

Conditional expectation

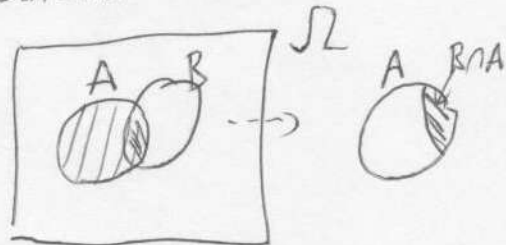
1

WARNING: This is HARD.

Introduction

Def: If $A, B \subset \Omega$ are events, $P(A) > 0$, then the conditional probability of B under the condition A is

$$P(B|A) = \frac{P(A \text{ and } B)}{P(A)}$$



[Philosophically: $\frac{\# \text{ favourable outcomes}}{\# \text{ all outcomes}}$]

Problem: What if $P(A) = 0$?

[Remark: A good question is why we are interested in that. We will come to that later.]

Bad news 1: If $P(A) = 0$, then $P(B|A) = \frac{P(A \cap B)}{P(A)}$ formally

makes no sense: division by zero.

Bad news 2: If $P(A) = 0$, then $P(B|A)$ REALLY

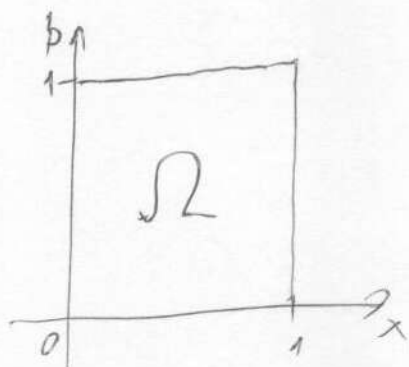
makes no sense: Not even intuitively!

Scary example:

2

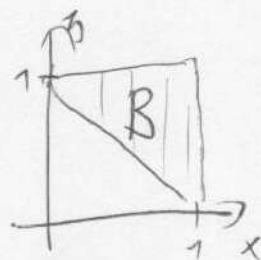
Consider the prob space $\Omega = [0,1]^2$ with $\mathbb{P} = \text{Leb}$,

so $\omega \in (x,y) \in \Omega$ is uniformly chosen from $[0,1]^2$

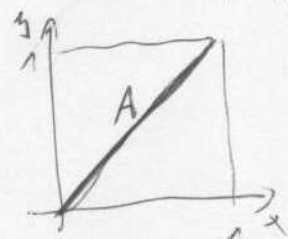


Now consider the events

$$B := \{(x,y) \in \Omega \mid x+y \geq 1\}$$



$$A := \{(x,y) \in \Omega \mid y=x\}$$



How can we give sense to $\mathbb{P}(B|A)$?

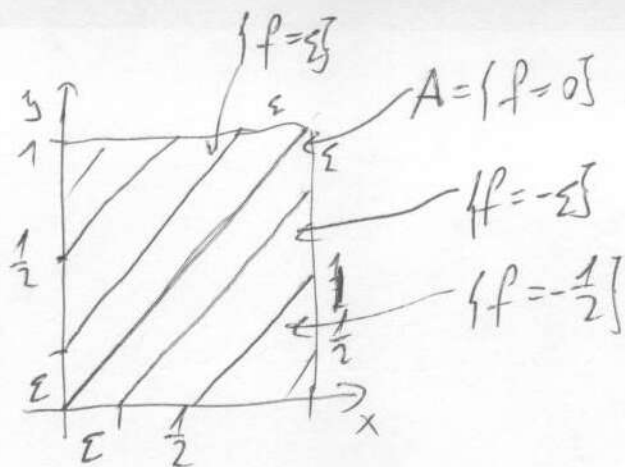
Idea 0: Approximate A with some A_ε such that $\mathbb{P}(A_\varepsilon) > 0$ (although small), so $\mathbb{P}(B|A_\varepsilon)$ makes sense, and $\approx \mathbb{P}(B|A)$ at least intuitively.

Q: How can we reasonably construct such an A_ε in a natural way?

Spoiler: I will give two equally natural constructions, leading to completely different results.

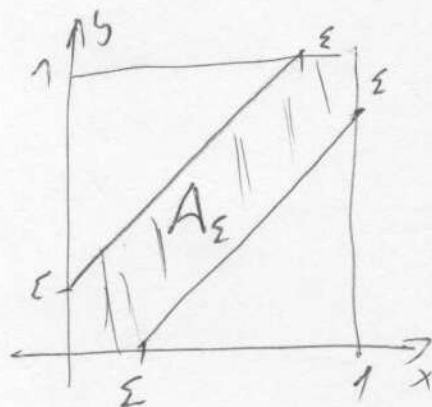
Idea 1: Let $f: \Omega \rightarrow \mathbb{R}$, $f(x,y) = y-x$.

So $A = \{(x,y) \in \Omega \mid f(x,y) = 0\}$ is a level set of f .
Other level sets are also line segments:



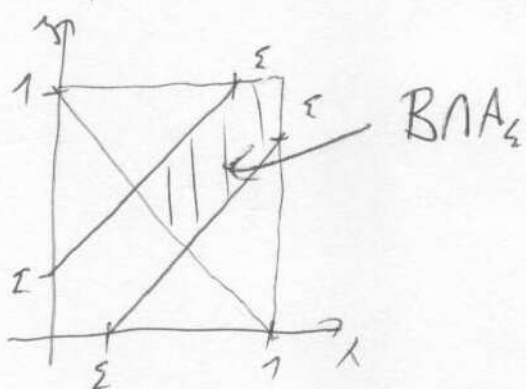
So let ~~$A := \{(x,y) \in \Omega \mid \varepsilon \leq f\}$~~ 3

$$A_\varepsilon := \{(x,y) \in \Omega \mid -\varepsilon \leq f(x,y) \leq \varepsilon\}$$



Easy to see:

$$P(B|A_\varepsilon) = \frac{1}{2} \text{ for } \forall \varepsilon > 0,$$



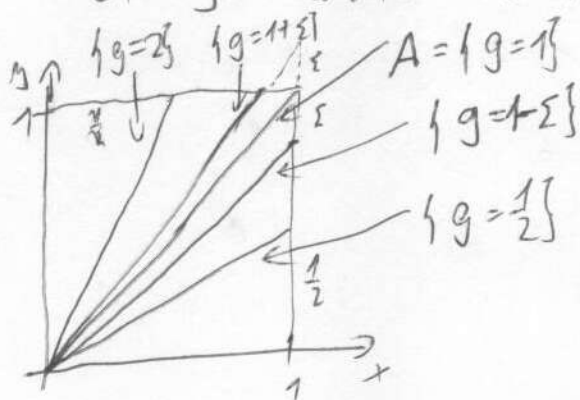
so

$$P(B|A) := \lim_{\varepsilon \rightarrow 0} P(B|A_\varepsilon) = \underline{\underline{\frac{1}{2}}}$$

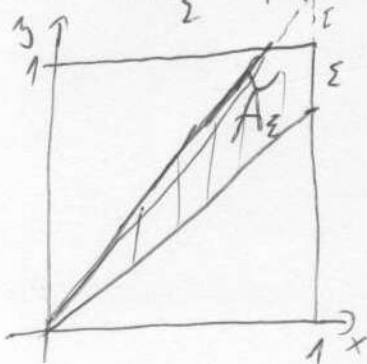
Idea 2: Let $g: \Omega \rightarrow \mathbb{R}$, $g(x,y) := \frac{y}{x}$ (defined \mathbb{P} -a.e.)

So $A = \{(x,y) \in \Omega \mid g(x,y) = 1\}$ is again a level set

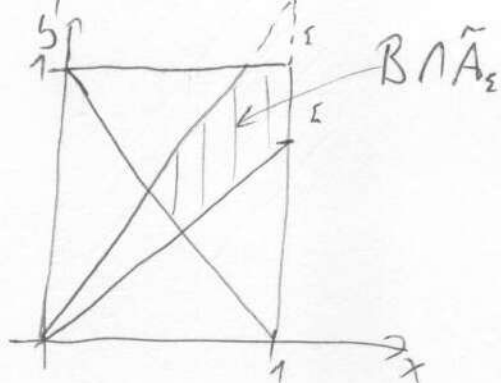
of g . Other level sets are also line segments:



So let $\tilde{A}_\varepsilon := \{(x,y) \in \Omega \mid 1-\varepsilon \leq g(x,y) \leq 1+\varepsilon\}$



Easy to see:



$$P(B|A_\epsilon) \approx \frac{3}{4}, \text{ so}$$

↑
 ≈ triangle
 cut in half

$$\tilde{P}(B|A) = \lim_{\epsilon \rightarrow 0} P(B|A_\epsilon) = \underline{\underline{\frac{3}{4}}}$$

Observation: The descriptions $\{f=0\} = A = \{g=1\}$

can be equally natural, depending on the context
 - whether A showed up while you studied $y-x$,
 or while you studied y/x .

~~Why we want to~~

Lesson to learn: For a single event $A \subset \Omega$ with $P(A) > 0$,
 $P(B|A)$ can not be made sensible.

So we will have to define $P(B|A_z^z)$ simultaneously
 for an entire family of events $A_z^z \subset \Omega$
 with $P(A_z^z) > 0$

[Say, in the above examples: $A_z^z := \{f=z\} \quad (z \in [-1, 1])$
 or $A_z^z := \{g=z\} \quad (z \in [0, 1])$]

Why we want to define $P(B|A)$ if $P(A)=0$

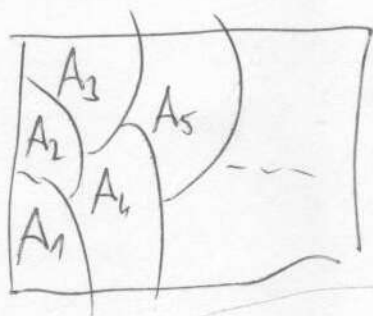
5

Q: What is conditional prob. good for?

A: It's ¹ always used to calculate total probabilities.
nearly

~~Thm~~ (~~Theorem of total~~ COUNTABLY MANY)

Def: The sequence of events $A_1, A_2, A_3, \dots \subset \Omega$ is called a partition if always exactly 1 of them occurs:



$$\bigcup_i A_i = \Omega \quad (\text{covering})$$

$$A_i \cap A_j = \emptyset \quad \forall i \neq j \quad (\text{pairwise disjoint})$$

Thm (Theorem of total probability) (countable)

Let $A_1, A_2, A_3, \dots \subset \Omega$ be a partition, and $B \subset \Omega$ an event

$$\text{Then } P(B) = \sum_i P(A_i) P(B|A_i).$$

T.T.P.

Remark Convention: $0 \cdot \text{undefined} = 0$, so

if $P(A_i) = 0$, then $P(A_i) P(B|A_i) \xrightarrow{\text{convention}} 0$.

This is natural for 2 reasons:

1.) $P(B|A_i)$ makes no sense, but if it would; then it would surely be $\in [0, 1]$ and $0 \cdot \text{anything in } [0, 1] = 0$.

2.) $P(A_i) P(B|A_i) \xrightarrow{\text{should be}} P(\underbrace{B \cap A_i}_{\subset A_i}) \leq P(A_i)$, so $= 0$ if $P(A_i) = 0$.

Conclusion: 1

It seems that we don't need to define $P(B|A_i)$

when $P(A_i) = 0$. it will/would be multiplied by 0 anyway.

Another tempting excuse:

$P(B|A)$ means: ~~what~~ what would be the chance of B occurring, if A would occur?

But who cares what would happen if A would occur?

If $P(A) = 0$, then A never occurs!

Problem with these:

1) This Theorem / Convention / Conclusion 1 only hold as the partition is countable!

Adding "more than countably many zeroes" is not so simple.

2) Zero probability events do occur all the time:

If X is a continuous random variable, then

$P(\{X = x\}) = 0$, but one of the (continuum many)

$\{X = x\}$ is sure to happen.

Conclusion 2:

If the theorem is the only important thing,
about conditional probability, then

Why not turn it into a definition?

axiomatic

Indeed: Our def. on conditional probability will essentially
be a generalization of the T.T.P. for uncountable
partitions

Idea of the construction

[WARNING again: this is a HARD construction/definition,
with many ingenious ideas.]

Idea 1: For given $B \subset \Omega$ and partition $\{A_i\}_{i \in I}$

we want to define all $P(B|A_i)$ simultaneously,

such that the T.T.P. holds:

$$P(B) = \sum_{i \in I} P(A_i) P(B|A_i).$$

However, this does not characterize $\{P(B|A_i)\}_{i \in I}$,

even when I is finite:

the equation $P(B) = \sum_{i \in I} P(A_i) x_i$ has many

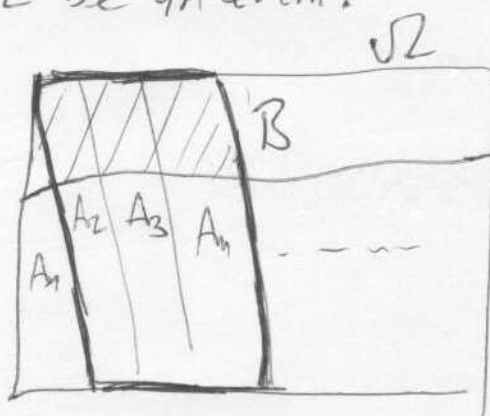
solutions: a single equation for many variables $\{x_i\}_{i \in I}$.

So consider the following easy corollary:

Thm (R.T.T.P. Extra): Let $\{A_i\}_{i \in I}$ be a countable partition. Let $B \subset \Omega$ be an event.

Furthermore, let $C \subset \Omega$

be an event like this: \rightarrow



$$C = \bigcup_{i \in J} A_i \text{ for some } J \subset I$$

or, equivalently,

$$C = A_2 \cup A_3 \cup A_4$$

for any $i \in I$ either $A_i \subset C$ or $A_i \cap C = \emptyset$.

$$\text{Then } P(B \cap C) = \sum_{i \in J} P(A_i) P(B|A_i)$$

[weight of the striped  region above]

This is what we will generalize.

Idea 2: Defining $P(B|A_i)$ simultaneously for all i means that the conditional probability is not a number, but a function:

Cond. prob: ~~A_i~~
 $\{ \text{partition elements} \} \rightarrow [0, 1]$
 $A_i \mapsto P(B|A_i)$

Q: How can we formulate such a function?

Answer 1: measure theorist's approach:

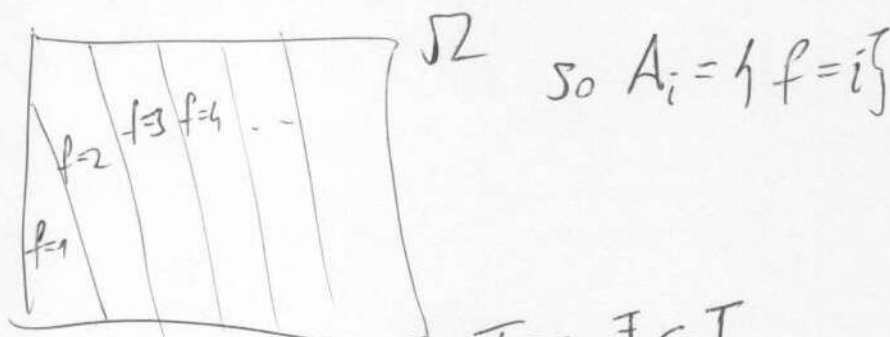
• Index your partition elements in your favourite way:

Say, indexing with $i \in I$ will do, and then

cond. prob. is a function $Y: I \rightarrow [0, 1]$
 $i \mapsto P(B|A_i)$

Equivalently: the A_i are level sets of the

function $f: \Omega \rightarrow I$ $f|_{\omega} = i$ for $\omega \in A_i$



Then the T.T.P. Extra reads: For $J \subset I$

$$P(B \mid f \in J) = \sum_{i \in J} P(A_i) Y(i).$$

This can then be generalized such that

$$\sum_{i \in J} P(A_i) Y(i) \text{ becomes } \int_J Y(i) d\mu(i)$$

↑
some measure
on I .

This is doable.

This is a standard construction
in measure theory, called conditional measure.

This is not what we will do.

Disadvantage: The indexing is quite arbitrary:

- ~~the~~ elements of the partition can be reordered
- equivalently: many functions $f: \Omega \rightarrow$ some set I
produce the same level sets
- brings in the extra object I , which has no real meaning.

Answer 2: Probabilist's approach

11

[This is what we will do] [slightly more general]
[than the previous]

Assigning numbers to partition elements $A_i \subset \Omega$

is ~~equivalent~~ to $Y: A_i \mapsto P(B|A_i)$

is equivalent to ~~def~~ giving a function Y

• defined on Ω , so $Y: \Omega \rightarrow [0, 1]$

• but constant on each A_i : $Y|_{A_i} = y_i = P(B|A_i)$

	A_2	A_3	...
A_1	$y = y_2$	$y = y_3$...
$y = y_1$			

Now Thm T.T.P. Extra reads: for $C = \bigcup_{i \in J} A_i$

$$P(B \cap C) = \sum_{i \in J} P(A_i) y_i = \sum_{i \in J} \int_{A_i} Y(\omega) dP(\omega) = \int_C Y dP$$

\uparrow
 $y = y_i$ on $A_i \subset \Omega$

Moreover: **GREAT OBSERVATION:**

Let $\mathcal{G} = \sigma(\{A_i\}_{i \in J})$ be the σ -algebra generated by the partition

12

Then • "Y is constant on each A_i " means Y is measurable w.r.t. \mathcal{G}

• " $C = \bigcup_{i \in I} A_i$ for some $I \subset I$ " means $C \in \mathcal{G}$

So Thm T.T.P. Extra reads:

Thm (Theorem of Total Probability Extra, equivalent 2nd version)

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a prob. space.

Let $\{A_i\}_{i \in I}$ be a countable partition of Ω . (into events $A_i \in \mathcal{F}$)

Let $\mathcal{G} = \sigma(\{A_i\}_{i \in I}) \subset \mathcal{F}$ be the sub- σ -algebra generated

Let $B \in \mathcal{F}$ be any event. ~~by the partition~~ by the partition.

Define $Y: \Omega \rightarrow [0, 1]$ as $Y(\omega) = \mathbb{P}(B|A_i)$ when $\omega \in A_i$

[Notice: The good old $\mathbb{P}(B|A_i) = \frac{\mathbb{P}(A_i \cap B)}{\mathbb{P}(A_i)}$ is well defined whenever $\mathbb{P}(A_i) > 0$. Since the partition is countable (by assumption), this $Y(\omega)$ is well defined for \mathbb{P} -almost every ω .]

Then • Y is \mathcal{G} -measurable, and

• for any $C \in \mathcal{G}$ $\mathbb{P}(B \cap C) = \int_C Y d\mathbb{P}$

This statement can be generalized to define the conditional probability for non-countable partitions.

This is ALMOST what we will do.

Last idea (p. quite obvious):

13

~~It's no extra diff~~ Without extra difficulty, we can immediately define the conditional probability expectation.
Then the cond. prob. will be a special case:

$$P(B|A) = E(\mathbb{1}_B | A).$$

Def: Let (Ω, \mathcal{F}, P) be a prob. space, $X: \Omega \rightarrow \mathbb{R}$ measurable, and assume that $E X$ exists. Let $A \in \mathcal{F}$, $P(A) > 0$.

Then the conditional expectation of X under the condition

A is $E(X|A) := \frac{\int_A X dP}{P(A)}$ [weighted average of values of X on A .]

Thm (Theorem of Total expectation) (T.T.E.)

Let (Ω, \mathcal{F}, P) be a prob. space, $X: \Omega \rightarrow \mathbb{R}$ measurable st. $E X$ exists.

Let $\{A_i\}_{i \in I}$ be a countable partition of Ω .

$$\text{Then } E X = \sum_{i \in I} P(A_i) E(X|A_i)$$

[with the convention 0-undefined := 0]

Proof: trivial from the definition.

As before, a trivial (equivalent) generalization:

Thm (T.T.E Extra): Under the same assumptions,

14

if $C = \bigcup_{i \in I} A_i$ for some $I \subset \mathcal{I}$, then

$$\int_C X dP = \sum_{i \in I} P(A_i) E(X|A_i).$$

And, just as before, $\sum_{i \in I} P(A_i) \cdot (\dots) \rightsquigarrow \int_C (\dots) dP$

Thm (Theorem of total expectation Extra, equivalent 2nd version)

Let (Ω, \mathcal{F}, P) be a prob. space.

Let $\{A_i\}_{i \in I}$ be a countable partition of Ω .

Let $\mathcal{G} = \sigma(\{A_i\}_{i \in I})$ be the generated σ -algebra.

Let $X: \Omega \rightarrow \mathbb{R}$ and assume that $E X$ exists.

Define $Y: \Omega \rightarrow [a_0, \infty]$ as $Y(\omega) := E(X|A_i)$ when $\omega \in A_i$.

[Again: this is well-defined for P -a.e. ω .]

Then \bullet Y is \mathcal{G} -measurable, and

$$\bullet \text{ for any } C \in \mathcal{G} \quad \int_C X dP = \int_C Y dP$$

Notice!! If $\int_C X dP = \int_C Y dP$ would hold for every $C \in \mathcal{F}$,

that would mean $Y = X$ ~~a.e.~~ almost surely.

But we only claim it for $\boxed{C \in \mathcal{G}}$!!

We are (finally) ready to define the conditional expectation ☺☺☺

Def: Conditional expectation with respect to a σ -algebra

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a prob. space, $\mathcal{G} \subset \mathcal{F}$ a sub- σ -algebra. Let $X: \Omega \rightarrow \mathbb{R}$ integrable

The function $Y: \Omega \rightarrow \mathbb{R}$ is called the conditional expectation of X w.r.t. \mathcal{G} , if

1) Y is \mathcal{G} -measurable, and

2) for any $A \in \mathcal{G}$ $\int_A X d\mathbb{P} = \int_A Y d\mathbb{P}$.

Notation: $Y = \mathbb{E}(X | \mathcal{G})$

Remark: Instead of "the conditional expectation" it would be fair to say "a version of the conditional expectation", since we haven't (yet) checked uniqueness.

Remark': If $Y = Z$ a.s., then $\int_A Y d\mathbb{P} = \int_A Z d\mathbb{P}$, so $\mathbb{E}(X | \mathcal{G})$ can only be unique up to modification on \mathbb{Q} -measure sets (respecting \mathcal{G} -measurability)

Example:

Let $\Omega = \{a, b, c\}$, $X = \mathbb{1}_{\{a, b\}}$, so $X(a) = 1$

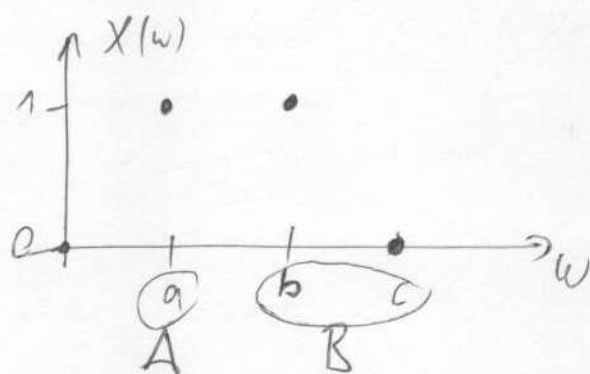
$P = \frac{1}{3} \mathcal{X}$: uniform prob. $X(b) = 1$

$P(\{a\}) = P(\{b\}) = P(\{c\}) = \frac{1}{3}$ $X(c) = 0$

$\mathcal{F} = 2^{\Omega}$, discrete σ -algebra.

$\mathcal{G} = \{\emptyset, \{a\}, \{b, c\}, \Omega\}$: this is generated by

the partition $\{\{a\}, \{b, c\}\}$



Set $A = \{a\}$

$B = \{b, c\}$

Then $Y = \Omega \rightarrow \mathbb{R}$ defined as

$$Y|_A := 1 \quad Y|_B := \frac{1}{2}$$

so $Y(a) = 1$
 $Y(b) = \frac{1}{2}$
 $Y(c) = \frac{1}{2}$

will do for $Y = E(X|\mathcal{G})$:

- $\int_{\Omega} X dP = 0 = \int_{\Omega} Y dP$ ✓
- $\int_A X dP = P(\{a\}) X(a) = \frac{1}{3} \cdot 1 = P(\{a\}) Y(a) = \int_A Y dP$ ✓
- $\int_B X dP = \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 0$; $\int_B Y dP = \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot \frac{1}{2}$ EQUAL ✓
- $\int_{\Omega} X dP = \frac{2}{3} = \int_{\Omega} Y dP$ ✓

Of course: we chose

$$\begin{aligned} Y|_A &:= \text{average of } X \text{ on } A = \{a\} \quad (\text{boring}) \\ Y|_B &:= \text{average of } Y \text{ on } B = \{b, d\} = \frac{1+0}{2} \end{aligned} \left. \begin{array}{l} \text{good old} \\ \text{conditional} \\ \text{expectation.} \end{array} \right\}$$

Thm (Theorem of total expectation extra, 3rd equivalent version)

Let (Ω, \mathcal{F}, P) be a prob. space.

Let $\{A_i\}_{i \in I}$ be a countable partition of Ω .

Let $\mathcal{G} = \sigma(\{A_i\}_{i \in I}) \subset \mathcal{F}$ be the generated sub- σ -algebra.

Let $X: \Omega \rightarrow \mathbb{R}$ integrable

Define $Y: \Omega \rightarrow \mathbb{R}$ as $Y(\omega) := E(X|A_i)$ when $\omega \in A_i$

[Where $E(X|A_i)$ is the good old $E(X|A_i) := \frac{\int_{A_i} X dP}{P(A_i)}$]
 [This is well defined a.s.]

Then $E(X|\mathcal{G}) = Y$.

So, indeed, $E(X|\mathcal{G})$ is the generalization of the old notion.

very Special case: If $A \in \mathcal{F}$, $P(A) \neq 0$ and

$\mathcal{G} = \{\emptyset, A, A^c, \Omega\}$, then $E(X|\mathcal{G})(\omega) = \begin{cases} E(X|A) & \text{if } \omega \in A \\ E(X|A^c) & \text{if } \omega \in A^c \end{cases}$

Thm: The conditional expectation is unique up to changes on a set of measure 0:

If $Y: \Omega \rightarrow \mathbb{R}$ and $Z: \Omega \rightarrow \mathbb{R}$ are both versions of $E(X|G)$, then $P(Z=Y) = 1$.

Proof: Y, Z are both G -measurable, so

$Y-Z$ is also G -measurable, so

$A := \{Y-Z > 0\} \in G$, so by definition of the conditional expectation

$$\int_A Y dP = \int_A X dP = \int_A Z dP \implies \int_A Y-Z dP = 0$$

$\underbrace{\hspace{10em}}_{> 0 \text{ on } A}$

$$P(Y > Z) = P(A) = 0$$

Similarly $P(Y < Z) = 0 \quad \square$

Existence of $E(X|G)$

Thm: The conditional expectation exists:

For any (Ω, \mathcal{F}, P) prob-space, $G \subset \mathcal{F}$, $X: \Omega \rightarrow \mathbb{R}$ integrable

$\exists Y: \Omega \rightarrow \mathbb{R}$ such that $Y = E(X|G)$

$$\left[\begin{array}{l} \text{meaning } Y \in \mathcal{G} \text{ (} G\text{-measurable) } \\ \text{and } \int_A X dP = \int_A Y dP \quad \forall A \in G \end{array} \right].$$

Proof: Assume 1st that $X \geq 0$.

19

Recall the Raden-Nikodym theorem: (with unusual notation)

~~If (Ω, \mathcal{G}) is a meas~~

If (Ω, \mathcal{G}) is a measurable space,

$\mu, \nu: \mathcal{G} \rightarrow \mathbb{R}$ σ -finite measures

and $\nu \ll \mu$,

then $\exists Y: \Omega \rightarrow \mathbb{R}^+$ \mathcal{G} -measurable

Such that $\nu(A) = \int_A Y d\mu$ for all $A \in \mathcal{G}$.

Of course Y is \mathcal{G} -measurable;

Of course $A \in \mathcal{G}$:

there's no other σ -algebra in the assumptions.

But this measurability is of key importance for us now.

[The first time during the course]

So consider the measurable space (Ω, \mathcal{G}) with the measures μ, ν on this space: $\mu, \nu: \mathcal{G} \rightarrow \mathbb{R}^+$

$\mu(A) := \mathbb{P}(A)$ for all $A \in \mathcal{G}$ - so $\mu = \mathbb{P}|_{\mathcal{G}}$

$\nu(A) := \int_A X d\mathbb{P}$ for $A \in \mathcal{G}$.

Lemma 1 These are both finite measures on (Ω, \mathcal{G}) .

Proof! easy HW - use the def. of the measure and basic properties of the integral.

In particular, $\mu(\Omega) = P(\Omega) = 1$

$$V(\Omega) = \int_{\Omega} X dP = EX < \infty.$$

Lemma 2 $V \ll \mu$.

Proof: ~~trivial~~ obvious: If $\mu(A) = P(A) = 0$, then

$$V(A) = \int_A X dP = 0 \quad \checkmark$$

So applying the Radon-Nikody theorem gives:

$\exists Y: \Omega \rightarrow \mathbb{R}$ \mathcal{G} -measurable s.t. $\forall A \in \mathcal{G}$ $\int_A Y d\mu = \int_A X dP$

$$V(A) = \int_A X dP$$

$$\int_A X dP$$

Lemma 3: For $A \in \mathcal{G} \subset \mathcal{F}$, $\mu = P|_{\mathcal{G}}$

$$\int_A Y d\mu = \int_A Y dP$$

integral on (Ω, \mathcal{G})

integral on (Ω, \mathcal{F})

[Formally they are not the same, so we would need to check that they are equal - DON'T DO IT!]

We are done when $X \geq 0$.

In the general case, write $X = X_+ - X_-$

\uparrow positive part \downarrow negative part,

so $X_+ \geq 0$, $X_- \geq 0$,

so ~~\mathbb{E}~~ $Y_+ := \mathbb{E}(X_+ | \mathcal{G})$ and $Y_- := \mathbb{E}(X_- | \mathcal{G})$ exist,

so $Y := Y_+ - Y_-$ will do (easy). □

Def. (Conditional probability, w.r.t. a σ -algebra)

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a prob. space, $A \in \mathcal{F}$, $\mathcal{G} \subset \mathcal{F}$ a sub-

σ -algebra. Then $\mathbb{P}(A | \mathcal{G}) \stackrel{\text{def}}{=} \mathbb{E}(\mathbb{1}_A | \mathcal{G})$

Properties of conditional expectation

1.) If $\mathcal{G} = \mathcal{F}$, then $\mathbb{E}(X | \mathcal{G}) = X$.

2.) More generally: If X is \mathcal{G} -measurable, then $\mathbb{E}(X | \mathcal{G}) = X$

Proof: obvious - check the definition.

3.) If $\mathcal{G} = \{\emptyset, \Omega\}$ (the indiscrete σ -algebra), then

$\mathbb{E}(X | \mathcal{G}) \equiv \mathbb{E}X$ (constant function)

Proof: check the def.

4.) If $a, b \in \mathbb{R}$ ~~then~~, X, Y are integrable then

$$\mathbb{E}(aX + bY | \mathcal{G}) = a\mathbb{E}(X | \mathcal{G}) + b\mathbb{E}(Y | \mathcal{G}) \quad : \text{linearity}$$

Proof: check the def.

5.) Tower rule: Let $\mathcal{G}_2 \subset \mathcal{G}_1 \subset \mathcal{F}$, so

\mathcal{G}_2 is a smaller / rougher σ -algebra than \mathcal{G}_1 .

Then $\mathbb{E}(\mathbb{E}(X | \mathcal{G}_1) | \mathcal{G}_2) = \mathbb{E}(X | \mathcal{G}_2)$,

but also $\mathbb{E}(\mathbb{E}(X | \mathcal{G}_2) | \mathcal{G}_1) = \mathbb{E}(X | \mathcal{G}_2)$

No typo:
the smaller /
rougher

σ -algebra wins.

Phenomenon: ~~a~~ smaller σ -algebra = less information.