

Title Journal

Abstract

Abstract of the Journal

Index Terms

Supervised Learning, Empirical Risk Minimization; Relative Entropy; Regularization; Gibbs Measure; Inductive Bias; Gibbs Algorithm; Sensitivity; and Generalization.

I. INTRODUCTION

II. FOLKLORE THEOREMS OF THE RADOM-NIKODYM DERIVATIVE

This section introduces relevant notational conventions alongside the Radon-Nikodym theorem. In particular, some equalities presented in this paper are valid almost surely with respect to a given measure. For clarity, given a measure space (Ω, \mathcal{F}, P) , the notation $\stackrel{a.s.}{P}$ is introduced and shall be read as “equal for all $x \in \Omega$ except on a negligible set with respect to P ”; or equivalently as “equal almost surely with respect to P ”. Moreover, given two measures P and Q on the same measurable space, the notation $P \ll Q$ stands for “the measure P is absolutely continuous with respect to Q ”. Using this notation, the Radon-Nikodym derivative is introduced by the following theorem.

Theorem 1 (Radon-Nikodym theorem, [?, Theorem 2.2.1]): Let P and Q be two measures on a given measurable space (Ω, \mathcal{F}) , such that Q is σ -finite and $P \ll Q$. Then, there exists a nonnegative Borel measurable function $g : \Omega \rightarrow \mathbb{R}$ such that for all $\mathcal{A} \in \mathcal{F}$,

$$P(\mathcal{A}) = \int_{\mathcal{A}} g(x) dQ(x). \quad (1)$$

Moreover, if another function h satisfies for all $\mathcal{A} \in \mathcal{F}$ that $P(\mathcal{A}) = \int_{\mathcal{A}} h(x) dQ(x)$, then $g(x) \stackrel{a.s.}{Q} h(x)$.

The function g in (1) is often referred to as the Radom-Nikodym derivative of P with respect to Q ; and is also written as $\frac{dP}{dQ}$, such that $g(x) = \frac{dP}{dQ}(x)$.

The Radon-Nikodym theorem is the foundational tool from which many folklore theorems in information theory originate. Some of these folklore theorems are thoroughly studied in the following sections.

III. BASIC FOLKLORE THEOREMS

This section focuses on basic folklore theorems, where “basic” denotes their well-established nature. One of the most common folklore theorems is often referred to as the “change of measure” theorem.

Theorem 2 (Change of Measure): Let P and Q be two measures on the measurable space (Ω, \mathcal{F}) with $P \ll Q$; and Q a σ -finite measure. Let $f : \Omega \rightarrow \mathbb{R}$ be a Borel measurable function such that the integral $\int_{\Omega} f(x) dP(x)$ exists. Then, for all $\mathcal{A} \in \mathcal{F}$,

$$\int_{\mathcal{A}} f(x) dP(x) = \int_{\mathcal{A}} f(x) \frac{dP}{dQ}(x) dQ(x). \quad (2)$$

Proof: The first part of the proof is developed under the assumption that the function f is simple. That is, for all $x \in \mathcal{X}$, $f(x) = \sum_{i=1}^m a_i \mathbb{1}_{\mathcal{A}_i}(x)$, for finite $m \in \mathbb{N}$, disjoint sets $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_m$ in \mathcal{F} and reals a_1, a_2, \dots, a_m . For all $\mathcal{A} \in \mathcal{F}$, and for all $i \in \{1, 2, \dots, m\}$, let $\mathcal{B}_i = \mathcal{A} \cap \mathcal{A}_i$, hence,

$$\int_{\mathcal{A}} f(x) \frac{dP}{dQ}(x) dQ(x) = \int_{\mathcal{A}} \frac{dP}{dQ}(x) \sum_{i=1}^m a_i \mathbb{1}_{\mathcal{A}_i}(x) dQ(x) \quad (3)$$

$$= \sum_{i=1}^m a_i \int_{\mathcal{B}_i} \frac{dP}{dQ}(x) dQ(x) \quad (4)$$

$$= \sum_{i=1}^m a_i P(\mathcal{B}_i), \quad (5)$$

where the equality in (4) follows from the linearity of the integral [?, Theorem 1.6.3]; and the equality in (5) follows from Theorem 1. On the other hand, for all $\mathcal{A} \in \mathcal{F}$

$$\int_{\mathcal{A}} f(x) dP(x) = \int_{\mathcal{A}} \sum_{i=1}^n a_i \mathbb{1}_{\mathcal{A}_i(x)} dP(x) \quad (6)$$

$$= \sum_{i=1}^n \int_{\mathcal{A}} a_i \mathbb{1}_{\mathcal{A}_i(x)} dP(x) \quad (7)$$

$$= \sum_{i=1}^n a_i \int_{\mathcal{B}_i} dP(x) = \sum_{i=1}^n a_i P(\mathcal{B}_i), \quad (8)$$

where the equality in (7) follows from the linearity of the integral [?, Theorem 1.6.3]. Hence, from Theorem 1, and equalities (5) and (8), it follows that when f is a simple function, the equality in (2) holds. This concludes the first part of the proof.

The second part of the proof proceeds by considering the following observations: (a) simple functions form a dense subset of the space of Borel measurable functions [?, Theorem 1.5.5(b)]; and (b) the integral is a continuous map from that space [?, Theorem 1.6.2]. Hence, from (a) and (b), it follows that (2) also holds for any Borel measurable function f . This completes the proof. ■

Another reputed folklore theorem, which is often referred to as the “proportional measures” theorem, establishes the explicit forms of the Radon-Nikodym derivatives between two measures, in which one is proportional to the other.

Theorem 3 (Proportional Measures): Let P and Q be two σ -finite measures on the measurable space (Ω, \mathcal{F}) , such that for all $\mathcal{A} \in \mathcal{F}$

$$Q(\mathcal{A}) = cP(\mathcal{A}), \quad (9)$$

with $c > 0$. Then, for all $x \in \Omega$

$$\frac{dP}{dQ}(x) \stackrel{a.s.}{=} \frac{1}{c}, \text{ and } \frac{dQ}{dP}(x) \stackrel{a.s.}{=} c. \quad (10)$$

Proof: First, note that (9) implies that the measures P and Q are mutually absolutely continuous. Hence, from Theorem 1, it follows that for all $\mathcal{A} \in \mathcal{F}$,

$$P(\mathcal{A}) = \int_{\mathcal{A}} dP(x) = \int_{\mathcal{A}} \frac{dP}{dQ}(x) dQ(x). \quad (11)$$

On the other hand, the equality (9) also implies

$$P(\mathcal{A}) = \frac{1}{c} Q(\mathcal{A}) = \int_{\mathcal{A}} \frac{1}{c} dQ(x). \quad (12)$$

Hence, it follows directly from Theorem 1 that the Radon-Nikodym derivative $\frac{dP}{dQ}$ is unique almost surely with respect to Q . Thus, $\frac{dP}{dQ}(x) \stackrel{a.s.}{=} \frac{1}{c}$. Using similar arguments and the fact that P and Q are mutually absolutely continuous, it is verified that $\frac{dQ}{dP}(x) \stackrel{a.s.}{=} c$. ■

In the case in which $c = 1$ in (10), measures P and Q are identical, thus, $\frac{dP}{dQ}(x) \stackrel{a.s.}{=} \frac{dQ}{dP}(x) \stackrel{a.s.}{=} 1$.

The following folklore theorem is often referred to as the “chain rule”.

Theorem 4 (Chain Rule): Let P , Q , and R be three measures on the measurable space (Ω, \mathcal{F}) such that $P \ll Q$; $Q \ll R$; and Q and R are σ -finite measures. Then,

$$\frac{dP}{dR}(x) \stackrel{a.s.}{=} \frac{dP}{dQ}(x) \frac{dQ}{dR}(x). \quad (13)$$

Proof: From the assumptions of the theorem, it follows that for all $\mathcal{A} \in \mathcal{F}$,

$$P(\mathcal{A}) = \int_{\mathcal{A}} dP(x) = \int_{\mathcal{A}} \frac{dP}{dQ}(x) dQ(x) \quad (14)$$

$$= \int_{\mathcal{A}} \frac{dP}{dQ}(x) \frac{dQ}{dR}(x) dR(x) \quad (15)$$

$$= \int_{\mathcal{A}} \frac{dP}{dR}(x) dR(x), \quad (16)$$

where the second equality in (14) follows from Theorem 1; and the equality in (15) follows from Theorem 2. The equality in (16) holds from Theorem 1 and by noticing that $P \ll R$. Therefore, the equalities in (15) and (16) together with Theorem 1 imply (13), which completes the proof. ■

The following folklore theorem shows the connection between the Radon-Nikodym derivative and its multiplicative inverse.

Theorem 5 (Multiplicative Inverse): Let P and Q be two mutually absolutely continuous measures on the measurable space (Ω, \mathcal{F}) ; and assume that for all $x \in \Omega$, $\frac{dQ}{dP}(x) > 0$. Then,

$$\frac{dP}{dQ}(x) \stackrel{a.s.}{=} \left(\frac{dQ}{dP}(x) \right)^{-1}. \quad (17)$$

Proof: From Theorem 4, it follows that

$$\frac{dP}{dQ}(x) \frac{dQ}{dP}(x) \stackrel{a.s.}{=} \frac{dQ}{dQ}(x) \stackrel{a.s.}{=} 1, \quad (18)$$

where the last equality follows from Theorem 3, with $c = 1$. This completes the proof. ■

The subsequent folklore theorem establishes the linearity of the Radon-Nikodym derivative.

Theorem 6 (Linearity): Let P be a σ -finite measure on (Ω, \mathcal{F}) and let also Q_1, Q_2, \dots, Q_n be finite measures on (Ω, \mathcal{F}) absolutely continuous with respect to P . Let c_1, c_2, \dots, c_n be positive reals; and let S be a finite measure on (Ω, \mathcal{F}) such that for all $\mathcal{A} \in \mathcal{F}$, $S(\mathcal{A}) = \sum_{t=1}^n c_t Q_t(\mathcal{A})$. Then,

$$\frac{dS}{dP}(x) \stackrel{a.s.}{=} \sum_{t=1}^n c_t \frac{dQ_t}{dP}(x). \quad (19)$$

Proof: The proof starts by noticing that, from the assumptions of the theorem, it holds that $S \ll P$. Hence, for all $\mathcal{A} \in \mathcal{F}$, it holds that

$$\int_{\mathcal{A}} \frac{dS}{dP}(x) dP(x) = \int_{\mathcal{A}} dS(x) = \sum_{t=1}^n c_t Q_t(\mathcal{A}) \quad (20)$$

$$= \sum_{t=1}^n \int_{\mathcal{A}} c_t dQ_t(x) = \sum_{t=1}^n \int_{\mathcal{A}} c_t \frac{dQ_t}{dP}(x) dP(x) \quad (21)$$

$$= \int_{\mathcal{A}} \sum_{t=1}^n c_t \frac{dQ_t}{dP}(x) dP(x), \quad (22)$$

where the first equality in (20) and the last equality in (21) follow from Theorem 2; and the equality (22) follows from the additivity property of the integral [?, Corollary 1.6.4]. The proof ends by using Theorem 1, which implies the equality in (19) from (22). ■

The following folklore theorem establishes the continuity of the Radon-Nikodym derivative.

Theorem 7 (Continuity): Let P be a σ -finite measure on (Ω, \mathcal{F}) , and let Q_1, Q_2, \dots be an infinite sequence of σ -finite measures on (Ω, \mathcal{F}) , converging to a measure Q . Suppose that for all $n \in \mathbb{N}$, $Q_n \ll P$. Then, $Q \ll P$ and

$$\lim_{n \rightarrow \infty} \frac{dQ_n}{dP}(x) \stackrel{a.s.}{=} \frac{dQ}{dP}(x). \quad (23)$$

Proof: From the assumptions of the theorem, for all $\mathcal{A} \in \mathcal{F}$, it holds that

$$Q(\mathcal{A}) = \lim_{n \rightarrow \infty} Q_n(\mathcal{A}) \quad (24)$$

$$= \lim_{n \rightarrow \infty} \int_{\mathcal{A}} \frac{dQ_n}{dP}(x) dP(x) \quad (25)$$

$$= \int_{\mathcal{A}} \lim_{n \rightarrow \infty} \frac{dQ_n}{dP}(x) dP(x), \quad (26)$$

where the equality in (25) follows from Theorem 2; and the equality in (26) follows from [?, Theorem 1.6.2]. The equality in (26) implies that $Q \ll P$. Hence, for all $\mathcal{A} \in \mathcal{F}$, it holds that

$$Q(\mathcal{A}) = \int_{\mathcal{A}} \frac{dQ}{dP}(x) dP(x). \quad (27)$$

Therefore, the equalities in (26) and (27) jointly with Theorem 1 imply equation (23), which completes the proof. ■

The ensuing folklore theorem establishes the relation between the Radon-Nikodym derivative of a product measure with respect to its component measures.

Theorem 8 (Product of Measures): For all $i \in \{1, 2\}$, let P_i and Q_i be a finite and a σ -finite measure on $(\Omega_i, \mathcal{F}_i)$, respectively; with $P_i \ll Q_i$. Let also $P_1 P_2$ and $Q_1 Q_2$ be the product measures on $(\Omega_1 \times \Omega_2, \mathcal{F}_1 \times \mathcal{F}_2)$ formed by P_1 and P_2 ; and Q_1 and Q_2 , respectively. Then,

$$\frac{dP_1 P_2}{dQ_1 Q_2}(x_1, x_2) \stackrel{a.s.}{=} \frac{dP_1}{dQ_1}(x_1) \frac{dP_2}{dQ_2}(x_2). \quad (28)$$

Proof: From the assumptions of the theorem, for all $\mathcal{A} \in (\Omega_1 \times \Omega_2)$,

$$P_1 P_2(\mathcal{A}) = \int_{\mathcal{A}} dP_1 P_2(x_1, x_2) \quad (29)$$

$$= \int \int_{\mathcal{A}_{x_2}} dP_1(x_1) dP_2(x_2) \quad (30)$$

$$= \int \int_{\mathcal{A}_{x_2}} \frac{dP_1(x_1)}{dQ_1} dQ_1(x_1) dP_2(x_2) \quad (31)$$

$$= \int \int_{\mathcal{A}_{x_2}} \frac{dP_1(x_1)}{dQ_1} \frac{dP_2(x_2)}{dQ_2} dQ_1(x_1) dQ_2(x_2) \quad (32)$$

$$= \int_{\mathcal{A}} \frac{dP_1}{dQ_1}(x_1) \frac{dP_2}{dQ_2}(x_2) dQ_1 Q_2(x_1, x_2), \quad (33)$$

where \mathcal{A}_{x_2} is the section of the set \mathcal{A} determined by x_2 , namely, $\mathcal{A}_{x_2} \triangleq \{x_1 \in \Omega_1 : (x_1, x_2) \in \mathcal{A}\}$; the equality in (29) arises from the definition of $P_1 P_2$ as the product of P_1 and P_2 ; the equality in (31) is a direct consequence of Theorem 2; the equality in (32) follows from Theorem 1; and finally, the equality in (33) is due to the construction of $Q_1 Q_2$ as the product measure of Q_1 and Q_2 .

The proof follows by observing that from the equality in (33), it holds that $P_1 P_2 \ll Q_1 Q_2$. Thus, for all $\mathcal{A} \in \mathcal{F}_1 \times \mathcal{F}_2$,

$$P_1 P_2(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_1 P_2}{dQ_1 Q_2}(x_1, x_2) dQ_1 Q_2(x_1, x_2). \quad (34)$$

The equalities in (33) and (34), together with Theorem 1, imply the equality in (28), which completes the proof. ■

IV. ADVANCED FOLKLORE THEOREMS

This section requires some additional notation. In particular, denote by $\Delta(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$, or simply $\Delta(\mathcal{X})$, the set of all probability measures on the measurable space $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$, where $\mathcal{F}_{\mathcal{X}}$ is a σ -algebra on \mathcal{X} . Using this notation, conditional probability measures can be defined as follows.

Definition 1 (Conditional Probability): A family $P_{Y|X} \triangleq (P_{Y|X=x})_{x \in \mathcal{X}}$ of elements of $\Delta(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ indexed by \mathcal{X} is said to be a conditional probability measure if, for all sets $\mathcal{A} \in \mathcal{F}_{\mathcal{Y}}$, the map

$$\mathcal{X} \rightarrow [0, 1] \quad (35)$$

$$x \mapsto P_{Y|X=x}(\mathcal{A}) \quad (36)$$

is Borel measurable. The set of such conditional probability measures is denoted by $\Delta(\mathcal{Y}|\mathcal{X})$.

A conditional probability $P_{Y|X} \in \Delta(\mathcal{Y}|\mathcal{X})$ and a probability measure $P_X \in \Delta(\mathcal{X})$ determine two unique probability measures in $\Delta(\mathcal{X} \times \mathcal{Y})$ and $\Delta(\mathcal{Y} \times \mathcal{X})$, respectively. These probability measures are denoted by P_{XY} and P_{YX} , respectively, and for all sets $\mathcal{A} \in \mathcal{F}_{\mathcal{X}} \times \mathcal{F}_{\mathcal{Y}}$, it follows that

$$P_{XY}(\mathcal{A}) = \int P_{Y|X=x}(\mathcal{A}_x) dP_X(x), \quad (37)$$

where \mathcal{A}_x is the section of the set \mathcal{A} determined by x , namely,

$$\mathcal{A}_x \triangleq \{y \in \mathcal{Y} : (x, y) \in \mathcal{A}\}. \quad (38)$$

Alternatively, for all sets $\mathcal{B} \in \mathcal{F}_{\mathcal{Y}} \times \mathcal{F}_{\mathcal{X}}$, it follows that

$$P_{YX}(\mathcal{B}) = \int P_{Y|X=x}(\mathcal{B}_x) dP_X(x), \quad (39)$$

where \mathcal{B}_x is the section of the set \mathcal{B} determined by x . For all sets $\mathcal{A} \in \mathcal{F}_{\mathcal{X}} \times \mathcal{F}_{\mathcal{Y}}$, let the set $\hat{\mathcal{A}} \in \mathcal{F}_{\mathcal{Y}} \times \mathcal{F}_{\mathcal{X}}$ be such that

$$\hat{\mathcal{A}} = \{(y, x) \in \mathcal{Y} \times \mathcal{X} : (x, y) \in \mathcal{A}\}. \quad (40)$$

Then, from (37) and (39), it holds that

$$P_{XY}(\mathcal{A}) = P_{YX}(\hat{\mathcal{A}}). \quad (41)$$

Using this notation, the notion of marginal probability measures can be introduced as follows.

Definition 2 (Marginals): Given two joint probability measures $P_{XY} \in \Delta(\mathcal{X} \times \mathcal{Y})$ and $P_{YX} \in \Delta(\mathcal{Y} \times \mathcal{X})$, satisfying (41), the marginal probability measures in $\Delta(\mathcal{X})$ and $\Delta(\mathcal{Y})$, denoted by P_X and P_Y , respectively satisfy for all sets $\mathcal{A} \in \mathcal{F}_{\mathcal{X}}$ and for all sets $\mathcal{B} \in \mathcal{F}_{\mathcal{Y}}$,

$$P_X(\mathcal{A}) \triangleq P_{XY}(\mathcal{A} \times \mathcal{Y}) = P_{YX}(\mathcal{Y} \times \mathcal{A}); \text{ and} \quad (42)$$

$$P_Y(\mathcal{B}) \triangleq P_{XY}(\mathcal{X} \times \mathcal{B}) = P_{YX}(\mathcal{B} \times \mathcal{X}). \quad (43)$$

From the total probability theorem [?, Theorem 4.5.2], it follows that for all $\mathcal{A} \in \mathcal{F}_{\mathcal{Y}}$,

$$P_Y(\mathcal{A}) = \int \int_{\mathcal{A}} dP_{Y|X=x}(y) dP_X(x); \quad (44)$$

and for all $\mathcal{B} \in \mathcal{F}_{\mathcal{X}}$,

$$P_X(\mathcal{B}) = \int \int_{\mathcal{B}} dP_{X|Y=y}(x) dP_Y(y). \quad (45)$$

The joint probability measures P_{XY} and P_{YX} can be described via the conditional probability measure $P_{Y|X}$ and the probability measure P_X as in (37) and in (39); or via the conditional probability measure $P_{X|Y} \in \Delta(\mathcal{X}|\mathcal{Y})$ and the marginal probability measure $P_Y \in \Delta(\mathcal{Y})$. More specifically, for all sets $\mathcal{A} \in \mathcal{F}_{\mathcal{X}} \times \mathcal{F}_{\mathcal{Y}}$, it follows that

$$P_{XY}(\mathcal{A}) = \int P_{X|Y=y}(\mathcal{A}_y) dP_Y(y), \quad (46)$$

where \mathcal{A}_y is the section of the set \mathcal{A} determined by y , namely,

$$\mathcal{A}_y \triangleq \{x \in \mathcal{X} : (x, y) \in \mathcal{A}\}. \quad (47)$$

Alternatively, for all sets $\mathcal{B} \in \mathcal{F}_{\mathcal{Y}} \times \mathcal{F}_{\mathcal{X}}$, it follows that

$$P_{YX}(\mathcal{B}) = \int P_{X|Y=y}(\mathcal{B}_y) dP_Y(y), \quad (48)$$

where \mathcal{B}_y is the section set of \mathcal{B} determined by y .

Within this context, the following folklore theorem highlights a property of conditional measures, which is reminiscent of the unit measure axiom in probability theory.

Theorem 9 (Unit Measure): Consider the conditional probability measures $P_{Y|X} \in \Delta(\mathcal{Y}|\mathcal{X})$ and $P_{X|Y} \in \Delta(\mathcal{X}|\mathcal{Y})$; the probability measures $P_Y \in \Delta(\mathcal{Y})$ and $P_X \in \Delta(\mathcal{X})$ that satisfy (44) and (45). Assume that for all $x \in \mathcal{X}$, the probability measure $P_{Y|X=x} \ll P_Y$. Then,

$$\int \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X(x) \stackrel{\text{a.s.}}{=} 1. \quad (49)$$

Proof: For all $\mathcal{A} \in \mathcal{F}_{\mathcal{Y}}$, from (44), it holds that

$$P_Y(\mathcal{A}) = \int \int_{\mathcal{A}} dP_{Y|X=x}(y) dP_X(x) \quad (50)$$

$$= \int \int_{\mathcal{A}} \frac{dP_{Y|X=x}}{dP_Y}(y) dP_Y(y) dP_X(x) \quad (51)$$

$$= \int_{\mathcal{A}} \int \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X(x) dP_Y(y), \quad (52)$$

where the equality in (51) follows from a change of measure (Theorem 2). Moreover, (52) is obtained using Fubini's theorem [?, Theorem 2.6.6]. The proof proceeds by noticing that $P_Y(\mathcal{A}) = \int_{\mathcal{A}} dP_Y(y)$, and thus from Theorem 1 and the equality in (52), the statement in (49) holds. ■

The following folklore theorem is reminiscent of the Bayes rule.

Theorem 10 (Bayes-like rule): Consider the conditional probability measures $P_{Y|X}$ and $P_{X|Y}$; the probability measures P_Y and P_X that satisfy (44) and (45); and the joint probability measures P_{YX} and P_{XY} in (39) and (46) respectively. Let also $P_X P_Y \in \Delta(\mathcal{X} \times \mathcal{Y}, \mathcal{F}_X \times \mathcal{F}_Y)$ and $P_Y P_X \in \Delta(\mathcal{Y} \times \mathcal{X}, \mathcal{F}_Y \times \mathcal{F}_X)$ be the measures formed by the product of the marginals P_X and P_Y . Assume that:

- (a) For all $x \in \mathcal{X}$, $P_{Y|X=x} \ll P_Y$; and
- (b) For all $y \in \mathcal{Y}$, $P_{X|Y=y} \ll P_X$.

Then,

$$\frac{dP_{XY}}{dP_X P_Y}(x, y) \stackrel{a.s.}{=} \frac{dP_{X|Y=y}}{dP_X}(x) \quad (53)$$

$$\stackrel{a.s.}{=} \frac{dP_{Y|X=x}}{dP_Y}(y) \quad (54)$$

$$\stackrel{a.s.}{=} \frac{dP_{YX}}{dP_Y P_X}(y, x). \quad (55)$$

Proof: Note that assumptions (a) and (b) are sufficient for the Radon-Nikodym derivatives of P_{XY} with respect to $P_X P_Y$ and P_{YX} with respect to $P_Y P_X$ to exist. Hence, it follows that for all sets $\mathcal{A} \in \mathcal{F}_X \times \mathcal{F}_Y$,

$$P_{XY}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{XY}}{dP_X P_Y}(x, y) dP_X P_Y(x, y), \quad (56)$$

which follows from Theorem 2. Note also that from (46), it follows that

$$P_{XY}(\mathcal{A}) = \int \int_{\mathcal{A}_y} dP_{X|Y=y}(x) dP_Y(y) \quad (57)$$

$$= \int \int_{\mathcal{A}_y} \frac{dP_{X|Y=y}}{dP_X}(x) dP_X(x) dP_Y(y) \quad (58)$$

$$= \int \int \mathbb{1}_{\mathcal{A}_y}(x) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X(x) dP_Y(y) \quad (59)$$

$$= \int \mathbb{1}_{\mathcal{A}}(x, y) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X P_Y(x, y) \quad (60)$$

$$= \int_{\mathcal{A}} \frac{dP_{X|Y=y}}{dP_X}(x) dP_X P_Y(x, y), \quad (61)$$

where, the set \mathcal{A}_y is defined in (47). Moreover, the equality in (58) follows from Assumption (b) and Theorem 1. Similarly, from (37), it follows that

$$P_{XY}(\mathcal{A}) = \int \int_{\mathcal{A}_x} dP_{Y|X=x}(y) dP_X(x) \quad (62)$$

$$= \int \int_{\mathcal{A}_x} \frac{dP_{Y|X=x}}{dP_Y}(y) dP_Y(y) dP_X(x) \quad (63)$$

$$= \int \int \mathbb{1}_{\mathcal{A}_x}(y) \frac{dP_{Y|X=x}}{dP_Y}(y) dP_Y(y) dP_X(x) \quad (64)$$

$$= \int \int \mathbb{1}_{\mathcal{A}_y}(x) \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X(x) dP_Y(y) \quad (65)$$

$$= \int \mathbb{1}_{\mathcal{A}}(x, y) \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X P_Y(x, y) \quad (66)$$

$$= \int_{\mathcal{A}} \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X P_Y(x, y), \quad (67)$$

where, the set \mathcal{A}_x is defined in (38). Moreover, the equality in (63) follows from Assumption (a) and Theorem 1; and the equality in (65) follows by exchanging the order of integration [?, Theorem 2.6.6]. Finally, from (41), it follows that

$$\begin{aligned} P_{XY}(\mathcal{A}) &= \int_{\hat{\mathcal{A}}} dP_{YX}(y, x) \\ &= \int_{\hat{\mathcal{A}}} \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_Y P_X(y, x) \end{aligned} \quad (68)$$

$$= \int \mathbb{1}_{\hat{\mathcal{A}}}(y, x) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_Y P_X(y, x) \quad (69)$$

$$= \int \int \mathbb{1}_{\hat{\mathcal{A}}_x}(y) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_Y(y) dP_X(x) \quad (70)$$

$$= \int \int \mathbb{1}_{\hat{\mathcal{A}}_y}(x) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_X(x) dP_Y(y) \quad (71)$$

$$= \int \mathbb{1}_{\hat{\mathcal{A}}}(x, y) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_X P_Y(x, y) \quad (72)$$

$$= \int_{\mathcal{A}} \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_X P_Y(x, y), \quad (73)$$

where the set $\hat{\mathcal{A}}$ is defined in (40). By performing a change of measure using Theorem 2 and assumption (a), the equality in (68) is obtained; and the equality in (71) follows by exchanging the order of integration [?, Theorem 2.6.6].

The proof is completed from Theorem 1 and by combining equations (56), (61), (67) and (73), which establish (53), (54) and (55). \blacksquare

Example 1:

From definition ??, we want to show that there exist an $x \in \mathcal{X}$ such that $P_{Y|X=x}(A) > 0$ and $P_X(\{x \in \mathcal{X} : P_{Y|X=x}(A) > 0\}) = 0$

For all $x \in \mathbb{R}$, for $\mathcal{A} \in \mathcal{F}_Y$, for $\mathcal{B} \in \mathcal{F}_Y$ and μ the Lebesgue measure, $P_{Y|X=x}(A)$ and $P_X(B)$ are defined as follows

$$P_{Y|X=x}(A) = \mu(A \cap (-x, x)) \quad (74)$$

$$P_X(B) = \mu(B) \quad (75)$$

Given $A \in \mathcal{B}(\mathbb{R})$ such that $P_Y(A) = 0$, then

$$P_Y(A) = \int P_{Y|X=x}(A) dP_X(x) \quad (76)$$

$$= \int \mu(A \cap (-x, x)) \mathbb{1}_{\{x \in \mathbb{N}\}} dP_X(x) \quad (77)$$

$$= \int_{\mathbb{N}} \mu(A \cap (-x, x)) dP_X(x) \quad (78)$$

$$= \sum_{t=0}^{\infty} \mu(A \cap (-t, t)) \mu(\{t\}) \quad (79)$$

$$= 0 \quad (80)$$

The equality in (77) follows from (44); the equality (80) follows from the definition of the Lebesgue measure of a singleton. Thus, this example constructs $P_{Y|X}$ and P_X such that for some $\mathcal{A} \in \mathcal{B}(\mathbb{R})$ $P_Y(A) = 0$ but $P_{Y|X=x}(A) > 0$ for some $x \in \mathcal{X}$, therefore the Radom-Nikodym derivative of $P_{Y|X} = x$ with respect to P_Y does not exist. This proves that assumption (a) in 10 cannot be weakened and is needed for all $x \in \mathcal{X}$. A example for assumption (b) in 10 can be constructed similarly.

Theorem 11 (Inverse Bayes-like Rule): Consider the conditional probability measures $P_{Y|X}$ and $P_{X|Y}$; and the probability measures P_Y and P_X that satisfy (44) and (45); and the joint probability measures P_{YX} and P_{XY} in (39) and (46) respectively. Assume that:

- (a) For all $x \in \mathcal{X}$, $P_Y \ll P_{Y|X=x}$; and
- (b) For all $y \in \mathcal{Y}$, $P_X \ll P_{X|Y=y}$.

Then,

$$\frac{dP_X P_Y}{dP_{XY}}(x, y) \stackrel{a.s.}{=} \frac{dP_X}{dP_{X|Y=y}}(x) \quad (81)$$

$$\stackrel{a.s.}{=} \frac{dP_Y}{dP_{Y|X=x}}(y) \quad (82)$$

$$\stackrel{a.s.}{=} \frac{dP_Y P_X}{dP_{YX}}(y, x). \quad (83)$$

Proof: The proof follows along the same lines as the proof of Theorem 10. The proof follows by observing that for all measurable sets $\mathcal{A} \in \mathcal{X} \times \mathcal{Y}$, the product measure $P_X P_Y \in \Delta(\mathcal{X} \times \mathcal{Y})$ satisfies

$$\begin{aligned} & P_X P_Y(\mathcal{A}) \\ &= \int_{\mathcal{A}} dP_X P_Y(x, y) \end{aligned} \quad (84)$$

$$= \int \int_{\mathcal{A}_y} dP_X(x) dP_Y(y) \quad (85)$$

$$= \int \int_{\mathcal{A}_y} \frac{dP_X}{dP_{X|Y=y}}(x) dP_{X|Y=y}(x) dP_Y(y) \quad (86)$$

$$= \int \int \frac{dP_X}{dP_{X|Y=y}}(x) \mathbb{1}_{\{x \in \mathcal{A}_y\}} dP_{X|Y=y}(x) dP_Y(y) \quad (87)$$

$$= \int \frac{dP_X}{dP_{X|Y=y}}(x) \mathbb{1}_{\{x \in \mathcal{A}_y\}} dP_{XY}(x, y) \quad (88)$$

$$= \int_{\mathcal{A}} \frac{dP_X}{dP_{X|Y=y}}(x) dP_{XY}(x, y), \quad (89)$$

where the set \mathcal{A}_y is defined in (47); the equality in (??) follows from Assumption (b) and [?, Theorem 2.2.3]; and the measure P_{XY} is defined in (46).

The proof proceeds by noticing that for all measurable sets $\mathcal{A} \in \mathcal{X} \times \mathcal{Y}$, the product measure $P_X P_Y \in \Delta(\mathcal{X} \times \mathcal{Y})$ also satisfies

$$\begin{aligned} & P_X P_Y(\mathcal{A}) \\ &= \int_{\mathcal{A}} dP_X P_Y(x, y) \end{aligned} \quad (90)$$

$$= \int \int_{\mathcal{A}_y} dP_X(x) dP_Y(y) \quad (91)$$

$$= \int \int_{\mathcal{A}_x} dP_Y(y) dP_X(x) \quad (92)$$

$$= \int \int_{\mathcal{A}_x} \frac{dP_Y}{dP_{Y|X=x}}(y) dP_{Y|X=x}(y) dP_X(x) \quad (93)$$

$$= \int \int \frac{dP_Y}{dP_{Y|X=x}}(y) \mathbb{1}_{\{y \in \mathcal{A}_x\}} dP_{Y|X=x}(y) dP_X(x) \quad (94)$$

$$= \int \frac{dP_Y}{dP_{Y|X=x}}(y) \mathbb{1}_{\{y \in \mathcal{A}_x\}} dP_{YX}(y, x) \quad (95)$$

$$= \int_{\hat{\mathcal{A}}} \frac{dP_Y}{dP_{Y|X=x}}(y) dP_{YX}(y, x) \quad (96)$$

$$= \int_{\mathcal{A}} \frac{dP_Y}{dP_{Y|X=x}}(y) dP_{XY}(x, y), \quad (97)$$

where the sets \mathcal{A}_x , \mathcal{A}_y , and $\hat{\mathcal{A}}$ are defined in (38), (47), and (40); and the measure P_{YX} is defined in (39). The equality in (??) follows by exchanging the order of the integrals [?, Theorem 2.6.6]; the equality in (??) follows from Assumption (a) and [?, Theorem 2.2.3].

The proof is completed by noticing that from (??) and (??), the following equalities hold:

$$P_{X,Y}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_X}{dP_{X|Y=y}}(x) dP_{XY}(x, y), \quad (98)$$

$$= \int_{\mathcal{A}} \frac{dP_Y}{dP_{Y|X=x}}(y) dP_{XY}(x, y), \quad (99)$$

which together with [?, Theorem 2.2.3] implies the equality in (??) almost surely with respect to the measure $P_{XY} \in \Delta(\mathcal{X} \times \mathcal{Y})$ in (37). This completes the proof. ■

Remark 1: An alternative proof for Theorem 9 can be written as follows, by combining Theorems 10 11 and 5.

$$\int \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X(x) \quad (100)$$

$$= \int \frac{dP_{Y|X=x}}{dP_Y}(y) \frac{dP_X}{dP_{X|Y=y}}(x) dP_{X|Y=y}(x) \quad (101)$$

$$= \int dP_{X|Y=y}(x) = 1 \quad (102)$$

Example 2:

This example shows that there exist a $z \in \mathcal{X}$ such that $P_{Y|X=z}(A) = 0$ and $P_Y(A) > 0$.

For all $x \in \mathbb{R}$, for $\mathcal{A} \in \mathcal{F}_Y$, for $\mathcal{B} \in \mathcal{F}_X$ and μ the Lebesgue measure, $P_{Y|X=x}(A)$ and $P_X(B)$ are defined as follows

$$P_{Y|X=x}(A) = \mathbf{1}_{\{A \cap \{(x,y) \in \mathbb{R}^2 : x=y\} \neq \emptyset\}} \quad (103)$$

$$P_X(B) = \mu(B) \quad (104)$$

Given $A \in \mathcal{B}(\mathbb{R})$ such that $P_Y(A) = 0$, then

$$P_Y(A) = \int P_{Y|X=x}(A) dP_X(x) \quad (105)$$

$$= \int \mathbf{1}_{\{A \cup \{(x,y) \in \mathbb{R}^2 : x=y\} \neq \emptyset\}} dP_X(x) \quad (106)$$

$$= \int \mathbf{1}_{\{A \cup \{(x,x)\} \neq \emptyset\}} dP_X(x) \quad (107)$$

$$= \int \mathbf{1}_{\{x \in A\}} dP_X(x) \quad (108)$$

$$= \int_A dP_X(x) \quad (109)$$

$$= \mu(A) \quad (110)$$

Let us define the infimum as follows,

$$\underline{x} = \inf_y \{x \in \mathbb{R} : (x, y) \in A\} \quad (111)$$

Then for all $z \in \mathbb{R}$ such that $z = \underline{x} - \epsilon$ for $\epsilon > 0$,

The equality in (77) follows from (44); the equality (80) follows from the definition of the Lebesgue measure of a singleton. Thus, this example constructs $P_{Y|X}$ and P_X such that for some $\mathcal{A} \in \mathcal{B}(\mathbb{R})$ but $P_{Y|X=x}(A) = 0$ but $P_Y(A) > 0$ for some $x \in \mathcal{X}$, therefore the Radon-Nikodym derivative of P_Y with respect to $P_{Y|X} = x$ does not exist. This proves that assumption (a) in 11 cannot be weakened and is needed for all $x \in \mathcal{X}$. A example for assumption (b) in 11 can be constructed similarly.

Example 3: Let $(\mathbb{R}, \mathcal{B}(\mathbb{R}), P_X)$ be a probability space, with P_X absolutely continuous with respect to the Lebesgue measure μ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. For all $x \in \mathbb{R}$ and for all $\mathcal{A} \in \mathcal{B}(\mathbb{R})$, let

$$P_{Y|X} \triangleq \{P_{Y|X=x} \in \Delta(\mathbb{R}, \mathcal{B}(\mathbb{R})) : x \in \mathbb{R}\} \quad (112)$$

such that:

$$P_{Y|X=x}(\mathcal{A}) = \mathbf{1}_{\{x \in \mathcal{A}\}} \quad (113)$$

Hence, there exists a unique measure on $(\mathbb{R}^2, \mathcal{B}(\mathbb{R}))$ denoted by P_{XY} such that:

$$P_{XY}(\mathcal{A}) = \int_{\mathbb{R}} \int_{\mathcal{A}_y} dP_{Y|X=x}(y) dP_X(x) \quad (114)$$

$$= \int_{\mathbb{R}} P_{Y|X=x}(\mathcal{A}_x) dP_X(x) \quad (115)$$

$$= \int_{\mathbb{R}} \mathbf{1}_{\{(x,x) \in \mathcal{A}\}} dP_X(x) \quad (116)$$

$$= P_X(\{t \in \mathbb{R} : (t, t) \in \mathcal{A}\}). \quad (117)$$

Assuming that $\mathcal{A} = \{(x, y) \in \mathbb{R}^2 : x = y\}$, it follows that

$$P_{XY} = P_X(\{t \in \mathbb{R} : (t, t) \in \mathcal{A}\}) \quad (118)$$

$$= P_X(\mathbb{R}) \quad (119)$$

$$= 1. \quad (120)$$

Moreover, let the measure P_Y on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ be such that for all $\mathcal{A} \in \mathcal{B}(\mathbb{R})$,

$$P_Y(\mathcal{A}) = P_{XY}(\mathbb{R} \times \mathcal{A}) \quad (121)$$

$$= P_X(\{t \in \mathbb{R} : (t, t) \in \mathcal{A}\}) = P_X(\mathcal{A}), \quad (122)$$

hence P_Y is absolutely continuous with respect to μ . Hence,

$$P_X P_Y(\{(x, y) \in \mathbb{R}^2 : x = y\}) \quad (123)$$

$$= \int_{\{(x, y) \in \mathbb{R}^2 : x = y\}} dP_X P_Y(x, y) \quad (124)$$

$$= \int_{\{(x, y) \in \mathbb{R}^2 : x = y\}} \frac{dP_X P_Y}{d\mu\mu}(x, y) d\mu\mu(x, y) \quad (125)$$

$$= \int_{\{(x, y) \in \mathbb{R}^2 : x = y\}} \frac{dP_X}{d\mu}(x) \frac{dP_Y}{d\mu}(y) d\mu\mu(x, y) \quad (126)$$

$$= \int_{\mathbb{R}} \int_{\mathbb{R}} \mathbf{1}_{\{x=y\}} \frac{dP_X}{d\mu}(x) \frac{dP_Y}{d\mu}(y) d\mu(x) d\mu(y) \quad (127)$$

$$= \int_{\mathbb{R}} \int_{\{y\}} \frac{dP_X}{d\mu}(x) \frac{dP_Y}{d\mu}(y) d\mu(x) d\mu(y) \quad (128)$$

$$= 0 \quad (129)$$

where the equality in (??) follows from the fact that for all $y \in \mathbb{R}$, $\mu(\{y\}) = 0$. Thus P_{XY} is not absolutely continuous with respect to $P_X P_Y$.

V. INFORMATION MEASURES

A. KL

The entropy of a given random variable is defined as follows.

Definition 3 (Entropy): Let X be a discrete random variable and assume it induces the probability measure $P_X \in \Delta(\mathcal{X}, \mathcal{F}_X)$ and a sigma finite probability measure $Q \in \Delta(\mathcal{X}, \mathcal{F}_X)$, with $P_X \ll Q$, on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Then, the entropy of X , denoted by $H(X)$, $H(\frac{dP_X}{dQ}(x))$ or $H(P_X)$, is:

$$H(X) \triangleq - \sum_{x \in \text{supp } P_X} \frac{dP_X}{dQ}(x) \log P_X(x). \quad (130)$$

When Q is the counting measure $\frac{dP_X}{dQ}$ is the probability mass function induced by P_X , and when Q is the Lebesgue measure $\frac{dP_X}{dQ}$ is the probability density function induced by P_X .

B. Mutual information

Let $P_{XY} \in \Delta(\mathcal{X} \times \mathcal{Y})$ and $P_{YX} \in \Delta(\mathcal{Y} \times \mathcal{X})$ be two joint probability measures, satisfying (41), such that $P_{XY} \triangleq P_{X|Y} P_X$ and $P_{YX} \triangleq P_{Y|X} P_Y$. Let the marginal probability measures in $\Delta(\mathcal{X})$ and $\Delta(\mathcal{Y})$, be denoted by P_X and P_Y and satisfy for all measurable sets $\mathcal{A} \subseteq \mathcal{X}$ and for all measurable sets $\mathcal{B} \subseteq \mathcal{Y}$, (42) and (43) respectively.. Let two product measures $P_X P_Y \in \Delta(\mathcal{Y} \times \mathcal{X})$ and $P_Y P_X \in \Delta(\mathcal{Y} \times \mathcal{X})$ be the product of the marginals. For all $x \in \mathcal{X}$, the probability measure $P_{Y|X=x}$ is absolutely continuous with respect to P_Y ; and for all $y \in \mathcal{Y}$, the probability measure $P_{X|Y=y}$ is absolutely continuous with respect to P_X .

$$I(P_{Y|X}; P_Y P_X) \triangleq D(P_{YX} \| P_Y P_X) \quad (131)$$

$$= \int \log \left(\frac{dP_{YX}}{dP_Y P_X}(x, y) \right) dP_{YX}(x, y) \quad (132)$$

$$= \int D(P_{Y|X=x} \| P_Y) dP_X(x), \quad (133)$$

$$I(P_{X|Y}; P_X P_Y) \triangleq D(P_{XY} \| P_X P_Y) \quad (134)$$

$$= \int \log \left(\frac{dP_{XY}}{dP_X P_Y}(x, y) \right) dP_{XY}(x, y) \quad (135)$$

$$= \int D(P_{X|Y=y} \| P_X) dP_Y(y). \quad (136)$$

$$D(P_{XY} \| P_X P_Y) \quad (137)$$

$$= \int \log \left(\frac{dP_{XY}}{dP_X P_Y}(x, y) \right) dP_{XY}(x, y) \quad (138)$$

$$= \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) dP_{XY}(x, y) \quad (139)$$

$$= \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) \frac{dP_{XY}}{dP_X P_Y}(x, y) dP_X P_Y(x, y) \quad (140)$$

$$= \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X P_Y(x, y) \quad (141)$$

$$= \int \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X(x) dP_Y(y). \quad (142)$$

Option 1:

$$\begin{aligned} & \int \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X(x) dP_Y(y) \\ &= \int \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) dP_{X|Y=y}(x) dP_Y(y) \end{aligned} \quad (143)$$

$$= \int D(P_{X|Y=y} \| P_X) dP_Y(y). \quad (144)$$

Option 2:

$$\begin{aligned} & \int \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X(x) dP_Y(y) \\ &= \int \int \log \left(\frac{dP_{Y|X=x}}{dP_Y}(y) \right) \frac{dP_{Y|X=x}}{dP_Y}(y) dP_X(x) dP_Y(y) \end{aligned} \quad (145)$$

$$= \int \int \log \left(\frac{dP_{Y|X=x}}{dP_Y}(y) \right) \frac{dP_{Y|X=x}}{dP_Y}(y) dP_Y(y) dP_X(x) \quad (146)$$

$$= \int \int \log \left(\frac{dP_{Y|X=x}}{dP_Y}(y) \right) dP_{Y|X=x}(y) dP_X(x) \quad (147)$$

$$= \int D(P_{Y|X=x} \| P_Y) dP_X(x). \quad (148)$$

Option 3:

$$\begin{aligned} & \int \int \log \left(\frac{dP_{X|Y=y}}{dP_X}(x) \right) \frac{dP_{X|Y=y}}{dP_X}(x) dP_X(x) dP_Y(y) \\ &= \int \int \log \left(\frac{dP_{YX}}{dP_Y P_X}(y, x) \right) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_X(x) dP_Y(y) \end{aligned} \quad (149)$$

$$= \int \int \log \left(\frac{dP_{YX}}{dP_Y P_X}(y, x) \right) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_Y(y) dP_X(x) \quad (150)$$

$$= \int \log \left(\frac{dP_{YX}}{dP_Y P_X}(y, x) \right) \frac{dP_{YX}}{dP_Y P_X}(y, x) dP_Y P_X(y, x) \quad (151)$$

$$= \int \log \left(\frac{dP_{YX}}{dP_Y P_X}(y, x) \right) dP_{YX}(y, x) \quad (152)$$

$$= D(P_{YX} \| P_Y P_X). \quad (153)$$

C. Lautum information

REFERENCES

- [1] R. B. Ash and C. A. Doleans-Dade, *Probability and Measure Theory*, 2nd ed. Burlington, MA, USA: Academic Press, 2000.
- [2] E. L. Lehmann, J. P. Romano, and G. Casella, *Testing Statistical Hypotheses*, 3rd ed. New York, NY, USA: Springer, 2005.