

NOVEL PARAMETER PRIORS FOR BAYESIAN SIGNAL IDENTIFICATION

Anthony Quinn

Department of Electronic and Electrical Engineering, University of Dublin, Trinity College, Dublin 2. Ireland.

Phone: +353-1-6081863/6081580; Fax: +353-1-6772442; Email: aquinn@ee.tcd.ie

ABSTRACT

The problem of eliciting priors on the parameter space of a signal hypothesis is considered in this paper, and two lesser-known approaches are emphasized. Each yields conservative priors appropriate for data-dominated Bayesian parameter inference. They are based, respectively, on the principles of (i) *a posteriori* transformation invariance, and (ii) *a priori* maximum entropy. Novel priors on a wide class of signal models are deduced. Their ability to regularize inference of the difference frequency between closely spaced tones is considered, and they are compared with the Ockham Prior which was studied in previous work.

1. INTRODUCTION

In the Bayesian approach, inferences concerning unknown parameters are drawn from the *a posteriori* density, which is formed, from Bayes' Theorem [1, 2], as the product of the *Likelihood Function (LF)* and the *prior*. It is the problem of assigning prior distributions in signal identification which is addressed in this paper. The prior not only quantifies any relevant information available before observations are made, but it can also be used to modify inferences so that they behave appropriately (Section 2.1).

The problem of prior elicitation remains the single most contentious area in all of mathematical inference. Many criteria have been proposed in the literature [1-7], some of which have been used in Bayesian signal identification:

- (i) **Bayes' Postulate:** this states that a uniform prior should be assigned when 'no' prior information is available. Inferences are then drawn exclusively from the LF.
- (ii) **Jeffreys' Priors:** $\Theta = \theta$ is a *location parameter* if the LF has the form

$$p(x | \theta) = h(x - \theta).$$

In this case, Jeffreys' Prior for Θ is uniform, as in (i). $\Sigma = \sigma$ is a *scale parameter* if

$$p(x | \sigma) = \frac{1}{\sigma} h\left(\frac{x}{\sigma}\right),$$

for which Jeffreys' Prior is $p(\sigma) \propto 1/\sigma$. The former is employed widely as the prior for 'signal' parameters,

the latter as the prior on the standard deviation of the additive 'noise' [1, 5, 8].

- (iii) **Conjugate Priors:** the prior, $q(\theta)$, is conjugate to the LF, $p(x | \theta)$, if the resulting posterior, $\pi(\theta | x)$, has the same form as $q(\cdot)$. This property is useful in ensuring tractability in problems of sequential inference [2].

The sequel defines the problem of assigning priors consistently (i.e. "elicitation"), and presents two strategies which are not widely known in the signal identification literature.

2. PRIOR ELICITATION

Consider a parametric hypothesis, \mathcal{I} , which seeks to explain a set of observations, $\mathbf{D} = \mathbf{d} \in \overline{\mathcal{C}}^N$, in terms of a set of unknown parameters, $\Theta = \theta \in \overline{\Theta} \subset \overline{\mathcal{C}}^k$, where $\overline{\mathcal{C}}$ denotes the set of complex numbers, $\overline{\Theta}$ is the space of θ , and notation of the type ' $\mathbf{X} = \mathbf{x}$ ' signifies a realization, \mathbf{x} , for the probabilistic parameter (p.p.) [9] vector, \mathbf{X} . \mathbf{d} may be viewed as a vector of N discrete samples from a random process. The direct probability of \mathbf{d} given candidate values of θ , i.e. $p(\mathbf{d} | \theta, \mathcal{I})$, may be interpreted as a function, $l(\theta | \mathbf{d}, \mathcal{I})$, of θ itself, being the *Likelihood Function (LF)*. In this paper, \mathcal{I} additively decomposes \mathbf{d} into a 'signal', $\mathbf{s} = \mathbf{g}(\theta)$, and 'noise', \mathbf{e} :

$$\mathbf{d} = \mathbf{g}(\theta) + \mathbf{e}, \quad (1)$$

where

$$\mathbf{e} \sim p_E(\mathbf{e} | \mathcal{I}). \quad (2)$$

' \sim ' denotes 'is distributed as', and $\mathbf{g}(\cdot)$ is the sampled signal function, parameterized by θ . It will be assumed that the analytical form of this function is known, so that the signal identification problem reduces to one of estimating Θ . The framework is readily extended to the case of model uncertainty [5, 9]. From (1,2), the LF is given by

$$p(\mathbf{d} | \theta, \mathcal{I}) \equiv l(\theta | \mathbf{d}, \mathcal{I}) = p_E(\mathbf{d} - \mathbf{g}(\theta) | \mathcal{I}). \quad (3)$$

Much of the *orthodox* signal identification literature is based on inferences from the LF. In contrast, the Bayesian approach recognizes that degrees of belief in a proposition are quantified by Probability. For example, the proposition that the (unknown) parameters, Θ , take on value θ —i.e. ' $\Theta = \theta$ '—has a belief level given by $p(\theta | \mathbf{d}, \mathcal{I})$ once \mathbf{d} is observed. This is the *a posteriori (AP)* pdf of Θ [1-3, 5].

Bayes' Rule inverts the direct probability (3) to yield the required AP pdf:

$$p(\theta | d, \mathcal{I}) \propto p(d | \theta, \mathcal{I})p(\theta | \mathcal{I}). \quad (4)$$

$p(\theta | \mathcal{I})$ is the *prior* on $\bar{\Theta}$. Data-dominated (i.e. objective) inference may be achieved using uniform or diffuse priors, as discussed in Section 1. Then, $p(\theta | d, \mathcal{I}) \propto l(\theta | d, \mathcal{I})$ from (4), but this fails to *regularize* the inference, as discussed next.

2.1. Prior Regularization

Normalizability of $p_E(\cdot)$ (2) implies that $p_E(e | \mathcal{I}) = f(\|e\|)$, where $f(\cdot)$ is monotonically decreasing and $\|e\|$ denotes some norm of e . For example, the (complex) Gaussian noise pdf—to be explored in this paper—is $p_E(e | \mathcal{I}) \propto \exp(-e^H \Sigma^{-1} e)$, where $e^H \Sigma^{-1} e$ is a Mahalanobis norm. The Maximum Likelihood (ML) estimate, from (3), is given by $\theta_{ML} = \arg. \min_{\theta} \|d - g(\theta)\|$. The criterion may overfit d with $g(\cdot)$, particularly if $k \ll N$. This manifests itself as thresholds in estimation and excessive order in model selection [5, 9]. Prior regularization can be accomplished via an appropriate non-uniform prior, to yield a Maximum *a Posteriori* (MAP) criterion of the type

$$\theta_{MAP} = \arg. \min_{\theta} [\|d - g(\theta)\| + \beta(\theta)], \quad (5)$$

where $\beta(\theta)$ is introduced via the prior (4). In tasks such as image segmentation [10], for example, Markovian priors may be adopted to encourage connectivity in the inferred label field. Two principles which yield novel regularizing priors for the signal identification task (1) are now presented.

3. TRANSFORMATION INVARIANT INFERENCE

Let $h: \bar{\Theta} \subset \bar{\mathcal{C}}^k \rightarrow \bar{\eta} \subset \bar{\mathcal{C}}^k | \eta = h(\theta)$ be a complex, bijective, holomorphic transformation with continuous first partial derivatives. The gradient matrix of the transformation is denoted by $\nabla_{\theta}(\eta^T)$, whose determinant is assumed to be non-zero, $\forall \theta \in \bar{\Theta}$. Denoting the AP inference for Θ by $p(\theta | d, \mathcal{I})$ (assumed positive $\forall \theta$), then the inference for η is [1, 3, 4]

$$q(\eta | d, \mathcal{I}) = \frac{1}{|\nabla_{\theta}(\eta^T)|^2} p[h^{-1}(\eta) | d, \mathcal{I}], \quad (6)$$

where $h^{-1}(\cdot)$ denotes the inverse transformation. The determinant term, $|\cdot|$, which is squared since the transformation is *complex*, is the hypervolume element expansion ratio for $h(\cdot)$, and arises in (6) because of the normalization requirement of probabilities.

One consequence of (6) is that *Bayesian* inferences—as opposed to inferences based on non-measure functions such as the LF [5]—are not, in general, invariant with respect to a transformation. For example, the Maximum *a Posteriori* (MAP) estimate will not satisfy $\eta_{MAP} = h(\theta_{MAP})$, which is desirable. Thus, the MAP estimate of f in the signal model $\sin(2\pi fn)$ will not be consistent with that of T in

$\sin(2\pi n/T)$, under the transformation $T = 1/f$. However, the following theorem, based on work by Jeffreys [4], provides a sufficient condition for transformation invariance to be met.

Theorem 1 Consider hypothesis \mathcal{I} , which explains observations $d \in \bar{\mathcal{C}}^N$ in terms of parameters $\Theta = \theta \in \bar{\Theta} \subset \bar{\mathcal{C}}^k$, via the LF, $l(\theta | d, \mathcal{I})$. Let $\eta = h(\theta)$ belong to the transformation class defined above. Then AP inferences are invariant with respect to $h(\cdot)$ if the following prior is used:

$$p(\theta | \mathcal{I}) \propto |J(\theta | \mathcal{I})|. \quad (7)$$

$J(\theta | \mathcal{I})$ is the Fisher information matrix under \mathcal{I} , whose elements are

$$J_{ij}(\theta | \mathcal{I}) = -\mathcal{E}_{\theta} \left[\frac{\partial^2 \ln l(\theta | d, \mathcal{I})}{\partial \theta_i \partial \theta_j} \right] \quad 1 \leq i, j \leq k. \quad \square \quad (8)$$

- A proof is given in [7]. The mild regularity conditions which must be met by $l(\theta | d, \mathcal{I})$ are given in [5].
- The appropriate invariance prior for *real* parameter transformations is $p(\theta | \mathcal{I}) \propto |J(\theta | \mathcal{I})|^{\frac{1}{2}}$.

3.1. Data-Translated Likelihood

If the LF can be expressed in the form

$$l(\theta | d, \mathcal{I}) = q[\phi(\theta) - f(d)], \quad (9)$$

then it is *data-translated* [1] in the transformed space, $\bar{\Phi}$, of ϕ , since d influences only the *location*, and not the *shape*, of the LF. While this property often holds in the single parameter case for asymptotically large samples, it is necessary to weaken the definition in the multiparameter case to one of *hypervolume invariance* of the LF. It has been shown [1] that the prior (7) is consistent with the assignment of a *uniform* prior in the space, $\bar{\Phi}$, where this hypervolume invariance is achieved. This confers another justification for choosing (7), though it is important to recognize the weakness of this invariance constraint [4].

3.2. Prior Stochastic Independence

If $\Theta = (\psi^T, \phi^T)^T$, and stochastic independence between ψ and ϕ is to be imposed *a priori*, then the prior factorization

$$p(\theta | k, l, \mathcal{I}) \propto |J_1(\psi | k, \mathcal{I})| |J_2(\phi | l, \mathcal{I})| \quad (10)$$

secures AP invariance under the decoupled transformations $\eta_1 = h_1(\psi)$, $\eta_2 = h_2(\phi)$, where $J_1(\psi | k, \mathcal{I})$ is calculated from $l(\psi | \phi = k, d, \mathcal{I})$ via (8), and $J_2(\phi | l, \mathcal{I})$ is calculated from $l(\phi | \psi = l, d, \mathcal{I})$. k and l are the necessary *hyperparameters* of the prior [5].

4. INVARIANCE PRIOR FOR A WIDE SIGNAL CLASS

Consider the hypothesis, \mathcal{I} , which analyzes observations, $d[n]$, in terms of m basis functions, $G_k[n, \omega]$, $\omega \in \bar{\mathcal{R}}^+$:

$$d[n] = \sum_{k=1}^m b_k G_k[n, \omega] + e[n], \quad n = 0, \dots, N-1. \quad (11)$$

b_k are the linear coefficients, and $e[n]$ are the residuals ('noise'). This rich class [5, 9, 11] embraces a wide range of signal analysis techniques as special cases, including Fourier and wavelet analyses, etc. Gathering the N observations into vector-matrix form, then

$$d = s + e = G(\omega)b + e, \quad (12)$$

where $G(\omega) \in \bar{\mathcal{C}}^{N \times m}$ with (n, k) th element $G_k[n, \omega]$, and $e \sim \mathcal{N}(0, \Sigma)$, being complex, zero-mean, Gaussian noise with known covariance matrix Σ (see Section 2.1). From (3,12):

$$l(\omega, b | d, \Sigma, \mathcal{I}) \propto \exp [-(d - Gb)^H \Sigma^{-1} (d - Gb)], \quad (13)$$

where $G \equiv G(\omega)$ for convenience. From (7,8,13), matrix calculus yields the following invariance prior for the signal parameter space [5]:

$$p(\omega, b | \Sigma, \mathcal{I}) \propto \Re \left[\begin{array}{ccc} b^H \frac{\partial G^H}{\partial \omega_i} \Sigma^{-1} \frac{\partial G}{\partial \omega_j} b & b^H \frac{\partial G^H}{\partial \omega_i} \Sigma^{-1} G & j b^H \frac{\partial G^H}{\partial \omega_i} \Sigma^{-1} G \\ G^H \Sigma^{-1} \frac{\partial G}{\partial \omega_j} b & G^H \Sigma^{-1} G & j G^H \Sigma^{-1} G \\ -j G^H \Sigma^{-1} \frac{\partial G}{\partial \omega_j} b & -j G^H \Sigma^{-1} G & G^H \Sigma^{-1} G \end{array} \right]^{\frac{1}{2}}. \quad (14)$$

Certain submatrices above are represented by their (i, j) th element, i th row or j th column as appropriate, for notational convenience. Techniques for marginalizing out Σ , when it is unknown *a priori*, are described in [1, 5]. Useful special cases include the following:

(i) If ω and b are independent *a priori*, then, from (10) [5]:

$$p(\omega, b | \Sigma, k, \mathcal{I}) \propto \Re \left[k^H \frac{\partial G^H}{\partial \omega_i} \Sigma^{-1} \frac{\partial G}{\partial \omega_j} k \right]^{\frac{1}{2}}, \quad (15)$$

where the information matrix on the right is denoted by its (i, j) th element.

(ii) If all ω_j and b_k are independent, then, from (10) [5]:

$$p(\omega, b | \Sigma, k, l, \mathcal{I}) \propto \left[\prod_{i=1}^r k^H \frac{\partial G_i^H}{\partial \omega_i} \Sigma^{-1} \frac{\partial G_i}{\partial \omega_i} k \right]^{\frac{1}{2}}. \quad (16)$$

G_i denotes G as a function of ω_i when all other ω_j are held constant at scalar values which are gathered together in the hyperparameter vector, l , on the lefthand side. Hence, (16) is a product of r scalar functions parameterized successively by ω_i , achieving the necessary *a priori* independence.

5. AN ENTROPIC PRIOR FOR SIGNAL IDENTIFICATION

The expected information (i.e. *entropy*)—in the Shannon sense—gained upon observing $\Theta \sim p(\theta | \mathcal{I})$ is given by

$$\mathcal{H}_{\bar{\Theta}} = - \int_{\bar{\Theta}} p(\theta | \mathcal{I}) \ln [p(\theta | \mathcal{I})] d\theta \quad (17)$$

(assuming Lebesgue measure on $\bar{\Theta}$ [3, 5]). An intuitively appealing prior assignment strategy is to maximize (17) with respect to $p(\theta | \mathcal{I})$, subject to any testable constraints on Θ . This yields the *minimum information prior* [2, 5-7]. A novel case is now described for the model structure (12).

From (12), the mean square value of the signal over the observation window is $\mathcal{T}(\omega, b) = \frac{1}{N} b^H G^H G b$, with ensemble average over $\bar{\Theta}$ given by

$$\mathcal{E}[\mathcal{T}(\omega, b)] = \gamma^2. \quad (18)$$

$\mathcal{E}[\cdot]$ is with respect to the prior $p(\omega, b | \mathcal{I})$, and γ^2 , the expected square value (i.e. 'power'), is testable via an ergodicity assumption. Lagrangian optimization of (17), constrained by (18), yields a density of the kind [5]

$$p(\omega, b | \alpha, \mathcal{I}) \propto \exp [-\alpha b^H G^H G b], \quad (19)$$

where α is a Lagrange multiplier. Standard manipulations yield $\alpha = m/N\gamma^2$ [5].

- It has been shown [5, 11] that (19) is similar to the *Ockham Prior* for the signal structure (12); i.e. it favours objectively 'simpler' models, thereby regularizing the AP inference for ω and b , in the sense explained by (5);
- (19) may be generalized to the case where there is testable information concerning signal autocorrelations to lag $N-1$, i.e.

$$\mathcal{E}[ss^H] = C. \quad (20)$$

$s = Gb$ (12) and $C \in \bar{\mathcal{C}}^{N \times N}$ is the known Hermitian correlation matrix for the signal. Maximization of (17), constrained by (20), yields the following entropic prior [5]:

$$p(\omega, b | C, \mathcal{I}) \propto \exp [-b^H G^H C^{-1} G b]. \quad (21)$$

- Prior independence is not engendered by the priors (19,21). If the second moments of b and ω are available independently *a priori*, then conventional independent Gaussian priors are obtained on the spaces, $\bar{\mathcal{B}}$ and $\bar{\mathcal{W}}$, of b and ω respectively, using the procedure above [5].

6. SUPERRESOLUTION OF CLOSELY SPACED FREQUENCIES

Consider the classic problem of inferring the difference frequency, ω , in the following signal model:

$$d[n] = e^{j\omega_1 n} [b_1 + b_2 e^{j\omega n}] + e[n], \quad n = 0, \dots, N-1, \quad (22)$$

where $E[n] = e[n] \sim \mathcal{N}(0, \sigma^2)$, $\forall n$. ω_1 is known *a priori* and b_1, b_2 and ω are to be inferred. In the sequel, *marginal a priori* densities for ω ,

$$p(\omega | \mathbf{x}, \mathcal{I}) = \int_{\mathbf{B}} p(\omega, \mathbf{b} | \mathbf{x}, \mathcal{I}) d\mathbf{b},$$

are compared, where \mathbf{x} denotes any necessary hyperparameters, and the integrand on the right is the joint parameter prior which will be assigned using the strategies outlined in this paper.

6.1. Invariance Prior

Assuming ω and \mathbf{b} are independent *a priori*, then, from (15) or (16), $p(\omega | \sigma, \mathbf{k}, \mathcal{I}) \propto p(\omega, \mathbf{b} | \sigma, \mathbf{k}, \mathcal{I})$ (assuming a range constraint on \mathbf{b}). Substituting (22) (via (11) and (12)) into (15) or (16), noting that $\omega = \omega$ (a scalar) in this case, and letting $\Sigma = \sigma^2 \mathbf{I}_N$ (i.e. diagonal), then $p(\omega | \sigma, \mathbf{k}, \mathcal{I}) \propto \text{const.}$ (see Fig. 1). In the Figure, one DFT bin is defined to be $2\pi/N$ rads/sample. If independence is relaxed *a priori*, then (14) must be employed, yielding a prior which is also illustrated in Fig. 1. It behaves similarly to the one proposed in [12], being approximately uniform for $\omega > 1$ DFT bin, and small for $\omega < 1$ bin, thereby eliminating the possibility of model rank deficiency *a priori*.

6.2. Entropic Prior

From (19), the marginal entropic prior for ω is [5]

$$\begin{aligned} p(\omega | \alpha, \mathcal{I}) &= \int_{\mathbf{B}} p(\omega, \mathbf{b} | \alpha, \mathcal{I}) d\mathbf{b} \propto |\mathbf{G}^H \mathbf{G}|^{-1} \\ &= \left(N^2 - \sin^2 \left(\frac{N\omega}{2} \right) / \sin^2 \left(\frac{\omega}{2} \right) \right)^{-1}, \quad (23) \end{aligned}$$

assuming $\mathbf{B} = \overline{\mathbf{C}}^m$. This is also plotted in Fig. 1. As a regularizing prior, it outperforms the others; i.e. it is large for $\omega < 1$ bin, and has minima at integer DFT bins, $\omega = k2\pi/N$, $k \in \mathbb{Z}^+$ (i.e. where the basis functions of (22) are orthogonal [5]). As such, this prior engenders Ockham characteristics in the inference: if data support for the hypothesis (22) is insufficient, (23) dominates the AP inference, $p(\omega | \mathbf{d}, \sigma, \mathcal{I})$, because of the singularity at $\omega = 0$, and this encourages the *rejection* of the hypothesis. If data support is sufficient, the LF will dominate the prior (23), and good estimates will result. As mentioned earlier, (19) and (23) have been deduced, from fundamentals, as the Ockham Prior, whose excellent regularization properties, in both Estimation and Model Selection, have been explored in simulations [5, 9, 11]. Robust, threshold-free inferences are possible over a wide range of Signal-to-Noise Ratios (SNRs) and observation window lengths (N).

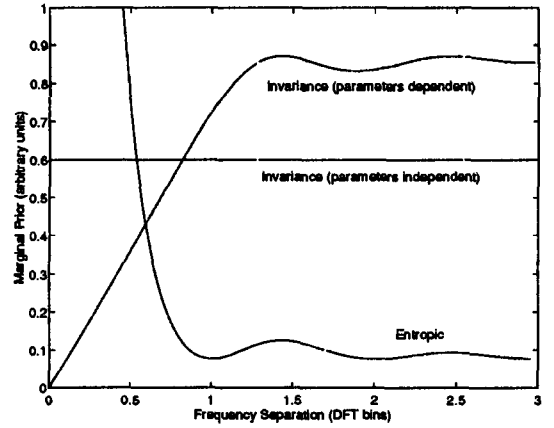


Figure 1: Invariance and Entropic Priors for ω in the 2-Cis Model (22).

References

- [1] G. E. P. Box and G. C. Tiao. *Bayesian Inference in Statistical Analysis*. Addison-Wesley, 1973.
- [2] J. O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag, 2nd edition, 1985.
- [3] J. M. Bernardo and A. F. M. Smith. *Bayesian Theory*. Wiley, 1994.
- [4] H. Jeffreys. *Theory of Probability*. Oxford Univ. Press, 3rd edition, 1961.
- [5] A. P. Quinn. *Bayesian Point Inference in Signal Processing*. PhD thesis, Cambridge University Engineering Dept., 1992.
- [6] R. D. Rosenkrantz, editor. *E. T. Jaynes: Papers on Probability, Statistics and Statistical Physics*. D. Reidel, Dordrecht-Holland, 1983.
- [7] A. Zellner. *An Introduction to Bayesian Inference in Econometrics*. John Wiley and Sons, 1971.
- [8] G. L. Bretthorst. *Bayesian Spectrum Analysis and Parameter Estimation*. Springer-Verlag, 1989.
- [9] A. P. Quinn. A consistent, numerically efficient Bayesian framework for combining the selection, detection and estimation tasks in model-based signal processing. In *Proc. IEEE Int. Conf. on Acoust., Sp. and Sig. Proc.*, Minneapolis, 1993.
- [10] O. Schwartz and A. Quinn. Fast and accurate texture-based image segmentation. In *Proc. 3rd IEEE Int. Conf. on Image Proc.*, Lausanne, 1996.
- [11] A. Quinn. A general complexity measure for signal identification using Bayesian inference. In *Proc. IEEE Workshop on Nonlinear Sig. and Image Process. (NSIP '95)*, Greece, 1995.
- [12] P. M. Djurić and H.-T. Li. Bayesian spectrum estimation of harmonic signals. *IEEE Sig. Proc. Letts.*, 2(11), 1995.