## PT = INTEGRALS & EXPECTATIONS

Consider a discrete random variable $X$ taking values $x_1, x_2, \dots$ with probability $p_i = \mathbb{P}(X = x_i)$. Suppose we perform $n$ independent replications of a random experiment, denoting $X_j$ the value in the $j$-th replication.

Put $n_i = \#\{j \leq n \mid X_j = x_i\}$
$\qquad$ = number of times the value $x_i$ is observed during the first $n$ experiments.

The average value after $n$ trials is

$$\overline{X}_n := \frac{1}{n} \sum_{j=1}^{n} X_j = \frac{1}{n} \sum_{j=1}^{n} \sum_{i} x_i \, \mathbb{1}(X_j = x_i)$$

*partition on the possible values of $X_j$*

$$= \frac{1}{n} \sum_{i} x_i \underbrace{\sum_{j} \mathbb{1}(X_j = x_i)}_{= n_i}$$

$$= \sum_{i} x_i \boxed{\frac{n_i}{n}} \qquad \cdot = n_i$$

*Frequency interpretation of probability gives*

$$\frac{n_i}{n} \approx \mathbb{P}(X = x_i)$$

$$\approx \sum_{i} x_i \, \mathbb{P}(X = x_i)$$

You probably remember this expression from previous courses, and used it as a definition of the expected value of a discrete random variable $X$. Likewise, if $X$ is AC with

density $f$, the expected value of $X$ is taken as $\int x f(x) \, dx$. However, this is more a computational rule, rather than a definition of expectation.

However, how would you make sense of the expression $\int x f(x) \, dx$ if $X$ is defined over $(\Omega, \mathcal{F}, \mathbb{P})$, with $\Omega$ = space of continuous functions? [ provided we can construct a suitable measure $\mathbb{P}$ on it ]

$\Rightarrow$ We need a more general expression.

## I - EXPECTED VALUE OF A RV

### ① General Definition

- Start with <u>indicators</u>.

Consider $(\Omega, \mathcal{F}, \mathbb{P})$ = probability space

For $A \in \mathcal{F}$, put $X = \mathbb{1}_A = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$

How to define $\mathbb{E} X$?

Again, making use of frequency interpretation,

$$\overline{X}_n = \frac{1}{n} \sum_{j=1}^{n} X_j = \frac{n_A}{n} \longrightarrow \mathbb{P}(A) \quad \text{as } n \to \infty$$

*same notation as on page 1.*

Thus, put $\boxed{\mathbb{E} X = \mathbb{P}(A) \quad \text{for} \quad X = \mathbb{1}_A}$.

- Next, we want to generalize $\mathbb{E} X$ for <u>simple RVs</u>.

To do so, we are guided once again by the frequency interpretation of probability. In particular, we would like to keep the property of linearity:

$$n \text{ trials} \nearrow \quad x_1, \ldots, x_n \rightarrow \text{mean value } \bar{x} = \frac{1}{n} \Sigma x_i$$
$$\searrow \quad y_1, \ldots, y_n \rightarrow \text{mean value } \bar{y} = \frac{1}{n} \Sigma y_i$$

$$\downarrow \text{ take the sum } \downarrow$$

$$x_1 + y_1 \qquad x_n + y_n$$

$$\downarrow$$

$$\text{mean value}$$

$$\overline{x+y} = \frac{1}{n} \Sigma (x_i + y_i) = \frac{1}{n} \Sigma x_i + \frac{1}{n} \Sigma y_i = \bar{x} + \bar{y}$$

so it would be desirable to construct $\mathbb{E}$ such that $\mathbb{E}(X+Y) \approx \overline{x+y} = \bar{x} + \bar{y} \approx \mathbb{E}X + \mathbb{E}Y$

$\Rightarrow$ For a simple RV
$$X = \sum_{k=1}^{n} \alpha_k \mathbb{1}_{A_k}, \quad \text{put}$$
$$\mathbb{E}X = \sum_{k=1}^{n} \alpha_k \mathbb{P}(A_k)$$

*not necessarily a partition of $\Omega$*

In particular, for a constant RV $X(\omega) = \alpha \mathbb{1}_{\Omega}$, this definition ensures that $\mathbb{E}X = \alpha$. Good.

Indeed,
$$\mathbb{E}X = \mathbb{E}\left( \sum_{k=1}^{n} \alpha_k \mathbb{1}_{A_k} \right) = \sum_{k=1}^{n} \alpha_k \mathbb{E}\mathbb{1}_{A_k} = \sum_{k=1}^{n} \alpha_k \mathbb{P}(A_k)$$

*enforcing linearity for indicators*

*expected value of an indicator (page 2)*

$\hookrightarrow$ <u>Consequences</u>. Linearity of expectation for simple RVs follows.

Indeed, take
$$X = \Sigma \alpha_i \mathbb{1}_{A_i}$$
$$Y = \Sigma \beta_j \mathbb{1}_{B_j}$$
} partitions of $\Omega$

Then
$$\mathbb{E}(\alpha X + \beta Y) = \mathbb{E} \sum_{i,j} (\alpha \alpha_i + \beta \beta_j) \mathbb{1}_{A_i \cap B_j}$$

$\in \mathbb{R}$

$$= \sum_{i,j} (\alpha \alpha_i + \beta \beta_j) \mathbb{P}(A_i \cap B_j)$$

$$= \alpha \sum_{i} \alpha_i \boxed{\sum_{j} \mathbb{P}(A_i \cap B_j)} = \mathbb{P}(A_i)$$

$$+ \beta \sum_{j} \beta_j \boxed{\sum_{i} \mathbb{P}(A_i \cap B_j)} = \mathbb{P}(B_j)$$

$$= \alpha \sum_{i} \alpha_i \mathbb{P}(A_i) + \beta \sum_{j} \beta_j \mathbb{P}(B_j)$$

$$= \alpha \mathbb{E}X + \beta \mathbb{E}Y$$

• For a simple RV $X \geq 0$, we have $\mathbb{E}X \geq 0$.

• Monotonicity follows as well: if $X \leq Y$, for $X$ and $Y$ simple RVs, then $Y - X \geq 0$, and so
$$\mathbb{E}(Y-X) \geq 0$$
(linearity) $\Rightarrow \mathbb{E}Y - \mathbb{E}X \geq 0$
$$\mathbb{E}Y \geq \mathbb{E}X.$$

Summarizing, for simple random variables $X = \sum_{k=1}^{n} \alpha_k \mathbb{1}_{A_k}$, defining $\mathbb{E}X = \sum_{k=1}^{n} \alpha_k \mathbb{P}(A_k)$ yields
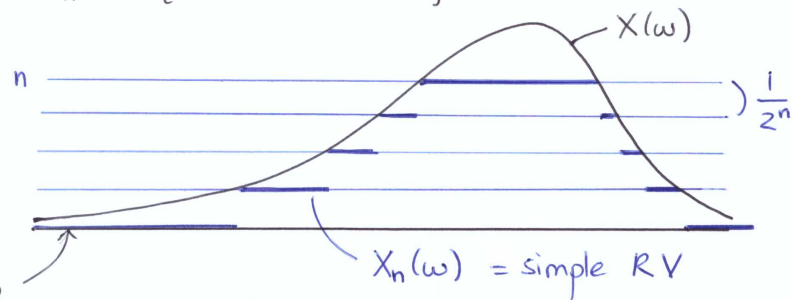
• <u>linearity</u> $\mathbb{E}(\alpha X + \beta Y) = \alpha \mathbb{E}X + \beta \mathbb{E}Y$
• <u>Constants</u> are expectations themselves $\mathbb{E}X = c$ for $X \equiv c$
• <u>Monotonicity</u>: $X \leq Y \Rightarrow \mathbb{E}X \leq \mathbb{E}Y$.

By the way, this definition is consistent. If

$$X = \sum_{k=1}^{n} \alpha_k \mathbb{1}_{A_k} = \sum_{k=1}^{n'} \alpha'_k \mathbb{1}_{A'_k}, \text{ then } \sum \alpha_k \mathbb{P}(A_k)$$
$$= \sum \alpha'_k \mathbb{P}(A'_k).$$

- Next, we consider general positive random variables $X \geqslant 0$. We know that there is a sequence of simple RVs $\{X_n\}$ such that

$$\forall \omega \in \Omega \quad X_n(\omega) \uparrow X(\omega) \text{ as } n \to \infty,$$

and this sequence can be constructed explicitly.

Set $A_{n,k} := \left\{ \omega \ \middle| \ \dfrac{k}{2^n} \leqslant X(\omega) < \dfrac{k+1}{2^n} \right\}$,

$$k = 0, 1, \ldots, n 2^n - 1$$

$$B_n := \{ \omega \mid X(\omega) \geqslant n \}$$



Put $X_n(\omega) = \begin{cases} 2^n/k & \text{for } \omega \in A_{n,k} \\ n & \text{for } \omega \in B_n \end{cases}$

Monotonicity of $\{X_n\}$ follows directly since

$$A_{n,k} = A_{n+1, 2k} \cup A_{n+1, 2k+1}, \quad k < n 2^n$$

$X_n$ and $X_{n+1}$ take the same value on this set

$X_{n+1}$ is larger than $X_n$ on this set (by $\frac{1}{2^{nn}}$)

Thus $X_n \leqslant X_{n+1}$.

---

Moreover, for a fixed $\omega$, $\qquad X_n(\omega) \leqslant X(\omega)$

for $X(\omega) < n \quad \to \quad X_n(\omega) \geqslant X(\omega) - 2^{-n}$

$$0 \leqslant X(\omega) - X_n(\omega) \leqslant 2^{-n}$$

Thus $\boxed{X_n(\omega) \uparrow X(\omega) \text{ as } n \to \infty.}$

↑ True for all $\omega$ such that $X(\omega) < \infty$, but if work for those $\omega$s for which $X(\omega) = \infty$ as well!

⟹ Consequence: Since $\{X_n\}$ is a sequence of simple RVs such that $X_1 \leqslant X_2 \leqslant X_3 \leqslant \ldots$, monotonicity implies that $\mathbb{E}X_1 \leqslant \mathbb{E}X_2 \leqslant \mathbb{E}X_3 \leqslant \ldots$
$$\underset{0}{\overset{\vee}{}}$$

⟹ the limit of this sequence exists, always! (but can be infinite).

Thus, $\boxed{\text{for a general } X \geqslant 0, \text{ we put } \mathbb{E}X := \lim_{n \to \infty} \mathbb{E}X_n}$

- Important remark: this is a consistent definition: the value of the limit does not depend on the choice of the sequence $\{X_n\}$. We only need that $X_n \uparrow X$.

Let's prove this.

⊗ Let $\{X_n\}$ and $\{\tilde{X}_n\}$ be two sequences of simple RVs such that $\forall \omega \in \Omega \quad X_n(\omega) \uparrow X(\omega)$, and $\tilde{X}_n(\omega) \uparrow X(\omega)$.

We want to show that $\lim_{n \to \infty} \mathbb{E}X_n = \lim_{n \to \infty} \mathbb{E}\tilde{X}_n$.

⊛ It suffices to prove that $\mathbb{E}\,\tilde{X}_k \leq \lim_{n\to\infty} \mathbb{E}\,X_n$ for any $k$, since this inequality implies that $\lim_{k\to\infty} \mathbb{E}\,\tilde{X}_k \leq \lim_{n\to\infty} \mathbb{E}\,X_n$. The inequality in the reverse direction follows by reversing the roles of $X_n$ and $\tilde{X}_n$. The two limits must then coincide.

⊛ Put $A_n := \{\omega \mid X_n(\omega) \geq \tilde{X}_k(\omega) - \varepsilon\}$, for some fixed value of $k$ and $\varepsilon > 0$. By definition of $A_n$, we have that $X_n \geq (\tilde{X}_k - \varepsilon)\,\mathbb{1}_{A_n}$

$$\mathbb{E}\,X_n \geq \mathbb{E}\{(\tilde{X}_k - \varepsilon)\,\mathbb{1}_{A_n}\} \quad \text{← monotonicity for simple RVs}$$
$$= \mathbb{E}\,\tilde{X}_k\,\mathbb{1}_{A_n} - \varepsilon\,\mathbb{E}\,\mathbb{1}_{A_n} \quad \text{← linearity}$$
$$= \mathbb{E}\,\tilde{X}_k(1 - \mathbb{1}_{A_n^c}) - \varepsilon\,\underbrace{\mathbb{P}(A_n)}_{\leq 1}$$
$$\geq \mathbb{E}\,\tilde{X}_k - \mathbb{E}\,\tilde{X}_k\,\mathbb{1}_{A_n^c} - \varepsilon$$
$$\geq \mathbb{E}\,\tilde{X}_k - \left[\max_{\omega\in\Omega} \tilde{X}_k(\omega)\right]\mathbb{P}(A_n^c) - \varepsilon$$

$\tilde{X}_k$ is a simple RV, so this value is $< \infty$.

Provided we show that $\mathbb{P}(A_n^c) \to 0$ as $n\to\infty$, we established that $\forall \varepsilon > 0$, $\lim_{n\to\infty} \mathbb{E}\,X_n \geq \mathbb{E}\,\tilde{X}_k - \varepsilon$

Since $\varepsilon$ is arbitrary, we must have $\lim_{n\to\infty} \mathbb{E}\,X_n \geq \mathbb{E}\,\tilde{X}_k$.

⊛ Proof that $\mathbb{P}(A_n^c) \to 0$ as $n\to\infty$.
Since $X_{n+1} \geq X_n$, we have $A_n \subset A_{n+1}$. Moreover, since $X_n \uparrow X \geq \tilde{X}_k > \tilde{X}_k - \varepsilon$, one has $A_n \uparrow \Omega$.
Then $\mathbb{P}(A_n) \to 1$ (page 13 Chp "Solid Foundations") & thus $\mathbb{P}(A_n^c) \to 0$

• For a general (not necessarily positive) RV $X$, write $X = X^+ - X^-$, where $X^+ = \max(0, X) \geq 0$, $X^- = -\min(0, X) \geq 0$

Note that in this notation, $|X| = X^+ + X^-$.

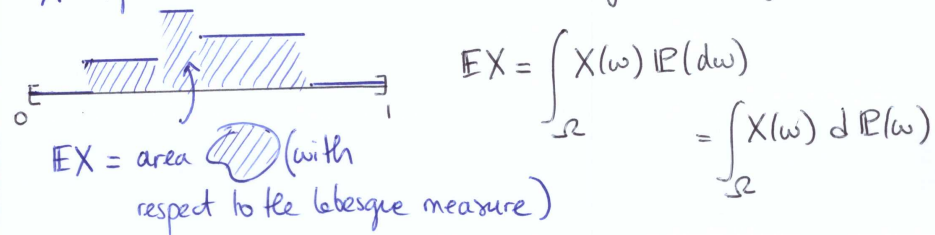A random variable $X$ is called __INTEGRABLE__ if $\mathbb{E}|X| < \infty$, that if both $\mathbb{E}\,X^+$ and $\mathbb{E}\,X^-$ are finite. Since both $X^+$ and $X^-$ are nonnegative, their expectation is well defined. Thus, for an integrable random variable $X$, put $\mathbb{E}\,X := \mathbb{E}\,X^+ - \mathbb{E}\,X^-$

If one of the $\mathbb{E}\,X^{+/-}$ is infinite, we can still make use of this definition, which will be $\pm\infty$ depending on which $\mathbb{E}\,X^{+/-}$ is infinite.

However, if both $\mathbb{E}\,X^{+/-} = \infty$, then $\mathbb{E}\,X$ is __undefined__. (what is $\infty - \infty$ ?)

This idea of approximating a function by piecewise constant functions should look very familiar (remember the Riemann integral ?). In fact, the above construction is nothing else than the __Lebesgue integral__ of $X$ with respect to the probability measure $\mathbb{P}$ defined on $(\Omega, \mathcal{F})$.

$X$ = simple RV on $[0,1]$ $\Rightarrow$ Notation of $\mathbb{E}$ using integrals:

$\mathbb{E}\,X = \int_\Omega X(\omega)\,\mathbb{P}(d\omega)$
$= \int_\Omega X(\omega)\,d\mathbb{P}(\omega)$

$\mathbb{E}\,X$ = area ⬭ (with respect to the Lebesgue measure)

$\Rightarrow$ $\mathbb{E}X$ inherits all properties of Lebesgue integrals, including

- **Monotonicity** : if $X \leq Y$ and $\mathbb{E}Y < \infty$, then $\mathbb{E}X \leq \mathbb{E}Y$

  First consider non-negative $X, Y$, and approximate them using simple RVs $X_n \uparrow X$ and $Y_n \uparrow Y$. Then $\mathbb{E}X_n \leq \mathbb{E}Y_n$ follows from $X_n \leq Y_n$. The result follows by letting $n \to \infty$.

  Next, drop the non-negativity assumption by considering separately $X^+, Y^+$ and $X^-, Y^-$.

- **Linearity** : if both $X$ and $Y$ are integrable, then for $a, b \in \mathbb{R}$, $\mathbb{E}(aX + bY) = a\mathbb{E}X + b\mathbb{E}Y$.

  First, note that $|aX + bY| \leq |a||X| + |b||Y|$, so that $\mathbb{E}|aX + bY| < \infty$ follows from $\mathbb{E}|X| < \infty$, and $\mathbb{E}|Y| < \infty$. The random variable $aX + bY$ is thus integrable.

  Next, proceed as before.

Remarks
- For $X$ integrable, we have that $|\mathbb{E}X| \leq \mathbb{E}|X|$.

  Indeed, since by definition $\mathbb{E}X = \mathbb{E}X^+ - \mathbb{E}X^-$, we have that
  $$|\mathbb{E}X| \leq |\mathbb{E}X^+| + |\mathbb{E}X^-| \quad \text{since } X^{+/-} \geq 0$$
  $$= \mathbb{E}X^+ + \mathbb{E}X^-$$
  $$= \mathbb{E}(X^+ + X^-) \quad \text{linearity}$$
  $$= \mathbb{E}|X|$$

- For $\mathbb{C}$-valued RVs, expectations are defined component-wise

---

For $Z = X + iY$, $X, Y \in \mathbb{R}$, put
$$\mathbb{E}Z := \mathbb{E}X + i\mathbb{E}Y$$
Same for random vectors, expectation is defined component wise.

- The integral $\mathbb{E}X = \int_\Omega X(\omega) \mathbb{P}(d\omega)$ make sense when integrating functions $X$ defined on more general spaces $\Omega$ than $\mathbb{R}$. For example, one may consider the space of continuous functions.

- A set $A \in \mathcal{F}$ is called a __null set__ (with respect to $\mathbb{P}$) if $\mathbb{P}(A) = 0$. We say that two random variables $X$ and $Y$ are equal almost surely (a.s.), and we write $X = Y$ a.s. if $\{\omega \in \Omega | X(\omega) \neq Y(\omega)\}$ is a null set.

  By construction of Lebesgue integral, it turns out that if $X = Y$ a.s, then $X$ is integrable if and only if $Y$ is integrable, and $\int X(\omega)\mathbb{P}(d\omega) = \int Y(\omega)\mathbb{P}(d\omega)$
  
  $\uparrow$
  *X and Y have the same expected value.*

  And this holds true as well for inequalities. Look: if $X \leq Y$ a.s. then $\mathbb{E}X \leq \mathbb{E}Y$
  $$\left( \int X(\omega)\mathbb{P}(d\omega) \leq \int Y(\omega)\mathbb{P}(d\omega) \right)$$

---

(i) Show that if $X \geq 0$ a.s and $\mathbb{E}X = 0$, then $X = 0$ a.s.

(ii) Show that if $\mathbb{E}X < \infty$, then $X < \infty$ a.s.

If instead of considering $(\Omega, \mathcal{F}, \mathbb{P})$ and $X: \Omega \to \mathbb{R}$ with induced measure $P_X$, we take $(\mathbb{R}, \mathcal{B}(\mathbb{R}), P_X)$ and set $X(x) = x$, using $x$ in place of the usual $\omega$, we see that

$$\int X(\omega)\, \mathbb{P}(d\omega) = \int x\, P_X(dx)$$

More generally, if $h = \underline{\text{measurable}}$ function, then $Y := h(X)$ has expected value given by

$$\mathbb{E}\,Y = \mathbb{E}\,h(X) = \int_{\Omega} h(X(\omega))\, \mathbb{P}(d\omega)$$
$$= \int_{\mathbb{R}} h(x)\, P_X(dx)$$

*as long as at least one of these two integrals make sense.*

This property is known as "Théorème du Transfert", in French.

To prove this theorem, the strategy is to start with indicator functions $h = \mathbb{1}_B$, $B \in \mathcal{B}(\mathbb{R})$, then simple functions, then positive measurable functions (need monotone convergence theorem), and finally measurable functions. We will not go through details. Good.

Consequences: • For a discrete random variable $X$, the law of $X$ is given by the discrete measure $P_X = \sum_{x} \mathbb{P}(X=x)\, \delta_x$

$X$ is integrable if and only if $\mathbb{E}|X| = \sum_{x} |x|\, \mathbb{P}(X=x) < \infty$.

In this case $\mathbb{E}\,X = \sum_{x} x\, \mathbb{P}(X=x)$, and we recover the formula on page **1**. Moreover, if $h$ is measurable, then

$h(X)$ is a discrete random variable; it is integrable if and only if $\sum_{x} |h(x)|\, \mathbb{P}(X=x) < \infty$. Its expected value is then $\mathbb{E}\big[h(X)\big] = \sum_{x} h(x)\, \mathbb{P}(X=x)$.

• If $X$ is AC with density $f$, the law of $X$ satisfies

$$\underline{P_X}(B) = \int_B f(x)\, dx \qquad \forall B \in \mathcal{B}(\mathbb{R}).$$

*To be more precise, AC $\equiv$ with respect to the Lebesgue measure.*

A consequence of the $\underline{\text{RADON-NIKODYM THEOREM}}$ (see below) yields $\mathbb{E}\,X = \int x\, f(x)\, dx$, provided $X$ is integrable, that is $\mathbb{E}|X| = \int |x|\, f(x)\, dx < \infty$.

Likewise, if $h = \mathbb{R} \to \mathbb{R}$ is measurable, $h(X)$ is integrable if $\int_{\mathbb{R}} |h(x)|\, f(x)\, dx < \infty$, and its expectation is $\mathbb{E}\big[h(X)\big] = \int_{\mathbb{R}} h(x)\, f(x)\, dx.$

$h(X)$ is not necessarily AC!

Let $X$ = AC random variable
$h$ = continuous function.
Construct $X$ and $h$ such that the random variable $h(X)$ is non-degenerate and discrete.

Remark: the Radon-Nikodym theorem. (RN)

Let $(\Omega, \mathcal{F})$ be a measurable space, endowed with two measures $\mathbb{P}$ and $Q$. We say that $\mathbb{P}$ is absolutely continuous (AC

with respect to $Q$ if $\forall A \in \mathcal{F}$, holds

$$Q(A) = 0 \Rightarrow \mathbb{P}(A) = 0,$$

and we write $\mathbb{P} \ll Q$

A null set for $Q$ is a null set for $\mathbb{P}$.

$\longrightarrow$ We have the following result (Radon-Nikodym).

If $\mathbb{P} \ll Q$, then there exists a measurable function $f : (\Omega, \mathcal{F}) \to [0, +\infty)$ such that $\forall A \in \mathcal{F}$,

$$\mathbb{P}(A) = \int_A f \, dQ$$

The function $f$ is called the Radon-Nikodym density of $\mathbb{P}$ with respect to $Q$. It is usually denoted

$$f = \frac{d\mathbb{P}}{dQ}.$$

Moreover, if $h$ is measurable, we have

$$\int_\Omega h \, d\mathbb{P} = \int_\Omega h f \, dQ$$ (as long as at least one of the two integrals make sense)

Not all distributions $P_X$ are Absolutely Continuous with respect to the Lebesgue measure $\lambda$ [defined such that $\lambda([a,b]) = b - a =$ length of the interval $[a,b]$, and extended to all Borel sets]

$\rightarrow$ No discrete distribution is AC with respect to $\lambda$ since $\mathbb{P}(X = x) > 0$ while $\lambda(\{x\}) = 0$.

$\rightarrow$ Those that are AC w.r.t. $\lambda$ are those with a density; that is the ones we encountered before.

Csq of RN theorem : $E h(X) = \int h(x) P_X(dx)$ $\Bigg\}$ $E h(X) =$

$P_X(B) = \int_B f(x) \lambda(dx)$ $\Bigg\}$ $\Rightarrow$ $\int h(x) f(x) \lambda(dx).$

---

Next, we give an alternative expression of the expected value of a non-negative random variable.

Theorem. Let $X \geq 0$ with distribution function $F_X$.
Then
$$\mathbb{E} X = \int_0^{+\infty} (1 - F_X(x)) \, dx$$
$$= \sum_{n \geq 1} n \, \mathbb{P}(X = n) = \sum_{n \geq 1} \mathbb{P}(X \geq n), \text{ if in addition}$$
$X$ is integer valued.

proof = For $X \geq 0$, one has $X = \int_0^X dx = \int_0^{+\infty} \mathbb{1}(X > x) \, dx$

Thus $\mathbb{E} X = \mathbb{E} \int_0^{+\infty} \mathbb{1}(X > x) \, dx$ $\rightarrow$ exchanging the order of integration (Tonelli)

$$= \int_0^{+\infty} \mathbb{E} \, \mathbb{1}(X > x) \, dx$$

$$= \int_0^{+\infty} \mathbb{P}(X > x) \, dx \xrightarrow[\text{case}]{\text{discrete}} \sum_{k=1}^{\infty} \int_{k-1}^{k} \mathbb{P}(X > x) \, dx$$

$$= \int_0^{+\infty} (1 - F_X(x)) \, dx$$

$$= \sum_{k \geq 1} \int_{k-1}^{k} \mathbb{P}(X \geq k) \, dx$$

$$= \sum_{k \geq 1} \mathbb{P}(X \geq k) \underbrace{\int_{k-1}^{k} dx}_{= 1}$$

$$= \sum_{k \geq 1} \mathbb{P}(X \geq k)$$

Theorem. If . $X_1$ and $X_2$ are independent RVs

. $g_1$ and $g_2$ such that $g_i(X_i)$ is integrable,

Then $\mathbb{E}\left[g_1(X_1)\, g_2(X_2)\right] = \mathbb{E}\left[g_1(X_1)\right]\,\mathbb{E}\left[g_2(X_2)\right]$

Sketch of proof: . First, consider $g_i(x) = \mathbb{1}(x \in B_i)$

$\longleftarrow \in B(\mathbb{R})$

Then $\mathbb{E}\left[g_1(X_1)\, g_2(X_2)\right] = \mathbb{E}\left[\mathbb{1}(X_1 \in B_1,\ X_2 \in B_2)\right]$

independence $\longrightarrow$

$= \mathbb{P}(X_1 \in B_1,\ X_2 \in B_2)$

$= \mathbb{P}(X_1 \in B_1)\, \mathbb{P}(X_2 \in B_2)$

$= \mathbb{E}\,\mathbb{1}(X_1 \in B_1)\ \mathbb{E}\,\mathbb{1}(X_2 \in B_2)$

$= \mathbb{E}\,g_1(X_1)\ \mathbb{E}\,g_2(X_2)$.

. Then, consider simple functions

. Use these to approximate general functions .

② Moments & Spaces $\mathscr{L}^P$.

Moments are special cases of $\mathbb{E}\,h(X)$, with $h(x) = x^P$, $p \geq 1$.

The p-th moment of $X$ is $\mathbb{E}\,X^P = \int x^P \, dF_X(x)$

The expected value of $X$ is the first moment.

provided it exists !

$\Rightarrow$ integrability condition apply !

$\longrightarrow$ For a RV $X$ and $p \geq 1$, let $\|X\|_P = \left(\mathbb{E}\,|X|^P\right)^{1/P}$.

A RV with $\|X\|_1 < \infty$ is called INTEGRABLE

$\|X\|_2 < \infty$ —"— SQUARE INTEGRABLE.

A random variable is called BOUNDED if there exists $K \in \mathbb{R}$ such that $|X| \leq K$ a.s. ; the quantity $\|X\|_\infty$ is by definition the smallest such $K$.

$\longrightarrow$ The spaces $\mathscr{L}^P := \{X : \Omega \to \mathbb{R} \mid \|X\|_p < \infty\}$ play a central role in functional analysis.

On $\mathscr{L}^P$, $\|.\|_p$ is almost a norm ( it is not a norm because $\|X\|_p = 0$ implies $X = 0$ a.s and not $X(\omega) = 0$ for all $\omega$ )

• $\mathscr{L}^P$ is almost a BANACH SPACE .

• $\mathscr{L}^2$ is almost a HILBERT SPACE .

Define $\langle X, Y \rangle_{\mathscr{L}^2} = \mathbb{E}(XY)$, for $X, Y \in \mathscr{L}^2$

Remark: $X, Y \in \mathscr{L}^2 \Rightarrow \mathbb{E}(XY) < \infty$ . Look:

$0 \leq (X \pm Y)^2 = X^2 + Y^2 \pm 2XY$

$\Rightarrow |XY| \leq \frac{1}{2}(X^2 + Y^2)$

Take $\mathbb{E}(\ldots)$ .

Result : $\langle \cdot, \cdot \rangle_{\mathscr{L}^2}$ is almost an inner product on $\mathscr{L}^2$.

(i) $\langle X + Y, Z \rangle_{\mathscr{L}^2} = \langle X, Z \rangle_{\mathscr{L}^2} + \langle Y, Z \rangle_{\mathscr{L}^2}$

(ii) $\langle \lambda X, Z \rangle_{\mathscr{L}^2} = \lambda \langle X, Z \rangle_{\mathscr{L}^2}$, $\lambda \in \mathbb{R}$

(iii) $\langle X, Z \rangle_{\mathscr{L}^2} = \langle Z, X \rangle_{\mathscr{L}^2}$

(iv) $\langle X, X \rangle_{\mathscr{L}^2} \geq 0$

However, $\langle X, X \rangle_{\mathscr{L}^2} = 0 \Rightarrow X = 0$ a.s.

In fact, we can modify the definition of $\mathscr{L}^2$ and regard two RVs in $\mathscr{L}^2$ as equal if they are equal $\mathbb{P}$-a.s, to make it a Hilbert space.
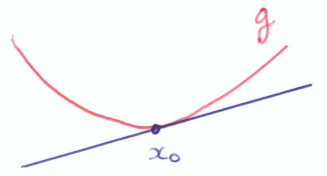
The induced norm is $\|X\|_2 = \sqrt{\langle X, X \rangle_{\mathscr{L}^2}} = \sqrt{\mathbb{E} X^2}$.

→ For $X \in \mathscr{L}^p$, the $p$-th **CENTRAL MOMENT** of $X$ is $\mathbb{E}(X - \mathbb{E}X)^p$. The **VARIANCE** is the second central moment $\text{Var}\, X = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \sigma_X^2$

$\sigma_X = $ **STANDARD DEVIATION** (same scale as $X$)

Next, we present two elementary inequalities:

> **Theorem** (JENSEN INEQUALITY)
> Let $X \in \mathscr{L}^1$ and $g : \mathbb{R} \to \mathbb{R}$ a convex function.
> Then $g(\mathbb{E}X) \leq \mathbb{E}(g(X))$

**proof**



For a convex function $g$,
$\forall x, x_0 \in \mathbb{R}$,

$$g(x) \geq g(x_0) + a(x - x_0),$$

where $a$ is the gradient of $g$ if $g$ is differentiable. If not, then $a$ is not unique.

Taking $x_0 = \mathbb{E}X$
$x = X$ , $g(X) \geq g(\mathbb{E}X) + a(X - \mathbb{E}X)$

$\Big($ Take $\mathbb{E}(\cdot)$ & use monotonicity + linearity.

$$\mathbb{E}[g(X)] \geq g(\mathbb{E}X)$$

> **Theorem** (CHEBYSHEV / MARKOV)
> If $g : \mathbb{R} \to \mathbb{R}$ is a positive non-decreasing function, then for any RV $X$ and $a \in \mathbb{R}$,
> $$\mathbb{P}(X > a) \leq \frac{\mathbb{E}[g(X)]}{g(a)}$$

**proof** : Since $g(a) \mathbb{1}(X \geq a) \leq g(X)$,

$$\mathbb{P}(X > a) = \mathbb{E}\,\mathbb{1}(X > a) \leq \mathbb{E}\left(\frac{g(X)}{g(a)}\right) = \frac{\mathbb{E}[g(X)]}{g(a)}$$

→ **Special cases:**

• $\mathbb{P}(|X| \geq a) \leq \dfrac{\mathbb{E}|X|^p}{a^p}$ for $p, a > 0$ $(X \in \mathscr{L}^p)$

  Take $g(x) = x^p$ and apply the inequality to $|X|$.

• $\mathbb{P}(|X - \mathbb{E}X| \geq a) \leq \dfrac{\text{Var}\, X}{a^2}$ for $a > 0$

  Apply the inequality to $|X - \mathbb{E}X|$ and with $g(x) = x^2$

  ↳ In particular, we obtain the "$3\sigma$ rule": the probability that a RV deviates from its mean by three $\sigma$ is small : Take $a = 3\sigma$,
  $$\mathbb{P}(|X - \mathbb{E}X| \geq 3\sigma) \leq 1/9 \quad \text{(crude band)}.$$

• $\mathbb{P}(X \geq a) \leq \dfrac{\mathbb{E}\, e^{tX}}{e^{ta}}$ for $t > 0$.

We present next further elementary results about the space $\mathscr{L}^p$.

**Proposition**

(i) $\mathcal{L}^p$ is a linear space: if $X, Y \in \mathcal{L}^p$, $\lambda \in \mathbb{R}$,
then $X + Y \in \mathcal{L}^p$ ; $\lambda X \in \mathcal{L}^p$

(ii) If $X \in \mathcal{L}^p$, $1 \leq q \leq p$, then $\|X\|_q \leq \|X\|_p$ ;
so that $\mathcal{L}^p \subset \mathcal{L}^q$    ↖ Lyapunov inequality

(iii) HÖLDER'S INEQUALITY

Let $p, q \geq 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$.
If $X \in \mathcal{L}^p$ and $Y \in \mathcal{L}^q$, then $|\mathbb{E} XY| \leq \|X\|_p \|Y\|_q$

(iv) MINKOWSKI'S INEQUALITY

If $X, Y \in \mathcal{L}^p$, then $\|X + Y\|_p \leq \|X\|_p + \|Y\|_p$

---

**proof** (i) We implicitly used this result on page 16 when
stating that $\mathcal{L}^p$ is a Banach space. It
follows from $|x + y|^p \leq \left[ 2 (\max(|x|, |y|)) \right]^p$
$$\leq 2^p (|x|^p + |y|^p).$$

(ii) This result implies that if the $p$-th moment is
finite, then the $q$-th moment is finite. In
particular, if the variance is finite, then the
mean is finite. (look back at the definition of
the $p$-th central moment on page 17).

It follows from Jensen's inequality with
$g(x) = x^{p/q}$ for $x \geq 0$ ; which is convex
for $p \geq q$. Putting $Y := |X|^q$,
$$g(\mathbb{E}Y) = (\mathbb{E}Y)^{\frac{p}{q}} \leq \mathbb{E}(g(Y)) = \mathbb{E}\left(Y^{\frac{p}{q}}\right)$$
$$\left(\mathbb{E}|X|^q\right)^{\frac{1}{q}} \leq \left(\mathbb{E}|X|^p\right)^{1/p}.$$

---

(iii) Making use of the convexity of exp, we
have that $\forall a, b$,
$$|ab| = \exp\left[ \frac{1}{p} \ln |a|^p + \frac{1}{q} \ln |b|^q \right]$$

sum to 1



$$\leq \left(\frac{1}{p}\right) \exp\left[ \ln |a|^p \right] + \left(\frac{1}{q}\right) \exp\left[ \ln |b|^q \right]$$

$\ln |a|^p \quad \ln |b|^q$

$$= \frac{1}{p} |a|^p + \frac{1}{q} |b|^q.$$

Applying this to $a = \dfrac{|X|}{\|X\|_p}$ ; $b = \dfrac{|Y|}{\|Y\|_q}$ , we obtain

$$\frac{|XY|}{\|X\|_p \|Y\|_q} \leq \frac{1}{p} \frac{|X|^p}{\|X\|_p^p} + \frac{1}{q} \frac{|Y|^q}{\|Y\|_q^q}$$

Taking $\mathbb{E}(\cdot)$

$$\frac{\mathbb{E}|XY|}{\|X\|_p \|Y\|_q} \leq \frac{1}{p} \left( \frac{\mathbb{E}|X|^p}{\|X\|_p^p} \right) + \frac{1}{q} \left( \frac{\mathbb{E}|Y|^q}{\|Y\|_q^q} \right)$$

$$= \frac{1}{p} + \frac{1}{q} = 1$$

(iv). Let $q^{-1} = 1 - p^{-1}$ $\left( \Leftrightarrow pq^{-1} = p - 1 \Leftrightarrow p = q(p-1) \right)$

• First, note that $|X + Y|^{p-1} \in \mathcal{L}^q$.
Indeed, $|X + Y|^p = |X + Y|^{p(q-1)}$
$\cap$
$\mathcal{L}^p$ thus integrable.

• $\mathbb{E}|X + Y|^p \leq \mathbb{E}|X||X+Y|^{p-1} + \mathbb{E}|Y||X+Y|^{p-1}$

Hölder ↳ $\leq (\|X\|_p + \|Y\|_p) \| |X+Y|^{p-1} \|_q$

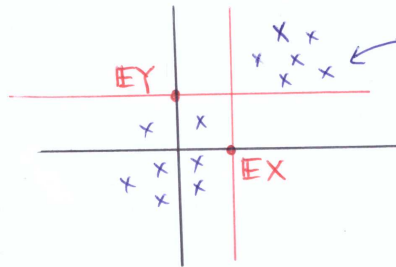$= (\|X\|_p + \|Y\|_p)(\mathbb{E}|X+Y|^p)^{\frac{1}{q}}$

- For $X, Y \in \mathcal{L}^2$, the <u>COVARIANCE</u> of $X$ and $Y$ is defined as $\text{Cov}(X, Y) = \mathbb{E}\big[(X - \mathbb{E}X)(Y - \mathbb{E}Y)\big]$
$$= \mathbb{E}(XY) - \mathbb{E}X\,\mathbb{E}Y$$

Remember, for $X, Y \in \mathcal{L}^2$, we have that $XY \in \mathcal{L}^1$



$\leftarrow$ If $X$ and $Y$ tend to be 'large' together or 'small' together, the covariance of $X$ and $Y$ is positive.

- The <u>CORRELATION</u> between $X$ and $Y$, provided $\text{Var}\,X$, $\text{Var}\,Y > 0$ is
$$\varrho := \text{Corr}(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}\,X\,\text{Var}\,Y}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

- <u>Properties of the correlation coefficient</u>

We made use on page 16 that for $X, Y \in \mathcal{L}^2$,
$$|XY| \leq \frac{1}{2}(X^2 + Y^2).$$

Applying this inequality with $X \to \dfrac{X}{\sqrt{\mathbb{E}X^2}}$

$Y \to \dfrac{Y}{\sqrt{\mathbb{E}Y^2}}$, we get

$$\left|\frac{XY}{\sqrt{\mathbb{E}X^2\,\mathbb{E}Y^2}}\right| \leq \frac{1}{2}\left(\frac{X^2}{\mathbb{E}X^2} + \frac{Y^2}{\mathbb{E}Y^2}\right)$$

Taking $\mathbb{E}(\cdot)$, this leads to the famous...

---

<u>CAUCHY - BUHYAKOVSKY INEQUALITY</u>

If $X, Y \in \mathcal{L}^2$, then $XY \in \mathcal{L}^1$, and $\mathbb{E}|XY| \leq \sqrt{\mathbb{E}X^2\,\mathbb{E}Y^2}$

$\uparrow$ This is just a special case of Hölder's inequality, by the way. Take $p = q = 2$.

Moreover, since $|\mathbb{E}XY| \leq \mathbb{E}|XY|$, (Jensen)
replace $X$ by $X - \mathbb{E}X$
$Y$ by $Y - \mathbb{E}Y$; & it follows from CB
inequality that $|\text{Cov}(X, Y)| \leq \sqrt{\text{Var}\,X\,\text{Var}\,Y}$,
and thus that $|\text{Corr}(X, Y)| \leq 1$

$\leftarrow$ When do we have an equality here?

$\hookrightarrow$ Assume that $\text{Corr}(X, Y) = +1$. Then, for the standardized RVs $X_1 = \dfrac{X - \mathbb{E}X}{\sqrt{\text{Var}\,X}}$ and

$Y_1 = \dfrac{Y - \mathbb{E}Y}{\sqrt{\text{Var}\,Y}}$, we have

$$\mathbb{E}(X_1 - Y_1)^2 = \underset{\underset{1}{\|}}{\mathbb{E}X_1^2} + \underset{\underset{1}{\|}}{\mathbb{E}Y_1^2} - 2\underbrace{\mathbb{E}(X_1 Y_1)}_{=\text{Corr}(X,Y)=1} = 0$$

Likewise, when $\text{Corr}(X, Y) = -1$, we have that
$\mathbb{E}(X_1 + Y_1)^2 = 0$

Either way, we have a RV $Z := (X_1 \pm Y_1)^2 \geq 0$ with $\mathbb{E}Z = 0$. We know from 📕 on page 10 that this implies that $Z = 0$ a.s. We just showed that:

$$\text{Corr}(X, Y) = \pm 1 \iff \mathbb{P}(X_1 \pm Y_1 = 0) = 1$$

Thus, with probability 1,

$$\frac{X - EX}{\sqrt{\text{Var } X}} \pm \frac{Y - EY}{\sqrt{\text{Var } Y}} = 0$$

$$\Rightarrow Y = aX + b$$

↑ Same sign as $\text{Corr}(X, Y)$.

$\Rightarrow$ Correlation between $X$ and $Y$ is $\pm 1$ where there is a perfect linear relationship between $X$ and $Y$.
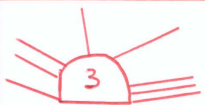
When $\text{Corr}(X, Y) = 0$, we say that $X$ and $Y$ are UNCORRELATED

⚠ Which is not the same as INDEPENDENCE.

If $X$ and $Y$ are independent, then $E(XY) = EX \, EY$, so that $\text{Cov}(X, Y) = 0$ and $X$ & $Y$ are uncorrelated. But the converse does not hold in full generality.

Summarizing

> Independence $\Rightarrow$ Uncorrelation
> ⇍

**3** Provide examples of uncorrelated random variables that are not independent

Remark: When dealing with random vectors $X = (X_1, \ldots, X_d)^t$, we use COVARIANCE MATRICES.

$$\Sigma := \begin{bmatrix} \text{Cov}(X_1, X_1) & \text{Cov}(X_1, X_2) & \cdots \\ \text{Cov}(X_2, X_1) & \text{Cov}(X_2, X_2) & \cdots \\ \vdots & \vdots & \end{bmatrix} = E\left[(X - EX)(X - EX)^t\right]$$

$(d \times d)$

---

Two key properties: [P1] $\Sigma$ is symmetric

[P2] $\Sigma$ is positive semi-definite:

$$\forall x \in \mathbb{R}^d, \quad x^t \Sigma x \geq 0$$

• [P1] is obvious

• To get [P2], Put $Y = x^t x \in \mathbb{R}$. Then

$$0 \leq \text{Var } Y = E(Y - EY)^2$$
$$= E(x^t x - E x^t x)^2$$
$$= E\left[(X - EX)^t x\right]^2$$
$$\quad \underset{\in \mathbb{R}}{}$$
$$= E\left[(X - EX)^t x \, (X - EX)^t x\right]^2$$
$$= E\left[\{(X - EX)^t x\}^t (X - EX)^t x\right]^2$$
$$= x^t E\left[(X - EX)(X - EX)^t\right] x$$
$$= x^t \Sigma x$$

In fact, any $(d \times d)$ matrix satisfying [P1]+[P2] is the covariance matrix of some distribution on $\mathbb{R}^d$ (why?)

• Geometrical Considerations.

Consider the space of centered square integrable random variables $\mathcal{L}_0^2 := \{X \in \mathcal{L}^2 \mid EX = 0\}$.

Let $X, Y \in \mathcal{L}_0^2$.

Then $\|X + Y\|_{\mathcal{L}^2}^2 = E(X + Y)^2$
$$= EX^2 + EY^2 + 2E(XY)$$

$$\Rightarrow \boxed{\|X + Y\|_{\mathcal{L}^2}^2 = \|X\|_{\mathcal{L}^2}^2 + \|Y\|_{\mathcal{L}^2}^2 + 2\langle X, Y \rangle.}$$

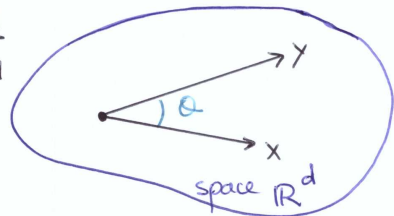$$\text{Var}(X + Y) = \text{Var } X + \text{Var } Y + 2\text{Cov}(X, Y)$$

Compare the last written equality with what happens in $\mathbb{R}^d$, endowed with the usual metric $\langle \cdot, \cdot \rangle_d$ (st. $\forall x, y \in \mathbb{R}^d$, $\langle x, y \rangle_d = x^t y$).

$$\|x+y\|_d^2 = \langle x+y, x+y \rangle_d$$
$$= \langle x, x \rangle_d + \langle y, y \rangle_d + 2\langle x, y \rangle_d$$

$$\Rightarrow \boxed{\|x+y\|_d^2 = \|x\|_d^2 + \|y\|_d^2 + 2\langle x, y \rangle_d .}$$

Moreover, the cosine of the angle between $x$ and $y$ is

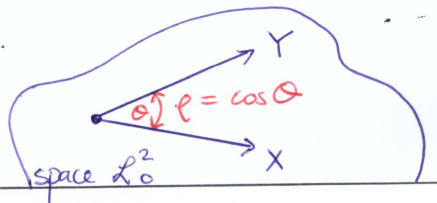$$\cos \theta = \frac{\langle x, y \rangle_d}{\|x\|_d \|y\|_d}$$


space $\mathbb{R}^d$

$\Rightarrow$ Therefore, for $X, Y \in \mathcal{L}_0^2$,

$$\rho = \text{Corr}(X, Y) = \frac{\text{Cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{\mathbb{E}(XY)}{\sigma_X \sigma_Y}$$

$$= \frac{\langle X, Y \rangle_d}{\|X\|_{\ell^2} \|Y\|_{\ell^2}} = \cos \theta$$

since
$$\sigma_X = \sqrt{\text{Var}\, X} = \sqrt{\mathbb{E}X^2} = \sqrt{\langle X, X \rangle_{\ell^2}} = \|X\|_{\ell^2}$$

$\Rightarrow$ The correlation coefficient represents the cosine of the angle between $X$ and $Y$ in the space $\mathcal{L}_0^2$.

In particular, when $\rho = 0$, $X$ and $Y$ are uncorrelated, $\theta = \frac{\pi}{2}$: $X$ and $Y$ are "perpendicular" and we write $X \perp Y$.


$\rho = \cos \theta$
space $\mathcal{L}_0^2$

---

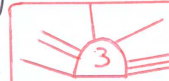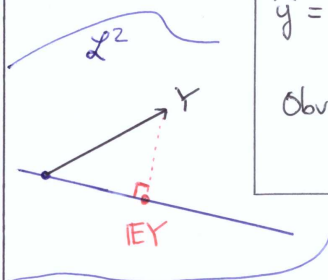# II - CONDITIONAL EXPECTATION.

## I) General definition.

Motivation : goal is to estimate the value of a variable of interest, denoted $Y$. If no further information other than observed values of $Y$ are provided, what is your best educated guess ?

$\hookrightarrow$ 'best' must be explicitly defined. We place ourselves in the space $\mathcal{L}^2$ of square integrable RVs, endowed with the inner product $\langle X, Y \rangle_{\mathcal{L}^2} = \mathbb{E}(XY)$. The induced metric is $d^2(X, Y) = \mathbb{E}[(X - Y)^2]$

• If we do not know anything about $Y$, you may want to estimate $Y$ using a constant value $\hat{y}$ such that $\hat{y}$ minimizes the distance $d^2(Y, y)$ :


$\mathcal{L}^2$
$Y$
$\mathbb{E}Y$

$$\hat{y} = \underset{y}{\text{argmin}}\ d^2(Y, y) = \underset{y}{\text{argmin}}\ \mathbb{E}[(Y - y)^2]$$

Obviously, $\hat{y} = \mathbb{E}Y$.


③
Compute $\tilde{y} = \underset{y}{\text{argmin}}\ \mathbb{E}|Y - y|$

• However, in a supervised learning context, we often know something about $Y$. For example, we may not know if a patient has some disease, but we might know the result of a medical test.

$\hookrightarrow$ Such extra information is commonly referred to as a 'predictor'; a 'feature'; or a 'covariate'.

- Suppose the only available extra piece of information is the answer to a 'yes/no' question : we know that some event $A$ occured ( for example, blood pressure is higher than some threshold ). What is your __best__ educated guess of $Y$ in this case.

 ↳ Proceed as before, and consider the minimization of $h(y) := \mathbb{E}\left[ (Y - y)^2 \mathbb{1}_A \right]$

 *Restrict the minimization on a subset of $\Omega$.*



*Intuitively, responds to the mean value of $Y$ after discarding all runs of the experiment where $A$ did not occur.*

$$h(y) = \mathbb{E}\left[ Y^2 \mathbb{1}_A \right] - 2y \mathbb{E}\left[ Y \mathbb{1}_A \right] + y^2 \mathbb{P}(A)$$

$$h'(y) = -2 \mathbb{E}\left[ Y \mathbb{1}_A \right] + 2\hat{y}\, \mathbb{P}(A) = 0$$

Gives $\hat{y} = \dfrac{\mathbb{E}\left[ Y \mathbb{1}_A \right]}{\mathbb{P}(A)}$

 __Notation:__ We use $\mathbb{E}(Y ; A)$ to denote $\mathbb{E}[Y \mathbb{1}_A]$

- We can do the same for $A^c$ instead of $A$. Our 'best' guess would be $\dfrac{\mathbb{E}[Y \mathbb{1}_{A^c}]}{\mathbb{P}(A^c)}$.

- In this simple situation, our predictor is the variable $X = \mathbb{1}_A$. Summarizing our findings, our best predictor of $Y$ given $X$, that we denote $\hat{Y}$, is

$$\hat{Y}(\omega) = \begin{cases} \dfrac{\mathbb{E}[Y ; A]}{\mathbb{P}(A)} & \text{if } \omega \in A \\[2mm] \dfrac{\mathbb{E}[Y ; A^c]}{\mathbb{P}(A^c)} & \text{if } \omega \notin A \end{cases}$$

*This object is a Random Variable !*

---

⇒ Rewriting $\hat{Y}$ slightly differently,

$$\hat{Y} = \begin{cases} \mathbb{E}[Y ; A] / \mathbb{P}(A) & \text{with probability } \mathbb{P}(A) \\[2mm] \mathbb{E}[Y ; A^c] / \mathbb{P}(A^c) & \text{—"— } \mathbb{P}(A^c) \end{cases}$$

The mean value of

$$\dots \quad \dfrac{\mathbb{E}[Y \mathbb{1}_A]}{\mathbb{P}(A)} \times \mathbb{P}(A) + \dfrac{\mathbb{E}[Y \mathbb{1}_{A^c}]}{\mathbb{P}(A^c)} \mathbb{P}(A^c) = \mathbb{E}(Y)$$

... equal to the expected value of $Y$. OK.

- Next, suppose that your predictor is a simple RV,
$X = \sum_{i=1}^n x_i \mathbb{1}_{A_i}$ , where $\{A_1, .., A_n\}$ is a partition of $\Omega$, and all $x_i$ are distinct

 ↳ If you observe $X = x_i$, then $\omega \in A_i$, and based on previous calculations, the 'best' forecast for $Y$ is



$$\hat{Y}(\omega) := \dfrac{\mathbb{E}[Y ; A_i]}{\mathbb{P}(A_i)} =: y_i \quad , \quad \omega \in A_i$$

 *← provided $\mathbb{P}(A_i) > 0$.*

 Thus

$$\hat{Y}(\omega) = \begin{cases} \dfrac{\mathbb{E}[Y ; A_1]}{\mathbb{P}(A_1)} = y_1, & \omega \in A_1 \\[2mm] \vdots & \\[2mm] \dfrac{\mathbb{E}[Y ; A_n]}{\mathbb{P}(A_n)} = y_n, & \omega \in A_n \end{cases}$$

*__Rk:__ If $\mathbb{P}(A_i) = 0$, define $\hat{Y}$ as you wish on $A_i$ ⇒ the resulting definition of $\hat{Y}$ is unique $\mathbb{P}$- almost-surely.*

*This object is a random variable !*

Since $A_i = \{\omega \mid X(\omega) = x_i\}$, we introduce (29)
a function $\varphi(x)$ by putting $\varphi(x_i) = y_i$. We
see that the new random variable $\hat{Y}$ is a function
of $X$ since in this notation, $\hat{Y} = \varphi(X)$.



$x_2$ ——— $X = $ simple RV
$y_2$ ——— $Y = $ constant on $A_i$
$x_3$ ———
$y_1$ ———
$x_1$ ———
$y_3$ ———
$A_1 \quad A_2 \quad A_3 \quad \leftarrow \quad \Omega$

Remarks: • In view of the expression of $\hat{Y}$, $\hat{Y}$ represents the
'conditional mean of $Y$ given $X$'. Instead of $\hat{Y}$,
we may use the notation $E(Y|X)$.

• The values of $X$ do not matter when defining
$E(Y|X)$. What matters is the partition
created by $X$. In other words, what matters
is $\sigma(X)$. Therefore, we may use the
notation $E(Y \mid \sigma(X))$ in place of $E(Y|X)$.

[If you are interested in the average weight
of inhabitants in a big city given their
postcode; does the postcode itself actually
matter, or the partition of the city it creates?]

• $Y$ and $\hat{Y}$ have the same average value over $A_i$:

$$E(\hat{Y}; A_i) = E(\hat{Y} \mathbb{1}_{A_i}) = E(y_i \mathbb{1}_{A_i})$$
$$= y_i \, \mathbb{P}(A_i)$$
$$= E(Y; A_i)$$

• $Y$ and $\hat{Y}$ also have the same average over (30)
arbitrary union of $A_i$; that is over any
element in $\sigma(X) = \{X^{-1}(B) \mid B \in \mathcal{B}(\mathbb{R})\}$.
Consider

$$X^{-1}(B) = \{X \in B\} = \bigcup_{k \mid x_k \in B} \{X = x_k\},$$

so that

$$\mathbb{1}_{\{X \in B\}} = \sum_{k \mid x_k \in B} \mathbb{1}_{A_k}.$$

Hence,

$$E(\hat{Y}; X \in B) = E(\hat{Y} \mathbb{1}_{\{X \in B\}})$$
$$= E\left(\hat{Y} \sum \mathbb{1}_{A_k}\right)$$
$$= \sum E(\hat{Y} \mathbb{1}_{A_k}) \quad \Bigg] \text{page } 29$$
$$= \sum E(Y \mathbb{1}_{A_k})$$
$$= E\left(Y \sum \mathbb{1}_{A_k}\right)$$
$$= E(Y \mathbb{1}_B)$$
$$= E(Y; X \in B), \text{ indeed.}$$

Picture =



$B$

$x_3$
$y_3$
$x_2$
$y_1$
$y_2$
$x_1$

original $Y$
$X = $ simple RV
$x_4$
$y_4$
$\hat{Y} = E(Y|X)$

$A_1 \quad A_2 \quad A_3 \quad A_4$

$X^{-1}(B)$
$\shortparallel$
$A_2 \cup A_3$

Average value of $Y$ over $A_2 \cup A_3$
$\shortparallel$
Average value of $\hat{Y}$ over $A_2 \cup A_3$

In summary, we see that the conditional expectation ③¹
$\hat{Y} = \mathbb{E}(Y \mid X)$ satisfies two fundamental properties.

[CE.1]  $\hat{Y}$ is a function of $X$ : $\hat{Y} = \varphi(X)$
[CE.2]  $\forall B \in \mathcal{B}(\mathbb{R})$,  $\mathbb{E}[\hat{Y}; X \in B] = \mathbb{E}[Y; X \in B]$
$\{ \forall A \in \sigma(X), \quad \mathbb{E}[\hat{Y}; A] = \mathbb{E}[Y; A] \}$

Remark: These two conditions uniquely specify $\hat{Y}$ in the case of simple RVs. Indeed, with $A = A_i$ in [CE.2], we have $\mathbb{E}[\hat{Y}; A_i] = \mathbb{E}[Y; A_i]$.
Using [CE.1], $\hat{Y} = \varphi(X)$, so that

$$\mathbb{E}[\varphi(X) \mathbb{1}_{A_i}] = \mathbb{E}[Y \mathbb{1}_{A_i}]$$
$$\| \qquad\qquad$$
$$\varphi(x_i) \mathbb{E}\, \mathbb{1}_{A_i} = \varphi(x_i) \mathbb{P}(A_i),$$

so that  $\varphi(x_i) = \dfrac{\mathbb{E}[Y \mathbb{1}_{A_i}]}{\mathbb{P}(A_i)}$ = value of $\hat{Y}$ on $A_i$

We formally define the Conditional Expectation (CE) for a general $X$ using [CE.1] and [CE.2]:

Theorem : Let $Y \in \mathcal{L}^1$ and $X$ be RVs defined on a common probability space.
Then there exists a unique random variable satisfying [CE.1] and [CE.2], and is called the CE of $Y$ given $X$ ; denoted $\mathbb{E}(Y \mid X)$. Moreover, $\mathbb{E}(Y \mid X) \in \mathcal{L}^1$.

Unique up to $\mathbb{P}$-null sets.

proof =
× Uniqueness = let $\hat{Y}'$ and $\hat{Y}''$ two RV in $\mathcal{L}^1$ such that $\forall A \in \sigma(X)$, $\mathbb{E}[\hat{Y}' \mathbb{1}_A] = \mathbb{E}[Y \mathbb{1}_A] = \mathbb{E}[\hat{Y}'' \mathbb{1}_A]$.
Take $A = \{\hat{Y}' > \hat{Y}''\} \in \sigma(X)$ since from [CE.1] both $\hat{Y}'$ and $\hat{Y}''$ are functions of $X$.

Then  $\mathbb{E}\Big[ \underbrace{(\hat{Y}' - \hat{Y}'') \mathbb{1}_{\{\hat{Y}' > \hat{Y}''\}}}_{\text{a positive random variable}} \Big] = 0$  ③²

$\Rightarrow$  $(\hat{Y}' - \hat{Y}'') \mathbb{1}_{\{\hat{Y}' > \hat{Y}''\}} = 0$  a.s.

Hence  $\hat{Y}' \leqslant \hat{Y}''$  a.s.
By symmetry, $\hat{Y}' \geqslant \hat{Y}''$ a.s, and we conclude that $\hat{Y}' = \hat{Y}''$ a.s.

× Existence = Suppose $Y \geqslant 0$.
Define a measure $Q$ on $(\Omega, \sigma(X))$ by
$\forall A \in \sigma(X)$,  $Q(A) = \mathbb{E}[X \mathbb{1}_A]$
(not a proba measure)

Since $\mathbb{P}$ is a probability measure defined on $(\Omega, \sigma(X))$, we see that $Q \ll \mathbb{P}$, by definition of $Q$ $\left( Q(A) = \displaystyle\int_A X(\omega)\, d\mathbb{P} \right)$.

$$\mathbb{P}(A) = 0 \Rightarrow Q(A) = 0$$

Radon-Nikodym theorem (page 13) ensures the existence of a measurable function
$\tilde{X} : (\Omega, \sigma(X)) \longrightarrow \mathbb{R}$ such that

$$Q(A) = \int_A \tilde{X}(\omega)\, d\mathbb{P} = \mathbb{E}[\tilde{X} \mathbb{1}_A]$$

Thus, $\tilde{X}$ is such that $\forall A \in \sigma(X)$, $\mathbb{E}[X \mathbb{1}_A]$
$\hookrightarrow \in \mathcal{L}^1$ : take $A = \Omega$
and thus verifies [CE.2]  $\overset{\|\sim}{\mathbb{E}[\tilde{X} \mathbb{1}_A]}$
· If $Y$ is not $\geqslant 0$, take $Y = Y^+ - Y^-$

Remark: Since $E(Y|X)$ does not depend on the value of $X$ but on the partition it creates, if $\varphi$ is a one-to-one function, then $E(Y|X) = E(Y|\varphi(X))$.

Ex: $E(Y|X) = E(Y|X^3) = E(Y|e^X) \cdots$

Example = Let $\quad X \sim \mathcal{P}(\lambda) \quad \mathbb{P}(X=x) = \frac{\lambda^x}{x!} e^{-\lambda}$
$\qquad\qquad\qquad Y \sim \mathcal{P}(\mu) \qquad\qquad\qquad x=0,1,\ldots$
$\qquad\qquad\qquad Z := X+Y \sim \mathcal{P}(\lambda+\mu)$

Show that $E(X|Z) = \frac{\lambda}{\lambda+\mu} Z$

[CE.1] is obvious since $\frac{\lambda}{\lambda+\mu}$ is a function of $Z$.

[CE.2] $E[X ; Z=k] = E[X \mathbb{1}(X+Y=k)]$

$\qquad = \sum_{i,j \geq 0} i \, \mathbb{1}(i+j=k) \, \mathbb{P}(X=i, Y=j)$
$\qquad\qquad\qquad \underbrace{\qquad\qquad}_{j=k-i \geq 0}$

$\qquad = \sum_{i=0}^{k} i \, \mathbb{P}(X=i, Y=k-i)$

$\qquad = \sum_{i=0}^{k} i \, \frac{\lambda^i}{i!} e^{-\lambda} \frac{\mu^{k-i}}{(k-i)!} e^{-\mu}$

$\qquad = e^{-(\lambda+\mu)} \sum_{i=1}^{k} \frac{i}{i!} \lambda^i \frac{\mu^{k-i}}{(k-i)!}$

$\qquad = e^{-(\lambda+\mu)} \lambda \sum_{i=1}^{k} \frac{1}{(i-1)!(k-i)!} \lambda^{i-1} \mu^{k-i}$

$\qquad = e^{-(\lambda+\mu)} \lambda \sum_{l=0}^{k-1} \frac{1}{l!(k-l-1)!} \lambda^{l} \mu^{k-l-1}$

$\boxed{l=i-1}$

$\qquad = e^{-(\lambda+\mu)} \frac{\lambda}{(k-1)!} \sum_{l=0}^{k-1} \frac{(k-1)!}{l!(k-l-1)!} \lambda^{l} \mu^{k-l-1}$
$\qquad\qquad\qquad\qquad\qquad \underbrace{\qquad\qquad}_{= \binom{k-1}{l}}$

$\qquad = e^{-(\lambda+\mu)} \frac{\lambda}{(k-1)!} (\lambda+\mu)^{k-1} \longleftarrow$ same!

On the other hand,

$E\left(\frac{\lambda}{\lambda+\mu} Z ; Z=k\right) = \frac{\lambda k}{\lambda+\mu} \mathbb{P}(Z=k) = \frac{\lambda k}{\lambda+\mu} \frac{(\lambda+\mu)^k}{k!} e^{-(\lambda+\mu)}$

$\longrightarrow$ What if the question was not to show that, but compute, not knowing the answer in advance?

General approach: find the UNDERLINED{CONDITIONAL DISTRIBUTION}, and then compute the expectation under it.

It remains to define conditional distribution of $Y$ given $X$.

① $\underline{X \text{ is discrete.}}$

Then we have already seen that $E(Y|X) = \varphi(X)$,

where $\varphi(x) = \dfrac{E[Y \mathbb{1}(X=x)]}{\mathbb{P}(X=x)}$

If $\mathbb{P}(X=x)=0$, $\varphi(x)$ is chosen arbitrarily.

$\qquad\qquad = \sum_{y} y \boxed{\dfrac{\mathbb{P}(X=x, Y=y)}{\mathbb{P}(X=x)}}$

define this guy as the conditional probability that $Y=y$ given $X=x$, and write

$\mathbb{P}(Y=y|X=x) := \dfrac{\mathbb{P}(X=x, Y=y)}{\mathbb{P}(X=x)}$

Compare with expression of expectation page 11.

First, we note that characterization [CE.2] can be generalized to

[CE.2'] For any random variable $Z$ which is $\sigma(X)$ – measurable, $E[\hat{Y} Z] = E[YZ]$

__Why?__ We know from [CE.2] that $\forall A \in \sigma(X)$,

$$E[\hat{Y} \mathbb{1}_A] = E[Y \mathbb{1}_A]$$

↖ instead $\mathbb{1}_A$, consider simple functions, then positive measurable functions, and finally any measurable function aka our random variable $Z$. We won't go into details.

So suppose that $(X, Y)$ has joint density $f(x, y)$, and $X$ has marginal distribution $f(x) = \int f(x, y) \, dy$.

Let $h : \mathbb{R} \to \mathbb{R}_+$ be a measurable function.

We compute $E[h(Y) | X]$ the following way =

$$E\left\{ E[h(Y) | X] \, z \right\} \overset{[CE.2']}{=} E\left\{ h(Y) \, g(X) \right\}$$

$Z$ is $\sigma(X)$ measurable
$\Rightarrow$ let's write it $g(X)$ for some measurable function $g$.

$$= \iint h(y) \, g(x) \, f(x, y) \, dx \, dy$$

$$= \int \left( \int h(y) f(x, y) \, dy \right) g(x) \, dx$$

$$= \int \left\{ \frac{\int h(y) f(x, y) \, dy}{f(x)} \right\} g(x) \, f(x) \, dx$$

↖ define this quantity as $\varphi(x)$

$$= \int \varphi(x) \, g(x) \, f(x) \, dx$$

$$= E\left[ \varphi(X) \, g(X) \right]$$

$\Rightarrow \forall$ measurable $h : \mathbb{R} \to \mathbb{R}_+$, we just showed that

$$E\left\{ E[h(Y) | X] \, z \right\} = E\left\{ \varphi(X) \, z \right\} \quad (\forall z)$$

We conclude that $\quad E[h(Y) | X] = \varphi(X)$

So that $\quad E[h(Y) | X = x] = \varphi(x)$

$$= \int h(y) \boxed{\frac{f(x, y)}{f(x)}} \, dy$$

define this quantity as the conditional density of $y$ given $x$, and write

$$f(y | x) = \frac{f(x, y)}{f(x)}.$$

In practice, to compute $E(Y | X)$, first compute $f(y | x)$, then compute $\varphi(x) = \int y \, f(y | x) \, dy$, and set $E(Y | X) = \varphi(x)$

[P1] Linearity. $\forall a, b \in \mathbb{R}$,

$$E(aX + bY \mid Z) = a E(X \mid Z) + b E(Y \mid Z)$$

We need to show that the Right Hand Side (RHS) satisfies [CE.1] and [CE.2].

[CE.1] is trivially satisfies, since a function of $Z$.

Next,

$$E[RHS \; ; \; Z \in B] = a E[E(X \mid Z) \; ; \; Z \in B]$$

$\underset{\text{linearity of } E(\cdot)}{\nearrow}$

$$\qquad\qquad\qquad + b E[E(Y \mid Z) \; ; \; Z \in B]$$

by definition of $\quad\nearrow\quad = a E[X \; ; \; Z \in B]$
$E(X \mid Z)$ and $E(Y \mid Z)$

$$\qquad\qquad\qquad + b E[Y \; ; \; Z \in B]$$

$$= E[aX + bY \; ; \; Z \in B]$$

so the RHS satisfies [CE.2]. 🔲

[P2] Monotonicity

$\quad$ If $X \leq Y$ a.s. then $E(X \mid Z) \leq E(Y \mid Z)$ a.s.

By contradiction, suppose the conclusion does not hold, so that

$$E(Y - X \mid Z) \underset{[P1]}{=} E(Y \mid Z) - E(X \mid Z) < 0$$

$\qquad\qquad\qquad\qquad\qquad$ with positive probability.

Put $\quad E(Y - X \mid Z) =: h(Z)$

[P1] + [P2] are analogous to the properties of $E(\cdot)$

---

Next,

$$\{w \in \Omega : E(Y - X \mid Z) < 0\} = \{w \in \Omega : h(Z) < 0\}$$
$$= \{w \in \Omega : Z \in B\}$$

$\qquad\qquad$ where we defined
$$B := \{x \mid h(x) < 0\}$$

By [CE.2],
$$E[h(Z) \; ; \; Z \in B] = E[Y - X \; ; \; Z \in B]$$

strictly negative on $B$ $\qquad$ positive by assumption

set of positive probability
$\Rightarrow$ contradiction 🔲

$\qquad\qquad$ measurable

[P3] If $\quad Y = g(Z)$ then $\quad E(XY \mid Z) = Y E(X \mid Z)$

The RHS is a function of $Z$ so [CE.1] is satisfied.

Regarding [CE.2], we consider first the case of indicators $Y = \mathbb{1}(Z \in C)$ for $C \in \mathcal{B}(\mathbb{R})$.

Then

$$\rightarrow E(XY \; ; \; Z \in B) = E(X \mathbb{1}(Z \in C) \mathbb{1}(Z \in B))$$

$$\rightarrow E(RHS \; ; \; Z \in B) = E(Y E(X \mid Z) \; ; \; Z \in B)$$

$$\underset{\substack{[CE.2] \\ \text{for } E(X \mid Z)}}{\searrow} = E(E(X \mid Z) \; ; \; Z \in B \times C)$$

$$= E(X \; ; \; Z \in B \times C),$$

which is the same as above.

Then move on to simple functions, their limits, etc. 🔲

**[P4]** If $X$ and $Y$ are independent, then
$$\mathbb{E}(Y|X) = \mathbb{E}(Y).$$

$\mathbb{E}(Y)$ is a constant, which is a trivial function of $X$, so that [CE.1] is verified. Next,

$$\mathbb{E}(Y \,;\, X \in B) = \mathbb{E}(Y \,\mathbb{1}(X \in B))$$

indpce

$$= \Big(\mathbb{E}\,Y\Big)\Big(\mathbb{E}\,\mathbb{1}(X \in B)\Big)$$

$$= \mathbb{E}\Big((\mathbb{E}\,Y)\,\mathbb{1}(X \in B)\Big)$$

$$= \mathbb{E}\Big(\mathbb{E}\,Y \,;\, X \in B\Big)$$

$$= \mathbb{E}\Big(\text{RHS} \,;\, X \in B\Big)$$

**[P5]** Double expectation law aka tower property.
$$\mathbb{E}\Big\{\,\underbrace{\mathbb{E}(Y|X_1, X_2)}_{\text{less crude}} \,|\, X_1\,\Big\} = \underbrace{\mathbb{E}(Y|X_1)}_{\text{more crude}}$$

Again, [CE.1] is immediate. For [CE.2],

$$\mathbb{E}\Big\{\mathbb{E}(Y|X_1, X_2) \,;\, X_1 \in B\Big\}$$

[CE.2] for $\mathbb{E}(Y|X_1,X_2)$
$$= \mathbb{E}\Big\{\mathbb{E}(Y|X_1, X_2) \,;\, (X_1, X_2) \in B \times \mathbb{R}\Big\}$$
$$= \mathbb{E}\Big\{\quad Y \quad \,;\, (X_1, X_2) \in B \times \mathbb{R}\Big\}$$

[CE.2] for $\mathbb{E}(Y|X_1)$
$$= \mathbb{E}\Big\{Y \,;\, X_1 \in B\Big\}$$
$$= \mathbb{E}\Big\{\mathbb{E}(Y|X_1) \,;\, X_1 \in B\Big\}$$
$$= \mathbb{E}\Big\{\text{RHS} \,;\, X_1 \in B\Big\}$$

In particular, taking $X_1 = \text{cst}$, $\mathbb{E}\Big\{\mathbb{E}(Y|X)\Big\} = \mathbb{E}\,Y$.

---

Property [P4] shows that if $Y$ is independent of $X$, then $\mathbb{E}(Y|X) = \mathbb{E}(Y)$. Show, however, that if for some $X$ and $Y$ such that $\mathbb{E}(Y|X) = \mathbb{E}(Y)$, this is not enough to conclude that $X$ and $Y$ are independent.

③ **Geometric insights.**

In the general definition of CE on page 31, the random variable $Y$ is assumed to belong to $\mathcal{L}^1$. If in addition it is square integrable, then it is possible to take a different route to define the notion of conditional expectation. In fact, for $Y \in \mathcal{L}^2$, we show next that $\hat{Y} = \mathbb{E}(Y|X)$ is the best forecast of $Y$ from $X$, under the metric $d^2(X, Y) = \mathbb{E}(Y - X)^2$.

Indeed, for a random variable $Z = \varphi(X)$ [your forecast should be a function of the predictor $X$]

$$d^2(Y, Z) = \mathbb{E}(Y - Z)^2 = \mathbb{E}\big((Y - \hat{Y}) + (\hat{Y} - Z)\big)^2$$

We are looking for the $Z = \varphi(X)$ which minimizes this quantity

$$= \mathbb{E}(Y - \hat{Y})^2$$
$$+ \mathbb{E}(\hat{Y} - Z)^2$$
$$+ 2\,\mathbb{E}(Y - \hat{Y})(\hat{Y} - Z)$$

$$\mathbb{E}(Y - \hat{Y})(\hat{Y} - Z) = \mathbb{E}\Big\{\mathbb{E}\big[(Y - \hat{Y})(\hat{Y} - Z)\,|\,X\big]\Big\}$$

[P5]

function of $X$
$\Rightarrow$ use [P3]

$$= \mathbb{E}\left\{ (\hat{Y}-Z)\, \boxed{\mathbb{E}[Y-\hat{Y}\mid X]} \right\}$$

(41)

$$\mathbb{E}(Y\mid X) - \mathbb{E}(\hat{Y}\mid X)$$
$$= \hat{Y} - \hat{Y} \qquad \uparrow \varphi(X)$$
$$= 0$$

Thus, $\mathbb{E}(Y-Z)^2 = \underbrace{\mathbb{E}(Y-\hat{Y})^2}_{\text{independent of } Z} + \underbrace{\mathbb{E}(\hat{Y}-Z)^2}_{\substack{\| \\ 0 \text{ if and only if} \\ Z=\hat{Y}}}$

We just showed that

$$\boxed{\hat{Y} = \mathbb{E}(Y\mid X) = \underset{Z=\varphi(X)}{\operatorname{argmin}}\left\{\mathbb{E}(Y-Z)^2\right\}}$$

↑ our best forecast is function of the covariate X

↑ Mean square error ≡ distance (square) between Y and its approximate value Z.

Picture:

$\mathcal{L}^2 =$ space of square int. functions

$Y \in \mathcal{L}^2$

$Y-\hat{Y} \perp \hat{Y}$ in $\mathcal{L}^2$ since

$\hat{Y} = \mathbb{E}(Y\mid X) = $ function of X

Subspace

$\mathcal{L}^2_X = \left\{ \varphi(X) \mid \mathbb{E}[\varphi(X)]^2 < \infty \right\}$

$\langle \hat{Y}, Y-\hat{Y} \rangle$
$= \mathbb{E}[\hat{Y}(Y-\hat{Y})]$  [P5]
$= \mathbb{E}\,\mathbb{E}[\hat{Y}(Y-\hat{Y})\mid X]$
$= \mathbb{E}[\hat{Y}\,\underbrace{\mathbb{E}(Y-\hat{Y}\mid X)}]$  [P3]

$\|$
$\mathbb{E}Y\mid X - \mathbb{E}\hat{Y}\mid X$
$\|$
$0$

Thus $\langle \hat{Y}, Y-\hat{Y}\rangle = 0$
i.e. $\hat{Y} \perp Y-\hat{Y}$.