

# Subjectivity in Application of the Principle of Maximum Entropy

Jan Peter Hessling

Measurement Technology, SP Technical Research Institute, Borås, Sweden  
Email: [peter.hessling@sp.se](mailto:peter.hessling@sp.se)

Received September 24, 2013; revised October 24, 2013; accepted November 1, 2013

Copyright © 2013 Jan Peter Hessling. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. In accordance of the Creative Commons Attribution License all Copyrights © 2013 are reserved for SCIRP and the owner of the intellectual property Jan Peter Hessling. All Copyright © 2013 are guarded by law and by SCIRP as a guardian.

## ABSTRACT

Complete prior statistical information is currently required in the majority of statistical evaluations of complex models. The principle of maximum entropy is often utilized in this context to fill in the missing pieces of available information and is normally claimed to be fair and objective. A rarely discussed aspect is that it relies upon *testable information*, which is never known but estimated, *i.e.* results from processing of raw data. The subjective choice of this processing strongly affects the result. Less conventional posterior completion of information is equally accurate but is computationally superior to prior, as much less information enters the analysis. Our recently proposed methods of lean deterministic sampling are examples of very few approaches that actively promote the use of minimal incomplete prior information. The inherited subjective character of maximum entropy distributions and the often critical implications of prior and posterior completion of information are here discussed and illustrated, from a novel perspective of consistency, rationality, computational efficiency and realism.

**Keywords:** Maximum; Entropy; Bayes; Monte Carlo; Uncertainty; Covariance; Deterministic; Sampling; Testable; Information; Model; Calculation; Simulation

## 1. Introduction

The principle of maximum entropy (PME) can be utilized to determine a probability density distribution function (pdf) from incomplete statistical information. The approach is not limited to determination of prior pdfs in Bayesian estimation, even though that is a common application. It is rather a general recipe how to make known but incomplete statistical information complete with the most ‘fair’ or the weakest possible hypotheses. As such it fits very well into our practice of statistics, where the lack of complete information is the rule rather than the exception. For instance, knowing only the mean and the variance of an uncertain parameter PME results in a normal distribution [1]. The known information must be well defined in a statistical sense, *i.e.* be formulated in terms of statistical expectations  $\langle f(\theta) \rangle$  of functions  $f(\theta)$  of the considered random quantity  $\theta$ . For instance,  $f(\theta) = \theta$  yields the mean and  $f = (\theta - \langle \theta \rangle)^q$  produces the  $q$ -th moment around the mean. Their expectations can be estimated from any set of observations

$\{\theta_k\}$  with estimators [2]  $\hat{f}$  to give *testable information*  $\phi \equiv \langle f(\theta) \rangle \approx \hat{f}(\{\theta_k\})$  [3].

Bayesian estimation generalizes traditional approaches limited to observations by inclusion of prior knowledge. Its claimed advantage or superiority relies heavily upon fair and truthful assignment of the prior pdf. If the applied method (like PME) to determine the prior pdf turns out to be subjective it would degrade the legitimacy of the approach. Indeed, PME is often motivated by its fairness or objectivity [1]. A minimum of supplementary (unknown) statistical information is imposed by, loosely speaking, maximizing the residual randomness as measured with the information entropy introduced by Shannon [4] and further explored by Jaynes [3]. In practice, the procedure does not provide a complete recipe. By starting with known testable information, PME avoids to define its the best form. There is a widely spread practice though, which probably originates from the ubiquitous use of Taylor expansions. Testable information is usually considered in a hierarchy starting from the mean, the covariance, the skewness, and the

kurtosis etc. Various statistical moments around the mean are tested, as if they were terms of a Taylor series. The functions  $g_q = (\theta - \langle \theta \rangle)^q$  are indeed the  $q$ -th order monomial related to a Taylor expansion around the mean. The expectation  $\langle g_q \rangle$  will contribute to the non-linear displacement  $\langle h(\theta) \rangle - h(\langle \theta \rangle)$ , or *scent* [5] of  $q$ -th order, of the model  $h$ . For any other statistic like e.g. covariance, no such simple direct relation holds [6]. Another line of reasoning might be that the mean describes the location of the distribution, the second moment the width, the third the lowest order asymmetry, while the fourth is the lowest order shape indicator etc. We might subjectively claim that these properties (location-width-asymmetry-shape) provide a hierarchy of testable information. This does not imply that the moments themselves drop in magnitude or relevance, with their order. On the contrary, the linearly scaled even moments  $M_{2q} = \langle g_{2q}(\theta) \rangle^{1/2q}$  usually *increase* with the order  $q$ . In the limit  $q \rightarrow \infty$ ,  $M_{2q} \rightarrow \max_k |\theta_k|$ , where  $\theta_k$  are the observations of  $\theta$ .

Given a set of observations it is not evident how they should be processed to provide testable information of the highest possible quality, *i.e.* with minimal residual uncertainty. For instance, which moments  $\langle g_q \rangle$  should be estimated with a quality (variance of estimator) that usually decrease with the order  $q$ ? As will be shown, the selection will directly influence the PME distribution function. It is also not evident why PME should be restricted to prior application, as usually practiced. Posterior utilization might in fact simplify the analysis greatly. After all, PME is a general method with no explicit reference to what the distribution describes, the input or the output of the analysis.

## 2. Method of Maximum Entropy

For a continuous sample space  $\Omega$ , the entropy [4] functional to be maximized in the method of maximum entropy is given by the Lebesgue integral,

$$S = - \int_{\bar{\Omega}} \bar{p}(\bar{\theta}) \log(\bar{p}(\bar{\theta})) d\bar{\theta}. \quad (1.1)$$

The bar-symbol restricts the integration over *distinguishable* outcomes  $\bar{\theta} \in \bar{\Omega} \subset \Omega \ni \theta$ . That accounts for our possible ignorance<sup>1</sup> of not distinguishing distinct outcomes. As there is an obvious contradiction of being aware of ignorance, it is better translated into irrelevance (for our stated problem). The integration over  $\bar{\Omega}$  can be extended to  $\Omega$ , by *locally* measuring the relative density of  $\bar{\Omega}$  with the Lebesgue measure  $m(\theta)$ ,  $d\bar{\theta} = m(\theta)d\theta$ . Consistency requires *transformation invariance* [1] of the probability  $dP$  in the interval  $\bar{\theta}, \bar{\theta} + d\bar{\theta}$ . That is,  $\bar{p}(\bar{\theta})d\bar{\theta} = p(\theta)d\theta$ , giving the

<sup>1</sup>Remark: "Ignorance" will here refer to the aspect of counting outcomes [1], not the concept of entropy [3].

substitution  $\bar{p}(\bar{\theta}) = p(\theta)/m(\theta)$ . This invariance is equivalent of requiring independence [3] of  $dP(\theta)$  on the subjectively chosen parameterization. In fact, that constraint provides a direct method of determining the Lebesgue measure: If  $\theta \rightarrow \psi$ ,  $m(\theta)d\theta = m(\psi)d\psi$ . For instance, ignorance to change of units correspond to the transformation  $\psi(\theta) = \beta\theta$ , where  $\beta$  describes the conversion factor between different units. It yields  $m(\theta)d\theta = m(\psi)d\psi = m(\beta\theta)d(\beta\theta) = \beta m(\beta\theta)d\theta$  or  $m(\theta) \propto 1/\theta$ .

Utilizing  $m(\theta)$  the integral in Equation (1.1) is converted into a Riemann integral,

$$S = - \int_{\Omega} p(\theta) \log \left( \frac{p(\theta)}{m(\theta)} \right) d\theta. \quad (1.2)$$

The optimization is subject to all known testable information  $\phi_i, i = 0, 1, 2, \dots, N$ ,

$$\phi_i \equiv \langle f_i(\theta) \rangle = \int \bar{p}(\bar{\theta}) f_i(\bar{\theta}) d\bar{\theta} = \int p(\theta) f_i(\theta) d\theta, \quad (1.3)$$

$$f_0 = 1, \phi_0 = 1. \quad (1.4)$$

The functions  $f_k$  will for convenience be denoted *test functions*. The mandatory zeroth constraint  $\phi_0$  (Equation (1.4)) is the normalization condition. As pointed out [3], for discrete sample spaces and no degeneracy ( $m = \text{const}$ ), there is a general expression for the maximum of Equation (1.1) in terms of a partition function  $Z$ . It can be generalized to non-constant measures  $m(\theta)$ , and continuous distributions using calculus of variations [7],

$$Z = \int_{\Omega} m(\theta) e^{-\sum_{i=1}^N \lambda_i f_i(\theta)} d\theta. \quad (1.5)$$

The Lagrange multipliers of optimization are implicitly given by the testable information,

$$-\frac{\partial}{\partial \lambda_i} \log Z + \phi_i = 0, \quad i = 1, 2, \dots, N, \quad (1.6)$$

$$\lambda_0 = \log Z. \quad (1.7)$$

The maximum entropy pdf is then given by,

$$p(\theta) = m(\theta) \cdot e^{-\sum_{i=0}^N \lambda_i f_i(\theta)}, \quad \theta \in \Omega. \quad (1.8)$$

### 2.1. Testable Information

The PME solution (Equation (1.8)) does not specify the test functions  $f_k, k \geq 1$ . They are of major importance though since the solution is directly expressed in them. As stated in the introduction there is a convention of setting  $f_1 = \theta$ ,  $f_2 = (\theta - \langle \theta \rangle)^2$  etc. That is a habit, not a prescription. The choice of test functions must consequently be considered to be at our free disposal.

The difficulty or accuracy of estimating  $\langle f_k(\theta) \rangle$  is

dependent on the explicit form of  $f_k$ . Also, the information contained in the observation set  $\{\theta_k\}$  is to a variable degree transferred to the estimate of  $\langle f_k(\theta) \rangle$ . For instance, if  $f(\theta) = \theta$  all observations are equally weighted, but if  $f_q = g_q$ ,

$$g_q(\theta) \equiv (\delta\theta)^q, \delta\theta \equiv \theta - \langle \theta \rangle, \quad (1.9)$$

the exponent  $q$  determines how much different observations contribute. In the asymptotic limit  $q \rightarrow \infty$ , the observation with the largest deviation is much more important than any other (which only contributes to the estimate of the mean). A lot of information is obviously disregarded. The estimator covariance will accordingly be large. Nevertheless, in many situations the range [6]  $\lim_{k \rightarrow \infty} \langle g_{2k}(\theta) \rangle^{(1/2k)} = \max_k |\theta - \langle \theta \rangle|$  is of much larger interest than any other information. The confidence interval is a more general statistic allowing for sample spaces without bound. From the perspective of objectivity, it thus appears difficult to prescribe any specific set  $f_k, k = 1, 2, \dots, N$ , or even their number ( $N$ ). Clearly, the larger amount of data or independent information that is available, the larger  $N$  is allowed without resulting in unacceptable estimator quality, as expressed with its bias and covariance.

To illustrate how subjectivity enters in practice, assume we have gathered a set of prior observations  $\{\theta_k\}$  of a phenomenological constant  $\theta$  contained in a computer model. To calibrate the model [8], *i.e.* determine the optimal values of its parameters and their uncertainties ( $\theta$  being one of them) from observations, Bayesian estimation is applied. That requires complete prior information, *i.e.* knowledge of the prior pdf  $p(\theta)$ . Applying PME starting from  $\{\theta_k\}$ , test functions  $f_k(\theta)$  must be selected. With the choice of  $f_1 = g_1(\theta)$ ,  $f_2 = g_q(\theta)$ , the partition function will according to Equation (1.5) read,

$$Z = \int_{-\infty}^{\infty} m(\theta) e^{-\lambda_1 g_1(\theta) - \lambda_2 g_q(\theta)} d\theta.$$

Since we are not aware of any ignorance, we set  $m(\theta) = 1$ . The condition Equation (1.6) on  $g_1(\theta)$  implies,

$$0 = Z^{-1} \int_{\Omega} \delta\theta e^{-\lambda_1 \delta\theta - \lambda_2 (\delta\theta)^q} d\theta.$$

A fair amount of subjective pragmatism now suggests  $q$  to be limited to even numbers and the support  $\Omega$  to extend to the whole real axis, since that allows us to determine  $\lambda_1$  with ease. The symmetry of the integrand then implies  $\lambda_1 = 0$ . If  $\theta < 0$  is prohibited,  $\lambda_1 \neq 0$ . The current assignment then translates into an approximation of the support  $\Omega$  as well as the factor  $e^{-\lambda_1(\theta - \langle \theta \rangle)} \approx 1$  of  $p(\theta)$ . The remaining Lagrange multiplier  $\lambda_2$  is obtained by re-scaling,

$$Z = \int_{-\infty}^{\infty} e^{-\lambda_2 (\delta\theta)^q} d\theta = \lambda_2^{-1/q} \int_{-\infty}^{\infty} e^{-t^q} dt,$$

$$\langle g_q(\theta) \rangle = -\frac{\partial \log Z}{\partial \lambda_2} = \frac{1}{\lambda_2 q}.$$

The resulting pdf,

$$p(\theta) \propto e^{-g_q(\theta/\Delta)/q}, \Delta \equiv \langle g_q(\theta) \rangle^{1/q}, \quad (1.10)$$

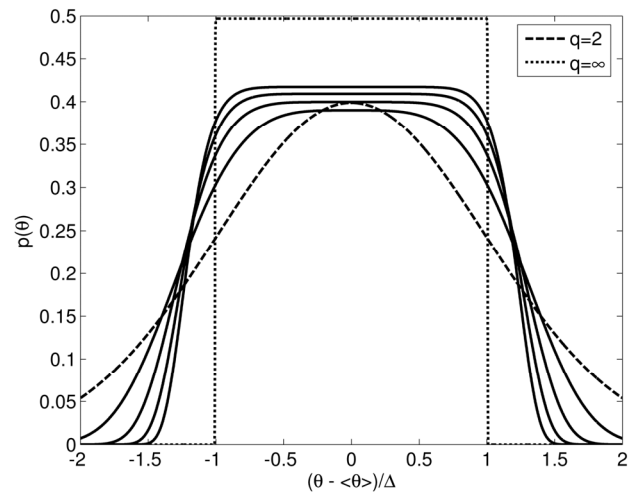
where  $\Delta$  describes the  $q$ -th order ‘‘width’’ of the pdf  $p(\theta)$ , is for different  $q$  displayed in **Figure 1**. Clearly, the PME pdf is to a significant extent controlled by our subjective choice of test function, *i.e.*  $q$ . The result varies between normal ( $q = 2$ ) and uniform ( $q = \infty$ ). More generally, with the understanding of the dependence on the test functions  $f_i$  (Equation (1.8)), almost any PME pdf can be generated. It is only a matter of formulating the question (selecting  $f_i$ ) to obtain the answer ( $p(\theta)$ ) on virtually any desired form.

### 2.2. Quality of Estimation

In practice, a testable piece of information  $\phi_i$  is estimated from raw observations with quality depending on the test function  $f_i$ . The maximum entropy pdf  $p(\theta)$  can be directly formulated in our observations  $\{\theta^{(k)}\}$ , using a specific estimator of the  $q$ -th moment around the mean. An obvious estimator is given by,

$$\hat{S}^{(q)} = B^{(q)} \sum_{k=1}^n \left( \theta^{(k)} - \frac{1}{n} \sum_{k=1}^n \theta^{(k)} \right)^q, \quad (1.11)$$

where  $B^{(q)}$  is an unknown constant which hopefully



**Figure 1.** The maximum entropy pdf  $p(\theta)$  for different kinds ( $q = 2, 4, 6, 8, 10, \infty$ ) of testable information  $\langle f(\theta) \rangle$ , see Equation (1.3). Case  $q$  corresponds to knowledge of the mean ( $f = \theta$ ) and the  $q$ -th statistical moment around the mean ( $f = g_q(\theta)$ , see Equation (1.9)).

can be selected to eliminate the bias, e.g.  $B^{(2)} = 1/(n-1)$ . To express the bias of  $\hat{S}^{(q)}$  in terms of  $B^{(q)}$  expand and calculate its expectation over observations  $\{\theta^{(k)}\}$ ,

$$\langle \hat{S}^{(q)} \rangle = B^{(q)} n \sum_{\substack{\sum_{j=0}^n k_j = q \\ \prod_{j=0}^n k_j!}} \frac{q!}{n} \left(-\frac{1}{n}\right)^{q-k_0} \times \langle \theta^{k_0+k_1} \rangle \prod_{j=2}^n \langle \theta^{k_j} \rangle. \quad (1.12)$$

On the other hand,

$$\langle S^{(q)} \rangle = \langle (\delta\theta)^q \rangle = \sum_{k=0}^q \frac{(-1)^{q-k} q!}{k!(q-k)!} \langle \theta^k \rangle \langle \theta \rangle^{q-k}. \quad (1.13)$$

Clearly, the number of terms of  $\langle S^{(q)} \rangle$  is much less than for  $\langle \hat{S}^{(q)} \rangle$  for  $q > 3$ . Elimination of bias by proper selection of  $B^{(q)}$  requires proportionality between  $\langle S^{(q)} \rangle$  and  $\langle \hat{S}^{(q)} \rangle$ . That cannot in general be achieved since one of them have much more terms than the other. Thus, no scaling of  $\langle \hat{S}^{(q)} \rangle$  can make it universally unbiased. Evaluating the first two coefficients of the expansion in Equation (1.12),

$$\begin{aligned} \langle \hat{S}^{(q)} \rangle &\equiv B^{(q)} n \left( c_q^{(q)} \langle \theta^q \rangle - q \cdot c_{q-1}^{(q)} \langle \theta^{q-1} \rangle \langle \theta \rangle \dots \right) \\ c_q^{(q)} = c_{q-1}^{(q)} &= \sum_{k=0}^{q-2} \frac{q!}{k!(q-k)!} \left(-\frac{1}{n}\right)^k + (q-1) \cdot \left(-\frac{1}{n}\right)^{q-1}. \end{aligned} \quad (1.14)$$

Surprisingly, these two terms satisfies the proportionality required to eliminate the bias. Bias can thus be eliminated by rescaling up to  $q=3$ , but not for  $q \geq 4$  (see estimation of kurtosis in [9]). This suggests the normalization,

$$B^{(q)} \rightarrow \left\{ n \cdot \left[ \sum_{k=0}^{q-2} \frac{q!}{k!(q-k)!} \left(-\frac{1}{n}\right)^k + (q-1) \cdot \left(-\frac{1}{n}\right)^{q-1} \right] \right\}^{-1}. \quad (1.15)$$

While  $B^{(2)} = 1/(n-1)$ ,  $B^{(4)} = 1/(n-4+6/n^2-3/n^3)$ . Clearly,  $B^{(4)}$  bears little resemblance to the conventional form  $1/(n-z)$  of  $B^{(2)}$ , where  $z$  is the number of degrees of freedom.

After failing to obtain an unbiased estimator of the  $q$ -th moment around an unknown mean, we may lower the ambition by assuming the mean  $\langle \theta \rangle$  is known. The corresponding estimator reads,

$$\hat{T}^{(q)} = B^{(q)} \sum_{k=1}^n (\theta_k - \langle \theta \rangle)^q, \quad (1.16)$$

This is a considerably simpler situation for which the normalization  $B^{(q)} = 1/n$  makes the estimator unbiased for all  $q$ . The expected precision of any estimator  $\hat{M}$  describing its typical relative variation may be defined by,

$$\Upsilon(\hat{M}) \equiv \frac{\sqrt{\text{var}(\hat{M})}}{\langle \hat{M} \rangle}. \quad (1.17)$$

The least possible variance of any unbiased estimator  $\tilde{S}^{(q)}$  of  $\langle g_q(\theta) \rangle$  is given by the Cramer-Rao lower bound [2]  $\Lambda^{(q)}$ , evaluated in Appendix A,

$$\Upsilon(\tilde{S}^{(q)}) \geq \frac{q}{\sqrt{2n}} \equiv \Lambda^{(q)}. \quad (1.18)$$

Since  $\Lambda^{(q)}$  increases with  $q$ , it is indeed more difficult to determine higher order moments accurately on an absolute scale, with *any* estimator. The efficiency  $\eta$  of our specific estimator  $\hat{T}^{(q)}$  measures its relative quality,

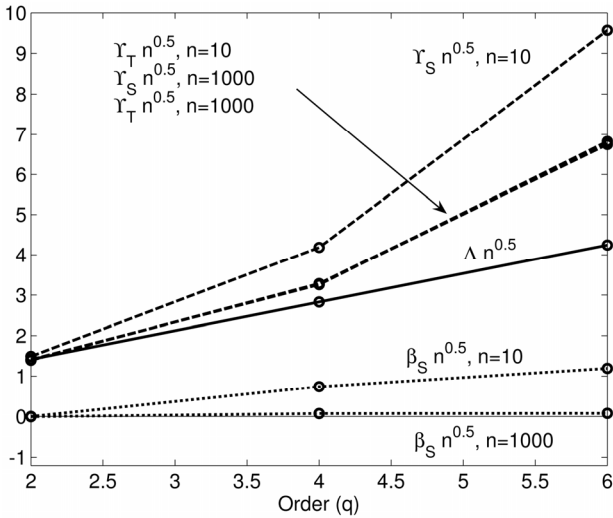
$$\eta^{(q)} \equiv \frac{\Lambda^{(q)}}{\Upsilon(\hat{T}^{(q)})} = \frac{q(q-1)!!}{\sqrt{2[(2q-1)!! - (q-1)!!]}} \leq 1, \quad (1.19)$$

where  $m!! = (m)(m-2)(m-4)\dots 1$ ,  $q$  is even and the precision  $\Upsilon(\hat{T}^{(q)})$  is evaluated in Appendix B. For an *efficient* [2] estimator,  $\eta = 1$ . A low value of  $\eta$  means that the potential for improvement of our estimator is large. Since  $\eta^{(q)}$  decreases rapidly with  $q$ ,  $\hat{T}^{(q)}$  also becomes worse on a relative scale as  $q$  increases. Nevertheless, the actual precision  $\Upsilon$  may or may not be satisfactory. The performance of the estimators  $\hat{S}, \hat{T}$  is illustrated and compared to the Cramer-Rao lower bound in **Figure 2**, by evaluating their relative variation  $\Upsilon$  and bias  $\beta$  numerically with multiple Monte Carlo ensembles.

To conclude, the determination of the PME pdf involved several subjective choices of processing the observations. The shape was here restricted by using monomial test functions  $(g_q)$ . The order  $q$  was limited to be even, to simplify the calculation of the integrals. Likewise, an infinite symmetric sample space  $\Omega$  was assigned. The difficulty of estimating various statistical moments around the mean increases rapidly with the order. If there are not more than a few samples it is in most cases impossible to reliably estimate any other moments than the first and the second. The application of PME thus relied upon rational selection rather than objective deduction.

### 2.3. Prior vs Posterior Application of PME

The PME constructs complete pdfs from incomplete testable information. It can be applied to the input (PME prior), or the output (PME post(-erior)) of the analysis. As no restriction or preference is stated in PME, both alternatives deserves to be considered and evaluated for efficiency and accuracy. The default method of linearization (LIN) [10] as well as the unscented Kalman filter (UKF) [11,12] and deterministic sampling (DS) [6]



**Figure 2.** The relative precision  $\Upsilon_S$  and  $\Upsilon_T$  (dashed, Equation (1.17)) for the estimators  $\hat{S}^{(q)}$  (Equation (1.11)) and  $\hat{T}^{(q)}$  (Equation (1.16)), respectively, compared to the Cramer-Rao lower bound  $\Lambda$  (solid, Equation (1.18)), for  $q = 2, 4, 6$  (connected by lines for clarity) and  $n = 10, 1000$  samples. The relative bias  $\beta^{(q)}$  of  $\hat{S}^{(q)}$  (dotted,  $B^{(q)}$  given by Equation (1.15)) is also shown. The calculations were made for  $10^7$  randomly generated samples of a normal distributed variable, equally split into ensembles of  $n$  samples. Scaling with  $n^{0.5}$  account for the main dependences on the number  $n$  of samples (Equation (1.18)).

in general propagate statistics like covariance which later can be expanded to information related to entire pdfs, e.g. confidence intervals. That is, LIN, UKF, and DS promotes PME post. Monte Carlo simulations [13] or random sampling (RS) must start with complete statistical information to determine its random generator, *i.e.* RS requires PME prior.

An elementary example illustrates the principal differences between PME prior and PME post. Assume it is of interest to calculate the 95% confidence interval  $[h_a, h_b]$  of an uncertain model  $h(\theta)$ , dependent on one parameter  $\theta$ . Let the model uncertainty [8] be given by the mean  $\langle \theta \rangle = 0$  and the variance  $\langle (\delta\theta)^2 \rangle = 1$ . PME prior determines the input pdf  $f(\theta)$  to be normal distributed. Repeated generation of RS ensembles estimates their variation as function of their size: To achieve a standard deviation of the input 95% confidence interval  $[\theta_a, \theta_b]$  for  $\theta$  of about 5%, no less than around 3000 randomly generated samples of  $\theta$  are required. In RS, the model is evaluated for all these samples of  $\theta$  and  $[h_a, h_b]$  is found by evaluating the results  $h(\theta^{(k)})$ ,  $k = 1, 2, \dots, 3000$ , ordering them, and extracting the 95% percentiles. The known statistics  $\langle \theta \rangle = 0$  and  $\langle (\delta\theta)^2 \rangle = 1$  can however be encoded

*exactly* in no more than two(!) deterministic (calculated with a rule) samples  $\theta^{(1)} = -1, \theta^{(2)} = 1$ . Assuming the model is close to linear-in-parameters, the model variance is mainly determined by the parameter variance, here represented by the parameter ensemble  $\theta^{(1)}, \theta^{(2)}$ . In DS, the model variance is then with the argument of consistency given by the second moment around the mean of the model ensemble  $h(\theta^{(k)})$ ,  $k = 1, 2$ , *i.e.*

$$\langle (\delta h)^2 \rangle \equiv \langle (h - \langle h \rangle)^2 \rangle = \left[ (\delta h)^2(\theta^{(1)}) + (\delta h)^2(\theta^{(2)}) \right] / 2.$$

Applying PME post to the result  $\langle h \rangle$  and  $\langle (\delta h)^2 \rangle$  implies that the resulting pdf is normal distributed, analogously to the input pdf in RS. That implies a coverage factor of  $k = 1.96$  resulting in a confidence interval,

$$\begin{aligned} & [h_a, h_b] \\ &= \langle h \rangle + k \sqrt{\frac{(\delta h)^2(\theta^{(1)}) + (\delta h)^2(\theta^{(2)})}{2}} [-1, 1]. \end{aligned} \quad (1.20)$$

For estimating confidence intervals, PME does not in fact need to be utilized at all. By using another DS technique of ours, sampling on confidence boundaries [5], the desired interval can be found directly as,

$$\begin{aligned} & [h_a, h_b] \\ &= \left[ h(\langle \theta \rangle - k \sqrt{\langle (\delta\theta)^2 \rangle}), h(\langle \theta \rangle + k \sqrt{\langle (\delta\theta)^2 \rangle}) \right]. \end{aligned} \quad (1.21)$$

Consistency in evaluating confidence intervals of models here lead to the related concept of confidence boundaries in parameter space. For multivariate models with non-linear dependencies on parameters, both these results needs to be properly generalized [5,6]. Propagating a PME prior probability density function thus typically require several thousands of model samples [13,14] or complex algorithms [15]. The number of model samples propagated with DS for final completion with PME post can be as few as  $n+1$  for  $n$  parameters and increases with the known statistical information, the complexity of the model and acceptable accuracy of evaluation [16].

PME prior requires much more information to be analyzed, than PME post. For the example above, it means 3000 or 2 model evaluations. Plain rationality thus strongly favors PME post. The reason for requesting complete prior information is likely an ambition to find unique unquestionable answers, even if that in practice requires an extensive amount of blind assignments or assumptions. Unique is not equivalent to accurate however, and the quality of any assignment should be critically judged. The loss in efficiency of PME prior compared to PME post is not compensated by superior accuracy. In a number of cases, DS produces the correct result without any error: Any moment of any model

given by a finite Taylor expansion can be exactly calculated with stratified DS [16]. Similarly, the confidence interval of any model  $h(\theta) = \eta(\beta^T \theta)$  of any number ( $n$ ) of parameters can be evaluated exactly [5], where  $\beta \equiv (\beta_1 \ \beta_2 \ \dots \ \beta_n)^T$ , are constants,  $\theta \equiv (\theta_1 \ \theta_2 \ \dots \ \theta_n)$  all parameters, and  $\eta(x)$  is any non-linear monotonic function. For comparison, we are not aware of any general analytic method to propagate a univariate pdf determined by PME prior through any non-linear model. In this case, RS is a numerical method that yield arbitrarily small errors, but at very high computational cost. For “genuinely” multivariate problems RS seizes to be accurate. Multivariate refers to non-trivial finite dependencies of any order, as required for optimal modeling [16]. As far as we know, higher order dependencies (beyond second moments and not normal distributions) can only be implemented in RS by excluding samples. Exceptionally dense sampling is then required to accurately represent sampling *density* over the entire  $n$ -dimensional sampling domain. Nevertheless, it is straight-forward to represent any arbitrary mixed moment in stratified DS [16]. The difficulty is just that the more requirements, the more samples are required. If not enough samples can be afforded it is possible to find the best approximating ensemble where different requirements are given different weights of importance. There are thus several reasons (as exemplified here) for PME post methods to not only be superior to PME prior approaches in efficiency, but also accuracy. The lean set of known input information in PME post is much simpler to analyze or propagate through any model than the complete information in PME prior.

With these observations, the traditional preference to prior instead of posterior completion of information with PME appears biased and strongly subjective. It is highly questionable if PME post methods like DS [6,11,16,17] can be claimed less accurate than state-of-the-art RS relying upon PME prior per se. As both practice *statistical sampling* once defined by Enrico Fermi [18], they are fully comparable. The choice is critical for complex models, as the low efficiency of RS easily render an impossible numerical task. Indeed, that is the main current obstacle for wide utilization of uncertainty quantification with RS.

Bayesian estimation is formulated for PME prior. The combination of the maximum likelihood function and the prior distribution makes Bayes approach superior to traditional approaches limited to the former. In practice though, two *functions* (PME prior) or two *discrete* sets of statistics (PME post) derived from observations are fused. Indeed, the common assumption of Gaussian noise and PME prior derived from the mean and covariance results in combination of covariance matrices (numbers) and not pdfs [8]. A non-trivial posterior distribution *function*

requires non-Gaussian and/or correlated noise and utilization of no less than an infinite set of testable *reliable* information in the PME prior. If not so, Bayes estimation can without any loss of accuracy and relevance be made with PME post instead of PME prior, e.g. with DS instead of RS [8].

PME post (DS) makes it evident that results seldom are unique, while PME prior (RS) conceals ambiguities in more or less dubious or blind assignments in the problem set up. The indefiniteness of the result is an unavoidable consequence of starting with incomplete information, not the analysis. For PME post analyses (DS) this can easily be illustrated, for instance using different ensembles [6]. It is considerably more difficult for PME prior analyses (RS) due to their prohibiting numerical complexity. The contradictory consequence is that PME prior generally are considered more credible than PME post analyses, even if the latter are more honest and realistic as their flaws easily can be illustrated.

## 2.4. Towards Consistency and Rationality

The PME test functions should be selected according to the evaluation of the model, its behavior, and the quality of estimating the corresponding testable information. If we are interested in propagating the covariance of a signal processing model [6], we should consequently use a distribution which has been tested for representation of covariance. That is a *consistent* rather than an objective choice. Objectivity is an impossible target—selecting *any* method is a subjective choice. Our ignorance of  $p(\theta)$  is rather reflecting relevance in perspective of our primary interest, or consistency throughout the analysis. It appears that consistency summarizes the goals of objectivity as well as our ignorance [1,3,4] and emphasizes the context. In contrast to objectivity, consistency expresses a relativity that can indeed be achieved in many situations, as well as measured, questioned and criticized.

The ubiquitous sparsity of statistical information in virtually all practical problems needs to be seriously addressed. If the realism of two approaches are comparable but their efficiency distinctively different, the choice should be based on *rationality* rather than consistency. Realism is distinct from resolution. Any complete set of assumptions will result in arbitrary high resolution, no matter how realistic the assumptions are. As the fidelity of the result never can be enhanced with blind assignments, it is questionable when statistical analyses should be made with distribution functions, rather than the testable information (statistics) these are derived from with PME. If desired, both approaches results in distribution functions. The completion of information is just made before (PME prior) or after (PME post) the analysis. The rational choice is PME post, as it propagates a

minute fraction of the statistical information used in PME prior.

### 3. Conclusions

The character of maximum entropy (PME) distribution functions has been discussed. There are two principal innovations of our study. The PME solutions can to some extent be controlled by how primary data is processed. The prevailing preference for applying PME to the input (prior) and not the output (posterior) of statistical analyses is difficult to justify, as their accuracy is comparable but the latter is computationally superior. The choice appears governed by the method of analysis. Hence, subjectivity enters into the processing of data as well how the analysis is made.

A non-trivial selection enters when observations are processed into testable information. A simple example illustrated common subjective choices, giving different results for identical observations.

Redirecting the focus from the treatment of known information to the targeted evaluation of the analysis emphasizes *consistency*, rather than objectivity. When consistency is indecisive, *rationality* or efficiency of the analysis provides obvious guidance. PME can be applied to find prior as well as posterior distributions. Its unconventional posterior application deserves to be seriously considered, as the analysis involves much less statistical information and is correspondingly more effective, than the prior. Consistency and rationality thus fundamentally questions the prevailing method of completing statistical information prior to the analysis, as in e.g. Monte Carlo simulations.

Maximal consistency and rationality are indeed primary goals of all our proposed methods of deterministic sampling. For complex models, such lean and customized approaches are often required to obtain any measure of modeling quality at all (within acceptable computational time). Without assessment of quality, any (modeling) result is of no value. These aspects are thus of paramount importance to our society where complex calculations (technical, physical, econometrical etc.) are rapidly increasing due to the fast development of computers.

### REFERENCES

- [1] D. S. Sivia and J. Skilling, "Data Analysis—A Bayesian Tutorial," Oxford University Press, Oxford, 2006.
- [2] S. M. Kay, "Fundamentals of Statistical Signal Processing, Estimation Theory," Volume 1, Prentice Hall, Upper Saddle River, 1993.
- [3] E. T. Jaynes, "Information Theory and Statistical Mechanics," *The Physical Review*, Vol. 106, No. 4, 1957, pp. 620-630. <http://dx.doi.org/10.1103/PhysRev.106.620>
- [4] C. E. Shannon, "A Mathematical Theory of Communication," *The Bell System Technical Journal*, Vol. 27, No. 3, 1948, pp. 379-423.
- [5] J. P. HESSLING and T. Svensson, "Propagation of Uncertainty by Sampling on Confidence Boundaries," *International Journal for Uncertainty Quantification*, Vol. 3, No. 5, 2013, pp. 421-444.
- [6] J. P. HESSLING, "Deterministic Sampling for Propagating Model Covariance," *SIAM/ASA Journal on Uncertainty Quantification*, Vol. 1, No. 1, 2013, pp. 297-318.
- [7] G. M. Ewing, "Calculus of Variations with Applications," Dover, New York, 1985.
- [8] J. P. HESSLING, "Identification of Complex Models," in Review, 2013.
- [9] L. Råde and B. Westergren, "Mathematics Handbook," 2 Edition, Studentlitteratur, Lund, 1990.
- [10] ISO GUM, "Guide to the Expression of Uncertainty in Measurement," Technical Report, International Organisation for Standardisation, Geneva, 1995.
- [11] S. Julier and J. Uhlmann, "Unscented Filtering and Non-linear Estimation," *Proceedings IEEE*, Vol. 92, No. 3, 2004, pp. 401-422.
- [12] S. Julier, J. Uhlmann and H. Durrant-Whyte, "A New Approach for Filtering Non-Linear Systems," *American Control Conference*, Seattle, 21-23 June 1995, pp. 1628-1632.
- [13] R. Y. Rubenstein and D. P. Kroese, "Simulation and the Monte Carlo Method," 2 Edition, John Wiley & Sons Inc., New York, 2007.
- [14] J. C. Helton and F. J. Davis, "Latin Hypercube Sampling and the Propagation of Uncertainty in Analyses of Complex Systems," *Reliability Engineering and System Safety*, Vol. 81, No. 1, 2003, pp. 23-69. [http://dx.doi.org/10.1016/S0951-8320\(03\)00058-9](http://dx.doi.org/10.1016/S0951-8320(03)00058-9)
- [15] T. Lovett, "Polynomial Chaos Simulation of Analog and Mixed-Signal Systems: Theory, Modeling Method, Application," Lambert Academic Publishing, Saarbrücken, 2006.
- [16] J. P. HESSLING, "Stratified Deterministic Sampling of Multivariate Statistics," in Preparation, 2013.
- [17] J. P. HESSLING, "Deterministic Sampling for Quantification of Modeling Uncertainty of Signals," In: F. P. G. Márquez and N. Zaman, Eds., *Digital Filters and Signal Processing*, INTECH, Rijeka, 2012.
- [18] N. Metropolis, "The Beginning of the Monte Carlo Method," *Los Alamos Science Special Issue*, Vol. 15, 1987, pp. 125-130.

### Appendix A Cramer-Rao Lower Bounds of Statistical Moments of a Normal Distributed Variable

Assume a random variable  $\theta$  with zero mean is normal distributed,  $\theta \sim N(0, \sigma)$ . By integrating by parts it can be shown that  $M^{(q)} = \sigma^q (q-1)!!$ , where  $(q-1)!! = (q-1)(q-3)(q-5)\dots 1$  is the semi-factorial function. The probability distribution function for  $\theta$  can then be expressed in  $M^{(q)}$ ,

$$p(\theta) \sim N(0, \sigma) \propto \sqrt{\frac{(q-1)!!}{M^{(q)}}} \exp\left\{-\theta^2 \frac{[(q-1)!!]^{2/q}}{2[M^{(q)}]^{2/q}}\right\}. \tag{1.22}$$

For  $n$  independent observations  $\{\theta_k\}, k = 1, 2, \dots, n$ , the likelihood function is given by,

$$p(\{\theta_k\} | M^{(q)}) \propto \prod_{k=1}^n \frac{1}{\sqrt{2\pi}} \sqrt{\frac{(q-1)!!}{M^{(q)}}} \exp\left\{-\theta_k^2 \frac{[(q-1)!!]^{2/q}}{2[M^{(q)}]^{2/q}}\right\} \tag{1.23}$$

$$\propto \left(\frac{(q-1)!!}{M^{(q)}}\right)^{(n/q)} \exp\left\{-\frac{[(q-1)!!]^{2/q}}{2[M^{(q)}]^{2/q}} \sum_{k=1}^n \theta_k^2\right\}.$$

Cramer-Rao lower bound ([2]) now states that for any estimator  $\hat{M}^{(q)}$  of  $M^{(q)}$ ,

$$\text{var}(\hat{M}^{(q)}) \geq F^{-1}(\theta), \tag{1.24}$$

where  $F(\theta)$  is the Fisher information matrix (scalar for one parameter),

$$F(\theta) = -\left\langle \frac{\partial^2 \ln p(\{\theta_k\} | M^{(q)})}{\partial \theta^2} \right\rangle = \frac{2n}{q^2 (M^{(q)})^2}. \tag{1.25}$$

The expected precision  $\Upsilon$  (Equation (1.17)) of the estimator  $M^{(q)}$  hence satisfies,

$$\Upsilon \equiv \frac{\sqrt{\text{var}(\hat{M}^{(q)})}}{\langle M^{(q)} \rangle} \geq \frac{q}{\sqrt{2n}}. \tag{1.26}$$

### Appendix B Variance of Estimator of Statistical Moments around Given Mean

An estimator of the  $q$ -th statistical moment  $M^{(q)} \equiv \langle (\theta - \langle \theta \rangle)^q \rangle$  of  $\theta$  around a known mean  $\langle \theta \rangle$  from a set of  $n$  independent observations  $\{\theta_k\}$  is given by,

$$\hat{M}^{(q)} = B \sum_{k=1}^n (\theta_k - \langle \theta \rangle)^q, \tag{1.27}$$

where the normalization constant  $B$  is chose to minimize its bias  $\langle \hat{M}^{(q)} \rangle - M^{(q)}$ . Since  $\langle \theta \rangle$  is a known constant it is trivially found that  $B = 1/n$  eliminates all bias, for all values of  $q$ . Its variance is found to be,

$$\text{var}(\hat{M}^{(q)}) = \langle [M^{(q)}]^2 \rangle - \langle M^{(q)} \rangle^2 = \frac{M^{(2q)}}{n} - \frac{[M^{(q)}]^2}{n}. \tag{1.28}$$

An explicitly value can be found for a normal distributed parameter,  $\theta \sim N(0, \sigma)$ . Then, by recursively integrating by parts it is found that  $M^{(q)} = \sigma^q (q-1)!!$ , where  $(q-1)!! = (q-1)(q-3)(q-5)\dots 1$  is the semi-factorial function, giving

$$\text{var}(\hat{M}^{(q)}) = \sigma^q \frac{(2q-1)!! - (q-1)!!}{n}. \tag{1.29}$$