# On information measures and prior distributions: a synthesis

Francisco Venegas-Martínez

## Abstract

This paper suggests a new approach to reconciling, in a systematic way, all inferential methods that maximize a specific criterion functional to produce *non-informative* and *informative* priors. In particular, Good's (1968) Minimax Evidence Priors (MEP), Zellner's (1971) Maximal Data Information Priors (MDIP) and Bernardo's (1979) Reference Priors (RP) are seen as special cases of maximizing a more general criterion functional. In a unifying approach Good-Bernardo-Zellner's priors are introduced and applied to a number of Bayesian inference problems, including the Kalman filter and Normal linear model. Moreover, the paper focuses, under plausible conditions, on the existence and uniqueness of the solutions of the derived optimization problems.

## 1  Introduction

The distinctive task in Bayesian inference of deriving priors, in such a way that the inferential content of the data is minimally affected in the posterior, has been of great interest for more than two centuries since the early work of Bayes (1763). More current approaches to this problem, based on the maximization of a specific criterion functional, have been suggested by Good (1968), Zellner (1971) and Bernardo (1979), among others. It is also important to mention that recent literature has included inference procedures to provide a posterior without having a prior, like the Bayesian method of moments (BMOM) introduced by Zellner (1996) and (1998).

   In Good's (1968) principle of maximum invariantized negative cross-entropy, the minimax evidence method of deriving priors was presented for the first time. In this approach the initial density is taken as the square root of Fisher's information. Zellner's (1971) book introduced a method to obtain priors through the maximization of the *total* information about the parameters provided by independent replications of an experiment (prior average information in the data minus the information in the prior). In Bernardo (1979) a procedure was proposed to produce *reference* priors by maximizing the expected information about the parameters provided by independent replications of an experiment (average information in the posterior minus the information in the prior). All of the above methods have comparative and absolute advantages in several respects and have been applied to a large number of inference problems:

   (*i*) While Zellner's method is based on an exact finite sample criterion functional, Good's approach uses a limiting criterion functional, and Bernardo's procedure lies in asymptotic results. In Bernardo's proposal a reference prior (posterior) is defined as the limit of a sequence of priors (posteriors) that maximize finite-sample criteria. In a pragmatic approach in which results are most important, many reference prior algorithms have been developed. For instance, Berger, Bernardo and Mendoza (1989), and Berger and Bernardo (1989), (1992a), (1992b), Bernardo and Smith (1994 , ch. 5), and Bernardo and Ramón (1997).

   (*ii*) The criterion functional used by Bernardo is a cross-entropy, which satisfies a number of remarkable properties, in particular, invariance with respect to one-to-one transformations of the parameters (Lindley 1956). In contrast, the total information functional employed by Zellner is invariant only for the location-scale family and under linear transformations of the parameters. To generate invariance under other relevant transformations, not necessarily one-to-one, side conditions could be needed, as suggested by Zellner (1971).

  (*iii*) These methods have been tested by seeing how well they perform in particular examples. The evaluation is often based on contrasting the derived priors with Jeffreys' (1961) priors, usually improper. Even though improper priors can be associated with unbounded measures consistent with Renyi's (1970) axioms

on probability measures, some technical difficulties remain, see: Box and Tiao (1973), p. 314; Akaike (1978), p. 58; and Berger and Bernardo (1992a), p. 37. It is also important to mention that Jeffreys' priors can lead to singularities producing inadequate results at certain values of the parameters; see Jeffreys (1967, p. 359). Of course, if MEP, MDIP, and RP priors were to be used to contrast the performance of other priors, the former priors could also produce unsatisfactory results under certain circumstances.

In this paper, we attempt to reconcile all inferential methods that maximize a criterion functional to produce *non-informative* and *informative* priors. In our general approach, Good's Minimax Evidence Priors (1968 and 1969), Zellner's Maximal Data Information Priors (1971, 1977, 1991, 1993, 1995, 1996a and 1996b) and Bernardo's Reference Priors (1979 and 1997) are seen as special cases of maximizing a more general indexed criterion functional. Thus, properties of the derived priors will depend on the choice of indexes from a wide range of possibilities, instead of on a few personal points of view with *ad hoc* modifications. In the spirit of Akaike (1978) and Smith (1979), we can say that this will look more like Mathematics than Psychology–without underestimating the importance of the latter in the Bayesian framework. This unified approach will enable us to explore a vast range of possibilities for constructing priors. It is worthwhile to note that our general method extends in a natural way Soofi's (1994) pyramid by adding more vertices and including their convex hull. In any event, a good choice will depend on the specific characteristics of the problem we are concerned with. Needless to say, the chosen method should also provide good predictions.

This work is organized as follows. In section 2, we will introduce an indexed family of information functionals. In section 3, on the basis of asymptotic normality, we will state a relationship between Bernardo's (1979) criterion functional and some members of the indexed family. In section 4, we will study a Bayesian inference problem associated with convex combinations of relevant members of the proposed indexed family. Here, we will introduce Good-Bernardo-Zellner's priors as well as their *controlled* versions as solutions of maximizing discounted entropy. We will pay special attention to the existence and uniqueness of the solution of the corresponding optimization problems. In section 5, we will study Good-Bernardo-Zellner's priors as Kalman Filtering priors. In section 6, we examine the relationship between Good-Bernardo-Zellner's

priors and the Normal linear model. Finally, in section 7, we will draw conclusions, acknowledge limitations, and make suggestions for further research.

## 2    An indexed family of information functionals

In this section, we define an indexed family of information functionals and study some distinguished members. For the sake of simplicity, we will remain in the single parameter case. The extension to the multi-dimensional parameter case will lead to conceptual complications. This is not surprising when dealing with information measures and priors; see Jeffreys (1961), Zellner (1971), Box and Tiao (1973), and Berger and Bernardo (1992a).

Suppose that we wish to make inferences about an unknown parameter $\theta \in \Theta \subseteq \mathbb{R}$ of a distribution $P_\theta$, from which there is available an observation, say, $X$. Assume that $P_\theta$ has density $f(x|\theta)$ (Radon-Nikodym derivative) with respect to some fixed dominating $\sigma$-finite measure $\lambda$ on $\mathbb{R}$ for all $\theta \in \Theta \subseteq \mathbb{R}$, that is, $dP_\theta/d\lambda = f(x|\theta)$ for all $\theta \in \Theta \subseteq \mathbb{R}$, thus $P_\theta(A) = \int_A f(x|\theta)d\lambda(x)$ for all Borel sets $A \subset \mathbb{R}$.

The Bayesian approach is to assume that there is a prior density, $\pi(\theta)$, describing initial knowledge about the likelihood of the values of the parameter, $\theta$. We will assume that $\pi(\theta)$ is a density with respect to some $\sigma$-finite measure $\mu$ on $\mathbb{R}$. Once a prior distribution, $\pi(\theta)$, has been prescribed, then the information provided by the data, $x$, about the parameter is used to modify the initial knowledge, as expressed in $\pi(\theta)$, *via* Bayes' theorem to obtain a posterior distribution of $\theta$, namely, $f(\theta|x) \propto f(x|\theta)\pi(\theta)$ for every $x \in \mathbb{R}$ (using $f$ generically to represent densities). The normalized posterior distribution is then used to make inferences about $\theta$. Let us define an infinite system of nested functionals:

$$(1) \qquad \mathcal{V}_{\gamma,\alpha,\delta}(\pi) = \frac{1}{1-\gamma} \int \pi(\theta) G(\mathcal{I}(\theta), \mathcal{F}(\theta), \gamma, \alpha, \delta) d\mu(\theta)$$

where

$$\begin{aligned} &G(\mathcal{I}(\theta), \mathcal{F}(\theta), \gamma, \alpha, \delta) \\ &= \log \left\{ \frac{\exp\{[\mathcal{F}(\theta)/\mathcal{I}(\theta)]^{1-\delta}[\mathcal{I}(\theta)]^{\frac{1-\gamma}{1+\alpha}} - \delta[\mathcal{I}(\theta)]^{1-\alpha}\}}{\pi(\theta)^{1-\gamma}} \right\}, \end{aligned}$$

$0 \leq \gamma < 1$, $\alpha \in \{0,1\}$, $\delta \in \{0,1\}$, and

(2)
$$\mathcal{I}(\theta) = \int \left( \frac{\partial}{\partial \theta} \log f(x|\theta) \right)^2 f(x|\theta) d\lambda(x)$$

is Fisher's information about $\theta$ provided by an observation $X$ with density $f(x|\theta)$, and

(3)
$$\mathcal{F}(\theta) = \int f(x|\theta) \log f(x|\theta) d\lambda(x)$$

is the negative Shannon's information of $f(x|\theta)$, provided $\mathcal{I}(\theta)$ and $\mathcal{F}(\theta)$ exist. In the case that $n$ independent observations of $X$ are drawn from $P_\theta$, say, $(X_1, X_2, ..., X_n)$, then $\mathcal{I}(\theta)$ and $\mathcal{F}(\theta)$ will still stand for the average Fisher's information and the average negative Shannon's information of $f(x|\theta)$ respectively. It is not unsual to deal with indexed functionals in inference problems about a distribution, as in Good (1968). It is worthwhile pointing out that for each triad $(\gamma, \alpha, \delta)$ taking values in $0 \leq \gamma < 1$, $\alpha \in \{0,1\}$, $\delta \in \{0,1\}$, then $\mathcal{V}_{\gamma,\alpha,\delta}(\pi)$ is a criterion functional that can be used to derive a prior $\pi(\theta)$, $\theta \in \Theta$, belonging to a feasible set $\mathcal{C}$. Usually, $\mathcal{C}$ is defined by constraints in terms of potential values of $\theta$.

Note now that for the location parameter family $f(x|\theta) = f(x - \theta)$, $\theta \in \mathbb{R}$, with the properties

$$\int [f'(x)]^2 / f(x) \ d\lambda(x) < \infty$$

and

$$\int f(x) \log f(x) \ d\lambda(x) < \infty,$$

where $\lambda = \mu$ stands for the Lebesgue measure, we have that both $\mathcal{I}(\theta)$ and $\mathcal{F}(\theta)$ are constant. Observe also that the scale parameter family $f(x|\theta) = (1/\theta)f(x/\theta)$, $\theta > 0$, with the above properties, satisfies the following relationship:

(4)
$$\mathcal{F}(\theta) = \tfrac{1}{2} \log \mathcal{I}(\theta) + \text{constant.}$$

The indexed family in which we will be concerned with is given by:

$$\mathcal{A} = conv[ \ \overline{\{\mathcal{V}_{\gamma,\alpha,\delta}(\pi)\}} \ ]$$
$$=\text{convex hull of the closure of the family} \{\mathcal{V}_{\gamma,\alpha,\delta}(\pi)\}.$$

We readily identify a number of distinguished members of $\mathcal{A}$:

($i$) Criterion for Maximum Entropy Priors (MAXENTP):

$$\mathcal{V}_{0,0,1}(\pi) = -\int \pi(\theta)\log\pi(\theta)d\mu(\theta),$$

which is just Shannon's information measure of a density $\pi(\theta)$, or Jaynes' (1957) criterion functional to derive maximum entropy priors. Notice also that (3) can be rewritten in a simpler way as $\mathcal{F}(\theta) = -\mathcal{V}_{0,0,1}(f(x|\theta))$.

($ii$) Criterion for Minimax Evidence Priors (MEP):

(5)  $\mathcal{V}_{1,1,1}(\pi) \overset{\text{def}}{=} \lim_{\gamma\to1} \mathcal{V}_{\gamma,1,1}(\pi) = -\int \pi(\theta)\log\frac{\pi(\theta)}{p(\theta)}d\mu(\theta) - \log C,$

which is Good's invariantized negative cross-entropy, taking as initial density $p(\theta) = C[\mathcal{I}(\theta)]^{\frac{1}{2}}$ with $C = \{\int[\mathcal{I}(\theta)]^{\frac{1}{2}}d\mu(\theta)\}^{-1}$, provided that $\int[\mathcal{I}(\theta)]^{\frac{1}{2}}d\mu(\theta) < \infty$. We can also write (5) as

(6)        $\mathcal{V}_{1,1,1}(\pi) - \mathcal{V}_{0,0,1}(\pi) = \int \pi(\theta)\log[\mathcal{I}(\theta)]^{\frac{1}{2}}d\mu(\theta).$

($iii$) Criterion for Maximal Data Information Priors (MDIP):

(7)        $\mathcal{V}_{0,0,0}(\pi) = \int\int f(x)f(\theta|x)\log\frac{\ell(\theta|x)}{\pi(\theta)}d\mu(\theta)d\lambda(x),$

which is Zellner's criterion functional in his MDIP approach. Here, as usual,

$$f(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{f(x)}, \qquad f(x) = \int f(x|\theta)\pi(\theta)d\mu(\theta),$$

and $\ell(\theta|x) = f(x|\theta)$ is the likelihood function. An alternative formulation of (7), which is often useful, is given by

(8)          $\mathcal{V}_{0,0,0}(\pi) - \mathcal{V}_{0,0,1}(\pi) = \int \pi(\theta)\mathcal{F}(\theta)d\mu(\theta).$

Some members of $\mathcal{A}$ define new criterion functionals in which the information provided by the sampling model, $\mathcal{I}(\theta)$, plays a role:

($iv$) Criterion for Maximal Modified Data Information Priors (MMDIP):

$$(9) \quad \mathcal{V}_{0,1,0}(\pi) = \int \int f(x)f(\theta|x) \log \frac{[\ell(\theta|x)]^{[\mathcal{I}(\theta)]^{\frac{1}{2}}}}{\pi(\theta)} d\mu(\theta)d\lambda(x),$$

which is the prior average information in the data *modified* by Fisher's information minus the information in the prior. Note that when $\mathcal{I}(\theta)$ is constant, (9) reduces to Zellner's criterion functional (up to a constant factor).

($v$) Criterion for Maximal Fisher Information Priors (MFIP):

$$(10) \qquad \mathcal{V}_{0,1,1}(\pi) = - \int \pi(\theta) \log \frac{\pi(\theta)}{\exp\{[\mathcal{I}(\theta)]^{\frac{1}{2}}\}} d\mu(\theta) - 1,$$

which is the prior average Fisher's information minus the information in the prior.

# 3   Revisiting Bernardo's reference priors

The maximization of Bernardo's (1979) criterion is usually a difficult problem to deal with. In order to get a simpler alternative procedure under specific conditions, we will derive a useful asymptotic approximation between Bernardo's criterion functional (or Lindley's information measure, 1956) and some members of the class $\mathcal{A}$. As stated in Bernardo (1979), the concept of reference prior is very general. However, in order to keep the analysis tractable, we will restrict ourselves to the continuous one-dimensional parameter case.

Suppose that there are available $n$ independent observations, say, $(X_1, X_2, \ldots, X_n)$, of a distribution $P_\theta$, $\theta \in \Theta \subseteq \mathbb{R}$. Accordingly, the random vector $(X_1, X_2, \ldots, X_n)$ has density

$$d\mathcal{P}_\theta/d\nu = f(\xi|\theta) = \prod_{k=1}^{n} f(x_k|\theta),$$

for all $\xi = (x_1, x_2, ..., x_n)$ and all $\theta \in \Theta \subseteq \mathbb{R}$, where

$$\mathcal{P}_\theta = \underbrace{P_\theta \otimes P_\theta \otimes \cdots \otimes P_\theta}_{n} \text{ and } \nu = \underbrace{\lambda \otimes \lambda \otimes \cdots \otimes \lambda}_{n}.$$

Following Lindley (1956), a measure of the expected information about $\theta$ of a sampling model $f(x|\theta)$ provided by a random sample of size $n$ when the prior distribution of $\theta$ is $\pi(\theta)$, is defined to be

$$(11) \qquad \mathcal{L}^{(n)}(\pi) = \int f(\xi) \int f(\theta|\xi) \log \frac{f(\theta|\xi)}{\pi(\theta)} d\mu(\theta) d\nu(\xi).$$

In order to obtain an asymptotic approximation of (11) in terms of $\mathcal{V}_{1,1,1}$ and $\mathcal{V}_{0,0,1}$, we state a limit theorem which justifies the passage of the limit under the integral signs in (11). The theorem rules out the possibility that the *essentials* of the statistical model, $f(\xi|\theta)$, change when samples grow in size. Let us rewrite (11) as:

$$\mathcal{L}^{(n)}(\pi) = \mathcal{V}_{\gamma,0,1}(\pi) + \log\sqrt{n}$$
$$(12) \qquad - \int \int \log\left(\int T_n(\omega) W_n(\omega) d\mu(\omega)\right) f(\xi|\theta)\pi(\theta) d\nu(\xi) d\mu(\theta),$$

where

$$(13) \qquad T_n(\omega) = \frac{f(X_1, X_2, ..., X_n|\theta + \frac{\omega}{\sqrt{n}})}{f(X_1, X_2, ..., X_n|\theta)}$$

and

$$(14) \qquad W_n(\omega) = \frac{\pi(\theta + \frac{\omega}{\sqrt{n}})}{\pi(\theta)}.$$

Throughout the paper, both $\lambda$ and $\mu$ will stand for the Lebesgue measure on $\mathbb{R}$. Also, we will assume that all densities involved are Lebesgue measurable in both arguments, $x$ and $\theta$.

**Theorem 3.1** *Assume that the following conditions hold:*

*(I) $\Theta$ is an open interval in $\mathbb{R}$;*

*(II) The function $\sqrt{f(x|\theta)}$ is absolutely continuous on $\theta$, and*

$$\{x|f(x|\theta) > 0\}$$

*is independent of $\theta$;*

*(III) If $\theta, \theta' \in \Theta$, then $\theta \neq \theta'$ implies $\lambda\{x|f(x|\theta) \neq f(x|\theta')\} > 0$;*

*(IV)* $\frac{\partial}{\partial\theta}\log f(x|\theta)$ *exists for all* $\theta \in \Theta$ *and every* $x$*;*

*(V)* $\mathcal{I}(\theta)$ *is a continuous and bounded function in* $\Theta$*;*

*(VI)* *For all* $\delta > 0$*, and all* $\theta \in \Theta$

$$\int_{B_\delta(\frac{\omega}{\sqrt{n}})} \left( \sqrt{f(x|\theta + \frac{\omega}{\sqrt{n}})} - \sqrt{f(x|\theta)} \right)^2 d\lambda(x) = o(\tfrac{1}{n}),$$

*where*

$$B_\delta(\tfrac{\omega}{\sqrt{n}}) = \{x : |\sqrt{f(x|\theta + \tfrac{\omega}{\sqrt{n}})} - \sqrt{f(x|\theta)}| > \delta\sqrt{f(x|\theta)} \};$$

*(VII)* *There exist* $c > 0$ *and* $\tau > 0$ *such that*

$$\int \left| \pi(\theta + u) - \pi(\theta) \right| d\mu(\theta) \leq c|u|^\tau;$$

*(VIII)* *For all* $\rho > 0$

$$\int_{|\omega| > n^\rho} \left( T_n(\omega)W_n(\omega) - T_n(\omega) \right) d\mu(\omega) \xrightarrow{P} 0;$$

*(IX)* *The sequence of random variables* $\{\log U_n\}_{n=1}^{\infty}$ *where*

$$U_n = \int T_n(\omega)W_n(\omega)d\mu(\omega)$$

*satisfies*

$$\lim_{\varepsilon \to \infty} \sup_{n \geq 1} \int_{|\log U_n| \geq \varepsilon} |\log U_n| dP = 0,$$

*where*

$$P\{\xi \in A, \ \theta \in B\} = \int_B \pi(\theta) \int_A f(\xi|\theta)d\nu(\xi)d\mu(\theta),$$

*for all* $A \in \mathbb{R}^n$ *and* $B \in \Theta$*.*

*Then, as* $n \to \infty$*,*

(15) $\qquad \mathcal{L}^{(n)}(\pi) - \mathcal{V}_{1,1,1}(\pi) = -\mathcal{V}_{0,0,1}(\varphi) + \log C\sqrt{n} + o(1),$

*where* $\varphi(z)$ *is the density of* $Z \sim \mathcal{N}(0,1)$*, and* $C$ *is taken as in (4).*

Some comments are in order: (I)-(IV) are standard regularity conditions, (V) states desirable properties for $\mathcal{I}(\theta)$, (VI) is a bounded variance condition, (VII) is a smoothness condition, (VIII) is a convergence condition, and (IX) says that the sequence $\{\log U_n\}_{n=1}^{\infty}$ is uniformly integrable with respect to $P$. It can be shown that (I)-(VI) lead to

$$(16) \qquad T_n(\omega) \xrightarrow{\mathcal{L}} \exp\big\{\omega\sqrt{\mathcal{I}(\theta)}\big[Z - \tfrac{1}{2}\omega\sqrt{\mathcal{I}(\theta)}\big]\big\},$$

where $Z \sim \mathcal{N}(0,1)$, and (16) along with (VII)-(IX) imply

$$\log U_n = \log \int T_n(\omega)W_n(\omega)d\mu(\omega) \xrightarrow{\mathcal{L}} \log\sqrt{2\pi/\mathcal{I}(\theta)} + \tfrac{1}{2}Z^2,$$

from where the conclusion of the theorem follows. Notice that the right-hand side of (3.5) is independent of $\pi$. Thus, if conditions (I)-(IX) are fulfilled, instead of maximizing $\mathcal{L}^{(\infty)}(\pi)$, which is usually a difficult problem, we have as an alternative procedure maximizing $\mathcal{V}_{1,1,1}(\pi)$, which is independent of $n$. Notice that for maximization purposes the right-hand side of (15) becomes a constant.

Finally, it is worthwhile to note that the location parameter family $f(x|\theta) = f(x - \theta)$, with $\sqrt{f(x)}$ absolutely continuous on $\mathbb{R}$, and $\int [f'(x)]^2/f(x)\ d\lambda(x) < \infty$, fully satisfies the conditions of Theorem 3.1.

## 4  Good-Bernardo-Zellner priors

In this section we introduce Good-Bernardo-Zellner's priors as solutions of convex combination of relevant members of the class $\mathcal{A}$. Very often, there exist priors for which entropy becomes infinite, specially when dealing with the non-informative case. In order to overcome this difficulty, we suggest the concept of discounted entropy. We also introduce Good-Bernardo-Zellner's *controlled* priors as solutions of maximizing discounted entropy. We emphasize the existence and uniqueness of the solutions of the corresponding variational and optimal control problems.

Throughout this section, we will be studying a number of Bayesian inferential problems related to convex combinations of distinctive elements of $\mathcal{A}$. Let

$$\mathcal{M}_\phi(\pi) \stackrel{\text{def}}{=} \phi\mathcal{V}_{1,1,1}(\pi) + (1-\phi)\mathcal{V}_{0,0,0}(\pi),$$

$0 \leq \phi \leq 1$. Plainly, $\mathcal{M}_\phi(\pi) \in \mathcal{A}$. To see that $\mathcal{M}_\phi(\pi)$ is concave w.r.t. $\pi$, it is enough to observe, as in Zellner (1991), that

$$\mathcal{V}_{0,0,0}(\pi(\theta)) = \mathcal{L}^{(1)}(\pi(\theta)) + \mathcal{V}_{0,0,1}(\pi(\theta)) - \mathcal{V}_{0,0,1}(f(x)),$$

is a sum of concave functions w.r.t. $\pi$ (up to the constant $\mathcal{V}_{0,0,1}(f(x))$). Since $\mathcal{V}_{1,1,1}(\pi)$ is concave w.r.t. $\pi$, $\mathcal{M}_\phi(\pi)$ is also concave w.r.t. $\pi$. Zellner (1996b) provides a criterion functional that agrees with $\mathcal{M}_\phi(\pi)$ given by

$$G_\phi[\pi(\theta)]$$
$$= \int \left[ \phi \mathcal{F}(\theta) + (1-\phi)\log[\mathcal{I}(\theta)]^{\frac{1}{2}} \right] \pi(\theta)d\mu(\theta) - \int \pi(\theta)\log\pi(\theta)d\mu(\theta).$$

Indeed, from (5), (6) and (8), we get

$$\begin{aligned} G_\phi[\pi(\theta)] =& \phi\left(\mathcal{V}_{0,0,0} - \mathcal{V}_{0,0,1}\right) + (1-\phi)\left(\mathcal{V}_{1,1,1} - \mathcal{V}_{0,0,1}\right) + \mathcal{V}_{0,0,1} \\ =& \phi\mathcal{V}_{0,0,0} + (1-\phi)\mathcal{V}_{1,1,1} - \mathcal{V}_{0,0,1} + \mathcal{V}_{0,0,1} \\ =& \mathcal{M}_\phi(\pi). \end{aligned}$$

Usually, in the absence of data *supplementary* information, in terms of expectations about the parameter, comes from additional knowledge of the experiment, or from the experience of the experimenter, namely,

$$(17) \qquad \int a_k(\theta)\pi(\theta)d\mu(\theta) = \bar{a}_k, \quad k = 1, 2, ..., s,$$

where both the functions $a_k$ and the constants $\bar{a}_k$, $k = 1, 2, ..., s$, are known. Hereafter, we will assume that (17) does not lead to any contradiction about $\pi(\theta)$. We will now concern with maximizing $\mathcal{M}_\phi(\pi)$ subject to *supplementary* information.

**Proposition 4.1** *Consider the Good-Bernardo-Zellner problem:*

*Maximize $\mathcal{M}_\phi(\pi)$ (with respect to $\pi$)*

*subject to $\mathcal{C} : \displaystyle\int a_k(\theta)\pi(\theta)d\mu(\theta) = \bar{a}_k, \quad k = 0, 1, 2, ..., s, \quad a_0 \equiv 1 = \bar{a}_0.$*

*Then a necessary condition for a maximum is*

$$(18) \qquad \pi_\phi^*(\theta) \propto [\mathcal{I}(\theta)]^{\frac{\phi}{2}} \exp\{(1-\phi)\mathcal{F}(\theta) + \sum_{k=0}^{s} \lambda_k a_k(\theta)\},$$

*where $\lambda_k$, $k = 0, 1, ..., s$, are the Lagrange multipliers associated with the constraints $\mathcal{C}$ (cf. Zellner 1995).*

Notice that when no supplementary information is available, $\pi_\phi^*(\theta)$ is appropiate for an unprejudiced experimenter, otherwise it will be suitable for an informed experimenter who is in favor of $\mathcal{C}$. Observe also that $\pi_1^*(\theta)$ is Good-Bernardo's prior, and $\pi_0^*(\theta)$ is Zellner's prior. Consider the binomial distribution for a single observation, $f(x|\theta) = \theta^x(1-\theta)^{1-x}$, $0 \le \theta \le 1$. In such a case, $\pi_1^*(\theta) \propto \theta^{-\frac{1}{2}}(1-\theta)^{-\frac{1}{2}}$ and $\pi_0^*(\theta) \propto \theta^\theta(1-\theta)^{1-\theta}$ for $\theta \in [0,1]$, which are quite different. Notice that $\pi_1^*(\theta)$ becomes infinite at $\theta = 0$ and $\theta = 1$. On the other hand $\pi_0^*(\theta)$ rises monotonically to 1.6186 at $\theta = 0$ and $\theta = 1$. Yet, another view in this regard (Geisser, 1993) states that when the sample size is fairly large it does not matter which prior is employed, and the uniform prior may as well be used for $\theta$.

**Corolary 4.1** *Consider the location and scale parameter families,*

$$f(x|\theta) = f(x - \theta), \ \theta \in \mathbb{R},$$

*and*

$$f(x|\theta) = (1/\theta)f(x/\theta), \ \theta > 0,$$

*respectively, both satisfying*

$$\int [f'(x)]^2/f(x) \ d\lambda(x) < \infty$$

*and $\int f(x)\log f(x) \ d\lambda(x) < \infty$. Then, Good-Bernardo's and Zellner's priors agree regardless of the value of $\phi \in (0,1)$.*

The proof of the above corollary for the scale parameter case follows from (4). It is important to point out that when there is no supplementary information, we require $\mu(\Theta) < \infty$. Of course, the parameter space $\Theta$ can have bounds as large as needed to consider where the likelihood for $\theta$ is appreciable.

Notice that Proposition 4.1 can be used recursively when there is more supplementary information to be added, say,

$$(19) \qquad \int a_k(\theta)\pi(\theta)d\mu(\theta) = \bar{a}_k, \quad k = s+1, s+2, ..., t.$$

In such a case, in a cross-entropy formulation (Kullback 1959), we take (18) as the initial density, and (19) as the additional information. Hence,

$$\pi_\phi^*(\theta) \propto [\mathcal{I}(\theta)]^{\frac{\phi}{2}} \exp\{(1-\phi)\mathcal{F}(\theta) + \sum_{k=0}^{s} \lambda_k a_k(\theta)\} \exp\{\sum_{k=s+1}^{t} \lambda_k a_k(\theta)\}$$

$$= [\mathcal{I}(\theta)]^{\frac{\phi}{2}} \exp\{(1-\phi)\mathcal{F}(\theta) + \sum_{k=0}^{t} \lambda_k a_k(\theta)\} .$$

To deal with the (local) uniqueness of the solution of the problem stated in Proposition 4.1, we rewrite the constraints, $\mathcal{C}$, as a function of the multipliers in the form $A(\Lambda) = [\int a_k(\theta)\pi_\phi^*(\theta)d\mu(\theta)]_{k=0}^{s} = \bar{A}$, where $\bar{A}^T = (\bar{a}_0, \bar{a}_1, ..., \bar{a}_s)$, and $\Lambda^T = (\lambda_0, \lambda_1, ..., \lambda_s)$ (the superindex $T$ denotes the usual vector or matrix transposing operation).

**Proposition 4.2** *Let $\pi_\phi^*(\theta)$ be as in (4.2), and suppose that $a_k$, $k = 0, 1, ..., s$, are linearly independent continuous functions in $L^2[\Theta, \pi_\phi^* d\mu]$ (the space of all $\pi_\phi^* d\mu$-measurable functions $a(\theta)$ defined on $\Theta$ such that $|a(\theta)|^2$ is $\pi_\phi^* d\mu$-integrable). Suppose that $A(\Lambda)$ is defined on an open set $\Delta \subset \mathbb{R}^{s+1}$, and let $\Lambda_o$ be a solution of $A(\Lambda) = \bar{A}$ for a fixed value of $\bar{A} = \bar{A}_o$. Then, there exists a neighborhood of $\Lambda_o$, $N(\Lambda_o)$, in which $\Lambda_o$ is the unique solution of $A(\Lambda) = \bar{A}_o$ in $N(\Lambda_o)$.*

The proof follows from the fact that $A(\Lambda)$ is continuously differentiable on $\Delta$, with nonsingular derivative

$$A'(\Lambda) = [\int a_j(\theta)a_\ell(\theta)\pi_\phi^*(\theta)d\mu(\theta)]_{0\leq j,\ell\leq s},$$

and from a straightforward application of inverse function theorem.

From (4.1) we may derive the following necessary condition, which is useful in practical situations.

**Proposition 4.3** *The multipliers $\Lambda^T = (\lambda_0, \lambda_1, ..., \lambda_s)$ appearing in (18) satisfy the following non-linear system of $s+1$ equations:*

$$1 = \lambda_0 + log\left\{\int [\mathcal{I}(\theta)]^{\frac{\phi}{2}} e^{(1-\phi)\mathcal{F}(\theta)} \prod_{k=1}^{s} e^{\lambda_k a_k(\theta)} d\mu(\theta)\right\},$$

$$1 = \lambda_0 - log\,\bar{a}_k + log\left\{\int a_k(\theta)[\mathcal{I}(\theta)]^{\frac{\phi}{2}} e^{(1-\phi)\mathcal{F}(\theta)} \prod_{u=1}^{s} e^{\lambda_u a_u(\theta)} d\mu(\theta)\right\},$$

$k = 1, 2, ..., s.$

*Moreover,*

(i) *if the integral in the first equality has a closed-form solution, then the rest of the multipliers can be found from the relations:*

$$\frac{\partial \lambda_0}{\partial \lambda_k} = \bar{a}_k, \qquad k = 1, 2, ..., s,$$

(ii) *the formula*

$$\phi \mathcal{V}_{1,1,1}(\pi_\phi^*) + (1 - \phi)[\mathcal{V}_{0,0,0}(\pi_\phi^*) - 2\mathcal{V}_{0,0,1}(\pi_\phi^*)] = 1 - \sum_{k=0}^{s} \lambda_k \bar{a}_k,$$

*holds for all* $0 \le \phi \le 1$.

Very often, experimenters are concerned with assigning weights $\bar{a}_k$, $k = 1, 2, ..., s$, to regions $A_k$, $k = 1, 2, ..., s$, to express, according to experience, how likely it is that $\theta$ belongs to each region. The following result, based on Proposition 4.3, characterizes Good-Bernardo-Zellner's priors when such a supplementary information comes in the form of quantiles, and both $\mathcal{I}(\theta)$ and $\mathcal{F}(\theta)$ are constant. Under such assumptions, the non-linear system of $s+1$ equations given in Proposition 4.3 is transformed into a homogeneous linear system of the same dimension as shown below:

**Proposition 4.4** *Suppose that the sets* $A_k = (b_k, b_{k+1}]$, $k = 1, 2, ..., s - 1$ *and* $A_s = (b_s, b_{s+1})$ *with* $b_1 < b_2 < \cdots < b_{s+1}$, $s \ge 2$, *constitute a partition of* $\Theta$, $0 < \mu(\Theta) < \infty$. *Suppose also that both* $\mathcal{I}(\theta)$ *and* $\mathcal{F}(\theta)$ *are constant. Let* $\bar{a}_1, \bar{a}_2, ..., \bar{a}_s > 0$ *be such that* $\sum_{k=1}^{s} \bar{a}_k = 1$, *and* $\int I_{A_k}(\theta)\pi(\theta)d\mu(\theta) = \bar{a}_k$, $k = 1, 2, ..., s$. *If we define new multipliers:* $\omega_0 = e^{1-\lambda_0}/D$ *where* $D = [\mathcal{I}(\theta)]^{\frac{\phi}{2}} e^{(1-\phi)\mathcal{F}(\theta)}$, *and* $\omega_k = e^{\lambda_k}$, $k = 1, 2, ..., s$. *Then,* $\Omega = (\omega_0, \omega_1, ..., \omega_s)$ *can be found from the following homogeneous linear system:*

(20)
$$\begin{pmatrix} -1 & u_1 & u_2 & \dots & u_s \\ -1 & v_1 & 0 & \dots & 0 \\ -1 & 0 & v_2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -1 & 0 & 0 & \dots & v_s \end{pmatrix} \begin{pmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \vdots \\ \omega_s \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

*where* $u_k = \mu(A_k)$, *and* $v_k = \bar{a}_k^{-1} u_k$, $k = 1, 2, ..., s$.

Observe that the determinant, $\Delta$, of the matrix in (20) is given by

$$\Delta = \left( \frac{\sum_{k=1}^s \bar{a}_k - 1}{\prod_{k=1}^s \bar{a}_k} \right) \prod_{k=1}^s u_k,$$

which guarantees that there exists a unique nontrivial solution since $\sum_{k=1}^s \bar{a}_k = 1$. In this case, the solution is $\Omega^{*T} = (1, v_1^{-1}, v_2^{-1}, ..., v_s^{-1})$, and $\pi_\phi^* = \sum_{k=1}^s v_k^{-1} I_{A_k}$.

The following proposition extends Good-Bernardo-Zellner's priors to a richer family by using the MMDIP and MFIP criteria:

**Proposition 4.5** *Let*

$$\mathcal{N}_{\phi,\psi}(\pi) \stackrel{def}{=} \phi \mathcal{V}_{1,1,1}(\pi) + (1-\phi)(1-\psi)\mathcal{V}_{0,0,0}(\pi) + (\psi(1-\phi)/2)[\mathcal{V}_{0,1,1} + \mathcal{V}_{0,1,0}],$$

$0 \le \phi, \psi \le 1$. *Then*

(i) $\mathcal{N}_{\phi,\psi}(\pi) \in \mathcal{A}$ *and is concave w.r.t.* $\pi$.

(ii) *A necessary condition for* $\pi$ *to be a maximum of the problem*

$$\text{Maximize} \quad \mathcal{N}_{\phi,\psi}(\pi)$$

$$\text{subject to} \quad \mathcal{C} : \int a_k(\theta)\pi(\theta)d\mu(\theta) = \bar{a}_k, \quad k = 0, 1, 2, ..., s,$$

*where* $a_0 \equiv 1 = \bar{a}_0$, *is given by*

$$\pi_{\phi,\psi}^*(\theta) \propto [\mathcal{I}(\theta)]^{\frac{\phi}{2}} \exp \left\{ (1-\phi)(1-\psi)\mathcal{F}(\theta) \right.$$

$$(21) \qquad + \frac{\psi(1-\phi)}{2} \left[ [\mathcal{I}(\theta)]^{\frac{1}{2}} + \frac{\mathcal{F}(\theta)}{[\mathcal{I}(\theta)]^{\frac{1}{2}}} \right] + \left. \sum_{k=0}^s \lambda_k a_k(\theta) \right\},$$

*where* $\lambda_k$, $k = 0, 1, ..., s$, *are the Lagrange multipliers associated with the constraints* $\mathcal{C}$.

The second term inside the exponential of (21) is the average between Fisher's information and the negative relative Shannon-Fisher's information. Notice that $\pi_{\phi,0}^*(\theta)$ is just Good-Bernardo-Zellner's prior.

In the following proposition, Good-Bernardo-Zellner type priors are derived as MAXENTP solutions by treating (5) and (8) as constraints (for the rationale of MAXENTP methods see Jaynes' 1982 seminal paper).

**Proposition 4.6** *Consider the Jaynes–Good–Bernardo–Zellner problem:*

$$\text{Maximize} \quad \mathcal{V}_{0,0,1}(\pi)$$

$$\text{subject to:} \begin{cases} \mathcal{V}_{1,1,1}(\pi) - \mathcal{V}_{0,0,1}(\pi) = \bar{b}_1, \\ \mathcal{V}_{0,0,0}(\pi) - \mathcal{V}_{0,0,1}(\pi) = \bar{b}_2, \\ \displaystyle\int a_k(\theta)\pi(\theta)d\mu(\theta) = \bar{a}_k, \quad k = 0, 1, 2, ..., s, \, a_0 \equiv 1 = \bar{a}_0. \end{cases}$$

*Then a necessary condition for a maximum is*

$$(22) \qquad \pi^*(\theta) \propto [\mathcal{I}(\theta)]^{\frac{\rho_1}{2}} \exp\{\rho_2 \mathcal{F}(\theta) + \sum_{k=0}^{s} \lambda_k a_k(\theta)\},$$

*where $\rho_j$, $j = 1, 2$, and $\lambda_k$, $k = 0, 1, ..., s$, are the Lagrange multipliers associated with the constraints.*

Unlike the coefficients $\phi$ and $1 - \phi$ appearing in (4.6), the multipliers $\rho_j$, $j = 1, 2$, do not necessarily add up to 1.

There typically exist priors for which Shannon-Jaynes entropy becomes infinite. One way to overcome this problem consists of discounting entropy at a constant rate $\nu > 0$. The following proposition introduces Good-Bernardo-Zellner's controlled priors as solutions of maximizing discounted entropy:

**Proposition 4.7** *Consider the discounted version of the problem stated in the preceding proposition:*

$$\text{Maximize} \quad -\int e^{-\nu\theta}\pi(\theta)\log\pi(\theta)d\mu(\theta),$$

*subject to:*

$$\begin{cases} \dfrac{1}{\pi(\theta)}\dfrac{dh_1(\theta)}{d\mu(\theta)} = \log[\mathcal{I}(\theta)]^{\frac{1}{2}}, \quad h_1(-\infty) = 0, \\ \qquad\qquad\qquad\qquad\quad h_1(\infty) = \mathcal{V}_{1,1,1}(\pi) - \mathcal{V}_{0,0,1}(\pi) < \infty, \\ \\ \dfrac{1}{\pi(\theta)}\dfrac{dh_2(\theta)}{d\mu(\theta)} = \mathcal{F}(\theta), \quad h_2(-\infty) = 0, \\ \qquad\qquad\qquad\qquad\quad h_2(\infty) = \mathcal{V}_{0,0,0}(\pi) - \mathcal{V}_{0,0,1}(\pi) < \infty, \\ \\ \dfrac{1}{\pi(\theta)}\dfrac{dg_k(\theta)}{d\mu(\theta)} = a_k(\theta), \, g_k(-\infty) = 0, \, g_k(\infty) < \infty, \, k = 0, 1, 2, ..., s \end{cases}$$

where $a_0 \equiv 1 = \bar{a}_0$. Then, a necessary condition for $\pi^*(\theta)$ to be an optimal control is given by

$$(23) \qquad \pi^*(\theta) \propto [\mathcal{I}(\theta)]^{\frac{\rho_1(\theta)}{2}} \exp\{\rho_2(\theta)\mathcal{F}(\theta) + \sum_{k=0}^{s} \lambda_k(\theta)a_k(\theta)\},$$

where $\rho_j(\theta) = \rho_{j0}e^{\nu\theta}$, $j = 1, 2$, and $\lambda_k(\theta) = \lambda_{k0}e^{\nu\theta}$, $k = 0, 1, ..., s$, are the costate variables associated with the state variables $h_j(\theta)$, $j = 1, 2$, and $g_k(\theta)$, $k = 0, 1, ..., s$, respectively. Furthermore, the constants $\rho_{j0}$, $j = 1, 2$, and $\lambda_{k0}$, $k = 0, 1, ..., s$, can be computed from the following non-linear system of $s + 3$ equations:

$$1 + \log h_1(\infty) = \log\left\{\int \log[\mathcal{I}(\theta)]^{\frac{1}{2}} m(\rho_{10}, \rho_{20}, \lambda_{00}, \lambda_{10}, ..., \lambda_{s0}; \theta)d\mu(\theta)\right\},$$

$$1 + \log h_2(\infty) = \log\left\{\int \mathcal{F}(\theta)m(\rho_{10}, \rho_{20}, \lambda_{00}, \lambda_{10}, ..., \lambda_{s0}; \theta)d\mu(\theta)\right\},$$

$$1 + \log g_k(\infty) = \log\left\{\int a_k(\theta)m(\rho_{10}, \rho_{20}, \lambda_{00}, \lambda_{10}, ..., \lambda_{s0}; \theta)d\mu(\theta)\right\},$$

$$k = 0, 1, 2, ..., s$$

where

$$m(\rho_{10}, \rho_{20}, \lambda_{00}, \lambda_{10}, ..., \lambda_{s0}; \theta)$$
$$= \left([\mathcal{I}(\theta)]^{\frac{\rho_{10}}{2}} e^{\rho_{20}\mathcal{F}(\theta)} e^{\lambda_{00}} \prod_{u=1}^{s} e^{\lambda_{u0}a_u(\theta)}\right)^{e^{\nu\theta}}.$$

## 5   Kalman filtering priors

In this section, we will study Good-Bernardo-Zellner's priors as Kalman Filtering priors (Kalman 1960, and Kalman and Bucy 1961). We will continue to work with the single parameter case, and focus our attention on both the location and scale parameter families.

Let $Y_1, Y_2, ..., Y_t$ be a set of indirect measurements, from a polling system or a sample survey, of an unobserved state variable $\beta_t$. The objective is to make inferences about $\beta_t$. The relationship between $Y_t$ and $\beta_t$ is specified by the measurement equation, sometimes also called the observation equation:

$$(24) \qquad\qquad\qquad Y_t = A_t\beta_t + \varepsilon_t,$$

where $A_t \neq 0$ is known, and $\varepsilon_t$ is the observation error distributed as $\mathcal{N}(0, \sigma^2_{\varepsilon_t})$ with $\sigma^2_{\varepsilon_t}$ known. Notice that the main difference between the

measurement equation and the linear model is that, in the former, the coefficient $\beta_t$ changes with time. Furthermore, we suppose that $\beta_t$ is driven by a first order autoregressive process, that is,

$$(25) \qquad \qquad \beta_t = Z_t \beta_{t-1} + \eta_{t-1},$$

where $Z_t \neq 0$ is known, and $\eta_t \sim \mathcal{N}(0, \sigma_{\eta t}^2)$ with $\sigma_{\eta t}^2$ known. In what follows, we will assume that $\beta_0$, $\varepsilon_t$, and $\eta_t$ are independent random variables. We might state nonlinear versions of (24) and (25), but this would not make any essential differences in the subsequent analysis.

Suppose now, that at time $t = 0$, supplementary information is given by $\widehat{\beta}_0$ and $\widehat{\sigma}_0^2$, the mean and variance of $\beta_0$ respectively. That is,

$$(26) \qquad \mathcal{C} \ : \ \begin{cases} \displaystyle\int_{-\infty}^{\infty} \pi(\beta_0) d\beta_0 = 1, \\[2ex] \displaystyle\int_{-\infty}^{\infty} \beta_0 \pi(\beta_0) d\beta_0 = \widehat{\beta}_0, \\[2ex] \displaystyle\int_{-\infty}^{\infty} (\beta_0 - \widehat{\beta}_0)^2 \pi(\beta_0) d\beta_0 = \widehat{\sigma}_0^2. \end{cases}$$

In this case, Good-Bernardo-Zellner's prior is given by

$$(27) \ \ \pi_\phi^*(\beta_0) \propto [\mathcal{I}(\beta_0)]^{\frac{\phi}{2}} \exp\{(1-\phi)\mathcal{F}(\beta_0) + \lambda_0 + \lambda_1 \beta_0 + \lambda_2 (\beta_0 - \widehat{\beta}_0)^2\},$$

where $\lambda_j$, $j = 0, 1, 2$, are Lagrange multipliers.

Suppose that, at time $t$, we wish to make inferences about the conditional state variable $\theta_t = \beta_t | I_t$, where $I_t = \{Y_1, Y_2, ..., Y_{t-1}\}$. To obtain a posterior distribution of $\theta_t$, the information provided by the measurement $Y_t$, with density $f(Y_t | \theta_t)$, is used to modify the initial knowledge in $\pi_\phi^*(\theta_t)$ according to Bayes' theorem:

$$(28) \qquad \qquad f(\theta_t | Y_t) \propto f(Y_t | \theta_t) \pi_\phi^*(\theta_t).$$

We are now in a position to state the Bayesian recursive updating procedure of the Kalman Filter (KF) for both the location and scale parameter families $f(Y_t | \theta) = f(Y_t - \theta)$, $\theta \in \mathbb{R}$, and $f(Y_t | \theta) = (1/\theta) f(Y_t / \theta)$, $\theta > 0$, respectively. To start off the KF procedure, we substitute (27) in (26), obtaining that Good-Bernardo-Zellner's prior at time $t = 0$, is given by $\mathcal{N}(\widehat{\beta}_0, \widehat{\sigma}_0^2)$, which is describing the initial knowledge of the system. Proceeding inductively, at time $t$, $\widehat{\beta}_{t-1}$ and

$\widehat{\sigma}_{t-1}^2$ become *supplementary* information, and therefore Good-Bernardo-Zellner's prior at time $t$ is represented by

$$(29) \qquad \theta_t = \beta_t | I_t \sim \mathcal{N}(Z_t \widehat{\beta}_{t-1}, M_t),$$

where

$$(30) \qquad M_t = Z_t^2 \widehat{\sigma}_{t-1}^2 + \sigma_{\eta_{t-1}}^2.$$

The sampling model (or likelihood function) is determined by

$$(31) \qquad Y_t | \theta_t \sim \mathcal{N}(A_t \beta_t, \sigma_{\varepsilon_t}^2).$$

The posterior distribution, at time $t$, is then obtained by substituting both (29) and (30) in (28), so

$$f(\theta_t | Y_t) \propto \exp\{-\tfrac{1}{2}[(A_t \beta_t - Y_t)^2 \sigma_{\varepsilon_t}^{-2} + (\beta_t - Z_t \widehat{\beta}_{t-1})^2 M_t^{-1}]\}.$$

Noting that $\pi_\phi^*(\theta_t)$ is a natural conjugate prior, we may complete the squares to get

$$\theta_t | Y_t \sim \mathcal{N}[Z_t \widehat{\beta}_{t-1} + K_t(Y_t - A_t Z_t \widehat{\beta}_{t-1}), M_t - K_t A_t M_t],$$

where

$$(32) \qquad K_t = M_t A_t (\sigma_{\varepsilon_t}^2 + A_t^2 M_t)^{-1}.$$

This, of course, means that

$$(33) \qquad \begin{cases} \widehat{\beta}_t = Z_t \widehat{\beta}_{t-1} + K_t(Y_t - A_t Z_t \widehat{\beta}_{t-1}), \\ \widehat{\sigma}_t^2 = M_t - K_t A_t M_t. \end{cases}$$

We then proceed with the next iteration. Equations (33), (30), and (32) are known in the literature as the KF.

The above analysis can be summarized in the following proposition:

**Proposition 5.1** *Consider the state-space representation:*

$$\begin{cases} Y_t = A_t \beta_t + \varepsilon_t, \\ \\ \beta_t = Z_t \beta_{t-1} + \eta_{t-1}, \end{cases}$$

*defined as in (24) and (25). Suppose that supplementary information on the mean and variance of $\beta_0$ is available. Let $\theta_t = \beta_t | I_t$, where*

$I_t = \{Y_1, Y_2, ..., Y_{t-1}\}$, *and consider the location and scale parameter families* $f(Y_t|\theta) = f(Y_t - \theta)$, $\theta \in \mathbb{R}$, *and* $f(Y_t|\theta) = (1/\theta)f(Y_t/\theta)$, $\theta > 0$, *respectively, along with the properties stated in Corollary 4.1. Then, under Good-Bernardo-Zellner's prior,* $\pi_\phi^*(\theta_t)$, *the posterior estimate of* $\beta_t$, $\widehat{\beta}_t$, *is given by*

$$\widehat{\beta}_t = \omega_t Z_t \widehat{\beta}_{t-1} + (1 - \omega_t)(Y_t/A_t),$$

*where* $\omega_t = \sigma_{\varepsilon_t}^2(\sigma_{\varepsilon_t}^2 + A_t^2 M_t)^{-1}$.

## 6    Revisiting the normal linear model

The results on Good-Bernardo-Zellner priors given so far can be easily extended to the multi-dimensional parameter case, namely, $\theta = (\theta_1, \theta_2, ..., \theta_m) \in \Theta \subseteq \mathbb{R}^m$, $m > 1$. Consider a vector of independent and identically distributed normal random variables $(X_1, X_2, ..., X_n)$ with common and known variance $\sigma^2$ satisfying

(34)        $\text{E}(X_k) = a_{k1}\theta_1 + a_{k2}\theta_2 + \cdots + a_{km}\theta_m, \qquad k = 1, 2, ..., n,$

where $A = (a_{ij})$ is a matrix of known coefficients for which $(A^T A)^{-1}$ exists.

Let $X$ and $\theta$ stand for the column vectors of variables $X_k$ and parameters $\theta_j$, respectively. Then (34) can be written in matrix notation as, $\text{E}(X) = A\theta$. In this case, we have

(35)            $f(\xi|\theta) = (\frac{1}{2\pi\sigma^2})^{\frac{n}{2}} \exp\{-\frac{1}{2\sigma^2}\|\xi - A\theta\|^2\},$

where $\xi = (x_1, x_2, ..., x_n)$. Since $\sigma^2$ has been assumed known, only the location parameter is unknown. The analogue of (2) is now given by the matrix:

$$\mathcal{I}_n(\theta) \equiv \left( \int \left( \frac{\partial}{\partial\theta_j} \log f(x|\theta) \right) \left( \frac{\partial}{\partial\theta_\ell} \log f(x|\theta) \right) f(x|\theta) d\lambda(x) \right)_{1 \leq j,\ell \leq m}$$

$$= \frac{1}{\sigma^2} A^T A,$$

and so $\det[\mathcal{I}_n(\theta)]$ is constant, which implies that the Good-Bernardo-Zellner prior distribution $\pi_\phi^*(\theta)$, describing a situation of vague information on $\theta$, must be a locally uniform prior distribution.

Let $\widehat{\theta}$ be the least squares estimate for $\theta$, then it is known that $A^T A \widehat{\theta} = A^T X$, $\mathrm{E}(\widehat{\theta}) = \theta$, and $\mathrm{Var}(\widehat{\theta}) = \sigma^2 (A^T A)^{-1}$. Noting from equation (35) that

$$ f(\xi|\theta) = (\tfrac{1}{2\pi\sigma^2})^{\frac{n}{2}} \exp\{-\tfrac{1}{2\sigma^2}(\|\xi - A\widehat{\theta}\|^2 + \langle A^T A(\theta - \widehat{\theta}), \theta - \widehat{\theta}\rangle)\}, $$

and applying Bayes' theorem, we get as the posterior distribution of $\theta$

$$ f(\theta|\xi) = (2\pi)^{-\frac{m}{2}}(\det[\tfrac{1}{\sigma^2}A^T A])^{\frac{1}{2}} \exp\{-\tfrac{1}{2}\langle\tfrac{1}{\sigma^2}A^T A(\theta - \widehat{\theta}), \theta - \widehat{\theta}\rangle\}. $$

If supplementary information in mean, $c$, and variance-covariance matrix, $D$, is now incorporated, then the (informative) Good-Bernardo-Zellner prior is given by

$$ \pi_\phi^*(\theta) = (2\pi)^{-\frac{m}{2}}(\det[D])^{-\frac{1}{2}} \exp\{-\tfrac{1}{2}\langle D^{-1}(\theta - c), \theta - c\rangle\}. $$

The posterior distribution is now

$$ \begin{aligned} f(\theta|\xi) = {}& (2\pi)^{-\frac{m}{2}}(\det[B])^{\frac{1}{2}} \\ &\times \exp\Big\{-\tfrac{1}{2}\langle B[\theta - ((DB)^{-1}c + \tfrac{1}{\sigma^2}B^{-1}A^T A\widehat{\theta})], \\ &\qquad \theta - ((DB)^{-1}c + \tfrac{1}{\sigma^2}B^{-1}A^T A\widehat{\theta})\rangle\Big\}, \end{aligned} $$

where $B = D^{-1} + \tfrac{1}{\sigma^2}A^T A$.

# 7 Summary and conclusions

We have presented, in a unified framework, a number of well-known methods that maximize a criterion functional to obtain non-informative and informative priors. Our general procedure is, by itself, capable of dealing with a range of interesting issues in Bayesian analysis. However, in this paper, we have limited our attention to Good-Bernardo-Zellner's priors as well as their application to some Bayesian inference problems, including the Kalman filter and the Normal linear model.

There exist priors for which Shannon-Jaynes entropy becomes infinite. In order to overcome this difficulty we proposed discounted entropy. We introduced Good-Bernardo-Zellner's controlled priors which maximize discounted entropy at a constant rate. Throughout the paper, we have emphasized the existence and uniqueness of the solutions of the corresponding variational and optimal control problems. There are, of

course, many other members of the class $\mathcal{A}$ that deserve much more attention than that we have attempted here. Needless to say, more work will be required in this direction. Results will be reported elsewhere.

### Acknowledgement

Francisco Venegas-Martínez
*Department of Finance*,
Tecnológico de Monterrey,
Calle del Puente 222,
14380 México D. F.
fvenegas@itesm.mx

# References

[1] Akaike H., *A new look at the Bayes procedure*, Biometrika **65** (1978), 53–59.

[2] Bayes T., *An essay towards solving a problem in the doctrine of chances*, Philos. Trans. R. Soc. London **53** (1763) 370–418. (Reprinted in Biometrika **45** (1958), 243–315.)

[3] Berger J. O.; Bernardo J. M., *Estimating a product of means: Bayesian analysis with reference priors*, J. Amer. Statist. Assoc. **84** (1989), 200–207.

[4] Berger J. O.; Bernardo J. M., *On the development of reference priors*, Bayesian Statistics **4** (1992a), 35–60.

[5] Berger J. O.; Bernardo J. M., *Ordered group reference priors with application to a multinomial problem*, Biometrika, **79** (1992b), 25–37.

[6] Berger J. O.; Bernardo J. M.; Mendoza M., *On priors that maximize expected information*, Recent Developments in Statistics and Their Applications (1989), 1–20.

[7] Bernardo J. M., *Noninformative priors do not exist*, J. Statist. Plann. Inference **B65** (1997), 177–189.

[8] Bernardo J. M., *Reference posterior distributions for Bayesian inference*, J. Roy. Statist. Soc. Ser. B Stat. Methodol. **41** (1979), 113–147.

[9] Bernardo J. M.; Ramón J. M., An Introduction to Bayesian Reference Analysis: Inference on the Ratio of Multinomial Parameters, Tech. Rep. 3–97, Universitat de València, Spain, 1997.

[10] Bernardo J. M.; Smith A. F. M., Bayesian Theory, John Wiley & Sons, Wiley Series in Probability and Mathematical Statistics, New York, 1994.

[11] Box G. E. P.; Tiao G. C.,  Bayesian Inference and Statistical Analysis, Addison-Wesley Series in Behavioral Science:  Quantitative Methods, Massachusetts, 1973.

[12] Geisser S., Predictive Inference:  An Introduction, Chapman & Hall, New York, 1993.

[13] Good I. J., *Utility of a distribution*, Nature **219** (1968), 1392.

[14] Good I. J., *What is the use of a distribution?*, Multivariate Analysis (Krishnaiah, ed.) Vol. II, Academic Press, New York (1969), 183–203.

[15] Jaynes E. T., *Information theory and statistical mechanics*, Phys. Rev. **106** (1957), 620–630.

[16] Jaynes E. T., *On the rationale of maximum-entropy methods*, Proc. of the IEEE **70** (1982), 939–952.

[17] Jeffreys H., Theory of Probability, 3rd. edition, Oxford University Press, Oxford, 1961.

[18] Kalman R. E., *A new approach to linear filtering and prediction problems*, Transactions ASME, Series D, J. of Basic Engineering **82** (1960), 35–45.

[19] Kalman R. E.; Bucy R., *New results in linear filtering and prediction theory*, Transactions ASME, Series D, J. of Basic Engineering **83** (1961), 95–108.

[20] Kullback S., Information Theory and Statistics, Wiley, New York, 1959.

[21] Lindley D. V., *On a measure of information provided by an experiment*, Ann. Math. Statist. **27** (1956), 986–1005.

[22] Rényi A., Foundations of Probability, Holdan-Day, San Francisco, 1970.

[23] Soofi E. S., *Capturing the intangible concept of information*, J. Amer. Statist. Assoc. **89** (1994), 1243–1254.

[24] Zellner A., An Introduction to Bayesian Inference in Econometrics, Wiley, New York, 1971.

[25] Zellner A., *Bayesian method of moments (BMOM) analysis of mean and regression models*, Bayesian Analysis in Econometrics and Statistics, Published by Edward Elgar Publishing Limited, UK (1997), 291–304. (Also available in Prediction and Modelling Honouring Seymor Geisser, Springer-Verlag, New York, Berlin, Heidelberg, (1996).)

[26] Zellner A., *Bayesian methods and entropy in economics and econometrics*, Maximum Entropy and Bayesian Methods, Dordrecht, Netherlands: Kluwer (1996), 17–31.

[27] Zellner A., *Maximal data information prior distributions*, New Developments in the Application of Bayesian Methods, Amsterdam: North-Holland (1977), 201–232.

[28] Zellner A., *Models, prior information and Bayesian analysis*, J. Econometrics **75** (1996a), 51–58.

[29] Zellner A., *Past and recent results on maximal data information priors*, J. Statist. Plann. Inference **49** (1996b), 3–8.

[30] Zellner A., *The finite sample properties of simultaneous equations' estimates and estimators: Bayesian and non-Bayesian approaches*, J. Econometrics **83** (1998), 185-212.