

PROBABILITY 2017 HANDOUT 3: CONDITIONING

GIOVANNI PISTONE

CONTENTS

1. Conditional expectation	1
2. Conditional distribution	4
References	5

Conditioning is one among the core concepts in reasoning about uncertainty in Probability, in Statistics, in Economics, in Machine Learning. In this notes we refer mainly to the textbook by D. Williams [2, Ch. 9]. A concise and fully rigorous review of the basic mathematics is in the monograph by C. Dellacherie and P.-A. Meyer [1, Ch. I-III].

1. CONDITIONAL EXPECTATION

1 (Definition). Let $(\Omega, \mathcal{F}, \mu)$ be a probability space, X a real random variable with finite expectation, $E_\mu(|X|) < +\infty$, \mathcal{G} a sub- σ -algebra of \mathcal{F} . A random variable \hat{X} is a *version of the conditional expectation of X given \mathcal{G}* if, and only if,

- (1) \hat{X} is integrable and \mathcal{G} -measurable;
- (2) for all bounded and \mathcal{G} -measurable random variable it holds

$$E_\mu(G\hat{X}) = E_\mu(GX) .$$

As the equation $E_\mu(G(\hat{X} - X)) = 0$, $G \in \mathcal{L}^\infty(\mathcal{G})$, is linear in G and continuous under bounded pointwise convergence, it is enough to check it for random variables of the form $\mathbf{1}_C$, $C \in \mathcal{C}$, \mathcal{C} π -system generating \mathcal{G} . [Monotone-Class Theorem [2, ¶3.14].]

2 (Almost sure equivalence). If \hat{X}_1, \hat{X}_2 , are two versions of the conditional expectation of X , then $E_\mu(G(\hat{X}_1 - \hat{X}_2)) = 0$ i.e. $\hat{X}_1 = \hat{X}_2$ μ -almost-surely. [Take $G = \text{sign}(\hat{X}_1 - \hat{X}_2)$ to get $E_\mu(|\hat{X}_1 - \hat{X}_2|) = 0$.] More generally, if $X_1 = X_2$ μ -almost-surely, then $\hat{X}_1 = \hat{X}_2$ μ -almost-surely. We write $E_\mu(X|\mathcal{G})$ to denote the μ -class of versions and, with abuse of notation, $\hat{X} = E_\mu(X|\mathcal{G})$. If $L^1(\mathcal{F}, \mu)$ is the vector space of classes μ -equivalent real random variables, there exists a mapping

$$L^1(\mathcal{F}, \mu) \ni X \mapsto E_\mu(X|\mathcal{G}) \in L^1(\mathcal{G}, \mu) .$$

3. The *existence issue* i.e., the fact that the previous mapping is actually defined on all of $L^1(\mathcal{F}, \mu)$, is discussed in [2, ¶9.5]. We skip this topic, together with a related issue namely, the notion of μ -complete σ -algebra. Many proofs of existence are actually available, either based on some result of Functional Analysis, or based on results from advanced Measure Theory such as the Radon-Nikodým Theorem. Here, we are mainly

focused on either *computing* a version of the conditional expectation of a given random variable, or *checking* that a random variable is a version of the conditional expectation of some random variable. We have defined the conditional expectation for integrable random variables. It is possible to define the conditional expectation for positive random variables, see the comments below about properties of the conditional expectation.

4. Projection property Let \mathcal{H} be a sub- σ -field of \mathcal{G} . Then $E_\mu(E_\mu(X|\mathcal{G})|\mathcal{H}) = E_\mu(X|\mathcal{H})$. In fact, the conditional expectation $X \mapsto E_\mu(X|\mathcal{F})$ operator is a projection operator on $L^1(\mathcal{F}, \mu)$. It is the transposed operator of the injection operator $\mathcal{L}^\infty(\mathcal{G}) \rightarrow \mathcal{L}^\infty(\mathcal{F})$.

5. Examples.

- (1) If $\mathcal{G} = \{\emptyset, \Omega\}$, then $E_\mu(X|\mathcal{G}) = E_\mu(X)$.
- (2) If $\mathcal{G} = \mathcal{F}$, then $E_\mu(X|\mathcal{G}) = X$.
- (3) Let $\{A_1, \dots, A_n\}$ be a measurable partition of Ω and let $\mathcal{G} = \sigma(A_1, \dots, A_n)$. Assume $\mu(A_j) \neq 0$, $j = 1, \dots, n$. It holds

$$E_\mu(X|\mathcal{G}) = \sum_{j=1}^n \frac{\int_{A_j} X d\mu}{\mu(A_j)} \mathbf{1}_{A_j} = \sum_{j=1}^n E_\mu(X|A_j) \mathbf{1}_{A_j} .$$

- (4) Let μ be a probability measure (Ω, \mathcal{F}) and $p \cdot \mu$ a μ -absolutely-continuous finite measure. Then the restriction of the measure to \mathcal{G} is $(p \cdot \mu)|_{\mathcal{G}} = E_\mu(p|\mathcal{G}) \cdot (\mu|_{\mathcal{G}})$.

6 (Conditioning to a random variable). Let (S, \mathcal{S}) be a measurable space, $Y: \Omega \rightarrow S$ a measurable mapping, and $\mathcal{Y} = \sigma(Y) = Y^{-1}(\mathcal{S})$. A real random variable is \mathcal{Y} -measurable if, and only if, it is of the form $\phi \circ Y$, where ϕ is a real random variable on (S, \mathcal{S}) . [The “if” part follows from $(\phi \circ Y)^{-1}(\mathcal{B}) \subset \mathcal{Y}$; the “only if” part is true because of the Monotone Class Theorem, see [2, ¶3.14].] In this situation, the definition of conditional expectation is rephrased as follows. A version of the conditional expectation of X given $\sigma(Y)$ is a μ -integrable real random variable of the form $\hat{\phi}_{\mu, X} \circ Y$ such that for all bounded measurable $\phi: S \rightarrow \mathbb{R}$ it holds $E_\mu(\phi(Y)\hat{\phi}_{\mu, X}(Y)) = E_\mu(\phi(Y)X)$. Notice that we could write this in terms of the joint distribution of the random variables X and Y as $\int \phi(y)\hat{\phi}_{\mu, X}(y) \mu_Y(dy) = \int \phi(y)x \mu_{X, Y}(dxdy)$. An imprecise, but widely used, notation is $\phi_{\mu, X}(y) = E_\mu(X|Y = y)$, which is called the *expected value of X , given $Y = y$* .

7. Special cases.

- (1) If $X \perp\!\!\!\perp Y$ then $E_\mu(X|\sigma(Y)) = E_\mu(X)$. in fact,

$$\int \phi(y)x \mu_{X, Y}(dxdy) = \int \phi(y) \left(\int x \mu_X(dx) \right) \mu_Y(dy) .$$

- (2) If $X \perp\!\!\!\perp Y$ then $E_\mu(f(X, Y)|\sigma(Y)) = \int f(x, Y) \mu_X(dx)$. In this case we have

$$\int \phi(y)f(x, y) \mu_X \otimes \mu_Y(dxdy) = \int \phi(y) \left(\int f(x, y) \mu_X(dx) \right) \mu_Y(dy) .$$

- (3) Let X, Y , be random variables in \mathbb{R}^m such that $(X - Y) \perp\!\!\!\perp Y$. Then

$$E_\mu(f(Y)|\sigma(Y)) = E_\mu(f((X - Y) + Y)|\sigma(Y)) = \int f(s, Y) \mu_{(X - Y)}(ds) .$$

Cf. the Gaussian case.

- (4) If $\mu_{X,Y}(dx, dy) = p_{X,Y} \cdot \nu_X \otimes \nu_Y$, then $\mu_Y = \left(\int p(x, y) \nu_X(dx) \right) \cdot \nu_Y(dy)$ and the characteristic equality becomes

$$\int \phi(y) \phi_X(y) \left(\int p(x, y) \nu_X(dx) \right) \cdot \nu_Y(dy) = \int \phi(y) \left(\int x p_{X,Y} \nu_X(dx) \right) \nu_Y(dy),$$

hence we can take

$$\hat{\phi}_X(y) = \int x p_{X|Y}(x|y) \nu_X(dx), \quad p_{X|Y}(x|y) = \frac{p_{X,Y}(x, y)}{p_X(x)}.$$

8. Properties. All random variables are defined on the probability space $(\Omega, \mathcal{F}, \mu)$ and \mathcal{G} is a sub- σ -algebra of \mathcal{F}

- (1) *Normalization.* $E_\mu(\mathbf{1}|\mathcal{G}) = \mathbf{1}$.
- (2) *\mathcal{G} -Linearity.* If $E_\mu(X|\mathcal{G}) = \hat{X}$ and $E_\mu(Y|\mathcal{G}) = \hat{Y}$, then $E_\mu(AX + BY|\mathcal{G}) = A\hat{X} + B\hat{Y}$ μ -almost-surely if $A, B \in \mathcal{L}^\infty(\mathcal{G})$.
- (3) *Positivity.* If $X \geq 0$ and $E_\mu(X|\mathcal{G}) = \hat{X}$, then $\hat{X} \geq 0$. Linearity and positivity together imply monotonicity. [Hint: take $G = \mathbf{1}_{\{\hat{X} \leq 0\}}$ in the characteristic property]
- (4) Normalization, linearity and monotonicity together imply *Jensen inequality*. Assume $\Phi: \mathbb{R} \rightarrow \mathbb{R}$ and assume both X and $\Phi(X)$ are integrable. Let $x \mapsto a + bx \leq \Phi(x)$. Then $a + b E_\mu(X|\mathcal{G}) \leq E_\mu(\Phi(X)|\mathcal{G})$. Chose a version $\hat{X} = E_\mu(X|\mathcal{G})$. Because of the convexity, for each $\omega \in \Omega$, there exists coefficients $a(\omega), b(\omega)$ such that $a(\omega) + b(\omega)\hat{X}(\omega) = \Phi(\hat{X}(\omega))$. We have shown that $\Phi(E_\mu(X|\mathcal{G})) \leq E_\mu(\Phi(X)|\mathcal{G})$. In particular, $E_\mu(|X|^\alpha|\mathcal{G}) \leq E_\mu(|X|^\alpha|\mathcal{G})$ if $\alpha \geq 1$.
- (5) *Monotone convergence.* If $0 \leq X_n \uparrow X$ and $\hat{X}_n = E_\mu(X_n|\mathcal{G})$, $n \in \mathbb{N}$, then random variable \hat{X} defined by $\hat{X}_n \uparrow \hat{X}$ is such that $E_\mu(G\hat{X}) = E_\mu(GX)$ if $0 \leq G \in \mathcal{L}^\infty(\mathcal{G})$. It follows immediatly from the monotone convergence for the expectation [Notice that here we are assuming each X_n to be 'integrable so that the conditional expectation is defined. This is not necessary if we define conditional expectation for non-negative random variable as it was for the expectation. We do not consider this generalization in this notes.] If moreover X happens to be integrable, then $\hat{X} = E_\mu(X|\mathcal{G})$.
- (6) *Fatou lemma.* If $0 \leq X_n$ and $\hat{X}_n = E_\mu(X_n|\mathcal{G})$, $n \in \mathbb{N}$, then $\wedge_{m \geq n} X_m \leq X_n$ if $m \geq n$, so that $E_\mu(\wedge_{m \geq n} X_m|\mathcal{G}) \leq \wedge_{m \geq n} E_\mu(X_m|\mathcal{G})$. From the monotone convergence it follows $E_\mu(G(\liminf_{n \rightarrow \infty} X_n)) \leq E_\mu(G(\liminf_{n \rightarrow \infty} E_\mu(X_n|\mathcal{G})))$ if $G \in \mathcal{L}^\infty(\mathcal{G})$ and $G \geq 0$. If $\liminf_{n \rightarrow \infty} X_n$ is integrable, then we can write $E_\mu(\liminf_{n \rightarrow \infty} X_n|\mathcal{G}) \leq \liminf_{n \rightarrow \infty} E_\mu(X_n|\mathcal{G})$.
- (7) *Dominated convergence.* If in the fatou lemma we assume that the sequence $(X_n)_{n \in \mathbb{N}}$ is dominated by the integrable random variable Y , by considering the non-negative sequence $(Y - X_n)_{n \in \mathbb{N}}$ we can obtain the inequality

$$E_\mu \left(\liminf_{n \rightarrow \infty} X_n \middle| \mathcal{G} \right) \leq \liminf_{n \rightarrow \infty} E_\mu(X_n|\mathcal{G}) \leq \limsup_{n \rightarrow \infty} E_\mu(X_n|\mathcal{G}) \leq E_\mu \left(\limsup_{n \rightarrow \infty} X_n \middle| \mathcal{G} \right).$$

If the sequence is convergent, then $\liminf_{n \rightarrow \infty} X_n = \lim_{n \rightarrow \infty} X_n = \limsup_{n \rightarrow \infty} X_n$ hence $\liminf_{n \rightarrow \infty} E_\mu(X_n|\mathcal{G}) = \limsup_{n \rightarrow \infty} E_\mu(X_n|\mathcal{G})$ and the sequence of conditional expectations is convergent to the expectation of the limit. The condition of positivity can be dropped by decomposing the positive and negative part of the sequence and the limit.

2. CONDITIONAL DISTRIBUTION

9 (Transition probability measure). Given a product measurable space $(\Omega_1 \times \Omega_2, \mathcal{F}_1 \otimes \mathcal{F}_2)$ a *transition* is a mapping $\mu_{1|2}: \mathcal{F}_1 \times \Omega_2$ such that

- (1) for each $x_2 \in \Omega_2$ the mapping $\mathcal{F}_1 \ni A_1 \mapsto \mu_{1|2}(A_1|x_2)$ is a probability measure on $(\Omega_1, \mathcal{F}_1)$ and
- (2) for each $A_1 \in \mathcal{F}_1$ the mapping $\Omega_2 \ni x_2 \mapsto \mu_{1|2}(A_1|x_2)$ is \mathcal{F}_2 -measurable.

10 (Integration of probability measures). Given a transition $\mu_{1|2}$ on $(\Omega_1 \times \Omega_2, \mathcal{F}_1 \otimes \mathcal{F}_2)$ and a probability measure μ_2 on $(\Omega_2, \mathcal{F}_2)$, there exists a unique probability measure $\mu = \int \mu_{1|2} d\mu_2$ on the product measurable space such that for each positive or μ -integrable function $f: \Omega_2 \times \Omega_2 \ni (x_1, x_2) \mapsto f(x_1, x_2)$ it holds

$$\int f d\mu = \int \left(\int f(x_1, x_2) \mu_{1|2}(dx_1|x_2) \right) \mu_2(dx_2) .$$

The measure μ is characterised on functions of the form $f(x_1, x_2) = f_1(x_1)f_2(x_2)$ by

$$\int f_1 f_2 d\mu = \int \left(\int f_1(x_1) \mu_{1|2}(dx_1|x_2) \right) f_2(x_2) \mu_2(dx_2) .$$

[The proof is a simple variation of the argument for Fubini theorem.]

11 (Transition densities). A simple case occurs when the transition has the form

$$\mu_{1|2}(A_1|x_2) = \int_{A_1} p_{1|2}(x_1|x_2) \nu_1(dx), \quad A_1 \in \mathcal{F}_1, x_2 \in \Omega_2$$

where $(x_1, x_2) \mapsto p_{1|2}(x_1|x_2)$ is measurable on the product space $(\Omega_1, \Omega_2, \mathcal{F}_1 \otimes \mathcal{F}_2)$ and $x_1 \mapsto p_{1|2}(x_1|x_2)$ is a ν_1 -probability density for each $x_2 \in \Omega_2$. In such a case,

$$\begin{aligned} \int \left(\int f_1(x_1) \mu_{1|2}(dx_1|x_2) \right) f_2(x_2) \mu_2(dx_2) &= \\ \int \left(\int f_1(x_1) p_{1|2}(x_1|x_2) \nu_1(dx_1) \right) f_2(x_2) \mu_2(dx_2) &= \\ \iint f_1(x_1) f_2(x_2) p_{1|2}(x_1|x_2) \nu_1(dx_1) \mu_2(dx_2) , & \end{aligned}$$

that is, $\mu = p_{1|2} \cdot \nu_1 \otimes \mu_2$. If moreover the second measure has itself a density, $\mu_2 = p_2 \cdot \nu_2$, then $\mu = (p_{1|2} \otimes p_2) \cdot \nu_1 \otimes \nu_2$

- 12** (Examples). (1) Let T_1, T_2 be independent and $\text{Exp}(1)$. Then the distribution of T_1 given $T_1 + T_2 = t$ is uniform on $]0, t[$.
- (2) If $(Y_1, Y_2) \sim N_{n_1+n_2}(0, \Sigma)$, $\det \Sigma \neq 0$, find the conditional distribution of Y_1 given Y_2 .
 - (3) If Y_1, Y_2 are independent and $N(0, 1)$, find the distribution of (y_1, Y_2) given $Y_1^2 + Y_2^2$.

13 (Regular version of the conditional expectation). With the notations above, denoting with X_1, X_2 the coordinate projection, the random variable $\hat{f}(X_2) = \int f(x_1, X_2) \mu_{1|2}(dx_1|X_2)$ is a version of the conditional expectation $E_\mu(f(X_1, X_2)|\sigma(X_2))$, namely a *regular version*. In fact,

$$E_\mu(f(X_1, X_2)g(X_2)) = \int \left(\int f(x_1, x_2) \mu_{1|2}(dx_1|x_2) \right) g(x_2) \mu_2(dx_2) = E_\mu \left(\hat{f}(X_2)g(X_2) \right) .$$

14 (Conditional independence (CI) of events). Conditional independence is a key property in Statistics e.g., in Graphical Models, in Stochastic Processes, in Markov processes, in Random Fields, in Machine Learning.

The nonzero events A, B, C are such that A and C are independent given B , $A \perp\!\!\!\perp C \mid B$, if each one of the following equivalent conditions are satisfied:

$$\begin{aligned} P(A \cap C \mid B) &= P(A \mid B) P(C \mid B) \quad (\text{CI}) \\ P(A \cap B \cap C) &= P(A \mid B) P(C \mid B) P(B) \quad (\text{J-CI}) \\ P(A \mid B \cap C) &= P(A \mid B) \quad (\text{M}) \\ P(A \cap B \cap C) &= P(A \mid B) P(B \mid C) P(C) \quad (\text{J-M}) \\ P(A \cap B \cap C) P(B) &= P(A \cap B) P(B \cap C) \quad (\text{A}) \end{aligned}$$

In fact, the first four equalities become the fifth one if the conditional probabilities are computed in terms of joint probabilities. The algebraic form (A) can be written in terms of indicator functions as

$$E(\mathbf{1}_A \mathbf{1}_B \mathbf{1}_C) E(\mathbf{1}_B) = E(\mathbf{1}_A \mathbf{1}_B) E(\mathbf{1}_B \mathbf{1}_C) .$$

which shows the bi-linearity in $\mathbf{1}_A$ and $\mathbf{1}_C$. For example, writing $\mathbf{1}_A = \mathbf{1} - \mathbf{1}_{A^c}$ one gets $A^c \perp\!\!\!\perp C \mid B$.

15 (Conditional Independence for random variables). *Random variables Y_1, Y_3 are conditionally independent given the random variable Y_2 , $Y_1 \perp\!\!\!\perp Y_3 \mid Y_2$* if each one of the following equivalent conditions are satisfied. If $f_i, i = 1, \dots, 3$, are bounded,

$$\begin{aligned} E(f_1(Y_1) f_3(Y_3) \mid Y_2) &= E(f_1(Y_1) \mid Y_2) E(f_3(Y_3) \mid Y_2) \\ E(f_1(Y_1) \mid Y_2, Y_3) &= E(f_1(Y_1) \mid Y_2) \end{aligned}$$

Let us prove the equivalence. The second one holds if, and only if, for all bounded $f_2(Y_2), f_3(Y_3)$

$$E(f_1(Y_1) f_2(Y_2) f_3(Y_3)) = E(E(f_1(Y_1) \mid Y_2) f_2(Y_2) f_3(Y_3))$$

The LHS is equal to

$$E(E(f_1(Y_1) f_3(Y_3) \mid Y_2) f_2(Y_2))$$

and the RHS is equal to

$$E(E(f_1(Y_1) \mid Y_2) f_2(Y_2) E(f_3(Y_3) \mid Y_2)) .$$

It follows that the first equation holds. By reversing the computation we get the other implication.

When a regular version of the conditional expectation given Y_2 is available, then conditional independence is equivalent to the product form of the transition.

16 (Markov process). A stochastic process Y_1, \dots, Y_N is a *Markov Process* if $(Y_1, \dots, Y_k) \perp\!\!\!\perp Y_{k+1}, \dots, Y_N \mid Y_k$, $k = 1, 2, \dots, N$. Equivalently,

$$E(f(Y_{k+1}, \dots, Y_N) \mid Y_1, \dots, Y_k) = E(f(Y_{k+1}, \dots, Y_N) \mid Y_k) .$$

REFERENCES

- [1] Claude Dellacherie and Paul-André Meyer, *Probabilities and potential*, North-Holland Mathematics Studies, vol. 29, North-Holland Publishing Co., Amsterdam-New York; North-Holland Publishing Co., Amsterdam-New York, 1978. MR 521810
- [2] David Williams, *Probability with martingales*, Cambridge Mathematical Textbooks, Cambridge University Press, Cambridge, 1991.

COLLEGIO CARLO ALBERTO

E-mail address: `giovanni.pistone@carloalberto.org`