

Conditional Expectation

- Random vectors
- Conditional Expectation - basic defn
- Conditional Expectation as an MMSE estimator
- Conditioning on a σ -field

Random Vectors

- Random vector \underline{X} of dimension n is a collection of n random variables $\underline{X} = (X_1, X_2, \dots, X_n)$

- CDF $F_{\underline{X}}(x_1, x_2, \dots, x_n) = \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n]$
↑ intersection of events

- If \underline{X} is discrete, then \underline{X} has a pmf
 $P_{\underline{X}}(x_1, x_2, \dots, x_n) = \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n]$

If \underline{X} is absolutely continuous, then it has a pdf $f_{\underline{X}}$
s.t. $F_{\underline{X}}(x_1, x_2, \dots, x_n) = \int_{-\infty}^{x_n} \int_{-\infty}^{x_{n-1}} \dots \int_{-\infty}^{x_1} f(z_1, z_2, \dots, z_n) dz_1 dz_2 \dots dz_n$

- For fn $g: \mathbb{R}^n \rightarrow \mathbb{R}$, its expectation is

$$\mathbb{E}[g(\underline{X})] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1, x_2, \dots, x_n) f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n$$

This inherits all the properties of $\mathbb{E}[\cdot]$ in one-dimension

- If $X_1 \perp\!\!\!\perp X_2 \perp\!\!\!\perp \dots \perp\!\!\!\perp X_n$ (mutually independent)

$$\text{then } F_{\underline{X}}(x_1, x_2, \dots, x_n) = F_{X_1}(x_1) F_{X_2}(x_2) \dots F_{X_n}(x_n)$$

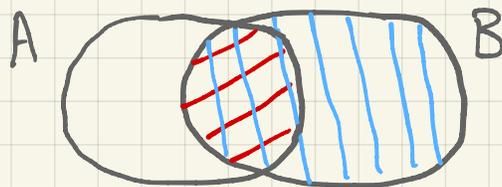
Basic Conditional Probability

(Revision of what you should have seen before)

- For $A, B \in \mathcal{F}$ st $P[B] > 0$,

$$P[A|B] = \frac{P[A \cap B]}{P[B]}$$

- Pictorially



- Similarly we can extend this to r.v.s conditioned on events

• For any r.v. X and event A

- conditional CDF $F_{X|A}(t) = P[X \leq t | A]$

- natural event - $A = \{Y = y\}$ for some r.v. Y

• For discrete r.v. X

$$P_{X|A}(t) = P[X=t | A], \quad E[X|A] = \sum_x x P_{X|A}(x)$$

• For continuous r.v. X, Y , and y st $f_Y(y) > 0$

$$f_{X|Y=y}(x) = \frac{f_{XY}(x,y)}{f_Y(y)}, \quad E[X|Y=y] = \int_{-\infty}^{\infty} x f_{X|Y=y}(x) dx$$

- Useful fact - for any r.v. X and event A

$$E[X \mathbb{1}_A] = E[X|A] P[A]$$

Conditional Expectation (via conditional probabilities)

- Now consider r.v. X and Y and let $\phi(y) \triangleq \mathbb{E}[X | Y=y]$.

Then the conditional expectation of X given Y is defined as

$$\mathbb{E}[X | Y] = \phi(Y)$$

Notes

i) $\mathbb{E}[X | Y]$ is a random variable!

ii) Sometimes denoted as $\mathbb{E}^Y[X]$ (see Brémaud)

Properties of $E[X|Y]$

We first look at some props of $E[X|Y]$, before trying to understand it in more detail.

$$\bullet \cdot E[\lambda_1 X_1 + \lambda_2 X_2 | Y] = \lambda_1 E[X_1 | Y] + \lambda_2 E[X_2 | Y]$$

(linearity)

$$- g_1(x) > g_2(x) \forall x \Rightarrow E[g_1(X) | Y] \geq E[g_2(X) | Y]$$

(monotonicity)

These follow from properties of $E[\cdot]$

Thm - $E[E[X|Y]] = E[X]$ (assuming $E[|X|] < \infty$)

Pf - $E[E[X|Y]] = \int_{-\infty}^{\infty} f_Y(y) E[X|Y=y] dy$ (tower rule)

$$= \int_{-\infty}^{\infty} f_Y(y) \left(\int_{-\infty}^{\infty} \frac{f_{XY}(x,y)}{f_Y(y)} x dx \right) dy$$

Fubini
(Assuming $E[|X|] < \infty$)

$$\Rightarrow \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} x f_{XY}(x,y) dy \right) dx$$

$$= \int_{-\infty}^{\infty} x f_X(x) dx = E[X]$$

Thm - $E[g(y) | Y] = g(y)$, and more

generally $E[g(y)h(x,y) | Y] = g(y)E[h(x,y) | Y]$
(pull-out property)

Pf - Again we will assume X, Y have a joint pdf f_{XY} . Now for any y st $f_Y(y) > 0$

$$\begin{aligned} E[g(y)h(x,y) | Y=y] &= \int_{-\infty}^{\infty} g(y)h(x,y) \frac{f_{XY}(x,y)}{f_Y(y)} dx \\ &= g(y) \int_{-\infty}^{\infty} h(x,y) \frac{f_{XY}(x,y)}{f_Y(y)} dx \\ &= g(y) E[h(x,y) | Y=y] \end{aligned}$$

Thus $E[g(y)h(x,y) | Y] = g(y)E[h(x,y) | Y]$

Thm - If $X \perp\!\!\!\perp Y$, $E[g(x) | Y] = E[g(x)]$
(independence & conditioning)

Pf - $E[g(x) | Y=y] = \int_{-\infty}^{\infty} g(x) \frac{f_{XY}(x,y)}{f_Y(y)} dx \quad \because X \perp\!\!\!\perp Y$

$$= \int_{-\infty}^{\infty} g(x) \frac{f_X(x)f_Y(y)}{f_Y(y)} dx = E[g(x)]$$

Conditional Expectation \equiv Estimation

- The best way to understand $E[X|Y]$ is in terms of estimation - In particular, suppose we have access to a random variable Y , and want to use it to approximate some other r.v. X as $\hat{X} = g(Y)$ for some fn g .

Claim - $g^*(Y) = E[X|Y]$ is the **MMSE** (minimum mean-squared error) approximation of X , i.e., it minimizes $E[(X - g(Y))^2]$ over all g s.t. $E[(g(Y))^2] < \infty$.

- This can actually be used to define $E[X|Y]$!
- You will see this in more detail in 6700; however, we will now see a brief proof of this.

Eg - For any r.v. X , suppose we want to approximate it by some constant $b \in \mathbb{R}$, such that we minimize the mean-squared error $E[(X-b)^2]$

Then we have

$$\begin{aligned}
 E[(X-b)^2] &= E\left[\underbrace{(X-E[X])}_{\text{Var}(X)} + (E[X]-b)\right]^2 \\
 &= E[(X-E[X])^2] + E[(E[X]-b)^2] \quad (\text{linearity of Expectation}) \\
 &= 2(E[X]-b) \underbrace{E[X-E[X]]}_{=0} + 2E[(X-E[X])(E[X]-b)] \\
 &= \text{Var}(X) + E[(b-E[X])^2]
 \end{aligned}$$

\therefore to minimize $E[(X-b)^2]$, we choose $b^* = E[X]$

• Now we can extend this to estimating X by $g(Y)$.

$$\text{We have } E[(X-g(Y))^2] = \int E[(X-g(Y))^2] f_Y(y) dy$$

$$\text{where } E[(X-g(Y))^2 | Y=y] = \int_{-\infty}^{\infty} (x-g(y)) \underbrace{f_{X|Y=y}(z)}_{= f_{XY}(z,y)/f_Y(y)} dx$$

From above, we know $E[(X-g(y))^2 | Y=y]$ is minimized

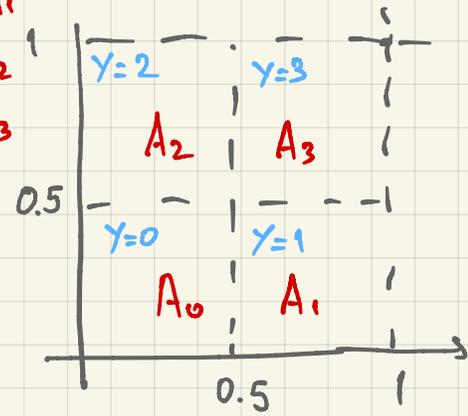
$$\text{by setting } g^*(y) = E[X|Y=y]$$

\Rightarrow The MMSE estimator $g^*(Y) = E[X|Y]$

Eg. 'Visualizing' $E[X|Y]$ 'distributed as'

Suppose $\Omega = [0, 1]^2$, $X = (X_1, X_2)$, $X_1, X_2 \sim \text{Unif}[0, 1]$
 $X_1 \perp\!\!\!\perp X_2$

and $Y = \begin{cases} 0 & ; \omega \in [0, 0.5] \times [0, 0.5] \equiv A_0 \\ 1 & ; \omega \in (0.5, 1] \times [0, 0.5] \equiv A_1 \\ 2 & ; \omega \in [0, 0.5] \times (0.5, 1] \equiv A_2 \\ 3 & ; \omega \in (0.5, 1] \times (0.5, 1] \equiv A_3 \end{cases}$
 (see fig on right)



- Recall we defined $\sigma(\mathcal{C})$ to be the smallest σ -field containing a given collection of sets \mathcal{C}

The σ -field $\sigma(\{A_0, A_1, A_2, A_3\})$ is referred to as the σ -field generated by Y , and denoted \mathcal{F}_Y , and

the conditional expectation $E[X|Y]$ can now be alternately written as $E[X|\mathcal{F}_Y]$ - this makes it clear that

$E[X|Y]$ is a function that associates a number with each set $A \in \mathcal{F}_Y$ (and such that the numbers obey Kolmogorov's axioms)

- In this case $E[X|Y=0] = \left(\int_0^{0.5} x_1 \cdot 2 dx_1, \int_0^{0.5} x_2 \cdot 2 dx_2 \right)$
 $= (0.25, 0.25)$

and $E[X|Y] = \begin{cases} (0.25, 0.25) & ; Y=0 \\ (0.75, 0.25) & ; Y=1 \\ (0.25, 0.75) & ; Y=2 \\ (0.75, 0.75) & ; Y=3 \end{cases}$

- Thus $E[X|Y] = E[X|\mathcal{F}_Y]$ essentially takes every set in \mathcal{F}_Y and associates the 'most likely' (in a mean-squared sense) number for X in that set.
 - You can think of this as a form of 'data compression' \equiv given some σ -field \mathcal{F}_Y (generated by Y), and a r.v. X , we 'smear' the information of X over \mathcal{F}_Y
-

• We next use this idea to give a more general definition of $E[X|Y]$, which covers the two definitions we have seen -

i) $E[X|Y] = g(Y)$, where $g(y) = E[X|Y=y]$

ii) $E[X|Y]$ is the (unique) fn $g(Y)$ with $E[|g(Y)|^2] < \infty$ which minimizes the mean-squared error $E[(X - g(Y))^2]$

Conditioning on a σ -field

We now see a more abstract defn of $E[X|Y]$ that generalizes the previous defns, and also the previous discussion.

- It is more general as it makes less assumptions (Note: for defn (i), we assumed X and Y have a pdf, for (ii), we needed $E[g(Y)]^2 < \infty$; in contrast we will now only need $E[|g(Y)|] < \infty$, which is weaker).
- It is more intuitive (even though more abstract!) once you get comfortable with the use of σ -fields
- It captures the idea of $E[X|Y]$ as a means of 'compressing information'.
- It will be important later when we talk about Markov chains & Martingales

We first need some defns. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a given probability space.

- i) A collection \mathcal{D} is a **sub σ -field** of \mathcal{F} if \mathcal{D} is a σ -field and $\mathcal{D} \subseteq \mathcal{F}$
- ii) A r.v. X is said to be **\mathcal{D} -measurable** or **adapted to \mathcal{D}** if $\{X \leq t\} \in \mathcal{D}$
- iii) For any collection of r.v. $\underline{Y} = \{Y_i; i \in I\}$, the **σ -field generated by \underline{Y}** (denoted as **$\sigma(\underline{Y})$ or $\mathcal{F}_{\underline{Y}}$**) is defined as the smallest sub- σ -field of \mathcal{F} containing all sets of the form $\{Y_i \leq t\}$, $i \in I$, $t \in \mathbb{R}$.
- iv) The σ -field $\mathcal{D} \triangleq \{\Omega, \emptyset\}$ is referred to as the **trivial σ -field**. The only r.v. which are measurable w.r.t \mathcal{D} are constants, i.e., $X(\omega) = c \quad \forall \omega \in \Omega$ (for some $c \in \mathbb{R}$)

Defn - Given prob space $(\Omega, \mathcal{F}, \mathbb{P})$, r.v.

X with $E[X] < \infty$, and sub σ -field $\mathcal{D} \subseteq \mathcal{F}$, the **conditional expectation of X given \mathcal{D}** , denoted $E[X|\mathcal{D}]$ is the (almost-sure) unique r.v. on $(\Omega, \mathcal{F}, \mathbb{P})$ s.t.

i) $E[X|\mathcal{D}]$ is \mathcal{D} -measurable

ii) $E[(X - E[X|\mathcal{D}]) \mathbb{1}_A] = 0 \quad \forall A \in \mathcal{D}$

Fact - \exists unique r.v. on $(\Omega, \mathcal{F}, \mathbb{P})$ satisfying the above
(See for example, Hajek Ch 10.1)

Note - $\because \Omega \in \mathcal{D}$ for any \mathcal{D} , we can set $A = \Omega$ in the above defn to get

$$E[(X - E[X|\mathcal{D}]) \mathbb{1}_\Omega] = E[X - E[X|\mathcal{D}]] = 0$$

$$\Rightarrow E[E[X|\mathcal{D}]] = E[X] \quad \forall \mathcal{D}$$

Indeed $E[X] = E[X|\{\Omega, \emptyset\}]$, i.e., the trivial σ -field.

• We can now re-state (and prove) properties of $E[X|\mathcal{D}]$. Below, we assume $E[X] < \infty \forall r.v.$

i) $E[ax+by|\mathcal{D}] = aE[X|\mathcal{D}] + bE[Y|\mathcal{D}]$ (linearity)

ii) If X is \mathcal{D} -measurable, then $E[g(X)|\mathcal{D}] = g(X)$
(more generally, $E[g(X)h(X,Y)|\mathcal{D}] = g(X)E[h(X,Y)|\mathcal{D}]$)
(pull-out property)

iii) If $\mathcal{A} \subset \mathcal{D} \subset \mathcal{F}$ (tower rule)
 $E[E[X|\mathcal{A}|\mathcal{D}]] = E[E[X|\mathcal{D}|\mathcal{A}]] = E[X|\mathcal{A}]$

• Note - The way to remember the tower rule is that if you condition \mathcal{A} on multiple σ -fields, then this is same as conditioning on the smallest (or coarsest) σ -field. This corresponds to the notion of 'conditioning as compression' - if you compress X to a coarse σ -field, then you cannot recover information!

• Pf of tower rule - Let $A \subset D \subset \mathcal{F}$. We want to show that $E[E[X|D]|A] = E[X|A]$. Note that both $E[X|A]$ and $E[E[X|D]|A]$ are A -measurable (by defn, since they are n.o. of the form $E[Y|A]$). Also note that we can write $X - E[E[X|D]|A] = (X - E[X|D]) - (E[E[X|D]|A] - E[X|D])$.

Now for any $A \in \mathcal{A}$, we have $A \in D$. By defn of $E[\cdot|A]$, we have $E[(X - E[X|D])\mathbb{1}_A] = 0$, $E[(E[E[X|D]|A] - E[X|D])\mathbb{1}_A] = 0$
 $\Rightarrow E[(X - E[E[X|D]|A])\mathbb{1}_A] = 0 \quad \forall A \in \mathcal{A}$

However, by the fact that $E[X|A]$ is a.s. unique, we must have $E[E[X|D]|A] = E[X|A]$

Note - Instead of defining $E[X|Y]$ in terms of $E[X|\mathcal{F}_Y]$ as above, we can directly define it as follows: given X, Y s.t. $E[|X|] < \infty$, then $E[X|Y]$ is the (unique) fn $g(Y)$ s.t. for every non-negative bounded fn ψ , we have $E[(X - g(Y))\psi(Y)] = 0$ a.s.

- See Brémaud Thm 2.3.15 for proof of existence & uniqueness