with

$$\Omega_1 \times \Omega_2 = \{ (\omega_1, \omega_2) \mid \omega_1 \in \Omega_1, \omega_2 \in \Omega_2 \}$$

$$\mathcal{A}_1 \otimes \mathcal{A}_2 = \{ A_1 \times A_2 \mid A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2 \}.$$

Consider two rvs
$$X_1 : (\Omega, \mathcal{A}, P) \to (\Omega_1, \mathcal{A}_1)$$
$$X_2 : (\Omega, \mathcal{A}, P) \to (\Omega_2, \mathcal{A}_2).$$

$$X := (X_1, X_2) , (\Omega, \mathcal{A}, P) \to (\Omega_1 \times \Omega_2, \mathcal{A}_1 \otimes \mathcal{A}_2)$$
$$(X_1, X_2)(\omega) = (X_1(\omega), X_2(\omega)).$$

The distribution $P_{(X_1, X_2)}$ on $(\Omega_1 \times \Omega_2, \mathcal{A}_1 \otimes \mathcal{A}_2)$ is called

the joint distribution of $X_1$ and $X_2$.

Example in ML: $(X, Y)$ where $X$ is the input data, $Y$ is
the label

Product measure: $(\Omega_1, \mathcal{A}_1, P_1)$, $(\Omega_2, \mathcal{A}_2, P_2)$ two

prob. spaces. We define the product measure $P_1 \otimes P_2$ on

the product space $(\Omega_1 \times \Omega_2, \mathcal{A}_1 \otimes \mathcal{A}_2)$ as

$$\left( P_1 \otimes P_2 \right) \left( A_1 \times A_2 \right) := P_1 \left( A_1 \right) \cdot P_2 \left( A_2 \right).$$

**Theorem**    Two rvs $X_1, X_2$ are independent if and only if their joint distribution coincides with the product distribution:

$$P_{(X_1, X_2)} = P_1 \otimes P_2 \, .$$

# Marginal distribution

Consider the joint distribution $P_{(X_1, X_2)}$ of two rvs $X := (X_1, X_2)$. The marginal distribution of $X$ wrt $X_1$ is the original distribution of $X_1$ on $(\Omega_1, \mathcal{A}_1)$, namely $P_{X_1}$. Similarly for $P_{X_2}$.

Example in the discrete case:

| $Y \backslash X$ | $x_1$ | $x_2$ | $x_3$ | $\Sigma$ |
|---|---|---|---|---|
| $Y_1$ | $p_1$ | $p_2$ | $p_3$ | $p_1 + p_2 + p_3 = P(Y = y_1)$ |
| $Y_2$ | $p_4$ | $p_5$ | $p_6$ | $p_4 + p_5 + p_6 = P(Y = y_2)$ |

$p_1 + p_4$
$= P(X = x_1)$ $\cdots$

marginal wrt $X$

$\hat{=}$ marginal distribution wrt $Y$.

## Marginal distributions in case of densities

$X, Y : (\Omega, \mathcal{A}, P) \to (\mathbb{R}, \mathcal{B}(\mathbb{R}))$, $Z := (X, Y)$. Assume that the joint distribution of $Z$ has a density $f$ on $\mathbb{R}^2$. Then the following statements hold:

(1) Both $X$ and $Y$ have densities on $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ given by

<span style="color:red">joint density</span>

$$f_X(x) = \int_{-\infty}^{\infty} f(x,y)\, dy \qquad \text{\color{red}sum over } Y$$

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

(2) $X$ and $Y$ are independent iff

$$f(x,y) = f_X(x) \cdot f_Y(y) \qquad a.s.$$

## Mixed cases

For example, consider $X$ a continuous rv with density and $Y$ a discrete rv.

<span style="color:red">So e.g., $X =$ income $\in \mathbb{R}$</span>

<span style="color:red">$Y =$ "yes" or "no", discrete</span>

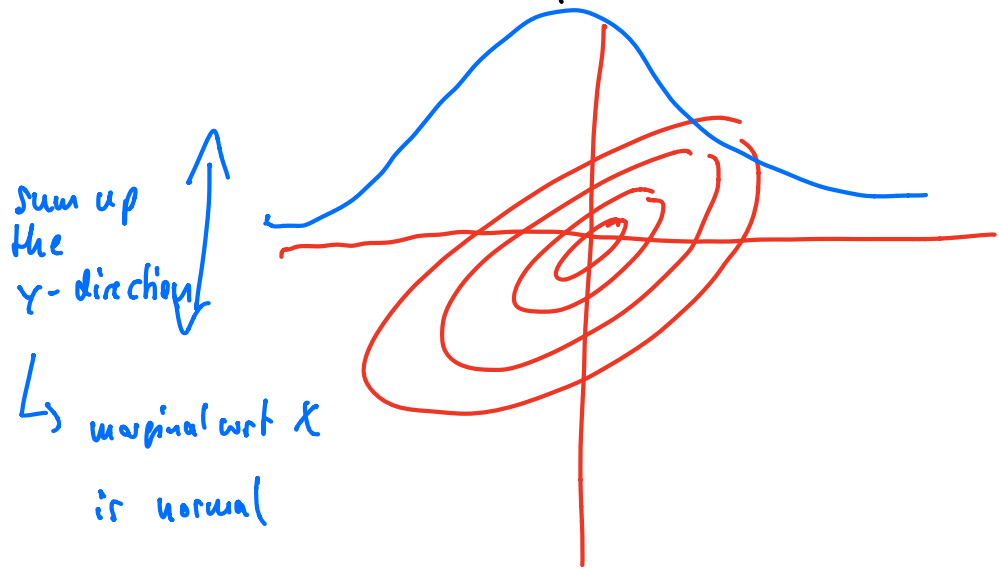## Special case: marginals of multivariate normal distributions

**2 dim** Consider a 2-dim normal rv $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ with mean

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \in \mathbb{R}^2 \quad \text{and} \quad \text{cov. } \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}.$$

Then the marginal distribution of $X$ wrt $X_1$ is again

a normal distribution with mean $\mu_1$ and var $\sigma_1^2$.



sum up
the
Y-direction

$\hookrightarrow$ marginal wrt $X$
is normal

__n-dim__

$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$ .   Group the variables:   $\left.\begin{matrix} x_1 \\ \vdots \\ x_k \end{matrix}\right\} \tilde{X} \in \mathbb{R}^k$

$\left.\begin{matrix} x_{k+1} \\ \vdots \\ x_n \end{matrix}\right\} \tilde{X}^{\#} \in \mathbb{R}^{n-k}$

Want to look at the marginal of $X$ wrt $\tilde{X}$,

$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$  mean,   $\tilde{\mu} := \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_k \end{pmatrix}$ ,  $\mu^{\#} = \begin{pmatrix} \mu_{k+1} \\ \vdots \\ \mu_n \end{pmatrix}$

$\underset{\substack{\uparrow \\ \mathbb{R}^{n \times n}}}{\Sigma} = \left(\begin{array}{c|c} \Sigma_{11} & \Sigma_{12} \\ \hline \Sigma_{21} & \Sigma_{22} \end{array}\right) \Big\} k$

$\underbrace{\phantom{mmmm}}_{k}$

Now the marginal of $X$ wrt $\tilde{X}$ is a normal distr. on $\mathbb{R}^k$

with mean $\tilde{\mu}$ and cov. $\Sigma_{11}$.

# Conditional distributions

### Discrete case:

Know conditional probabilities: $P(A \mid B)$
defined for events $A, B \in \mathcal{A}$, and $P(B) > 0$.

Let $X, Y : (\Omega, \mathcal{A}, P) \to \mathbb{R}$ be discrete rv, $y \in \mathbb{R}$ such that
$P(Y = y) > 0$. Then we can define the conditional probability
measure $P_{X \mid Y = y}$ : $A \longmapsto P(X \in A \mid Y = y)$.

This is a probability measure.

### For general rv this is surprisingly complicated!

~> "regular conditional probabilities"   ~> skipped

### Conditional distributions in case of densities

Assume $Z := (X, Y)$ has a joint density $f : \mathbb{R}^2 \to \mathbb{R}$,
and marginal densities $f_X, f_Y : \mathbb{R} \to \mathbb{R}$. Then the function

$$ f_{X \mid Y = y}(x) := \frac{f(x, y)}{f_Y(y)} $$

is then also a density on $\mathbb{R}$, called the conditional density of
$X$ given $Y = y$.

## Example: normal distributions

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_u \end{pmatrix} \begin{matrix} \} \tilde{\mu} \\ \} \mu^{\#} \end{matrix} \qquad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

If $X = \begin{pmatrix} x_1 \\ \vdots \\ x_u \end{pmatrix} \sim N(\mu, \Sigma)$, then the conditional distributions

of $\tilde{X} = \begin{pmatrix} x_1 \\ \vdots \\ x_u \end{pmatrix}$ wrt $x^{\#} = \begin{pmatrix} x_{u+1} \\ \vdots \\ x_u \end{pmatrix}$ is given by

$$P_{\tilde{X} \mid x^{\#}} \sim N\left( \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x^{\#} - \tilde{\mu}), \right.$$

$$\left. \Sigma_{22} - \Sigma_{12}^{t} \Sigma_{11}^{-1} \Sigma_{12} \right).$$



$x_1 = 1$

# Conditional expectation

<u>Def</u> (discrete case)  $X, Y : (\Omega, \mathcal{A}, P) \to \mathbb{R}$

assume $X$ takes finitely (countably) many values
$x_1, \ldots, x_n \in \mathbb{R}$,  $Y$ takes finitely (countably) many values
$y_1 \ldots y_m \in \mathbb{R}$, always with a positive probability.

$$E(Y \mid X = x_i) := \sum_{j=1}^{m} y_j \underbrace{P(Y = y_j \mid X = x_i)}_{\text{well defined}}$$

<u>Example</u> :  two dice, $X =$ first one, $Y =$ second one, independent

$$E(\text{sum} \mid X = 1) = \sum_{i=1}^{12} i \cdot P(\text{sum} = i \mid X = 1)$$

$$= \sum_{k=1}^{6} (1+k) \cdot P(Y = k \mid X = 1)$$

$$= \sum_{k=1}^{6} (1+k) \, P(Y = k) = \sum_{k=1}^{6} (1+k) \cdot \frac{1}{6} = 4.5$$

So far we defined $E(Y \mid X = x_i)$, but often we want to
consider the "function" $E(Y \mid X)_{(\omega)}$.   This is a rv:

$E(Y \mid X) : (\Omega, \mathcal{A}, P) \to (\mathbb{R}, \mathcal{B})$.

Leads to the following:

<u>Def</u> (discrete case) $X, Y$ as before. Then the conditional expectation is defined as follows:

$$E(Y|X) := f(X) \quad \text{with}$$

$$f(x) = \begin{cases} E(Y|X=x) & \text{if} \quad P(X=x) > 0 \\ \\ \text{arbitrary, say } 0 & \text{otherwise} \end{cases}$$
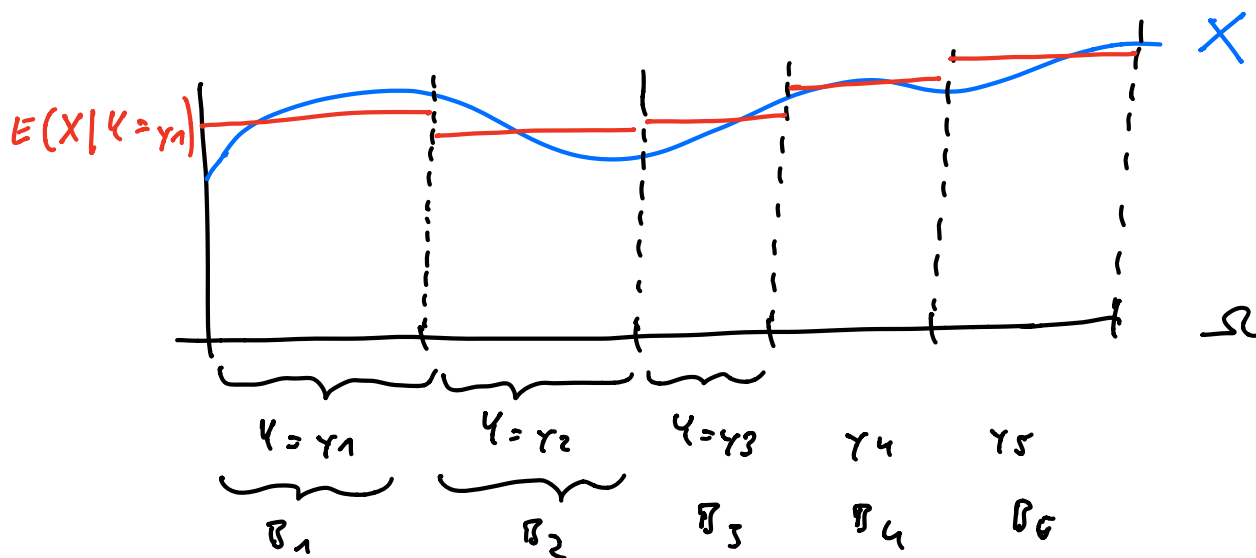
⚠ $E(Y|X)$ is only defined a.s.

Now we want to move to the more general case.

Sketch: $X$ continuous rv

$Y$ discrete rv $\leadsto Y_1, \ldots, Y_5$

Want to look at $E(X|Y)$



Want to "define" $E(X|Y) := \sum_{i=1}^{5} E(X|Y=Y_i) \cdot \mathbb{1}_{B_i}(\omega)$

But need to make sure that it is measurable wrt $\sigma(Y)$.
("the blur")

**Def** (Conditional expectation on $L_1$)

Consider rv $X: (\Omega, \mathcal{F}_0, P) \to \mathbb{R}$, $X \in L_1(\Omega, \mathcal{F}_0, P)$.

Let $\mathcal{F}$ be a sub-$\sigma$-algebra of $\mathcal{F}_0$. (Intuition: $\mathcal{F}_0$ will be the $\sigma$-alg. generated by the variable $Y$ we want to condition on).
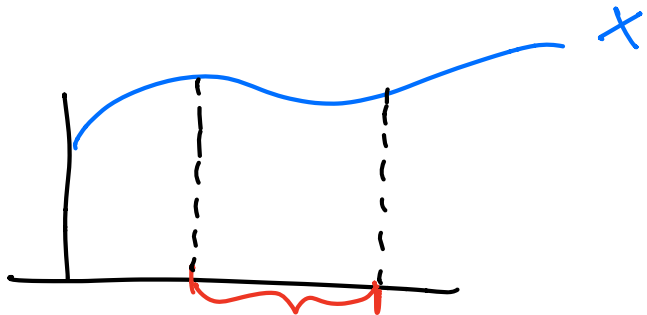
We now define the <u>cond. exp. of $X$ given $\mathcal{F}$</u>
$E(X | \mathcal{F})$ as any random variable $Z$ that satisfies

(1)   $Z$ is measurable w.r.t $\mathcal{F}$

(2)   For all $A \in \mathcal{F}$ we have

$$\int_A X \, dP = \int_A Z \, dP .$$



- Existence of $E(X | \mathcal{F})$ is not clear $^A$ a priori, it needs to be proved.

- $E(X | Y) := E(X | \sigma(Y))$

<u>Examples</u> ( extreme cases)

- $X = Y$. Then $E(X | Y) = X$   (a.s.)

- $X \perp\!\!\!\perp Y$.   $E(X | Y) = E(X)$   (a.s.)

# Case of joint densities

$X, Z: \Omega \to \mathbb{R}$ have a joint density $f(x,z)$.

Let $g: \mathbb{R} \to \mathbb{R}$ bounded, put $Y := g(Z)$. Assume we want to compute $E(Y|X) = E(\underbrace{g(Z)}_{Y}|X)$.

Recall $X$ has density $f_X(x) = \int f(x,z)\, dz$.

The conditional density of $Z$ given $X = x$ is

$$f_{X=x}(z) = \frac{f(x,z)}{f_X(x)} \qquad (\text{if } f_X(x) \neq 0)$$

Now consider $h(x) := \int \underbrace{g(z)}_{Y} f_{X=x}(z)\, dz$, now define

$$E(Y|X) = h(x).$$