

PROBABILITY 2018 HANDOUT 3: CONDITIONING

GIOVANNI PISTONE

CONTENTS

1. Conditional expectation	1
2. Conditional distribution	6
3. Conditioning of jointly Gaussian vectors	7
4. Conditional independence	9
References	11

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

Exercise 1. Let X be a measurable function from (Ω, \mathcal{F}) to (S, \mathcal{S}) . Let \mathcal{G} be the σ -algebra generated by X i.e., $\mathcal{G} = X^{-1}\mathcal{S}$. Every \mathcal{G} -measurable real random variable Y is of the form $Y = f \circ X$, where f is a real random variable on (S, \mathcal{S}) . [Hint: If Y is simple, $Y = \sum_{j=1}^n y_j \mathbf{1}_{B_j}$, with $B_j \in \mathcal{G}$, then $B_j = X^{-1}(A_j)$, $A_j \in \mathcal{S}$. It follows that $Y = \sum_{j=1}^n y_j \mathbf{1}_{X^{-1}(A_j)} = \sum_{j=1}^n y_j \mathbf{1}_{A_j} \circ X$, hence $f = \sum_{j=1}^n y_j \mathbf{1}_{A_j}$. If X is non-negative, take an increasing sequence of simple random variable ...]

Definition 1. Let $(\Omega, \mathcal{F}, \mu)$ be a probability space, X a real random variable with finite expectation, $\mathbb{E}_\mu[|X|] < +\infty$, \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

$$\mathbb{E}_\mu[G\hat{X}] = \mathbb{E}_\mu[GX] .$$

The sub- μ in the notation is there to remember that the conditional expectation *depends on the probability*. The conditions (1) and (2) in the definition provide actual equations to compute the conditional expectation, as the following examples show.

Exercise 2. If X is a real random variable with a positive density p , let \mathcal{G} be the σ -algebra generated by $|X|$. That is, the absolute value only, not the sign, is observed. In this case the conditional expectation of X given $\mathcal{G} = \sigma(|X|)$, briefly, given $|X|$, is a random variable of the form $\hat{X} = \hat{f}(|X|)$ (condition (1)) such that $\mathbb{E}(\hat{X}G) = \mathbb{E}(XG)$

for all $G = g(|X|)$, g bounded (condition (2)). As a density is given, we write the defining equation

$$\int \hat{f}(|x|)g(|x|)p(x) dx = \int xg(|x|)p(x) dx .$$

[Hint: To compute \hat{f} , split $\int = \int_{-\infty}^0 + \int_0^{+\infty}$ and change the variable $x \rightarrow -x$ in the first integral to get

$$\int_0^{+\infty} \hat{f}(|x|)g(|x|)(p(x) + p(-x)) dx = \int_0^{+\infty} g(|x|)(xp(x) - xp(-x)) dx ,$$

hence

$$\hat{f}(|x|)(p(x) + p(-x)) = xp(x) - xp(-x) .$$

Finally, notice that $\frac{xp(x) - xp(-x)}{p(x) + p(-x)}$ is symmetric.]

Exercise 3. Let S_1, S_2 be independent and exponential with mean 1. The joint density is $p_{S_1, S_2}(x_1, x_2) = e^{-(x_1+x_2)}(x_1, x_2 > 0)$. We want to compute the conditional expectation of S_1 given $S_1 + S_2$. We need to find \hat{f} such that for all bounded g we have

$$\iint_0^{\infty} \hat{f}(x_1 + x_2)g(x_1 + x_2)e^{-(x_1+x_2)} dx_1 dx_2 = \iint_0^{\infty} x_1 g(x_1 + x_2)e^{-(x_1+x_2)} dx_1 dx_2 .$$

[Hint. Let us make the transformation $y = x_1 + x_2, z = x_1$. The inverse transformation is $x_1 = z, x_2 = y - z$ with determinant -1 . We have

$$(x_1, x_2 > 0) = (z > 0)(y - z > 0) = (0 < z < y)$$

then the equation becomes

$$\iint_{\{0 < z < y\}} \hat{f}(y)g(y)e^{-y} dy dz = \iint_{\{0 < z < y\}} zg(y)e^{-y} dz dy .$$

Computing the dz integrals on both sides we get

$$\int_0^{\infty} \hat{f}(y)g(y)ye^{-y} dy = \int_0^{\infty} g(y)\frac{y^2}{2}e^{-y} dy ,$$

hence $\hat{f}(y) = \frac{y}{2}$.]

Exercise 4. Let $Z = (Z_1, Z_2) \sim N_2(0, I)$ and define $X = Z_1, Y = Z_1 + Z_2, \mathcal{G} = \sigma(Y)$. To compute a version of $E(X|\mathcal{G})$ we look for a function \hat{f} such that $\hat{f}(Y)$ satisfies

$$\mathbb{E}(Xg(Y)) = \mathbb{E}(\hat{f}(Y)g(Y)) \quad \text{for all bounded } g .$$

[Hint: As

$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} Z_1 \\ Z_1 + Z_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix}$$

we have $(X, Y) \sim N_2\left(0, \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}\right)$ and $Y \sim N_1(0, 2)$. We have $\det\left(\begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}\right) = 1$ and

$\begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}^{-1} = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix}$ so that the density of (X, Y)

$$p_{X,Y}(x, y) = (2\pi)^{-1} \exp\left(-\frac{1}{2}(2x^2 - 2xy + y^2)\right) .$$

We want

$$\begin{aligned} \iint xg(y) (2\pi)^{-1} \exp\left(-\frac{1}{2}(2x^2 - 2xy + y^2)\right) dx dy = \\ \int \hat{f}(y)g(y) (2\pi \cdot 2)^{-1/2} \exp\left(-\frac{1}{2 \cdot 2}y^2\right) dy \end{aligned}$$

Let us perform first the dx integration in the RHS:

$$\begin{aligned} \int x \exp\left(-\frac{1}{2}(2x^2 - 2xy + y^2)\right) dx &= \int x \exp\left(-\left(x^2 - xy + \frac{1}{2}y^2\right)\right) dx = \\ &= \int x \exp\left(-\left(x - \frac{1}{2}y\right)^2 - \frac{1}{4}y^2\right) dx = \\ &= \exp\left(-\frac{1}{4}y^2\right) \int \pi^{1/2}x \pi^{-1/2} \exp\left(-\left(x - \frac{1}{2}y\right)^2\right) dx = \\ &= \frac{\pi^{1/2}}{2}y \exp\left(-\frac{1}{4}y^2\right) . \end{aligned}$$

The defining equality becomes

$$\begin{aligned} \int g(y) (2\pi)^{-1} \frac{\pi^{1/2}}{2}y \exp\left(-\frac{1}{4}y^2\right) dy = \\ \int f(y)g(y) (2\pi \cdot 2)^{-1/2} \exp\left(-\frac{1}{2 \cdot 2}y^2\right) dy \end{aligned}$$

so that, g being generic, $\hat{f}(y) = y/2$. (We are going to see below a simpler and more principled way to do this computation.)]

1. As the equation $\mathbb{E}_\mu \left[G(\hat{X} - X) \right] = 0$, $G \in \mathcal{L}^\infty(\mathcal{G})$, 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 $\mathbb{E}_\mu \left[G(\hat{X}_1 - \hat{X}_2) \right] = 0$ i.e. $\hat{X}_1 = \hat{X}_2$ μ -almost-surely. [Take $G = \text{sign}(\hat{X}_1 - \hat{X}_2)$ to get $\mathbb{E}_\mu \left[|\hat{X}_1 - \hat{X}_2| \right] = 0$.] More generally, if $X_1 = X_2$ μ -almost-surely, then $\hat{X}_1 = \hat{X}_2$ μ -almost-surely. We write $\mathbb{E}_\mu(X|\mathcal{G})$ to denote the μ -class of versions and, with abuse of notation, $\hat{X} = \mathbb{E}_\mu(X|\mathcal{G})$. If $L^1(\mathcal{F}, \mu)$ is the vector space of classes μ -equivalent real random variables, there exists a mapping

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

3 (Existence). 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 discussion, 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 (existence of orthogonal projection), or based on results from advanced Measure Theory such as the Radon-Nikodým Theorem (see below). 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 (Image of a density). On the measurable space (Ω, \mathcal{F}) , consider the probability measure μ and the probability density P . If Φ is measurable from (Ω, \mathcal{F}) to (S, \mathcal{S}) , consider the image of the probability measure $p \cdot \mu$ under Φ . The image $\nu = \Phi_{\#}(p \cdot \mu)$ is characterized by

$$\int_S g(y) \nu(dy) = \int_{\Omega} g \circ \Phi(x) p(x) \mu(dx), \quad g \in \mathcal{L}^{\infty}(S, \mathcal{S}).$$

Now, $g \circ \Phi$ is the generic bounded $\sigma(\Phi)$ -measurable random variable, then

$$\int_{\Omega} g \circ \Phi(x) p(x) \mu(dx) = \int_{\Omega} g \circ \Phi(x) \hat{p} \circ \Phi(x) \mu(dx),$$

where $\hat{p} \circ \Phi$ is a version of the conditional expectation of p given $\sigma(\Phi)$. Now apply again the definition of image to the RHS to get

$$\int_S g(y) \Phi_{\#}(p \cdot \mu)(dy) = \int_S g(y) \hat{p}(y) \Phi_{\#}(\mu)(dy).$$

We have found the density of the image measure.

5 (Projection property). Let \mathcal{H} be a sub- σ -field of \mathcal{G} . It is easy to check that

$$\mathbb{E}_{\mu}(\mathbb{E}_{\mu}(X|\mathcal{G})|\mathcal{H}) = \mathbb{E}_{\mu}(X|\mathcal{H}).$$

In particular, the conditional expectation operator $X \mapsto \mathbb{E}_{\mu}(X|\mathcal{F})$ is a projection operator on $L^1(\mathcal{F}, \mu)$. [One could say that it is the transposed operator of the injection operator $\mathcal{L}^{\infty}(\mathcal{G}) \rightarrow \mathcal{L}^{\infty}(\mathcal{F})$.

6 (Orthogonal projection). The conditioning operator is an *orthogonal projection*. Assume Y in $L^2(\Omega, \mathcal{F}, \mu)$ that is, $\mathbb{E}(Y^2) < \infty$. If $\hat{Y} = \mathbb{E}(Y|\mathcal{G})$, then $\hat{Y} \in L^2(\Omega, \mathcal{G}, \mu)$ and

$$\mathbb{E}\left((Y - \hat{Y})Z\right) = 0, \quad z \in L^2(\Omega, \mathcal{G}, \mu).$$

This property should not be confused with *linear regression*. Let be given $Y \in L^2$ and let $X_1, \dots, X_m \in L^2$ be *explanatory variables*. We want a vector $\theta = (\theta_0, \theta_1, \dots, \theta_d) \in \mathbb{R}^{d+1}$ such that

$$\text{quadratic error} = \mathbb{E}\left(\left(Y - \theta_0 - \sum_{j=1}^d \theta_j X_j\right)^2\right)$$

be minimum. As a function of θ the quadratic error is a convex function then the minimum is obtained by imposing the gradient to be zero.

Exercise 5. Check all detail of the previous paragraph.

Exercise 6 (Examples). (1) If $\mathcal{G} = \{\emptyset, \Omega\}$, then $\mathbb{E}_{\mu}(X|\mathcal{G}) = \mathbb{E}_{\mu}[X]$.

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

(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

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

7 (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}) . 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 $\mathbb{E}_{\mu} \left[\phi(Y) \hat{\phi}_{\mu, X}(Y) \right] = \mathbb{E}_{\mu} [\phi(Y)X]$. 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 $\phi_{\mu, X}(y) = \mathbb{E}_{\mu}(X|Y = y)$, which is called the *expected value of X , given $Y = y$* .

8 (Special cases). (1) If $X \perp\!\!\!\perp Y$ then $\mathbb{E}_{\mu}(X|\sigma(Y)) = \mathbb{E}_{\mu}[X]$. 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 $\mathbb{E}_{\mu}(f(X, Y)|\sigma(Y)) = \int f(x, Y) \mu_X(dx)$. 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

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

Cf. the Gaussian case below.

(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

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

9 (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.* $\mathbb{E}_{\mu}(\mathbf{1}|\mathcal{G}) = \mathbf{1}$.

(2) *\mathcal{G} -Linearity.* If $\mathbb{E}_{\mu}(X|\mathcal{G}) = \hat{X}$ and $\mathbb{E}_{\mu}(Y|\mathcal{G}) = \hat{Y}$, then $\mathbb{E}_{\mu}(AX + BY|\mathcal{G}) = A\hat{X} + B\hat{Y}$ μ -almost-surely if $A, B \in \mathcal{L}^{\infty}(\mathcal{G})$.

(3) *Positivity.* If $X \geq 0$ and $\mathbb{E}_{\mu}(X|\mathcal{G}) = \hat{X}$, 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 $a + b\mathbb{E}_{\mu}(X|\mathcal{G}) \leq \mathbb{E}_{\mu}(\Phi(X)|\mathcal{G})$. Chose a version $\hat{X} = \mathbb{E}_{\mu}(X|\mathcal{G})$. 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 $\Phi(\mathbb{E}_{\mu}(X|\mathcal{G})) \leq \mathbb{E}_{\mu}(\Phi(X)|\mathcal{G})$. In particular, $\mathbb{E}_{\mu}(|X|\mathcal{G})^{\alpha} \leq \mathbb{E}_{\mu}(|X|^{\alpha}|\mathcal{G})$ if $\alpha \geq 1$.

- (5) *Monotone convergence.* If $0 \leq X_n \uparrow X$ and $\hat{X}_n = \mathbb{E}_\mu(X_n|\mathcal{G})$, $n \in \mathbb{N}$, then random variable \hat{X} defined by $\hat{X}_n \uparrow \hat{X}$ is such that $\mathbb{E}_\mu[G\hat{X}] = \mathbb{E}_\mu[GX]$ 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 $\hat{X} = \mathbb{E}_\mu(X|\mathcal{G})$.
- (6) *Fatou lemma.* If $0 \leq X_n$ and $\hat{X}_n = \mathbb{E}_\mu(X_n|\mathcal{G})$, $n \in \mathbb{N}$, then $\wedge_{m \geq n} X_m \leq X_n$ if $m \geq n$, so that $\mathbb{E}_\mu(\wedge_{m \geq n} X_m|\mathcal{G}) \leq \wedge_{m \geq n} \mathbb{E}_\mu(X_m|\mathcal{G})$. From the monotone convergence it follows $\mathbb{E}_\mu[G(\liminf_{n \rightarrow \infty} X_n)] \leq \mathbb{E}_\mu[G(\liminf_{n \rightarrow \infty} \mathbb{E}_\mu(X_n|\mathcal{G}))]$ 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 $\mathbb{E}_\mu(\liminf_{n \rightarrow \infty} X_n|\mathcal{G}) \leq \liminf_{n \rightarrow \infty} \mathbb{E}_\mu(X_n|\mathcal{G})$.
- (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

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

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 $\liminf_{n \rightarrow \infty} \mathbb{E}_\mu(X_n|\mathcal{G}) = \limsup_{n \rightarrow \infty} \mathbb{E}_\mu(X_n|\mathcal{G})$ 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

10 (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.

11 (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.]

12 (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 $\mu = (p_{1|2} \otimes p_2) \cdot \nu_1 \otimes \nu_2$

Exercise 7 (Examples).

- (1) Let X be a real random variable with positive density p . The conditional distribution of X given $|X|$ is
- (2) 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 $]0, t[$.
- (3) 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 .
- (4) If Y_1, Y_2 are independent and $N_1(0, 1)$, 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 $\mathbb{E}_\mu(f(X_1, X_2)|\sigma(X_2))$, namely a *regular version*. In fact,

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

3. CONDITIONING OF JOINTLY GAUSSIAN VECTORS

Exercise 8. Recall that for each $\Sigma \in \text{Sym}_+(n)$ there exists an orthogonal $U \in O(n)$ and a non-negative diagonal $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ such that $\Sigma = U\Lambda U^*$. By discarding the zero eigen-values, we can write $\Sigma = SDS^*$ with $S \in \text{Mat}(n \times r)$, $S^*S = I_r$, and D positive diagonal, where r is the rank of Σ . If $D = \text{diag}(\lambda_1, \dots, \lambda_r)$, we define $D^{-1} = \text{diag}(\lambda_1^{-1}, \dots, \lambda_r^{-1})$ and $\Sigma^+ = SD^{-1}S^*$. It follows that

$$\Sigma^+\Sigma = SD^{-1}S^*SDS^* = SS^* \quad \text{and} \quad \Sigma\Sigma^* = SDS^*SD^{-1}S^* = SS^* .$$

We have $\Pi = SS^* \in \text{Sym}_+(n)$ and $\Pi^2 = \Pi$. The matrix Π is the orthogonal projector onto the image of Σ . In fact, for all $x \in \mathbb{R}^n$,

$$\Pi x = SS^*x = SDS^*SD^{-1}S^*x = \Sigma SD^{-1}S^*x .$$

Moreover, for each $x, y \in \mathbb{R}^n$

$$\begin{aligned} (x - \Pi x) \cdot (\Sigma y) &= \\ (x - \Pi x)^*(\Sigma y) &= [(I - SS^*)x]^*(SDS^*y) = x^*(I - SS^*)SDS^*y = \\ &= x^*(SDS^* - SS^*SDS^*) = 0 \end{aligned}$$

Proposition 1.

(1) *The Gaussian random vector with components*

$$\begin{aligned}\tilde{Y}_1 &= Y_1 - (b_1 + L_{12}(Y_2 - b_2)), \quad L_{12} = \Sigma_{12}\Sigma_{22}^+ \\ \tilde{Y}_2 &= Y_2 - b_2\end{aligned}$$

is such that $\mathbb{E}(\tilde{Y}_1) = 0$, $\text{Var}(\tilde{Y}_1) = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^+\Sigma_{21}$, and $\tilde{Y}_1 \perp\!\!\!\perp \tilde{Y}_2$. It follows

$$\mathbb{E}(Y_1|Y_2) = b_1 + L_{12}(Y_2 - b_2)$$

(2) *The conditional distribution of Y_1 given $Y_2 = y_2$ is Gaussian with*

$$Y_1|Y_2 = y_2 \sim N_{n_1}(b_1 + L_{12}(y_2 - b_2), \Sigma_{11} - L_{12}\Sigma_{21})$$

(3) *The conditional density of Y_1 given $Y_2 = y_2$ in terms of the partitioned concentration is*

$$\begin{aligned}p_{Y_1|Y_2}(y_1|y_2) &= (2\pi)^{-\frac{n_1}{2}} \det(K_{1|2})^{\frac{1}{2}} \times \\ &\exp\left(-\frac{1}{2}(y_1 - b_1 - K_{11}^{-1}K_{12}(y_2 - b_2))^T K_{11}(y_1 - b_1 - K_{11}^{-1}K_{12}(y_2 - b_2))\right)\end{aligned}$$

Proof. (1) We have

$$\begin{bmatrix} \tilde{Y}_1 \\ \tilde{Y}_2 \end{bmatrix} = \begin{bmatrix} I & -\Sigma_{12}\Sigma_{22}^+ \\ 0 & I \end{bmatrix} \begin{bmatrix} Y_1 - b_1 \\ Y_2 - b_2 \end{bmatrix} \sim N_{n_1+n_2}\left(0, \begin{bmatrix} \Sigma_{1|2} & 0 \\ 0 & \Sigma_{22} \end{bmatrix}\right)$$

It follows

$$\mathbb{E}(Y_1|Y_2) = \mathbb{E}\left(\tilde{Y}_1 + b_1 + L_{12}(Y_2 - b_2) \middle| Y_2\right) = \mathbb{E}(\tilde{Y}_1) + b_1 + L_{12}(Y_2 - b_2)$$

(2) The conditional distribution of Y_1 given Y_2 is a transition probability $\mu_{Y_1|Y_2} : \mathcal{B}(\mathbb{R}^{n_1}) \times \mathbb{R}^{n_2}$ such that for all bounded $f : \mathbb{R}^{n_1}$

$$\mathbb{E}(f(Y_1)|Y_2) = \int f(y_1) \mu_{Y_1|Y_2}(dy_1|Y_2).$$

We have

$$\mathbb{E}(f(Y_1)|Y_2) = \mathbb{E}\left(f(\tilde{Y}_1 + \mathbb{E}(Y_1|Y_2)) \middle| Y_2\right) = \int f(x + \mathbb{E}(Y_1|Y_2)) \gamma(dx; 0, \Sigma_{1|2})$$

where $\gamma(dx; 0, \Sigma_{1|2})$ is the measure of $N_{n_1}(0, \Sigma_{1|2})$. We obtain the statement by considering the effect on the distribution $N_{n_1}(0, \Sigma_{1|2})$ of the translation $x \mapsto x + (b_1 + L_{12}(y_2 - b_2))$.

(3) A further application of the Schur complement gives

$$\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} I & \Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma_{1|2} & 0 \\ 0 & \Sigma_{22} \end{bmatrix} \begin{bmatrix} I & 0 \\ \Sigma_{22}^{-1}\Sigma_{21} & I \end{bmatrix}$$

whose inverse is

$$\begin{aligned}\begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix} &= \begin{bmatrix} I & 0 \\ -\Sigma_{22}^{-1}\Sigma_{21} & I \end{bmatrix} \begin{bmatrix} \Sigma_{1|2}^{-1} & 0 \\ 0 & \Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} I & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I \end{bmatrix} \\ &= \begin{bmatrix} \Sigma_{1|2}^{-1} & 0 \\ -\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{1|2}^{-1} & \Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} I & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I \end{bmatrix} \\ &= \begin{bmatrix} \Sigma_{1|2}^{-1} & -\Sigma_{1|2}^{-1}\Sigma_{12}\Sigma_{22}^{-1} \\ -\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{1|2}^{-1} & \Sigma_{22}^{-1}\Sigma_{21}\Sigma_{1|2}^{-1}\Sigma_{12}\Sigma_{22}^{-1} + \Sigma_{22}^{-1} \end{bmatrix}\end{aligned}$$

In particular, we have $K_{11} = \Sigma_{1|2}^{-1}$ and $K_{11}^{-1}K_{12} = -\Sigma_{12}\Sigma_{22}^{-1}$, hence

$$Y_1|Y_2 = y_2 \sim N_{n_1}(b_1 - K^{-1}K_{12}(y_2 - b_2), K_{11}^{-1})$$

so that the exponent of the Gaussian density has the factor

$$(y_1 - b_1 + K_{11}^{-1}K_{12}(y_2 - b_2))^T K_{11}(y_1 - b_1 + K_{11}^{-1}K_{12}(y_2 - b_2))$$

□

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