

On the loss of Fisher information in some multi-object tracking observation models

J. Houssineau*, A. Jasra† and S. S. Singh‡

Abstract

The concept of Fisher information can be useful even in cases where the probability distributions of interest are not absolutely continuous with respect to the natural reference measure on the underlying space. Practical examples where this extension is useful are provided in the context of multi-object tracking statistical models. Upon defining the Fisher information without introducing a reference measure, we provide remarkably concise proofs of the loss of Fisher information in some widely used multi-object tracking observation models.

1 Introduction

The Fisher information is a fundamental concept in Statistics and Information Theory (Rissanen 1996), e.g. it features in Jeffreys prior (Jeffreys 1946), the Cramér-Rao lower bound (Cramér 1946, Rao 1992) and in the analysis of the asymptotics of maximum-likelihood estimators (Le Cam 1986, Douc et al. 2004, 2011). Although different generalisations have been proposed, see e.g. (Lutwak et al. 2005, 2012), the standard formulation of the Fisher information often involves a parametric family of probability measures which are all absolutely continuous with respect to a common reference measure in order to define the corresponding probability density functions. This though can be a restrictive assumption for some statistical models.

Let $\Theta \subseteq \mathbb{R}$ be a given open set of parameters and let $\{P_\theta\}_{\theta \in \Theta}$ be a parametric family of probability measures on a Polish space E equipped with its Borel σ -algebra $\mathcal{B}(E)$ and with a reference measure λ . Most often, E is a subset of \mathbb{R}^d for some $d > 0$ and λ is the Lebesgue measure, although Haar measures can be considered more generally for locally-compact topological groups. We will consider the former since the main practical limitation with the usual definition of Fisher information does not come from the lack of natural reference measure but instead from the irregularity of the probability distributions of interest. The usual setting is to assume that for all $\theta \in \Theta$ it holds that P_θ is absolutely continuous with respect to λ , denoted $P_\theta \ll \lambda$. In this case, the probability density function p_θ can be defined as the Radon-Nikodym derivative

$$p_\theta = \frac{dP_\theta}{d\lambda}$$

that is, as the function on E defined uniquely up to a λ -null set by

$$P_\theta(A) = \int \mathbf{1}_A(x) p_\theta(x) \lambda(dx)$$

for all $A \in \mathcal{B}(E)$. In this situation, assuming that p_θ is differentiable with respect to θ , the *score* is defined as $\frac{\partial}{\partial \theta} \log p_\theta(x)$ or indeed $\frac{\partial}{\partial \theta} p_\theta(x) / p_\theta(x)$. Under the final assumption that the score is square

*Department of Statistics and Applied Probability, National University of Singapore. Email:stahje@nus.edu.sg

†Department of Statistics and Applied Probability, National University of Singapore Email:staja@nus.edu.sg

‡Department of Engineering, University of Cambridge and the Alan Turing Institute. Email:sss40@cam.ac.uk

integrable, the Fisher information (Lehmann & Casella 1998) is defined as

$$\mathcal{I}(\theta) = \int \left(\frac{\frac{\partial}{\partial \theta} p_\theta(x)}{p_\theta(x)} \right)^2 p_\theta(x) \lambda(dx). \quad (1)$$

The objective in this article is twofold. For some applications, it is necessary to relax the requirement that $P_\theta \ll \lambda$ holds for all $\theta \in \Theta$, or indeed any θ , and an appropriate definition of $\mathcal{I}(\theta)$ is needed in these cases. Upon addressing this issue, our second objective is then to study the Fisher information of some observation models frequently used in multi-object tracking. Our starting point is the following generalisation of the score $\frac{\partial}{\partial \theta} p_\theta(x)/p_\theta(x)$ given in Heidergott & Vázquez-Abad (2008),

$$\frac{dP'_\theta}{dP_\theta}(x), \quad (2)$$

where P'_θ is the (yet to be formally defined) derivative of the probability measure P_θ with respect to θ and the ratio in eq. (2) is the Radon-Nikodym derivative of P'_θ with respect to P_θ . Heidergott & Vázquez-Abad (2008) introduced this definition of the score in the context of sensitivity analysis for performance measures of Markov chains (Rubinstein & Shapiro 1993). We define the Fisher information using this expression for the score and then study the loss of information in the context of some statistical estimation problems arising in Engineering (see section 2.) Indeed, as shown in proposition 3, when the family P_θ have differentiable densities with respect to the Lebesgue measure, the Fisher information defined using the score in eq. (2) coincides with eq. (1).

The first problem studied in section 2.1 concerns fitting a parametric model to random vectors which are observed through a sensor that randomly permutes the components of the vector. This problem arises in the context of multi-object tracking (Houssineau et al. 2017) where the random vector corresponds to recorded measurements from distinct objects (e.g. vehicles) being tracked using a radar. The radar is able to provide (noisy) measurements of the locations of these object but without knowledge of the association of recorded measurements to the objects themselves. Our analysis involves studying a parametric model that does not have a common dominating measure and through the proposed definition of the Fisher information we provide a simple proof that association uncertainty results in a loss of information. This fact is surprisingly undocumented in the literature despite the numerous articles in Engineering on statistical inference for these types of models.

Multi-object observation models often also include thinning and clutter. Clutter are spurious observations, unrelated to the objects being tracked, generated by radar reflections from non-targets. Thinning is the random deletion of target generated measurements which models the occasional obscuring of targets by obstacles. The augmented set of thinned and spurious observations can be modelled as a spatial point process and section 2.2 concerns fitting a parametric model to a spatial point process that is observed under thinning and superposition. Like random permutation, thinning and superposition results in a loss of information, which is easily shown using the Fisher information defined via eq. (2) and its associated properties. These properties are invoked in the proofs in section 2 but are formally stated and proven in the final section, section 3.

2 Motivating examples

2.1 Random permutation of a random vector

Consider a parametric probability measure P_θ , $\theta \in \Theta \subseteq \mathbb{R}$. For each θ , P_θ is the law of a random vector (X_1, \dots, X_n) where each X_i are in \mathbb{R}^d , i.e. P_θ is a probability measure on $(\mathbb{R}^{dn}, \mathcal{B}(\mathbb{R}^{dn}))$. Assume $n, d \in \mathbb{N}$ are fixed. Let $(X'_1, \dots, X'_n) = (X_{\varsigma(1)}, \dots, X_{\varsigma(n)})$, a random permutation of (X_1, \dots, X_n) , where ς is a random variable with values in the set $\text{Sym}(n)$ of permutations of $\{1, \dots, n\}$. Throughout this section, $x_{1:n}$ denotes the vector (x_1, \dots, x_n) .

In multi-object tracking, each X_i corresponds to a measurement of a distinct object being tracked; there are n of them. The sensor acquiring (X_1, \dots, X_n) , e.g. a radar, returns the vector but with the

association of observations to the n targets lost, which can be modelled as (X'_1, \dots, X'_n) . Filtering for such models has spawned an entire family of algorithms. e.g. see Blackman (1986), Bar-Shalom (1987).

The following theorem shows that the Fisher information $\mathcal{I}'(\theta)$ of the law of $X'_{1:n}$, i.e. after the random permutation, is smaller than the Fisher information $\mathcal{I}(\theta)$ of P_θ . The concept of weak-differentiability will be defined formally in the next section.

Theorem 1. *Assume the family $\{P_\theta\}_{\theta \in \Theta}$ is weakly-differentiable. Then any random permutation of $X_{1:n}$ that is independent of θ incurs a loss of information, that is $\mathcal{I}'(\theta) \leq \mathcal{I}(\theta)$.*

Proof. Let π be the probability distribution of ς on $\text{Sym}(n)$, then a version of the conditional law of $X'_{1:n}$ given $X_{1:n}$ is

$$Q(B_1 \times \dots \times B_n \mid X_{1:n}) = \sum_{\sigma \in \text{Sym}(n)} \pi(\sigma) \prod_{i=1}^n \delta_{X_{\sigma(i)}}(B_i),$$

for any $B_1 \times \dots \times B_n \in \mathcal{B}(\mathbb{R}^{dn})$. The fact that Q does not depend on θ follows from the independence of the random permutation ς from the parameter. From lemma 1 and corollary 1 the score corresponding to the extended model $(X_{1:n}, X'_{1:n})$ can then be expressed as

$$\frac{d(P_\theta \times Q)'}{dP_\theta \times Q}(x_{1:n}, x'_{1:n}) = \frac{dP'_\theta \times Q}{dP_\theta \times Q}(x_{1:n}, x'_{1:n}) \quad (3a)$$

$$= \frac{dP'_\theta}{dP_\theta}(x_{1:n}) \quad (3b)$$

for all $x_{1:n}$ and all $x'_{1:n}$ in \mathbb{R}^{dn} . Note that $(X_{1:n}, X'_{1:n})$ is not absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^{2dn} even when P_θ has a density with respect to the Lebesgue measure. Using the extension of the Fisher identity (see proposition 4), it follows that

$$\frac{d\hat{P}'_\theta}{d\hat{P}_\theta}(X'_{1:n}) = \mathbb{E}_\theta \left(\frac{d(P_\theta \times Q)'}{dP_\theta \times Q}(X_{1:n}, X'_{1:n}) \mid X'_{1:n} \right) \quad (4a)$$

$$= \mathbb{E}_\theta \left(\frac{dP'_\theta}{dP_\theta}(X_{1:n}) \mid X'_{1:n} \right) \quad \text{almost surely,} \quad (4b)$$

with \hat{P}_θ the marginal law of $X'_{1:n}$. Applying Jensen's inequality to the function $y \mapsto y^2$, we conclude that

$$\mathcal{I}'(\theta) \doteq \mathbb{E}_\theta \left(\left(\frac{d\hat{P}'_\theta}{d\hat{P}_\theta}(X'_{1:n}) \right)^2 \right) \leq \mathbb{E}_\theta \left(\left(\frac{dP'_\theta}{dP_\theta}(X_{1:n}) \right)^2 \right) = \mathcal{I}(\theta),$$

which concludes the proof of the theorem. \square

Remark 1. A different proof of this result has been proposed in Houssineau et al. (2017) using the standard formulation of Fisher information. However the proof presented here is remarkably concise and less tedious thanks to the possibility of defining in eq. (3) the score of the extended parametric model $(X_{1:n}, X'_{1:n})$ which does not have a common dominating measure. The final result then follows from the identity in eq. (4) and Jensen's inequality.

It is not possible to establish a strict information loss in general, e.g. if P_θ is symmetrical or if θ is related to some summary statistics that is not affected by random permutation. Additional assumption that guarantee a strict loss are given in Houssineau et al. (2017).

2.2 Thinning and superposition of point processes

Spatial point processes are important in numerous applications (Baddeley et al. 2006), e.g. Forestry (Stoyan & Penttinen 2000) and Epidemiology (Elliot et al. 2000). In addition, point process models are widely used in formulating multi-object tracking problems (Mahler 2007) as they naturally account for an unknown number of objects which are observed indirectly without association and under thinning and superposition. We adopt the approach of the previous section but now characterise the Fisher information of a family of point process parametrized by $\theta \in \Theta$ observed under thinning and superposition. (Note the loss of Fisher information in the presence of association uncertainty has already been established in section 2.1.)

Let Φ denote a point process on \mathbb{R}^d with parametrised distribution P_θ on $E = \bigcup_{n \geq 0} \mathbb{R}^{dn}$, with \mathbb{R}^0 denotes an arbitrary isolated point representing the absence of points in the process. A realisation from P_θ is a random vector (x_1, \dots, x_n) where both the number of points n and their locations $x_i \in \mathbb{R}^d$ are random. However, point-process distributions on \mathbb{R}^d are not always absolutely continuous with respect to the corresponding Lebesgue measure. In particular, the distribution of a non-simple point process, which is a point process such that there is a positive probability of two or more points of its realisation, say x_i and x_j of (x_1, \dots, x_n) , being identical; see Schoenberg (2006) for a discussion about non-simple point processes and examples, e.g. by duplicating the points in a realisation as discussed further below. Assuming that the family $\{P_\theta\}_{\theta \in \Theta}$ is weakly-differentiable, the Fisher information $\mathcal{I}_\Phi(\theta)$ corresponding to the parametrised distribution of Φ can then be expressed as

$$\mathcal{I}_\Phi(\theta) = \sum_{n \geq 0} \pi_\theta(n) \int \left(\frac{dP'_\theta}{dP_\theta}(x_1, \dots, x_n) \right)^2 P_\theta(d(x_1, \dots, x_n) | n), \quad (5)$$

where π_θ is a probability mass function on \mathbb{N}_0 characterising the number of points N in Φ and where $P_\theta(\cdot | n)$ is the conditional distribution of the location of the points in Φ given that the number of points is n (which is supported by \mathbb{R}^{dn}). A straightforward example is when Φ is an independently identically distributed point process. Its distribution factorises as

$$P_\theta(B_1 \times \dots \times B_n) = \pi_\theta(n) \prod_{i=1}^n \mu_\theta(B_i)$$

for any $B_1, \dots, B_n \in \mathcal{B}(\mathbb{R}^d)$ and any $n \in \mathbb{N}_0$, where μ_θ is a probability measure on \mathbb{R}^d . Using the product rule of corollary 1 the expression of the Fisher information simplifies in the independently identically distributed case to

$$\begin{aligned} \mathcal{I}_\Phi(\theta) &= \mathcal{I}_N(\theta) + \sum_{n \geq 0} n^2 \pi_\theta(n) \int \left(\frac{d\mu'_\theta}{d\mu_\theta}(x) \right)^2 \mu_\theta(dx) \\ &= \mathcal{I}_N(\theta) + \mathbb{E}(N^2) \mathcal{I}_X(\theta) \end{aligned} \quad (6)$$

where X is a random variables with distribution μ_θ .

Example 1. A trivial construction of a non-simple point process can be obtained from an independently identically distributed point process Φ by duplicating its realisation. The resulting point process, denoted Φ_2 , has each point of Φ present twice. The Fisher information of Φ_2 can be expressed with the proposed formulation in spite of the lack of absolute continuity with respect to to the reference measure on E . Indeed, the law P_θ^+ of the point process Φ_2 is

$$P_\theta^+(B_1 \times \dots \times B_{2n}) = \pi_\theta(n) \sum_{\sigma \in \text{Sym}(2n)} \prod_{i=1}^n \bar{\mu}_\theta(B_{\sigma(2i-1)} \times B_{\sigma(2i)})$$

and $P_\theta^+(\mathbb{R}^{d(2n+1)}) = 0$, where $\bar{\mu}_\theta$ a probability measure supported by the diagonal of $\mathbb{R}^d \times \mathbb{R}^d$ such that $\bar{\mu}_\theta(B \times B') = \mu_\theta(B \cap B')$ for any $B, B' \in \mathcal{B}(\mathbb{R}^d)$. One can verify that $\bar{\mu}'_\theta(B \times B') = \mu'_\theta(B \cap B')$

so that

$$\frac{d\bar{\mu}'_{\theta}}{d\bar{\mu}_{\theta}}(x, x') = \begin{cases} \frac{d\mu'_{\theta}}{d\mu_{\theta}}(x) & \text{if } x = x' \\ 0 & \text{otherwise,} \end{cases}$$

from which it follows that $\mathcal{I}_{\Phi_2}(\theta) = \mathcal{I}_{\Phi}(\theta)$, that is, duplicating each point in the point process Φ does not change the Fisher information. In the context of parameter inference, this is in agreement with the natural approach of removing the duplicate points before estimating θ .

Returning now to a general point process Φ which is not necessarily independently identically distributed. For each $\alpha \in [0, 1]$, let Φ_{α} denote the thinned version of Φ where each point of its realisation is retained independently of the other points with probability α . In multi-object tracking, an independently thinned point processes arises because a radar can fail to return a credible observation for an object in its surveillance region.

Theorem 2. *Let Φ be a point process characterised by a weakly-differentiable family of probability distributions parametrised by Θ , then $\mathcal{I}_{\Phi}(\theta) \geq \mathcal{I}_{\Phi_{\alpha}}(\theta)$ holds for any $\alpha \in [0, 1]$. If $\mathcal{I}_{\Phi}(\theta) > 0$ then the inequality is strict when $\alpha < 1$.*

Proof. The probability distribution Q_{α} of the thinned point process Φ_{α} given Φ can be expressed as

$$Q_{\alpha}(B_1 \times \cdots \times B_k \mid x_1, \dots, x_n) = \sum_{I \subseteq \{1, \dots, n\}: |I|=k} \alpha^k (1 - \alpha)^{n-k} \prod_{i \in I} \delta_{x_i}(B_{s(i)})$$

for any $B_1, \dots, B_k \in \mathcal{B}(\mathbb{R}^d)$, any $x_1, \dots, x_n \in \mathbb{R}^d$ and any integers n, k such that $k \leq n$, with $s(i) = |\{1, \dots, i\} \cap I|$ so that i is the $s(i)$ th element of I . We obtain from the Fisher identity that the score associated with the point process Φ_{α} with law $P_{\theta, \alpha}$ verifies

$$\begin{aligned} \frac{dP'_{\theta, \alpha}}{dP_{\theta, \alpha}}(x_1, \dots, x_k) &= \mathbb{E} \left(\frac{dP'_{\theta} \times Q_{\alpha}}{dP_{\theta} \times Q_{\alpha}}(\Phi, \Phi_{\alpha}) \mid \Phi_{\alpha} = (x_1, \dots, x_k) \right) \\ &= \mathbb{E} \left(\frac{dP'_{\theta}}{dP_{\theta}}(\Phi) \mid \Phi_{\alpha} = (x_1, \dots, x_k) \right), \end{aligned}$$

where the use of Φ as an argument of point-process distributions is possible because of the irrelevance of the points' ordering. The proof of $\mathcal{I}_{\Phi}(\theta) \geq \mathcal{I}_{\Phi_{\alpha}}(\theta)$ can now be concluded using the decomposition in (5) and invoking Jensen's inequality as in theorem 1. The proof of the strict inequality is deferred to the Appendix. \square

The decrease of the Fisher information demonstrated in theorem 2 can be quantified in the special case of an independently identically distributed point process as follows.

Proposition 1. *Let Φ be an independently identically distributed point process characterised by a weakly-differentiable family of probability distributions parametrised by $\theta \in \Theta$ and assume its cardinality distribution $\pi_{\theta} = \{\pi_{\theta}(n) : n \in \mathbb{N}_0\}$ does not depend on θ , then*

$$(\mathcal{I}_{\Phi_{\alpha}}(\theta) - \mathcal{I}_{\Phi_{\alpha'}}(\theta)) / \mathcal{I}_X(\theta) = ((\alpha - \alpha') - (\alpha^2 - \alpha'^2))\mathbb{E}(N) + (\alpha^2 - \alpha'^2)\mathbb{E}(N^2) \geq 0$$

for any $0 \leq \alpha' \leq \alpha \leq 1$.

Proof. The parameter θ of the distribution π_{θ} is omitted in this proof as a consequence of the assumption of independence. Additionally, thinning does not affect the common distribution of the points in Φ so that, from (6), both point processes have $\mathcal{I}_N(\theta) = 0$ and their $\mathcal{I}_X(\theta)$ terms are equal. Thus, denoting N_{α} the random number of points in Φ_{α} , the objective is to show that $\mathbb{E}(N_{\alpha}^2)$ is greater than $\mathbb{E}(N_{\alpha'}^2)$. It holds that the distribution π_{α} of N_{α} verifies

$$\pi_{\alpha}(n) = \sum_{k \geq n} \pi(k) \binom{k}{n} \alpha^n (1 - \alpha)^{k-n},$$

for any $n \geq 0$, so that

$$\mathbb{E}(N_\alpha^2) = \sum_{k \geq 0} \pi(k) \sum_{n=0}^k n^2 \binom{k}{n} \alpha^n (1-\alpha)^{k-n}.$$

The second sum in the right hand side can be recognised to be the second moment of Bernoulli random variable so that

$$\begin{aligned} \mathbb{E}(N_\alpha^2) &= \sum_{k \geq n} \pi(k) k \alpha ((k-1)\alpha + 1) \\ &= (\alpha - \alpha^2) \mathbb{E}(N) + \alpha^2 \mathbb{E}(N^2), \end{aligned}$$

from which the result follows. \square

Proposition 1 sheds light on the source of the information loss when applying independent thinning to a point process: the quantity $(\mathcal{I}_{\Phi_\alpha}(\theta) - \mathcal{I}_{\Phi_{\alpha'}}(\theta))/\mathcal{I}_X(\theta)$, which can be seen as a relative loss of Fisher information, is shown to be related to the first and second moments of the random variable associated with the number of points in the process. This is because the operation of thinning applied to the considered type of independently identically distributed point process incurs a loss of information only through the decrease of the number of points.

The focus is now on how information evolves when the points of Φ are augmented with that of another point process which has a distribution not depending on θ . In the context of multi-object observation models, the point process being augmented to Φ are spurious observations called clutter which is unrelated to the objects being tracked, e.g. generated by radar reflections from non-targets. This, combined with the fact that the number of clutter points received is *a priori* unknown, shows that treating clutter as a θ -independent point process is appropriate. Superposition is less straightforward than thinning since the resulting augmented point process will have an altered spatial distribution and cardinality distribution. However, the operation of superposition can be expressed as a Markov kernel that transforms Φ to a new point process Φ' and this Markov kernel is independent of θ . Thus the same approach as in theorem 2 can be applied to show that superposition (in general) also leads to a loss of Fisher information. In the following proposition, $\Phi + \tilde{\Phi}$ stands for the point process resulting from the superposition of Φ with another point process $\tilde{\Phi}$.

Proposition 2. *Let Φ be a point process characterised by a weakly-differentiable family of probability distributions parametrised by Θ and let $\tilde{\Phi}$ be another point process whose conditional distribution given Φ does not depend on θ . Then $\mathcal{I}_\Phi(\theta) \geq \mathcal{I}_{\Phi + \tilde{\Phi}}(\theta)$.*

Proof. Let $\tilde{P}(\cdot | \Phi)$ be the conditional law of $\tilde{\Phi}$ given Φ , then the law of the point process $\Phi + \tilde{\Phi}$ given a realisation (x_1, \dots, x_k) of Φ is

$$Q(B_1 \times \dots \times B_n | x_1, \dots, x_k) = \frac{1}{n!} \sum_{\sigma \in \text{Sym}(n)} \mathbf{1}_{B_{\sigma(1)} \times \dots \times B_{\sigma(k)}}(x_1, \dots, x_k) \tilde{P}(B_{\sigma(k+1)} \times \dots \times B_{\sigma(n)} | x_1, \dots, x_k)$$

for any $B_1, \dots, B_n \in \mathcal{B}(\mathbb{R}^d)$. The desired can be now established by proceeding as in the proof of theorem 2; details are omitted. \square

3 Fisher information via the weak derivative

To start with, the derivative P'_θ has to be defined formally. For this purpose, we consider the following weak form of measure-valued differentiation (Pflug 1992), where the notation $\mu(f)$ is used to denote the integral $\int f(x)\mu(dx)$. Henceforth, the set E will be assumed to be Polish with $\mathcal{B}(E)$ its Borel σ -algebra.

Definition 1. Let $\{\mu_\theta\}_{\theta \in \Theta}$ be a parametric family of finite measures on $(E, \mathcal{B}(E))$, then $\theta \rightarrow \mu_\theta$ is said to be *weakly differentiable* at $\theta \in \Theta$ if there exists a signed finite measure μ'_θ on $(E, \mathcal{B}(E))$ such that

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} (\mu_{\theta+\epsilon}(f) - \mu_\theta(f)) = \mu'_\theta(f)$$

holds for all bounded continuous functions f on E .

Although the signed measure μ'_θ is only characterised by the mass it gives to bounded continuous functions, one can show that this characterisation is sufficient to define μ'_θ on the whole Borel σ -algebra $\mathcal{B}(E)$, see lemma 2 in the Appendix.

Assuming that $\theta \mapsto P_\theta$ has a derivative at $\theta \in \Theta$, that P'_θ is absolutely continuous with respect to P_θ , and that the square of the score is integrable, the Fisher information is defined to be

$$\mathcal{I}(\theta) = \int \left(\frac{dP'_\theta}{dP_\theta}(x) \right)^2 P_\theta(dx).$$

Simple cases where this more versatile definition of Fisher information is useful can be given using Dirac measures on the real line as in the following examples.

Example 2. Consider $\Theta = [0, 1]$, $P_\theta = \theta\delta_{-x} + (1 - \theta)\delta_x$ for some given $x \in E = \mathbb{R}$. Indeed, in this case, P_θ is not absolutely continuous with respect to the natural reference measure on the real line, the Lebesgue measure λ . However,

$$P'_\theta = \delta_{-x} - \delta_x,$$

which is a signed measure and

$$\frac{dP'_\theta}{dP_\theta} = \frac{1}{\theta} \mathbf{1}_{\{-x\}} - \frac{1}{1-\theta} \mathbf{1}_{\{x\}},$$

where the Radon-Nikodym derivative is assumed without loss of generality to be equal to 0 everywhere it is not uniquely defined. It follows from basic calculations that

$$\mathcal{I}(\theta) = \frac{1}{\theta(1-\theta)}.$$

This unsurprisingly is the Fisher information of a Bernoulli experiment with probability of success equal to θ . Example 2 is meant to be an illustrative calculation executing the definition of $\mathcal{I}(\theta)$: indeed the same result can be recovered by simply restricting the domain of definition of P_θ to the set $\{-x, x\}$ for all $\theta \in \Theta$. The following result illustrates a usual setting one would expect both definitions of the Fisher information to coincide.

Proposition 3. For some dominating measure λ , assume $P_\theta \ll \lambda$ for all $\theta \in \Theta$ and let p_θ denote its density. For each x , assume $p_\theta(x)$ is differentiable w.r.t. θ and

$$\left| \frac{\partial}{\partial \theta} p_\theta(x) \right| \leq g(x) \tag{9}$$

for all $\theta \in \Theta$ and λ -almost all $x \in E$ where g is some integrable function on E . Then $\mathcal{I}(\theta) = \mathcal{I}(\theta)$.

Remark 2. The assumption of eq. (9) is often invoked in the analysis of maximum likelihood estimation (Douc et al. 2004, Dean et al. 2014) to interchange the order of integration and differentiation, and thus not unique to us. An alternative to assumption in eq. (9) is to assume that the mapping $\theta \rightarrow \int \left| \frac{\partial}{\partial \theta} p_\theta(x) \right| \lambda(dx) < \infty$ is a continuous function of θ . This will imply

$$\lim_{\epsilon \rightarrow 0} \int \left| \frac{p_{\theta+\epsilon} - p_\theta}{\epsilon} - \frac{\partial}{\partial \theta} p_\theta \right| \lambda(dx) = 0 \tag{10}$$

and thus preserving the conclusion of proposition 3. The proof of eq. (10) follows similarly to that of (Van der Vaart 1998, lemma 7.6).

Proof. Recalling that the probability density function p_θ of P_θ with respect to λ is defined as

$$P_\theta(A) = \int \mathbf{1}_A(x) p_\theta(x) \lambda(dx)$$

for all $A \in \mathcal{B}(E)$, it follows from Leibniz's rule that

$$\begin{aligned} P'_\theta(f) &= \lim_{\epsilon \rightarrow \infty} \frac{1}{\epsilon} \int f(x) (p_{\theta+\epsilon}(x) - p_\theta(x)) \lambda(dx), \\ &= \int f(x) \frac{\partial}{\partial \theta} p_\theta(x) \lambda(dx), \end{aligned}$$

for any bounded continuous mappings f on E , and we conclude that $\frac{\partial}{\partial \theta} p_\theta$ is the Radon-Nikodym derivative of P'_θ with respect to λ . Rewriting the Fisher information $\mathcal{I}(\theta)$ as

$$\mathcal{I}(\theta) = \int \left(\frac{\frac{dP'_\theta}{d\lambda}(x)}{p_\theta(x)} \right)^2 P_\theta(dx) = \int \left(\frac{\frac{\partial}{\partial \theta} p_\theta(x)}{p_\theta(x)} \right)^2 p_\theta(x) \lambda(dx) = \mathcal{I}(\theta).$$

concludes the proof of the proposition. \square

The proposed expression of Fisher information can be easily extended to cases where the parameter θ is vector-valued: each component of the Fisher information matrix can be simply defined based on the partial version of the weak differentiation introduced in definition 1.

Another Polish space F is now considered in order to study the Fisher information for probability measures on product spaces. A function Q on $E \times \mathcal{B}(F)$ is said to be a *signed kernel* from E to F if $Q(x, \cdot)$ is a signed finite measure for all $x \in E$ and if $Q(\cdot, B)$ is measurable for all $B \in \mathcal{B}(F)$ (with \mathbb{R} equipped with the Borel σ -algebra, which will be considered by default). If, in particular, $Q(x, \cdot)$ is a probability measure for all $x \in E$ then Q is said to be a *Markov kernel*. If P is a probability measure on E then we denote by $P \times Q$ the probability measure on $(E \times F, \mathcal{B}(E) \otimes \mathcal{B}(F))$ characterised by $P \times Q(A \times B) = \int \mathbf{1}_A(x) Q(x, B) P(dx)$ for all $A \times B$ in the product σ -algebra $\mathcal{B}(E) \otimes \mathcal{B}(F)$. A family $\{Q_\theta\}_{\theta \in \Theta}$ of Markov kernels from E to F is said to be weakly-differentiable if the measure $Q_\theta(x, \cdot)$ is weakly-differentiable for all $x \in E$ and for all $\theta \in \Theta$; it is additionally said to be *bounded weakly-differentiable* if

$$\sup_g \left| \int g(y) Q'_\theta(x, dy) \right| < \infty,$$

where the supremum is taken over all bounded continuous functions. If the latter condition is satisfied, then Q'_θ is itself a signed kernel (see (Heidergott et al. 2008, theorem 1)). Some technical results are first required.

A formal approach to the weak differentiability of product measures has been considered in Heidergott & Leahu (2010) and we consider here an easily-proved corollary of (Heidergott & Leahu 2010, theorem 6.1).

Corollary 1. *Let $\{P_\theta\}_{\theta \in \Theta}$ be a weakly-differentiable parametric family of probability measures on E and let $\{Q_\theta\}_{\theta \in \Theta}$ be a bounded weakly-differentiable parametric family of Markov kernels from E to F , then*

$$(P_\theta \times Q_\theta)' = P'_\theta \times Q_\theta + P_\theta \times Q'_\theta.$$

Corollary 1 was used at several occasions in the examples of section 2 for the special case where the kernel does not depend on θ , that is $(P_\theta \times Q)' = P'_\theta \times Q$. In these examples, the key argument was the simplification of terms that appear both in the numerator and denominator of the score function, using the following lemma.

Lemma 1. *Let μ and τ be finite signed measures on $(E, \mathcal{B}(E))$ such that $\mu \ll \tau$ and let ν and η be signed kernels from E to F such that $\nu(x, \cdot) \ll \eta(x, \cdot)$ for all $x \in E$, then*

$$\frac{d\mu \times \nu}{d\mu \times \eta}(x, y) = \frac{d\nu(x, \cdot)}{d\eta(x, \cdot)}(y), \quad \frac{d\tau \times \eta}{d\mu \times \eta}(x, y) = \frac{d\tau}{d\mu}(x)$$

for $(\mu \times \eta)$ -almost every $(x, y) \in E \times F$.

Proof. Denoting f the Radon-Nikodym derivative of $\mu \times \nu$ by $\mu \times \eta$, it holds by definition that

$$\mu \times \nu(A \times B) = \int \mathbf{1}_{A \times B}(x, y) f(x, y) \mu \times \eta(d(x, y))$$

for all $A \times B \in \mathcal{B}(E) \otimes \mathcal{B}(F)$, so that

$$\int \mathbf{1}_A(x) \nu(x, B) \mu(dx) = \int \mathbf{1}_A(x) \int \mathbf{1}_B(y) f(x, y) \eta(x, dy) \mu(dx)$$

which implies that, for all $B \in \mathcal{B}(F)$, it holds that

$$\nu(x, B) = \int \mathbf{1}_B(y) f(x, y) \eta(x, dy) \quad (12)$$

for μ -almost every $x \in E$. Since F is a Polish space, there exists a countable collection \mathcal{G} of subsets of F that is a π -system and that is generating $\mathcal{B}(F)$. Equation (12) implies that for all $B \in \mathcal{G}$, there exists a subset E_B of E with full μ -measure such that $\nu(x, B) = \int \mathbf{1}_B(y) f(x, y) \eta(x, dy)$ is true for all $x \in E_B$. Considering the countable intersection $E_{\mathcal{G}} = \bigcap_{B \in \mathcal{G}} E_B$, it follows that the statement of interest is true for all $x \in E_{\mathcal{G}}$ and all $B \in \mathcal{G}$. To prove the equality of the measures defined on each side of eq. (12) it is sufficient to prove their equality on a π -system as demonstrated. We conclude that $f(x, \cdot)$ is also the Radon-Nikodym derivative of $\nu(x, \cdot)$ by $\eta(x, \cdot)$ for μ -almost every x , which proves the first result. The second result can be proved in a similar but simpler way. \square

Now assuming that the interest is in the marginal law \hat{P}_θ of $P_\theta \times Q_\theta$ on $(F, \mathcal{B}(F))$, it is often easier to express \hat{P}_θ as

$$\hat{P}_\theta(B) = P_\theta Q_\theta(B) \doteq \int \mathbf{1}_B(y) Q_\theta(x, dy) P_\theta(dx),$$

for any $B \in \mathcal{B}(F)$. In this case, the score can be computed as in the following proposition.

Proposition 4 (Fisher identity). *Let \hat{P}_θ be the law of a random variable Y from $(\Omega, \Sigma, \mathbb{P})$ to $(F, \mathcal{B}(F))$ defined as the marginal of the law $P_\theta \times Q_\theta$ of (X, Y) on $(E \times F, \mathcal{B}(E) \otimes \mathcal{B}(F))$, and let $\{P_\theta\}_{\theta \in \Theta}$ and $\{Q_\theta\}_{\theta \in \Theta}$ be respectively weakly-differentiable and bounded weakly-differentiable, then*

$$\frac{d\hat{P}'_\theta}{d\hat{P}_\theta}(Y) = \mathbb{E}_\theta \left(\frac{d(P_\theta \times Q_\theta)'}{dP_\theta \times Q_\theta}(X, Y) \middle| Y \right) \quad \text{almost surely} \quad (13)$$

with $\mathbb{E}_\theta(\cdot | Y)$ the conditional expectation for a given $\theta \in \Theta$.

Proof. For any $\theta \in \Theta$, the marginal \hat{P}_θ is simply the probability measure $B \mapsto P_\theta \times Q_\theta(E \times B)$, so that the family $\{\hat{P}_\theta\}_{\theta \in \Theta}$ inherits weak-differentiability from $\{P_\theta\}_{\theta \in \Theta}$ and $\{Q_\theta\}_{\theta \in \Theta}$. The derivative \hat{P}'_θ can then be characterised for all $B \in \mathcal{B}(E)$ by

$$\begin{aligned} \hat{P}'_\theta(B) &= (P_\theta \times Q_\theta)'(E \times B) \\ &= \int \frac{d(P_\theta \times Q_\theta)'}{dP_\theta \times Q_\theta}(x, y) \mathbf{1}_{E \times B}(x, y) P_\theta \times Q_\theta(d(x, y)) \\ &= \int \mathbb{E}_\theta \left(\frac{d(P_\theta \times Q_\theta)'}{dP_\theta \times Q_\theta}(X, Y) \middle| Y = y \right) P_\theta Q_\theta(dy) \end{aligned}$$

Recalling that $\hat{P}_\theta = P_\theta Q_\theta$ concludes the proof of the proposition. \square

The Fisher identity is particularly important when the interest is in the Fisher information with respect to the successive observations of a state space model (Douc et al. 2004, Dean et al. 2014), in which case it is defined as the limit

$$\mathcal{I}(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \int \left(\frac{d\bar{P}'_\theta}{d\bar{P}_\theta}(y_1, \dots, y_n) \right)^2 \bar{P}_\theta(d(y_1, \dots, y_n)),$$

where n refers to the time horizon and where \bar{P}_θ is the stationary distribution of the observation process.

The results of corollary 1 and lemma 1 also lead to the following extension of a known property of Fisher information, involving the Fisher information $\mathcal{I}_{Y|X}(\theta)$ of a random variable Y calculated with respect to the conditional law of Y given another random variable X , defined as

$$\mathcal{I}_{Y|X}(\theta) = \int \mathcal{I}_Y(\theta; x) P(dx),$$

where P is the law of X and where

$$x \mapsto \mathcal{I}_Y(\theta; x) = \int \left(\frac{dQ'_\theta(x, \cdot)}{dQ_\theta(x, \cdot)}(y) \right)^2 Q_\theta(x, dy)$$

is assumed to be a measurable mapping, with Q_θ a Markov kernel identified with the conditional law of Y given X . Note that making the law of X dependent on the parameter θ does not induce any difficulties.

Proposition 5. *Let X and Y be random variables on a common probability space $(\Omega, \Sigma, \mathbb{P})$ whose laws are parametrised by $\theta \in \Theta$, let the family of laws of X be weakly-differentiable, and let the family of laws of Y given X be bounded and weakly-differentiable, then the Fisher information $\mathcal{I}_{X,Y}(\theta)$ corresponding to the law of (X, Y) can be expressed as*

$$\mathcal{I}_{X,Y}(\theta) = \mathcal{I}_{Y|X}(\theta) + \mathcal{I}_X(\theta)$$

where $\mathcal{I}_{Y|X}(\theta)$ and $\mathcal{I}_X(\theta)$ correspond to the random variables $Y|X$ and X respectively.

Proof. Let $\{P_\theta\}_{\theta \in \Theta}$ be the (weakly-differentiable) parametric family of laws of X and let $\{Q_\theta\}_{\theta \in \Theta}$ be the (bounded weakly-differentiable) parametric family of conditional laws of Y given X , then

$$\mathcal{I}_{X,Y}(\theta) = \int \left(\frac{d(P_\theta \times Q_\theta)'}{dP_\theta \times Q_\theta}(x, y) \right)^2 P_\theta \times Q_\theta(d(x, y)).$$

Using corollary 1 and lemma 1, it follows that

$$\mathcal{I}_{X,Y}(\theta) = \int \left(\frac{dP'_\theta}{dP_\theta}(x) + \frac{dQ'_\theta(x, \cdot)}{dQ_\theta(x, \cdot)}(y) \right)^2 P_\theta \times Q_\theta(d(x, y))$$

which concludes the proof of the proposition. \square

A straightforward corollary of proposition 5 can be stated as follows: if X and Y are independent random variables, then $\mathcal{I}_{X,Y}(\theta) = \mathcal{I}_X(\theta) + \mathcal{I}_Y(\theta)$. Note that proposition 5 could also be used to prove theorem 1.

Acknowledgements

S.S. Singh would like to thank Prof. Ioannis Kontoyiannis for helpful remarks. All authors were supported by Singapore Ministry of Education AcRF tier 1 grant R-155-000-182-114. AJ is affiliated with the Risk Management Institute, OR and analytics cluster and the Center for Quantitative Finance at NUS.

Appendix

Proofs and technical details

Proof of strict inequality in theorem 2. Jensen's inequality is strict unless it is applied to a non-strictly-convex function or to a degenerate random variable. In the context of theorem 2, the involved function is $y \mapsto y^2$ so that we only have to verify that the random variable

$$S_\theta(\Phi) = \frac{dP'_\theta}{dP_\theta}(\Phi)$$

is not $\sigma(\Phi_\alpha)$ -measurable:

1. We can rule out $S_\theta(\Phi) = c$ (for some constant c) almost surely as follows: since $\mathbb{E}(S_\theta(\Phi)) = 0$, it follows that $c = 0$. But this violates the assumption that $\mathcal{I}_\Phi(\theta) > 0$.

2. Since $S_\theta(\Phi)$ is not a constant almost surely, there exists a set $A \in \mathcal{B}(\mathbb{R})$ such that $1 > \mathbb{E}(\mathbb{I}_A(S_\theta(\Phi))) > 0$. Then

$$\mathbb{E}(\mathbb{I}_A(S_\theta(\Phi)) \mathbb{I}_{\mathbb{R}^0}(\Phi_\alpha)) = \mathbb{E}(\mathbb{I}_A(S_\theta(\Phi)) \mathbb{E}(\mathbb{I}_{\mathbb{R}^0}(\Phi_\alpha) \mid \Phi)) \quad (15a)$$

$$= \mathbb{E}(\mathbb{I}_A(S_\theta(\Phi))(1 - \alpha)^{|\Phi|}) > 0 \quad (15b)$$

since $(1 - \alpha)^{|\Phi|} > 0$ almost surely where $|\Phi|$ denotes the number of points in Φ and where $\mathbb{I}_{\mathbb{R}^0}(\Phi_\alpha)$ is the indicator of the event $|\Phi_\alpha| = 0$. We can similarly show that eq. (15) holds with A replaced with A^c . Thus

$$\mathbb{E}(\mathbb{I}_{\mathbb{R}^0}(\Phi_\alpha)) > \mathbb{E}(\mathbb{I}_A(S_\theta(\Phi)) \mathbb{I}_{\mathbb{R}^0}(\Phi_\alpha)) > 0$$

which violates the following fact: Let X and Y be integrable random variables, assume $Y = c$ is an atom of $\sigma(Y)$ and $Y = c$ has positive probability. If X is $\sigma(Y)$ measurable then $\mathbb{E}(\mathbb{I}_A(X) \mathbb{I}_{\{c\}}(Y))$ is either 0 or equal to $\mathbb{E}(\mathbb{I}_{\{c\}}(Y))$.

□

Lemma 2. *If μ be a finite signed measure on a metric space E characterised by the value of $\mu(f)$ for all bounded continuous mappings f on E . Then μ is uniquely defined on $\mathcal{B}(E)$.*

Proof. Let τ be another finite signed measure that is characterised by $\tau(f) = \mu(f)$ for all bounded continuous functions f on E . We first prove that μ and τ agree on the closed subsets of E . Let $\rho(x, y)$ be the metric on E and let $\rho(x, C)$ denote the usual distance between a point x and set C . Let f_ϵ be the continuous function $f_\epsilon(x) = (1 - \rho(x, C)/\epsilon)^+$ for some some closed set C and some $\epsilon > 0$ where g^+ denotes the positive part of a function g . Note that $f_\epsilon(x)$ is a continuous function that approximates $\mathbf{1}_C(x)$ and

$$\mathbf{1}_C(x) \leq f_\epsilon(x) \leq \mathbf{1}_{C^\epsilon}(x)$$

with C^ϵ the ϵ -neighbourhood of C , so that $\eta(f_\epsilon)$ tends to $\eta(C)$ when $\epsilon \rightarrow 0$ for any finite signed measure η . It follows from that relation $\tau(f_\epsilon) = \mu(f_\epsilon)$ that $\tau(C) = \mu(C)$. This result can be extended to $\mu = \tau$ as follows. Noticing that the set $\mathcal{G} = \{B \in \mathcal{B}(E) : \mu(B) = \tau(B)\}$ is a λ -system that contains the closed sets and that the set of closed sets are themselves a π -system (which generates $\mathcal{B}(E)$), we conclude by the π - λ theorem that $\mathcal{B}(E)$ is contained in \mathcal{G} . Thus $\mathcal{B}(E) = \mathcal{G}$ and therefore $\mu(B) = \tau(B)$ for all $B \in \mathcal{B}(E)$. □

References

Baddeley, A., Gregori, P., Mateu, J., Stoica, R. & Stoyan, D. (2006), *Case studies in spatial point process modeling*, Springer.

- Bar-Shalom, Y. (1987), *Tracking and data association*, Academic Press Professional, Inc.
- Blackman, S. S. (1986), *Multiple-target tracking with radar applications*, Dedham, MA, Artech House, Inc.
- Cramér, H. (1946), *Mathematical Methods of Statistics*, Vol. 9, Princeton university press.
- Dean, T. A., Singh, S. S., Jasra, A. & Peters, G. W. (2014), ‘Parameter estimation for hidden Markov models with intractable likelihoods’, *Scandinavian Journal of Statistics* **41**(4), 970–987.
- Douc, R., Moulines, E., Olsson, J. & Van Handel, R. (2011), ‘Consistency of the maximum likelihood estimator for general hidden Markov models’, *The Annals of Statistics* **39**(1), 474–513.
- Douc, R., Moulines, E. & Ryden, T. (2004), ‘Asymptotic properties of the maximum likelihood estimator in autoregressive models with Markov regime’, *The Annals of Statistics* **32**(5), 2254–2304.
- Elliot, P., Wakefield, J. C., Best, N. G. & Briggs, D. J. (2000), *Spatial epidemiology: methods and applications*, Oxford University Press.
- Heidergott, B., Hordijk, A. & Weisshaupt, H. (2008), ‘Derivatives of Markov kernels and their Jordan decomposition’, *Journal of Applied Analysis* **14**(1), 13–26.
- Heidergott, B. & Leahu, H. (2010), ‘Weak differentiability of product measures’, *Mathematics of Operations Research* **35**(1), 27–51.
- Heidergott, B. & Vázquez-Abad, F. J. (2008), ‘Measure-valued differentiation for Markov chains’, *Journal of Optimization Theory and Applications* **136**(2), 187–209.
- Houssineau, J., Singh, S. S. & Jasra, A. (2017), ‘Identification of multi-object dynamical systems: consistency and Fisher information’, arXiv preprint arXiv:1707.04371.
- Jeffreys, H. (1946), An invariant form for the prior probability in estimation problems, in ‘Proceedings of the Royal Society of London. Series A: Mathematical and Physical Sciences’, Vol. 186, The Royal Society, pp. 453–461.
- Le Cam, L. (1986), *Asymptotic methods in statistical decision theory*, Springer-Verlag.
- Lehmann, E. L. & Casella, G. (1998), *Theory of point estimation*, Springer.
- Lutwak, E., Lv, S., Yang, D. & Zhang, G. (2012), ‘Extensions of Fisher information and Stam’s inequality’, *IEEE Transactions on Information Theory* **58**(3), 1319–1327.
- Lutwak, E., Yang, D. & Zhang, G. (2005), ‘Cramér-Rao and moment-entropy inequalities for Renyi entropy and generalized Fisher information’, *IEEE Transactions on Information Theory* **51**(2), 473–478.
- Mahler, R. P. (2007), *Statistical multisource-multitarget information fusion*, Artech House, Inc.
- Pflug, G. C. (1992), ‘Gradient estimates for the performance of Markov chains and discrete event processes’, *Annals of Operations Research* **39**(1), 173–194.
- Rao, C. R. (1992), Information and the accuracy attainable in the estimation of statistical parameters, in ‘Breakthroughs in Statistics’, Springer, pp. 235–247.
- Rissanen, J. J. (1996), ‘Fisher information and stochastic complexity’, *IEEE Transactions on Information Theory* **42**(1), 40–47.
- Rubinstein, R. Y. & Shapiro, A. (1993), *Discrete event systems: Sensitivity analysis and stochastic optimization by the score function method*, John Wiley & Sons Inc.

Schoenberg, F. P. (2006), 'On non-simple marked point processes', *Annals of the Institute of Statistical Mathematics* **58**(2), 223–233.

Stoyan, D. & Penttinen, A. (2000), 'Recent applications of point process methods in forestry statistics', *Statistical Science* pp. 61–78.

Van der Vaart, A. W. (1998), *Asymptotic statistics*, Vol. 3, Cambridge university press.