

Prior near ignorance for inferences in the k -parameter exponential family

A. Benavoli* and M. Zaffalon

IDSIA, Galleria 2, CH-6928 Manno (Lugano), Switzerland

This paper proposes a model of prior ignorance about a multivariate variable based on a set of distributions \mathcal{M} . In particular, we discuss four minimal properties that a model of prior ignorance should satisfy: invariance, near-ignorance, learning and convergence. Near-ignorance and invariance ensure that our prior model behaves as a vacuous model with respect to some statistical inferences (e.g., mean, credible intervals, etc.) and some transformation of the parameter space. Learning and convergence ensure that our prior model can learn from data and, in particular, that the influence of \mathcal{M} on the posterior inferences vanishes with increasing numbers of observations. We show that these four properties can all be satisfied by a set of conjugate priors in the multivariate exponential families if the set \mathcal{M} includes finitely additive probabilities obtained as limits of truncated exponential functions. The obtained set \mathcal{M} is a model of prior ignorance with respect to the functions (queries) that are commonly used for statistical inferences and, because of conjugacy, it is tractable and easy to elicit. Applications of the model to some practical statistical problems show the effectiveness of the approach.

Keywords: Prior ignorance, exponential families, lower and upper expectations.

AMS Subject Classification: 62C10; 62F15.

1. Introduction

In performing a Bayesian analysis concerning the unknown value w of a variable $W \in \mathcal{W}$, one is required to express prior beliefs about w in the form of a prior distribution. However, we are often in a condition of prior ignorance about w and, thus, the problem is how to choose a prior distribution representing “ignorance”. There are two main avenues to represent ignorance. The first assumes that ignorance can be modelled satisfactorily by a so-called noninformative prior, by which it is meant a prior that contains no information about w . The problem is how to define the meaning of “containing no prior information”. It is common to call noninformative a prior that is invariant under certain re-parametrizations of the parameter space. For instance, when w is a location parameter, i.e., $\mathcal{W} = \mathbb{R}^k$, the prior should be invariant to translations. The justification is that, since the choice of an origin for a unit of measurement is often arbitrary, the noninformative prior should be independent of this choice. The uniform density on \mathbb{R}^k , i.e., $p(w) = 1$, is the improper prior that is invariant to translations. Other invariance properties that are often imposed to derive noninformative priors are: scale invariance when $\mathcal{W} = \mathbb{R}^+$, permutation invariance, symmetry, etc. A discussion about different ways to derive noninformative priors can be found in [1], [2, Ch. 6] and [3, Sec. 5.6.2]. Much criticism has been raised concerning the use of noninformative priors [2, Sec. 3.4.5]; for instance that they are often improper. A way to solve this problem is to use an approximation of these improper priors. For instance, in case of the uniform on \mathbb{R} we can restrict the support of the prior to be a large closed interval, so that the resulting uniform truncated on this interval is proper. An issue of this approach is how to choose this interval; to do that we need some prior information that is not always available. Another possible way to justify the use of improper priors has been proposed in [4] and [5]. Noninformative priors can be regarded as limits of proper informative priors (a limit of proper truncated

priors when the truncation vanishes). The advantage of this interpretation is that it allows us to compute prior inferences (expectations) w.r.t. noninformative priors by computing this limit. This is not possible when we directly use improper priors, since the integrals diverge.

It can be shown that the probability measures that we obtain by this limit are finitely additive; this provides a justification of improper priors from a Bayesian perspective. The advantage of this interpretation of improper priors is that (i) we do not need to define a finite truncation; (ii) we can compute the prior expectation by calculating limits. A-posteriori, if the distribution that we obtain by combining the limit improper prior and the likelihood is proper, we can compute posterior expectations by moving the limit under the integral sign. In other words, we can compute first the limit which gives us the improper prior and then use this prior to compute the posterior. The presence of finitely additive measures is hidden in this case.

Other two criticisms to noninformative priors are that they are not unique, see for instance [2, Sec. 3.3.3], and that they do not really model lack of prior information, but only invariance [6]. To solve these issues, many authors [6–11] suggest that lack of prior information should be expressed in terms of a family \mathcal{M} consisting of all prior distributions that are compatible with the available prior information. Inferences should then be carried out by considering the whole family \mathcal{M} . This approach is known as *Bayesian robustness* or *Bayesian sensitivity analysis*—examples are ε -contamination models [2, 12], restricted ε -contamination models [9], intervals of measures [10, 12], the density ratio class [6, 10], etc.

In case almost no prior information is available on W , \mathcal{M} should be as large as possible in order to describe this state of prior ignorance. The natural candidate for \mathcal{M} to represent complete ignorance is the set of all distributions. However, it turns out that the posterior inferences obtained from this set are *vacuous* [6, Sec. 7.3.7], i.e., the posterior set of distributions coincides with the prior set of distributions. This means that our prior beliefs do not change with experience, or, in other words, that there is no learning from data. Therefore, the vacuous prior model is not a statistically useful way to model our prior ignorance. There is then a compromise to be made. Walley suggests, as an alternative, the use of an almost vacuous model which he calls a “near-ignorance prior” [6, 11]. This is a model that behaves a priori as a vacuous model for some basic inferences (e.g., prior mean, prior credible regions) but it always provides non-vacuous posterior inferences. Near-ignorance models allow us to start a statistical analysis with very weak assumptions about the problem of interest and, thus, to implement a genuine objective-minded approach to inference. Near-ignorance models have been presented in [6, 11] to address inference problems in which: (i) the observations are discrete and the unknown W is a vector of variables with bounded components (the so-called *imprecise Dirichlet model*, or IDM); (ii) the observations are continuous and the unknown W is an unbounded but scalar variable (Gaussian or double exponential family models [11]). Near-ignorance often naturally leads to invariance. For instance, the IDM is known to satisfy permutation invariance and invariance to coarsening/refinement of the categories [6]. When this is not the case, invariance properties can be imposed to near-ignorance priors as suggested in [11].

In a recent paper [13] we have made a first proposal to generalize near-ignorance to the one-parameter exponential families of distributions. Considering the case that the likelihood model is a density belonging to the one-parameter exponential families, i.e., $w \in \mathcal{W} \subseteq \mathbb{R}$, and \mathcal{M} includes the corresponding *proper conjugate exponential priors*; we have shown that there exists a choice of the parametrization of \mathcal{M} which satisfies near-ignorance without leading to vacuous posterior inferences, and which is unique up to the choice of its size which determines the degree of robustness (or caution) of the inferences. Stated differently, we have proven that the set of priors \mathcal{M} satisfying near-ignorance without leading to vacuous posterior inferences can be uniquely obtained by letting the parameters of the conjugate exponential prior vary jointly in suitable sets. This work allowed us to derive again the imprecise Beta model, which is the univariate version of the IDM, and the set of Gaussian priors discussed in [11] from a more general model of near-ignorance. Furthermore, we were able to derive new models of near-ignorance that were not available before such as, for instance, a model based on a set of Gamma priors.

In this paper, we extend the work in [13] in two main directions. First, we propose a way to reconcile the two approaches to design a prior model in case of lack of prior information, i.e., invariance and Walley's prior near-ignorance. We define four *minimal properties* that a prior ignorance model \mathcal{M} should satisfy: prior-ignorance, invariance, learning from data and asymptotic convergence. To apply these properties to define a prior model, the analysts must choose: (i) w.r.t. which functions (e.g., mean, variance, credible intervals, etc.) she/he wants to be prior ignorant; (ii) w.r.t. which group transformations she/he wants to be invariant; (iii) w.r.t. which functions she/he wants to satisfy learning and convergence. These are general principles that allow the analyst to elicit a prior model in a condition of ignorance. Then we specialize these properties to the k -dimensional exponential families $\mathcal{W} = \mathbb{R}^k$, which is the main contribution of this work. For the multivariate exponential family in $\mathcal{W} = \mathbb{R}^k$, it is natural to impose: (i) prior-ignorance w.r.t. the function $b(w)$ (the mean of the observation variable); (ii) invariance to translations (w is a location parameter) and, depending on the exponential family we are considering, invariance to permutations, symmetry, invariance to coarsening/refinement, etc.; (iii) learning and convergence w.r.t. any (measurable) function.

Moreover, we show that translation invariance, near-ignorance and learning can all be satisfied by a set of conjugate priors \mathcal{M} if it includes finitely additive probabilities obtained as limits of truncated increasing/decreasing exponential functions. As discussed previously, when the set of posteriors we obtain by combining the likelihood and the improper increasing/decreasing exponential functions includes only proper distributions, we can move the limit under the integral sign and compute inferences as we normally do with improper priors. When this is not the case, we must deal with this limit to obtain the lower and upper expectations. We will show that for the most important inferences in statistical analysis (mean, credible interval, etc.) computing this limit is easy. We discuss how the obtained near-ignorance models compare with other models that have been proposed for near-ignorance: the imprecise Dirichlet model [6], the bounded derivative model [14] and the nonparametric predictive inference model [15].

Finally, we show the application of our near-ignorance model to some practical statistical problems. In particular, to state the practical effectiveness of our prior ignorance model, we employ it to implement a prior near-ignorance version of the one-sample Bayesian t-test. Our new test has the peculiarity that it can suspend the judgement, which means that it stays indeterminate with respect to the choice of the null or the alternative hypothesis, in case it deems that there is not enough information in the data to take a reliable decision. By means of numerical simulations, we compare our test with the frequentist t-test, a Bayesian t-test based on an improper prior, a hierarchical Bayesian t-test based on an improper hyper-prior and show that our test is more robust. This means that if we take a random decision, with 50/50 chance, on the instances where our test is indeterminate (and take the recommended decision on the others), then we obtain a test that has the same performance (power) of the frequentist t-test and the Bayesian t-tests. This is an important result from a practical point of view, because it shows that using near-ignorance leads to isolating instances on which traditional tests have no clue. Finally, as another example, we apply our prior near-ignorance model to real data: the USA 2004 election poll.

Notation x, w, x_0, w_0, y, y_0 denote vectors in \mathbb{R}^k . Their components, for $i = 1, \dots, k$, are denoted by x_i, w_i for variables that do not have a subscript and by x_{0i}, w_{0i} for variables that have a subscript. g denotes a bounded real-valued function on a subset of \mathbb{R}^k . $E[\cdot]$ denotes the expectation operator, while $\underline{E}[\cdot]$ and $\overline{E}[\cdot]$ denote the lower and, respectively, upper expectation operators.

2. Prior Near-Ignorance and Invariance

The aim of this section is to define which minimal properties the set of priors \mathcal{M} should satisfy in case there is (almost) no prior information about a random variable W taking values in $\mathcal{W} \subseteq \mathbb{R}^k$.

We start by defining prior ignorance.

Definition 2.1. A state of ignorance about a real-valued bounded function (RVBF) g of W can be modelled as follows: $\underline{E}[g] = \inf g$ and $\overline{E}[g] = \sup g$, where $\underline{E}, \overline{E}$ denote the lower and upper bounds of the expectation of g computed w.r.t. the set of priors \mathcal{M} .

It can be observed that the range of $E[g]$ under the set of priors \mathcal{M} is the same as the original range of g . In other words, by specifying the set of priors \mathcal{M} , we are not giving any information on the value of the expectation of g . It corresponds to stating that we are in a state of ignorance.

In case of lack of prior information, another property that is desirable for \mathcal{M} is that of satisfying some invariance principle. Consider for instance the case in which $\mathcal{W} = \mathbb{R}$. Then, in case of lack of prior information, it seems to be reasonable that, for all subsets of reals A and for all real numbers a , the prior probabilities of the set A and the shifted set $A + a$ should coincide. This property is called *translation invariance*. The principle of invariance to transformations can be formalized as follows [6, Sec. 3.5] (hereafter, we use \mathcal{F} to denote a group of transformations of \mathcal{W}).¹

Definition 2.2. A lower expectation \underline{E} on $\mathcal{L}(\mathcal{W})$, where $\mathcal{L}(\mathcal{W})$ is the linear space of all RVBFs on \mathcal{W} , is called \mathcal{F} -invariant in $\mathcal{G} \subseteq \mathcal{L}(\mathcal{W})$ if $\underline{E}[g(f)] = \underline{E}[g]$ whenever $g \in \mathcal{G}$, $f \in \mathcal{F}$ and $g(f) \in \mathcal{G}$.

For example, let $\mathcal{W} = \mathbb{R}$ and define the translations f_a by $f_a(w) = w + a$. The group of transformations $\mathcal{F} = \{f_a : a \in \mathbb{R}\}$ is called the translation group on \mathbb{R} . Since $g(f_a(w)) = g(w + a)$, \mathcal{F} -invariance is translation invariance in this case.² If $g = I_{\{A\}}$,³ with $A \subseteq \mathcal{W}$ measurable, the above equality $\underline{E}[g(f)] = \underline{E}[g]$ means that the set A and the shifted set $A + a$ should have the same lower probability for any value of a .

The only set of priors which satisfies prior ignorance, i.e., $\underline{E}[g] = \inf g$ and $\overline{E}[g] = \sup g$, for all $g \in \mathcal{L}(\mathcal{W})$ and which is \mathcal{F} -invariant to all possible \mathcal{F} is the set of all probabilities (vacuous prior model). However, it turns out that the posterior inferences obtained from this set are *vacuous* [6, Sec. 7.3.7], i.e., the posterior set of probabilities coincides with the prior set of probabilities. This means that our prior beliefs do not change with experience (i.e., there is no learning from data). Therefore, the vacuous prior model is not a useful way to model our prior ignorance in statistics. There is then a compromise to be made: imposing prior ignorance and invariance only on a subset of the possible RVBFs and group of transformations. The set of minimal properties that \mathcal{M} or, equivalently, the lower and upper expectations it generates, should satisfy to be an invariant model of prior ignorance and produce consistent and meaningful posterior inferences are formalized hereafter.

- (A1) **$\{\mathcal{F}_i\}$ -prior invariance.** For a chosen set of RVBFs \mathcal{G}_1 and a class of groups of transformations $\{\mathcal{F}_i\}$, the prior upper and lower expectations are $\{\mathcal{F}_i\}$ -invariant, i.e., $\underline{E}[g(f)] = \underline{E}[g]$ ⁴ for any $\mathcal{F}_i \in \{\mathcal{F}_i\}$, $g \in \mathcal{G}_1$, $f \in \mathcal{F}_i$ and $g(f) \in \mathcal{G}_1$.
- (A2) **\mathcal{G}_0 -prior ignorance.** The prior upper and lower expectations of some suitable set of RVBFs \mathcal{G}_0 under \mathcal{M} are vacuous, i.e., $\underline{E}[g] = \inf g$ and $\overline{E}[g] = \sup g$ for all $g \in \mathcal{G}_0$.
- (A3) **\mathcal{G} -learning.** For a chosen set of RVBFs $\mathcal{G} \supseteq \mathcal{G}_0$ and for each $g \in \mathcal{G}$ satisfying $\overline{E}[g] - \underline{E}[g] > 0$, there exists a finite $\delta > 0$ (possibly dependent on g) such that for each $n \geq \delta$ and sequence of observations $y^n = (y_1, \dots, y_n)$, at least one of these two conditions is

¹More precisely, \mathcal{F} denotes a semigroup of transformations of \mathcal{W} . That is, each $f \in \mathcal{F}$ maps \mathcal{W} into itself, and the composition $f_1 f_2$ defined by $f_1(f_2(w))$ is in \mathcal{F} whenever $f_1, f_2 \in \mathcal{F}$. The semigroup \mathcal{F} is Abelian if $f_1 f_2 = f_2 f_1$ whenever $f_1, f_2 \in \mathcal{F}$.

²In this paper we mainly focus on translation invariance. However, for multivariate models, we will impose other invariance properties: invariance to permutations and invariance to representation.

³Note that $I_{\{A\}}$ is the indicator function of set A , i.e., $I_{\{A\}}(x) = 1$ if $x \in A$ and zero otherwise.

⁴Equivalently, if $-g$ and $-g(f)$ belong to \mathcal{G}_1 then $\underline{E}[-g(f)] = \underline{E}[-g]$, which implies that $\overline{E}[g(f)] = \overline{E}[g]$ being $\overline{E}[g] = -\underline{E}[-g]$ for any g .

satisfied:

$$\underline{E}[g|y^n] \neq \underline{E}[g], \quad \overline{E}[g|y^n] \neq \overline{E}[g], \quad (1)$$

where $\underline{E}[\cdot|y^n]$ and $\overline{E}[\cdot|y^n]$ denote the posterior lower and upper expectations of g .

(A4) Convergence. For each RVBF $g \in \mathcal{G}$ and sequence of observations y^n , the following conditions are satisfied for $n \rightarrow \infty$:

$$|E^*[g|y^n] - \underline{E}[g|y^n]| \rightarrow 0, \quad |E^*[g|y^n] - \overline{E}[g|y^n]| \rightarrow 0, \quad (2)$$

where $E^*[g|y^n]$ is the posterior expectation w.r.t. the posterior density derived, via Bayes' rule, from the likelihood model and some prior on \mathcal{W} that is dominated by the likelihood.¹

Property (A1) states that \mathcal{M} should be prior invariant w.r.t. some class of groups of transformations $\{\mathcal{F}_i\}$ and RVBFs \mathcal{G}_1 . In general we expect \mathcal{G}_1 to be very large and coincide with \mathcal{G} defined in (A3). Again in case \mathcal{M} includes all possible distributions then (A1) holds for any \mathcal{F}_i [6, Sec. 3.5.6]. Here, conversely, we require that (A1) is satisfied for some specific choice of $\{\mathcal{F}_i\}$. This choice depends on the aims and on the application of the set of priors \mathcal{M} .

Property (A2) states that \mathcal{M} should be vacuous a priori w.r.t. some set of RVBFs \mathcal{G}_0 , i.e., the lower and upper expectations of $g \in \mathcal{G}_0$ respectively coincide with the infimum and the supremum of g . In case \mathcal{M} includes all possible distributions then (A1) holds for any function g . Here, conversely, we require that (A1) is satisfied for some subset of RVBFs \mathcal{G}_0 . The subset of RVBFs \mathcal{G}_0 used in (A1) should include the RVBFs g w.r.t. which we state our condition of prior near-ignorance.

The sets \mathcal{G}_0 in (A2) and $\{\mathcal{F}_i\}$ in (A1) should be also as large as possible to guarantee that also \mathcal{M} is as large as possible, but no too large to be incompatible with the requirement (A3) of learning. In fact, property (A3) states that \mathcal{M} should be non-vacuous a-posteriori for any RVBF $g \in \mathcal{G} \supseteq \mathcal{G}_0$, which is a condition for learning from the observations. The set of RVBFs \mathcal{G} used in (A3) should consist of the RVBFs g w.r.t. which we are interested in computing expectations (i.e., making inferences). The fact that \mathcal{G} must include \mathcal{G}_0 is the only constraint on \mathcal{G} , meaning that (A3) requires that \mathcal{M} is not vacuous w.r.t. all these RVBFs for which the prior near-ignorance has been imposed. Since \mathcal{M} is a model of prior near-ignorance, it is also desirable that the influence of \mathcal{M} on the posterior inferences vanishes with increasing number of observations n . This is a sort of ultimate agreement by the accumulation of evidence.

Property (A4) states that, for $n \rightarrow \infty$, \mathcal{M} should give the same lower and upper expectations of $g \in \mathcal{G}$ as those obtained from some sufficiently regular prior on \mathcal{W} .

In order to better understand properties (A1)–(A4), in Section 4 we show their instantiation in the case of the exponential families. Before discussing these results, in the next section we introduce the exponential families of densities and review their main properties [3, 17–19].

3. Exponential Families

Consider a sampling model where i.i.d. samples of a k -dimensional random vector Y are taken from a sample space \mathcal{Y} . For $w \in \mathbb{R}^k$ define $b(w) = \ln \int_{\mathcal{Y}} u(y) \exp(y^T w) dy$, (a sum in case \mathcal{Y} is discrete), where $u(\cdot) : \mathcal{Y} \rightarrow \mathbb{R}^+$ and T denotes transpose. Let \mathcal{W} be $\{w : b(w) < \infty\}$ [20]; it can be shown that \mathcal{W} is a convex set and it is called the *natural parameter space*. It is further

¹We point the reader to [16, Ch. 20] for a general discussion about dominated priors. When the likelihood belongs to the exponential families (the focus of this paper), as dominated prior we may consider any proper conjugate prior, the improper uniform or other sufficiently regular priors. The posterior becomes asymptotically Normal in these cases.

assumed that \mathcal{W} is a non-empty open set in \mathbb{R}^k , i.e., we restrict attention to regular exponential families [20]. The probability density

$$p(y|w) = u(y) \exp(y^T w - b(w)), \quad y \in \mathcal{Y}, \quad (3)$$

is called the canonical form of representation of the exponential family; $y \in \mathcal{Y}$ is a so-called sufficient statistic; $w \in \mathcal{W}$ is a so-called natural parameter. The functions $u(\cdot) : \mathcal{Y} \rightarrow \mathbb{R}^+$ and $b(\cdot) : \mathcal{W} \rightarrow \mathbb{R}$ characterize each particular family of the exponential families of distributions. The canonical form has some useful properties. First, the mean and variance of Y are given by

$$E[Y|w] = \nabla b(w), \quad \text{Var}[Y|w] = \nabla^2 b(w), \quad (4)$$

where $\nabla b = [\frac{\partial b}{\partial w_1}, \frac{\partial b}{\partial w_2}, \dots, \frac{\partial b}{\partial w_k}]^T$ is the gradient (a column vector), w_i denotes the i -th component of the vector and $\nabla^2 b(w) > 0$ for each $w \in \mathcal{W}$ is the Hessian matrix. Second, in the case of n i.i.d. observations y_i , it follows that

$$\prod_{i=1}^n p(y_i|w) = \exp(n(\hat{y}_n^T w - b(w))) \prod_{i=1}^n u(y_i), \quad w \in \mathcal{W}, \quad (5)$$

where $\hat{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$ is the sample mean of y_i which, together with n , is a sufficient statistic of y_1, \dots, y_n for inference about W under the i.i.d. assumption.

Furthermore, by interpreting the density function in (5) as a likelihood function, we can define the corresponding conjugate prior. A probability density $p(w|n_0, y_0)$, parametrized by $n_0 \in \mathbb{R}$ and $y_0 \in \mathbb{R}^k$, is said to be the canonical conjugate prior of (3) if

$$p(w|n_0, y_0) \propto \exp(n_0(y_0^T w - b(w))), \quad w \in \mathcal{W}, \quad (6)$$

where n_0 is the so-called number of pseudo-observations and y_0 is the so-called pseudo-observation. Consider the convex hull of the support set of the measure $u(dy)$, and then let \mathcal{Y}_0 be the interior of this convex set. From (4), it follows that \mathcal{Y}_0 is the set of possible values of ∇b [20].

Lemma 3.1. *Consider the prior $p(w|n_0, y_0) \propto \exp(n_0(y_0^T w - b(w)))$. If $y_0 \in \mathcal{Y}_0$ and $0 < n_0 < \infty$, then the density is proper, i.e., there exists a normalization constant $k(n_0, y_0) > 0$ such that*

$$\int k(n_0, y_0) \exp(n_0(y_0^T w - b(w))) dw = 1.$$

The lemma is proven in [20, Th. 1]. Lemma 3.1 states that if $n_0 \in \mathbb{R}^+$ and $y_0 \in \mathcal{Y}_0$ then the kernel $\exp(n_0(y_0^T w - b(w)))$ is integrable and, thus, proper.

The likelihood and conjugate prior pair in the canonical exponential families satisfy a set of interesting properties, most of them are particularly useful to represent the nature of the Bayesian learning process. A list of such properties is given in the following lemmas.

Lemma 3.2. *For a pair of likelihood and conjugate prior in the canonical exponential family, it holds that the posterior density for w is:*

$$p(w|n_p, y_p) = k(n_p, y_p) \exp(n_p(y_p^T w - b(w))), \quad w \in \mathcal{W}, \quad (7)$$

where $n_p = n + n_0$ and $y_p = \frac{n_0 y_0 + n \hat{y}_n}{n + n_0}$, which is always proper if $n_0 > 0$ and $y_0 \in \mathcal{Y}_0$.

This follows by conjugacy, see [3, Ch. 5]. In (4), it has been shown that ∇b is the mean of Y conditional on w . Hence, ∇b is the quantity about which we will have prior beliefs before seeing

the data y and posterior beliefs after observing the data. Hence, it is interesting to compute the prior and posterior mean of ∇b . For sampling models such that $\nabla b(w) = x$, these expectations give also a closed formula for the prior and posterior expectations of x .

Lemma 3.3. *The prior mean of the function ∇b is $E[\nabla b|n_0, y_0] = y_0$ if $y_0 \in \mathcal{Y}_0$ and $n_0 > 0$ and the posterior mean: $E[\nabla b|n_p, y_p] = y_p = \frac{n_0 y_0 + n \hat{y}_n}{n + n_0}$.*

These results follow from the properties of exponential families [3, Ch. 5], [20]. Examples of conjugate exponential families can be found for instance in [3, Ch. 5].

4. Near-Ignorance For k -Parameter Exponential Families

Consider the problem of statistical inference about the real-valued parameter $w \in \mathcal{W} = \mathbb{R}^k$ from noisy measurements y_1, \dots, y_n and assume that the likelihood is completely described by an exponential family probability density function (PDF) (or probability mass function (PMF) in the discrete case):

$$\prod_{i=1}^n p(y_i|w) = \exp(n(\hat{y}_n^T w - b(w))) \prod_{i=1}^n u(y_i), \quad (8)$$

where the parameters of the likelihood, i.e., sample mean $\hat{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$ and $n \in \mathbb{N}$, are known (the likelihood is precisely known, a single PDF/PMF specified by the parameters \hat{y}_n and n). By conjugacy and following a Bayesian approach, as prior for w we may consider the proper PDF $p(w|n_0, y_0)$ defined in (6) for a given value of the parameter $y_0 \in \mathcal{Y}_0$ and $n_0 \in \mathbb{R}^+$. In the case there is not enough information about the value w of W to uniquely determine the values of the parameters y_0 and n_0 , we can consider the family \mathcal{M} of prior densities which represent prior-ignorance and satisfy some invariance principles, but without leading to vacuous posterior inferences. In other words, the question is how to construct a family \mathcal{M} of prior densities which satisfies properties (A1)–(A4) discussed in Section 2. For tractability of inferences, we restrict this set of priors only to include densities such that the resulting set of posteriors is conjugate with the likelihood model. Before answering this question for the general multivariate case, $\mathcal{W} = \mathbb{R}^k$, we focus on the univariate case, $\mathcal{W} = \mathbb{R}$. We start proving the following lemma.

Lemma 4.1. *The lower expectation model¹*

$$\underline{E}[g] = \liminf_{r \rightarrow \infty} \int g(w) p(w) dw, \quad (9)$$

with $p(w) = \frac{\ell}{\exp(\ell r)} \exp(\ell w) I_{\mathcal{W}_r}(w)$ for $\ell \neq 0$,

$$\mathcal{W}_r = \begin{cases} (-\infty, r] & \text{if } \ell > 0, \\ [-r, \infty) & \text{if } \ell < 0, \end{cases} \quad (10)$$

$r \in \mathbb{R}$ and $p(w) = \frac{1}{2r} I_{\mathcal{W}_r}(w)$ with $\mathcal{W} = [-r, r]$ for the case $\ell = 0$, satisfies $\underline{E}[g(f_a)] = \underline{E}[g]$ with $f_a(w) = w + a$ (translation invariance) for any g in the set of Borel-measurable RVBFs \mathcal{K} and $a \in \mathbb{R}$.

¹Let the least term ϕ of a sequence be a term which is smaller than all but a finite number of the terms which are equal to ϕ . Then ϕ is called the lower limit of the sequence.

Proof. For $\ell = 0$, $p(w) = \frac{1}{2r}I_{\mathscr{W}_r}(w)$ coincides with the uniform distribution on $\mathscr{W} = [-r, r]$, then

$$\begin{aligned}\underline{E}[g(f_a)] &= \liminf_{r \rightarrow \infty} \frac{1}{2r} \int_{-r}^r g(w+a)dw = \liminf_{r \rightarrow \infty} \frac{1}{2r} \int_{-r}^{r+a} g(w)dw \\ &= \liminf_{r \rightarrow \infty} \frac{1}{2r} \left(\int_{-r}^r g(w)dw + \int_r^{r+a} g(w)dw - \int_{-r}^{-r+a} g(w)dw \right) \\ &= \underline{E}[g] + \liminf_{r \rightarrow \infty} \frac{1}{2r} \left(\int_r^{r+a} g(w)dw - \int_{-r}^{-r+a} g(w)dw \right).\end{aligned}$$

Hence, being g bounded (it is a RVBF), it follows that both $\int_r^{r+a} g(w)dw$ and $\int_{-r}^{-r+a} g(w)dw$ are bounded for any r and, thus, the limit in the last equation above goes to zero for $r \rightarrow \infty$. Therefore, it results that $\underline{E}[g(f_a)] = \underline{E}[g]$. We have thus proven that the limit for $r \rightarrow \infty$ of the uniform distribution on $\mathscr{W} = [-r, r]$ satisfies translation invariance.

For $\ell \neq 0$ and $\mathscr{W}_r = (-\infty, r]$ (the case $\mathscr{W}_r = [-r, \infty)$ is similar),

$$\underline{E}[g(f_a)] = \inf_p \int g(w+a)p(w)dw = \inf_p \int g(w)p(w-a)dw,$$

where \inf_p means $\liminf_{r \rightarrow \infty}$ with respect to $p(w) = \frac{\ell}{\exp(\ell r)} \exp(\ell w)I_{(-\infty, r]}(w)$ or, equivalently, to $p(w-a) = \frac{\ell}{\exp(\ell r)} \exp(\ell(w-a))I_{(-\infty, r]}(w-a)$. Define $r_1 = r+a$, then $p(w-a) = \frac{\ell}{\exp(\ell r_1)} \exp(\ell w)I_{(-\infty, r_1]}(w)$ and, thus,

$$\underline{E}[g(f_a)] = \liminf_{r_1 \rightarrow \infty} \frac{\ell}{\exp(\ell r_1)} \int g(w) \exp(\ell w) I_{(-\infty, r_1]}(w) dw,$$

which is equal to $\underline{E}[g]$ and, thus, satisfies translation invariance. \square

Consider the kernel $\exp(\ell w)$; transformed back to the original parameter space, it becomes $\exp(\ell x)$ in the Normal case, $\theta^{\ell-1}(1-\theta)^{-\ell-1}$ in the Beta case and $\lambda^{\ell-1}$ in the Gamma case.¹ We can see these priors are generalizations of the translation invariant improper priors that are normally used in Bayesian analysis. These priors can be obtained by setting $\ell = 0$: they are respectively 1 (improper uniform), $\theta^{-1}(1-\theta)^{-1}$ (Haldane prior) and λ^{-1} (Jeffreys prior). All these priors are improper. Improper priors can be defined as a limit of truncations of proper priors as in (9). It can be shown that the probability measures that we obtain by this limit are finitely additive; this provides a justification of improper priors from a Bayesian perspective. We discuss this connection with more details in Appendix A. A-posteriori, if the distribution that we obtain by combining the limit improper prior and the likelihood is proper, we can compute posterior expectations by moving the limit under the integral sign. Hereafter, for notational convenience, we will often interchange the limiting procedure in $\underline{E}[g]$ with its limit kernel $\exp(\ell w)$, since they are in most of the cases equivalent in terms of posterior inferences. However, the reader should be aware that, for instance, in referring to the prior expected value of g w.r.t. $\exp(\ell w)$ we actually mean (9). Furthermore, again for notational convenience, we will denote the prior with $p(w) = \frac{\ell}{\exp(\ell r)} \exp(\ell w)I_{\mathscr{W}_r}(w)$ even in the case $\ell = 0$ where it should be $p(w) = \frac{1}{2r}I_{[-r, r]}(w)$ (this because the lower and upper expectations of RVBFs hereafter considered are obtained by priors of this form $p(w) = \frac{\ell}{\exp(\ell r)} \exp(\ell w)I_{\mathscr{W}_r}(w)$).

Theorem 4.2. *Assume that $\mathscr{W} = \mathbb{R}$ and consider properties (A1)–(A4), with $\mathscr{G} = \mathscr{G}_1$ including sufficiently smooth RVBFs,² $\mathscr{G}_0 = \{b'\}$, where b' is the first derivative of b , and $\mathscr{F}_i = \mathscr{F}$, where*

¹The differences are due to the Jacobians of the transformations.

²With sufficiently smooth RVBFs, we mean integrable w.r.t. the kernel $\exp(n(\hat{y}_n^T w - b(w))) \exp(\ell w)$ for any $\ell \in [-c, c]$, $n \in \mathbb{N}$ and $\hat{y}_n \in CI(\mathscr{G}_0)$, with support in \mathscr{W} and continuous on a neighbourhood of the point where the posterior relative to the improper uniform prior concentrates for $n \rightarrow \infty$.

$\mathcal{F} = \{f_a : a \in \mathbb{R}\}$ and $f_a(w) = w + a$, i.e., a translation of the parameter space. Given the parameter $c \in \mathbb{R}^+$, the following set of priors:

$$\mathcal{M} \propto \{\exp(\ell w), \ell \in [-c, c]\}, \quad (11)$$

satisfies (A1)–(A4) and conjugacy between likelihood and posterior.

Proof. Consider first property (A1). Define $\underline{E}[g] = \inf\{\underline{E}_\ell[g] : \ell \in [-c, c]\}$, where $\underline{E}_\ell[g] = \liminf_{r \rightarrow \infty} \frac{\ell}{\exp(\ell r)} \int g(w) \exp(\ell w) I_{\mathcal{W}_r}(w) dw$. From Lemma 4.1, it follows that $\underline{E}_\ell[g(f_a)] = \underline{E}_\ell[g]$ for any ℓ and, thus, translation invariance holds for any kernel $\exp(\ell w)$. Thus, it holds also for the kernel ℓ which attains the infimum in $\underline{E}[g]$. Therefore, $\underline{E}[g(f_a)] = \underline{E}[g]$ for any $g \in \mathcal{G}_1$ and, thus, the set of priors \mathcal{M} is invariant under translations. Furthermore, it can be seen that the elements of \mathcal{M} satisfy conjugacy between posteriors and the likelihood model, in fact: $\exp(n(\hat{y}_n w - b(w))) \exp(\ell w) = \exp(n(\frac{n\hat{y}_n + \ell}{n} w - b(w)))$.

Consider now (A2). From Section 3 it follows that $\nabla b = b' \in \mathcal{Y}_0$ (we are in the scalar case, $\mathcal{W} = \mathbb{R}$) and the first derivative of b is always increasing in \mathcal{W} since the second derivative b'' is positive. Therefore, since for $\ell > 0$, $\exp(\ell w)$ is always increasing in \mathcal{W} , it follows that:

$$\bar{E}[b'] = \sup_{\ell \in [-c, c]} \limsup_{r \rightarrow \infty} \frac{\ell}{\exp(\ell r)} \int b'(w) \exp(\ell w) I_{\mathcal{W}_r}(w) = \sup \mathcal{Y}_0,$$

where $\mathcal{W}_r = (-\infty, r]$. The last equality can be easily derived by the fact that $\exp(\ell w)$ has a maximum in $w = r$ and thus the mass of $\exp(\ell w)/\exp(\ell r)$ is concentrated on a small neighborhood of $w = r$. In fact, given r for any $r_1 < r$ it holds that

$$\begin{aligned} & \frac{\ell}{\exp(\ell r)} \int_{-\infty}^r b'(w) \exp(\ell w) \\ &= \frac{\ell}{\exp(\ell(r-r_1)) \exp(\ell r_1)} \int_{-\infty}^{r_1} b'(w) \exp(\ell w) + \frac{\ell}{\exp(\ell r)} \int_{r_1}^r b'(w) \exp(\ell w) \\ &> \frac{\ell}{\exp(\ell(r-r_1)) \exp(\ell r_1)} \int_{-\infty}^{r_1} b'(w) \exp(\ell w) + \frac{\exp(\ell r) - \exp(\ell r_1)}{\exp(\ell r)} b'(r_1) \stackrel{r \rightarrow \infty}{=} b'(r_1). \end{aligned}$$

Similarly, it can be shown that $\lim_{r \rightarrow \infty} \frac{\ell}{\exp(\ell r)} \int_{-\infty}^r b'(w) < \lim_{r \rightarrow \infty} b'(r)$. Hence, for any finite r_1 ,

$$b'(r_1) < \lim_{r \rightarrow \infty} \frac{\ell}{\exp(\ell r)} \int_{-\infty}^r b'(w) \exp(\ell w) < \lim_{r \rightarrow \infty} b'(r),$$

which by definition of supremum implies that $\bar{E}[b'] = \sup b'(r) = \sup \mathcal{Y}_0$.¹ For $r \rightarrow \infty$, the mass goes to infinity and, thus, concentrates on the value of w which gives the maximum of b' . In a similar way, since for $\ell < 0$, $\exp(\ell w)$ is always decreasing in \mathcal{W} , one has:

$$\underline{E}[b'] = \inf_{\ell \in [-c, c]} \liminf_{r \rightarrow \infty} \frac{\ell}{\exp(\ell r)} \int b'(w) \exp(\ell w) I_{[-r, \infty)}(w) = \inf \mathcal{Y}_0.$$

¹This holds for any $\ell \in (0, c]$. All $\ell \in (0, c]$ are equivalent w.r.t. this property, since all $\exp(\ell w)$ are increasing in \mathcal{W} for $\ell > 0$.

Thus, (A2) is also satisfied. Consider property (A3). The set of posteriors resulting from (11) is:

$$\mathcal{M}_p \propto \left\{ \exp \left(n \left(\frac{\ell + n\hat{y}_n}{n} w - b(w) \right) \right), \ell \in [-c, c] \right\}. \quad (12)$$

Since \mathcal{Y}_0 is an open convex subset of \mathbb{R} [20] (and thus dense in \mathbb{R}) and since $\hat{y}_n \in Cl(\mathcal{Y}_0)$, there exists a $\delta > 0$ such that for $n > \delta$ either $\frac{\ell + n\hat{y}_n}{n} \in \mathcal{Y}_0$ or $\frac{-\ell + n\hat{y}_n}{n} \in \mathcal{Y}_0$ for $\ell \in (0, c]$. Thus, for any \hat{y}_n and for a sufficiently large n , there are always probabilities in \mathcal{M}_p that are countably additive (they have proper PDFs) for either $\ell \in (0, c]$ or $\ell \in [-c, 0)$. By applying Proposition 3.3 with $n_0 y_0 = \ell$ and $n_0 \rightarrow 0$ to the proper PDFs in \mathcal{M}_p , this implies that $\underline{E}[b'|n, \hat{y}_n] = \frac{-c + n\hat{y}_n}{n}$ and/or $\bar{E}[b'|n, \hat{y}_n] = \frac{c + n\hat{y}_n}{n}$ for $n \geq \delta$ (in case $\hat{y}_n \in \mathcal{Y}_0$ but $\hat{y}_n \notin Cl(\mathcal{Y}_0) \setminus \mathcal{Y}_0$ all probabilities in \mathcal{M}_p are countably additive for all $\ell \in [-c, c]$ for a suitably large n). This means that (A3) holds w.r.t. the lower and/or the upper for the RVBF $g = b'$ for which we have stated our condition of prior ignorance (A2). Furthermore, since for $n \rightarrow \infty$, $\frac{\ell + n\hat{y}_n}{n} \rightarrow \hat{y}_n$ and \mathcal{M}_p reduces to the single posterior obtained from the improper prior $p(w) = 1$, it follows that the lower and upper posterior expectations of any RVBF $g \in \mathcal{G}$ converge to the posterior expectation obtained from the improper uniform prior $p(w) = 1$. This proves (A4) but also (A3); in fact since, for $n \rightarrow \infty$, $\underline{E}[g|n, \hat{y}_n] = \bar{E}[g|n, \hat{y}_n]$, this implies that either $\underline{E}[g|n, \hat{y}_n] \neq \underline{E}[g]$ or $\bar{E}[g|n, \hat{y}_n] \neq \bar{E}[g]$ for all g such that $\bar{E}[g] - \underline{E}[g] > 0$. \square

Summarizing, Theorem 4.2 states that:

- Translation invariance and prior ignorance are satisfied by the set of finitely additive probabilities obtained from the kernels in (11): that is, as lower limit for $r \rightarrow \infty$ of the truncated densities $\exp(\ell w) I_{\mathcal{W}_r}(w)$, see Appendix A.
- A-posteriori, for any $n > 0$, there are always probabilities in \mathcal{M}_p that are countably additive and, at the increase of n these probabilities give either the lower or upper expectation of g or both (since either $\frac{c + n\hat{y}_n}{n} \in \mathcal{Y}_0$ or $\frac{-c + n\hat{y}_n}{n} \in \mathcal{Y}_0$, while if they are both in \mathcal{Y}_0 then all the probabilities in \mathcal{M}_p are countably additive). This explains why also (A3)–(A4) are satisfied.
- The result in Theorem 4.2 is different from the one we have obtained in [13], because of the additional requisite of translation invariance that is not satisfied by the set of conjugate priors in [13]. The advantage of imposing translation invariance is twofold. First, it implies that the lower and upper posterior inferences are invariant to translations of the sample mean—this aspect will be clarified in Section 5. Second, it allows us to reduce to one (the constant c) the number of parameters that specify the set of priors, while in [13] there are two parameters. We compare the two models with more details in Section 7.1.

A consequence of Theorem 4.2 is given in the following corollary.

Corollary 4.3. *The lower and upper posterior expectation of b' obtained from the set of posteriors (12) is:*

$$\underline{E}[b'|n, \hat{y}_n] = \max \left(\inf \mathcal{Y}_0, \frac{n\hat{y}_n - c}{n} \right), \quad \bar{E}[b'|n, \hat{y}_n] = \min \left(\sup \mathcal{Y}_0, \frac{n\hat{y}_n + c}{n} \right), \quad (13)$$

and there exists $\delta > 0$ such that, for any $n > \delta$ either the lower or the upper is different from $\inf \mathcal{Y}_0$ or, respectively, $\sup \mathcal{Y}_0$.¹

Proof. In case $\frac{n\hat{y}_n + \ell}{n} \in \mathcal{Y}_0$ for any $\ell \in [-c, c]$ and any n , any posterior in (12) is proper and, thus, the result follows by Proposition 3.3 with $n_0 y_0 = \ell$ and $n_0 \rightarrow 0$. Conversely, when this is not the

¹Notice that this behavior in general is not monotone and depends on how \hat{y}_n converges with the number of observations.

case, observe that:

$$\frac{d \exp(n(y_p w - b(w)))}{dw} = n(y_p - b'(w)) \exp(n(y_p w - b(w))),$$

where $y_p = \frac{n\hat{y}_n + \ell}{n}$. Since $b''(w) > 0$ for each $w \in \mathcal{W}$, the function b' is always increasing in \mathcal{W} and obtains its infimum $\inf \mathcal{Y}_0$ and supremum $\sup \mathcal{Y}_0$ for $w \rightarrow \inf \mathcal{W}$ and, respectively, $w \rightarrow \sup \mathcal{W}$. Since $y_p \in Cl(\mathcal{Y}_0)$ and $b' \in \mathcal{Y}_0$, then if $y_p \leq \inf \mathcal{Y}_0$, it follows that $y_p - b'(w) \leq 0$ for any $w \in \mathcal{W}$ and, thus, the kernel $\exp(n(y_p w - b(w)))$ is always decreasing in w . Similarly it holds that if $y_p \geq \sup \mathcal{Y}_0$ then $\exp(n(y_p w - b(w)))$ is always increasing in w . Therefore, at the limit $r \rightarrow \infty$, the posterior mean computed w.r.t. the prior $\exp(\ell w)I_{(-\infty, r]}(w)$ if $\ell > 0$ or, $\exp(\ell w)I_{[-r, \infty)}(w)$ if $\ell < 0$, will be equal to $\sup \mathcal{Y}_0$ in the case $y_p > \sup \mathcal{Y}_0$ or, respectively, $\inf \mathcal{Y}_0$ in the case $y_p < \inf \mathcal{Y}_0$. For the second part of the corollary, as in the proof of Theorem 4.2, observe that, in case $\hat{y}_n \in Cl(\mathcal{Y}_0) \setminus \mathcal{Y}_0$, for sufficiently large n either $\frac{\ell + n\hat{y}_n}{n} \in \mathcal{Y}_0$ or $\frac{-\ell + n\hat{y}_n}{n} \in \mathcal{Y}_0$ for $\ell \in (0, c]$. This means that, for a sufficiently large n , if $\underline{E}[b'|n, \hat{y}_n] = \inf \mathcal{Y}_0$ then $\bar{E}[b'|n, \hat{y}_n] = \frac{n\hat{y}_n + c}{n}$ while if $\bar{E}[b'|n, \hat{y}_n] = \sup \mathcal{Y}_0$ then $\underline{E}[b'|n, \hat{y}_n] = \frac{n\hat{y}_n - c}{n}$. \square

We illustrate the results in Theorem 4.2 and Corollary 4.3 with the following example.

Example 4.4. In the Poisson case, the likelihood is $\prod_{i=1}^n p(y_i|\theta) = \prod_{i=1}^n P(y_i|\lambda) = \lambda^{\sum_{i=1}^n y_i} \exp(-n\lambda)$. Assume we do not have prior information about the unknown parameter λ . To model this condition of prior ignorance, we can employ the set of priors (11), that transformed to the original parameter space reduces to $\mathcal{M} \propto \{\lambda^{\ell-1}, \ell \in [-c, c]\}$, where λ^{-1} is the determinant of the Jacobian. If $c > 1$, it includes Jeffreys' prior $\lambda^{-\frac{1}{2}}$ and the improper prior λ^{-1} and it satisfies the prior-ignorance property $\underline{E}[\lambda] = 0$ and $\bar{E}[\lambda] = \infty$. The resulting set of posteriors obtained by combining priors and likelihood is:

$$\mathcal{M}_p \propto \left\{ \lambda^{\sum_{i=1}^n y_i + \ell - 1} \exp(-n\lambda), \ell \in [-c, c] \right\}. \quad (14)$$

It can be noticed that if $\sum_{i=1}^n y_i + \ell > 0$ for any ℓ , \mathcal{M}_p includes proper Gamma densities (it can thus be normalized), otherwise \mathcal{M}_p includes densities that again concentrate on $\lambda = 0$ for $\ell < 0$. Thus, from Corollary 4.3 it follows that a-posteriori:

$$\underline{E}[\lambda|n, \hat{y}_n] = \max\left(0, \frac{-c + \sum_{i=1}^n y_i}{n}\right), \quad \bar{E}[\lambda|n, \hat{y}_n] = \frac{c + \sum_{i=1}^n y_i}{n}.$$

\square

We are now ready to extend the results of Theorem 4.2 to the multivariate case. Before doing that, we state the following lemma that will be exploited in the theorem.

Lemma 4.5. *Let \mathbb{L} be a bounded closed convex subset of \mathbb{R}^k strictly including the origin (i.e., the origin is in the interior of \mathbb{L}), and let $\ell = [\ell_{\cdot 1}, \dots, \ell_{\cdot k}]^T$ belong to \mathbb{L} ; then it holds that*

$$\min_{\ell \in \mathbb{L}} \ell_{\cdot i} < 0, \quad \max_{\ell \in \mathbb{L}} \ell_{\cdot i} > 0. \quad (15)$$

Proof. This results follows straightforwardly by considering any ball $\mathcal{B}(\mathbb{L}) \subset \mathbb{L}$ centred at the origin and observing that (15) holds for $\ell \in \mathcal{B}(\mathbb{L})$. \square

Theorem 4.6. *Assume that $\mathcal{W} = \mathbb{R}^k$ and consider properties (A1)–(A4), with $\mathcal{G}_0 = \left\{ \frac{\partial b}{\partial w_{\cdot 1}}, \dots, \frac{\partial b}{\partial w_{\cdot k}} \right\}$, $\mathcal{G} = \mathcal{G}_1$ including sufficiently smooth RVBFs and $\mathcal{F}_i = \mathcal{F}$, where $\mathcal{F} = \{f_a : a \in$*

\mathbb{R}^k and $f_a(w) = w + a$. The following set of priors:

$$\mathcal{M} \propto \left\{ \exp(\ell^T w), \ell = [\ell_{.1}, \dots, \ell_{.k}]^T \in \mathbb{L} \right\}, \quad (16)$$

where \mathbb{L} is a bounded closed convex subset of \mathbb{R}^k strictly including the origin, satisfies (A1)–(A4) and conjugacy.

Proof. Consider first property (A1). Define $\underline{E}[g] = \inf\{\underline{E}_\ell[g] : \ell \in \mathbb{L}\}$, where

$$\underline{E}_\ell[g] = \liminf_{r_{.1} \rightarrow \infty} \dots \liminf_{r_{.k} \rightarrow \infty} \int \dots \int g(w) \exp(\ell^T w) \prod_{i=1}^k I_{\mathcal{W}_{r_{.i}}}(w_{.i}) \frac{\ell_{.i}}{\exp(\ell_{.i} r_{.i})} dw_{.i},$$

and $\mathcal{W}_{r_{.i}}$ defined in (10). Observe that Proposition 4.1 also holds for $\exp(\ell^T w)$ in the multivariate case (it is enough to proceed component-wise). In fact, since $g(f_a) = g(w + a)$ then after the change of variables $w' = w - a$,

$$\underline{E}_\ell[g(f_a)] = \liminf_{r_{.1} \rightarrow \infty} \dots \liminf_{r_{.k} \rightarrow \infty} \int \dots \int g(w') \exp(\ell^T w') \prod_{i=1}^k I_{\{\mathcal{W}_{r_{.i}+a_{.i}}\}}(w'_{.i}) \frac{\ell_{.i}}{\exp(\ell_{.i}(r_{.i}+a_{.i}))} dw'_{.i}.$$

Thus, by redefining $r_{.i} = r_{.i} + a_{.i}$ and $\mathcal{W}_{r_{.i}+a_{.i}} = \mathcal{W}_{r_{.i}}$, it follows that $\underline{E}_\ell[g(f_a)] = \underline{E}_\ell[g]$ for any ℓ and, thus, translation invariance holds for any kernel $\exp(\ell^T w)$. Then it holds also for the kernel ℓ which attains the infimum in $\underline{E}[g]$. Therefore, $\underline{E}[g(f_a)] = \underline{E}[g]$ for any $g \in \mathcal{G}_1$ and, thus, the set of priors \mathcal{M} is invariant under translations.

Property (A2) follows from Theorem 4.2 by proceeding component-by-component for each $\frac{\partial b}{\partial w_{.i}}$ for $i = 1, \dots, k$. In fact,

$$\begin{aligned} \underline{E}_\ell \left[\frac{\partial b}{\partial w_{.i}} \right] &= \liminf_{r_{.1} \rightarrow \infty} \dots \liminf_{r_{.k} \rightarrow \infty} \int \dots \int \frac{\partial b}{\partial w_{.i}} \exp(\ell^T w) \prod_{i=1}^k I_{\mathcal{W}_{r_{.i}}}(w_{.i}) \frac{\ell_{.i}}{\exp(\ell_{.i} r_{.i})} dw_{.i} \\ &= \liminf_{r_{.i} \rightarrow \infty} \int \frac{\partial b}{\partial w_{.i}} \exp(\ell_{.i} w_{.i}) I_{\mathcal{W}_{r_{.i}}}(w_{.i}) \frac{\ell_{.i}}{\exp(\ell_{.i} r_{.i})} dw_{.i}. \end{aligned}$$

Since \mathbb{L} includes the origin, then from Proposition 4.5 it follows that $\ell_{.i}$ can assume both positive and negative values. Thus, we are back to Theorem 4.2 (univariate case) in which $\ell_{.i}$ is free to vary in a closed interval including the origin. Thus, the infimum and supremum of $\frac{\partial b}{\partial w_{.i}}$ are obtained for $w_{.i} \rightarrow -\infty$ and, respectively, $w_{.i} \rightarrow +\infty$. By combining likelihood and the priors in (16), one obtains the posterior kernels:

$$\mathcal{M}_p = \left\{ \exp \left(n \left(\left(\frac{\ell + n \hat{y}_n}{n} \right)^T w - b(w) \right) \right), \ell \in \mathbb{L} \right\}. \quad (17)$$

Observe that $\frac{\ell + n \hat{y}_n}{n} = \hat{y}_n + \frac{\ell}{n}$. Since \mathbb{L} is bounded, the effect of $\frac{\ell}{n}$ vanishes for $n \rightarrow \infty$, which proves (A4) and thus (A3). \square

Theorem 4.6 requires \mathbb{L} to be bounded and to include the origin in order to satisfy (A1)–(A4). For symmetry reasons, we restrict \mathbb{L} to be symmetric w.r.t. the origin so to reduce to the set $[-c, c]$ in the $k = 1$ case. It is clear that while in the one-dimensional case, convexity, boundedness and symmetry are satisfied only by the interval $[-c, c]$ in the ($k > 1$)-dimensional case there are many subsets of \mathbb{R}^k that satisfy these three characteristics. In other words, Theorem 4.6 does not specify uniquely \mathbb{L} and, thus, the modeller can add further characteristics to \mathbb{L} in order to satisfy additional properties, for instance to determine the shape of the set of posterior means.

Corollary 4.7. *The expectation of ∇b calculated from the set of posteriors (17) satisfies:*

$$E[\nabla b|n, \hat{y}_n] \in \left\{ \frac{\ell + n\hat{y}_n}{n}, \ell \in \mathbb{L} \right\}, \quad (18)$$

provided that all the posteriors in (17) are proper.

Proof. The result follows directly from Lemma 3.3 assuming that the set \mathcal{M}_p contains all proper distributions. This always holds for a large enough n because (17) satisfies (A4). \square

Corollary 4.7 is weaker than Corollary 4.3, because it provides the shape of the posterior set of means only in the case all the posteriors in (17) are proper. In the multivariate case it is difficult to characterize in general the shape of the set of posterior means when (17) can include some improper distribution. This difficulty is due to the interaction among the components of w , which is determinate by the function $b(w)$, and changes depending on the exponential family we are considering. Although we cannot provide a general characterization, we can extend Corollary 4.7 case-by-case, as it is shown in the next examples about the multivariate Normal and categorical distributions.

4.1 Multivariate Normal with Known Variance

In the multivariate Normal case with known variance, $\prod_{i=1}^n p(y_i|x) = \prod_{i=1}^n N(y_i; x, V)$ with $x, y_i \in \mathbb{R}^d$, $0 < V \in \mathbb{R}^{d \times d}$ and $k = d$. Let us assume that $\mathbb{L} = \{\ell \in \mathbb{R}^d : \ell_i \in [-c_i, c_i], c_i > 0, i = 1, \dots, k\}$, this models a sort of ‘‘uncorrelated’’ prior uncertainty (i.e., ℓ_i is free to vary in $[-c_i, c_i]$ independently of the values that ℓ_j with $j \neq i$ assumes). Notice in fact that \mathbb{L} corresponds to a box in \mathbb{R}^d . In this case, the set of posteriors

$$\mathcal{M}_p = \left\{ N\left(x; \frac{V\ell + n\hat{y}_n}{n}, \frac{1}{n}V\right), \ell \in \mathbb{L} \right\}$$

includes Normal PDFs with variance V/n and means free to vary in a hyper-rectangle whose vertices are the extremes of $\frac{\ell + n\hat{y}_n}{n}$ for $\ell_i \in [-c_i, c_i]$, which implies:

$$E[X_i|n, \hat{y}_n] \in \left[\frac{-c_i + n\hat{y}_{ni}}{n}, \frac{c_i + n\hat{y}_{ni}}{n} \right]. \quad (19)$$

Observe that, provided that all the components of x are observed, \mathcal{M}_p includes only proper priors. Thus (19) follows by Corollary 4.7. Conversely, if some component of x was not observed, because of marginalization invariance of the Normal distribution, (19) would hold for all observed components, while for the unobserved ones $E[X_i|n, \hat{y}_n] \in (-\infty, \infty)$. We are prior ignorant until we observe x_i .

Notice that the set of constraints $\ell_i \in [-c_i, c_i]$ for $i = 1, \dots, k$ are equivalent to $\frac{1}{c_i}|\ell_i| \leq 1$ for $i = 1, \dots, k$ and can be replaced by the unique constraint $\max_i \frac{1}{c_i}|\ell_i| \leq 1$ or, equivalently, $\|C^{-1}\ell\|_\infty \leq 1$, where $\ell = [\ell_1, \dots, \ell_k]^T$ and $C = \text{diag}(c_1, \dots, c_k)$. From Theorem 4.6, it follows that prior ignorance and translation invariance are also satisfied in the case $\|C^{-1}\ell\|_\infty \leq 1$ is replaced by $\|C^{-1}\ell\|_1 \leq 1$ or $\|C^{-1}\ell\|_2 \leq 1$. It can also be noticed that, in the case the components of x have the same ‘‘physical’’ meaning, there is no reason to use different values of the parameters c_i , which can be thus set all equal, $c_i = c$. Furthermore, when $c_i = c$ and $\|C^{-1}\ell\|_p \leq 1$ for some p -norm, the resulting set of priors is also invariant to permutations of the components of x . This means that all the components have the same posterior imprecision, i.e., $\overline{E}[x_i|n, \hat{y}_n] - \underline{E}[x_i|n, \hat{y}_n] = 2c/n$.

4.2 Categorical Distribution

In the multinomial case, $\prod_{i=1}^n p(y_i|x) = \prod_{i=1}^n Ca(y_i;x)$ with $x \in (0,1)^d$, $\sum_i x_i = 1$ and $y_{ij} \in \{0,1\}$, the set of posteriors obtained based on the parametrization [17, Ch.1] assuming $\mathbb{L} = \{\ell \in \mathbb{R}^d : \ell_i \in [-c_i, c_i], c_i > 0, i = 1, \dots, k\}$ is:

$$\mathcal{M}_p = \left\{ \theta_{.1}^{n\hat{y}_{n1} + \ell_{.1} - 1} \theta_{.2}^{n\hat{y}_{n2} + \ell_{.2} - 1} \dots \theta_{.d}^{n(1 - \sum_{i=1}^{d-1} \hat{y}_{ni}) - \sum_{i=1}^{d-1} \ell_{.i} - 1}, \ell_i \in [-c_i, c_i] \right\}. \quad (20)$$

Observe that in this case Theorem 4.6 can be applied since $\mathcal{W} = \mathbb{R}^{d-1}$ and thus $\mathcal{W} = \mathbb{R}^k$ with $k = d - 1$. It can be noticed that in (20) the exponents of the components θ_i with $i \neq d$ can vary between $n\hat{y}_{ni} - c_i - 1$ and $n\hat{y}_{ni} + c_i - 1$, while the exponent of the component θ_d can vary between $n(1 - \sum_{i=1}^{d-1} \hat{y}_{ni}) - \sum_{i=1}^{d-1} c_i - 1$ and $n(1 - \sum_{i=1}^{d-1} \hat{y}_{ni}) + \sum_{i=1}^{d-1} c_i - 1$. This means that the imprecision on θ_d is larger than the imprecision on the other components even if $c_i = c$ for $i = 1, \dots, d - 1$. This asymmetry is due to the choice of the parameter transformation [17, Ch. 1], a different parameter transformation would result in asymmetry for a different component of θ . This is not desirable for a model of prior ignorance in which, because of lack of prior information, we expect the ignorance to be *symmetric*.

A way to overcome this problem is that of imposing to the set of priors not only to satisfy translation invariance but also permutation invariance: that is, that the set of priors \mathcal{M} is invariant to permutations of the components of θ .

Corollary 4.8. *Consider Theorem 4.6 in the multinomial case, the set of posteriors satisfy (A1)–(A4) and permutation invariance of the components of θ provided that \mathbb{L} is chosen as follows:*

$$\mathbb{L} = \left\{ \ell \in \mathbb{R}^d : \ell_i \in [-c, c], \sum_{i=1}^{d-1} \ell_i \in [-c, c], i = 1, \dots, k \right\}.$$

This follows by Theorem 4.6 imposing the additional constraint of permutation invariance. The above choice of \mathbb{L} implies that the set of priors is:

$$\mathcal{M}_{SYM} = \left\{ \theta_{.1}^{\ell_{.1} - 1} \theta_{.2}^{\ell_{.2} - 1} \dots \theta_{.d}^{-\sum_{i=1}^{d-1} \ell_{.i} - 1}, \ell_i \in [-c, c], \sum_{i=1}^{d-1} \ell_i \in [-c, c] \right\}, \quad (21)$$

which satisfies permutation invariance of the components of θ . For instance in the case $d = 3$, the PDFs with extreme exponents that belong to (21) are of the following type: $\theta_{.1}^{c-1} \theta_{.2}^{-c-1} \theta_{.3}^{-1}$, $\theta_{.1}^{c-1} \theta_{.2}^{-1} \theta_{.3}^{c-1}$, $\theta_{.1}^{-c-1} \theta_{.2}^{c-1} \theta_{.3}^{-1}$, $\theta_{.1}^{-c-1} \theta_{.2}^{-1} \theta_{.3}^{c-1}$, $\theta_{.1}^{-1} \theta_{.2}^{-c-1} \theta_{.3}^{c-1}$ and $\theta_{.1}^{-1} \theta_{.2}^{c-1} \theta_{.3}^{-c-1}$. For the case $d = 4$, one has $\theta_{.1}^{c-1} \theta_{.2}^{-c-1} \theta_{.3}^{c-1} \theta_{.4}^{-c-1}$, $\theta_{.1}^{c-1} \theta_{.2}^{-1} \theta_{.3}^{-c-1} \theta_{.4}^{c-1}$, etc. This means that $\sum_{i=1}^d s_i c = 0$, where s_i is the sign of c as it appears in the exponent of θ_i .

Corollary 4.9. *The lower and upper posterior means resulting from (21) are:*

$$\underline{E}[\theta_i | n, \hat{y}_n] = \max \left(0, \frac{-c + n\hat{y}_{ni}}{n} \right), \quad \bar{E}[\theta_i | n, \hat{y}_n] = \min \left(1, \frac{c + n\hat{y}_{ni}}{n} \right) \quad \forall i.$$

In general, the lower and upper mean of $\sum_{i \in J} \theta_i$, where J is a subset of $\{1, \dots, d\}$, is equal to:

$$\begin{aligned} \underline{E}[\sum_{i \in J} \theta_i | n, \hat{y}_n] &= \max \left(0, \frac{1}{n} \left[\sum_{i \in J} n\hat{y}_{ni} + \max(-c|J|, -c(d - |J|)) \right] \right), \\ \bar{E}[\sum_{i \in J} \theta_i | n, \hat{y}_n] &= \min \left(1, \frac{1}{n} \left[\sum_{i \in J} n\hat{y}_{ni} + \min(c|J|, c(d - |J|)) \right] \right), \end{aligned} \quad (22)$$

with $|J|$ denoting the cardinality of J .

Proof. Consider the general form (22), the posterior kernel is:

$$p(\theta|n, \hat{y}_n) \propto \theta_{.1}^{n\hat{y}_{n1}+\ell_{.1}-1} \theta_{.2}^{n\hat{y}_{n2}+\ell_{.2}-1} \dots \theta_{.d}^{n(1-\sum \hat{y}_{ni})-\sum_{i=1}^{d-1} \ell_{.i}-1}.$$

We prove it for $J = \{1, 2\}$ and $d = 4$ but the steps can be repeated in general. We need to compute

$$\frac{1}{K} \int_0^1 \int_0^{1-\theta_1} (\theta_{.1} + \theta_{.2}) \theta_{.1}^{n\hat{y}_{n1}+\ell_{.1}-1} \theta_{.2}^{n\hat{y}_{n2}+\ell_{.2}-1} \int_0^{1-\theta_1-\theta_2} \theta_{.3}^{n\hat{y}_{n3}+\ell_{.3}-1} (1-\theta_{.1}-\theta_{.2}-\theta_{.3})^{n(1-\sum_{i=1}^3 \hat{y}_{ni})-\sum_{i=1}^3 \ell_{.i}-1} d\theta_{.1} d\theta_{.2} \theta_{.3},$$

where K is the normalization constant. Now consider the transformation of coordinate $\theta'_{.3} = \theta_{.3}/(1-\theta_{.1}-\theta_{.2})$ (in general $\theta'_{.k} = \theta_{.k}/(1-\sum_{i \in J} \theta_{.i})$ for $k \notin J$):

$$\begin{aligned} & \frac{1}{K} \int_0^1 \int_0^{1-\theta_1} (\theta_{.1} + \theta_{.2}) \theta_{.1}^{n\hat{y}_{n1}+\ell_{.1}-1} \theta_{.2}^{n\hat{y}_{n2}+\ell_{.2}-1} \\ & (1-\theta_{.1}-\theta_{.2})^{n(1-\hat{y}_{n1}-\hat{y}_{n2})-\sum_{i=1}^2 \ell_{.i}-1} \int_0^1 \theta'_{.3}^{n\hat{y}_{n3}+\ell_{.3}-1} (1-\theta_{.3})^{n(1-\sum \hat{y}_{ni})-\sum_{i=1}^3 \ell_{.i}-1} d\theta_{.1} d\theta_{.2} \theta_{.3} \\ & = \frac{1}{K} \int_0^1 \int_0^{1-\theta_1} (\theta_{.1} + \theta_{.2}) \theta_{.1}^{n\hat{y}_{n1}+\ell_{.1}-1} \theta_{.2}^{n\hat{y}_{n2}+\ell_{.2}-1} (1-\theta_{.1}-\theta_{.2})^{n(1-\hat{y}_{n1}-\hat{y}_{n2})-\sum_{i=1}^2 \ell_{.i}-1} d\theta_{.1} d\theta_{.2}. \end{aligned}$$

Now consider the transformation of coordinates $\theta_{.1} = uv$, $\theta_{.2} = u(1-v)$ with $0 < u, v < 1$ (in general $\theta_{.i} = uv_i \forall i \in J$ and $1 - \sum_{i \in J} \theta_{.i} = u(1 - \sum_{i \in J} v_i)$):

$$\begin{aligned} & \frac{1}{K} \int_0^1 \int_0^{1-\theta_1} (\theta_{.1} + \theta_{.2}) \theta_{.1}^{n\hat{y}_{n1}+\ell_{.1}-1} \theta_{.2}^{n\hat{y}_{n2}+\ell_{.2}-1} (1-\theta_{.1}-\theta_{.2})^{n(1-\hat{y}_{n1}-\hat{y}_{n2})-\sum_{i=1}^2 \ell_{.i}-1} d\theta_{.1} d\theta_{.2} \\ & = \frac{1}{K} \int_0^1 \int_0^1 u^{n\hat{y}_{n1}+\ell_{.1}+n\hat{y}_{n2}+\ell_{.2}-1} (1-u)^{n(1-\hat{y}_{n1}-\hat{y}_{n2})-\sum_{i=1}^2 \ell_{.i}-1} v^{n\hat{y}_{n1}+\ell_{.1}-1} (1-v)^{n\hat{y}_{n2}+\ell_{.2}-1} dudv \\ & = \frac{1}{K} \int_0^1 u^{n\hat{y}_{n1}+\ell_{.1}+n\hat{y}_{n2}+\ell_{.2}-1} (1-u)^{n(1-\hat{y}_{n1}-\hat{y}_{n2})-\sum_{i=1}^2 \ell_{.i}-1} du. \end{aligned}$$

Note that we are reduced to a univariate expectation w.r.t. a Beta distribution. We can therefore apply Corollary 4.3, with the constraints $\ell_{.i} \in [-c, c]$, $\sum_{i=1}^{d-1} \ell_{.i} \in [-c, c]$, to obtain (22). \square

It can be noticed that, fixed J , the posterior imprecision of $\sum_{i \in J} \theta_{.i}$ increases with d , i.e., the components of θ . For some problem this can be undesirable, especially because the choice of the number of components of θ is often arbitrary. We would like to make the posterior imprecision invariant to the number of components of θ .

Corollary 4.10. *Consider Theorem 4.6. The set of posteriors satisfy (A1)–(A4), permutation invariance of the components of θ and invariance of the posterior imprecision w.r.t. the number of components of θ provided that \mathbb{L} is chosen as follows:*

$$\mathbb{L} = \left\{ \ell \in \mathbb{R}^d : \|\ell\|_1 \leq 2c, \sum_{i=1}^{d-1} \ell_{.i} \in [-c, c] \right\},$$

where $\ell = [\ell_{.1}, \dots, \ell_{.d-1}]^T$.

This follows by Theorem 4.6 imposing the additional constraint of permutation invariance and invariance w.r.t. the number of components of θ . This implies that the set of priors in this case

is:

$$\mathcal{M}_{RIP} = \left\{ \theta_1^{\ell_1-1} \theta_2^{\ell_2-1} \dots \theta_d^{-\sum_{i=1}^{d-1} \ell_i-1}, \|\ell\|_1 \leq 2c, \sum_{i=1}^{d-1} \ell_i \in [-c, c] \right\}. \quad (23)$$

For the case $d = 2$ and $d = 3$, (23) coincides with (21), but in the case $d = 4$, the PDFs with extreme exponents are $\theta_1^{c-1} \theta_2^{-c-1} \theta_3^{-1} \theta_4^{-1}$, $\theta_1^{c-1} \theta_2^{-1} \theta_3^{-c-1} \theta_4^{-1}$, $\theta_1^{c-1} \theta_2^{-1} \theta_3^{-1} \theta_4^{-c-1}$, $\theta_1^{-1} \theta_2^{c-1} \theta_3^{-c-1} \theta_4^{-1}$ etc. In this case, there are at most two components whose exponents assume value $\pm c - 1$. With this constraint, the set of posteriors satisfies the one-step ahead representation invariance principle (RIP) [6], i.e., the posterior upper and lower probabilities assigned to the event “the next observation belongs to the categories in J ” should not depend on the sample space in which the event and the previous observations are represented. In fact, one has that:

$$\underline{E}[\sum_{i \in J} \theta_i | n, \hat{y}_n] = \max \left(0, \frac{1}{n} \left[\sum_{i \in J} n \hat{y}_{ni} - c \right] \right), \quad \bar{E}[\sum_{i \in J} \theta_i | n, \hat{y}_n] = \min \left(1, \frac{1}{n} \left[\sum_{i \in J} n \hat{y}_{ni} + c \right] \right), \quad (24)$$

and, thus, it can be noticed that the difference between the upper and lower expectation does not depend on the number of categories in J .

4.3 How to Choose Parameters c_i

The aim of this section is to give guidelines for the choice of the parameters c_i . The parameters c_i determine the precision of posterior inferences and, thus, their robustness. To understand the meaning of c_i , we can compare the expression of the posterior expectation of $\partial b / \partial w_i$ obtained by considering a standard Bayesian analysis using the conjugate prior $\exp(n_0(y_0^T w - b(w)))$ and the lower and upper posterior expectations obtained using our model, i.e.,

$$\frac{n \hat{y}_{ni} + n_0 y_{0i}}{n + n_0} \longleftrightarrow \frac{n \hat{y}_{ni} \pm c_i}{n}. \quad (25)$$

Looking at the expression on the left, it is evident that the prior parameter n_0 attached to the prior mean y_{0i} plays the same role to the sample size n attached to the data mean \hat{y}_{ni} . Since $n \hat{y}_{ni} = \sum_j y_{ji}$, we can also interpret $n_0 y_{0i}$ as the sum of n_0 pseudo-observations that are added to $\sum_j y_{ji}$. Matching the two expectations in (25), it implies that $|n_0 y_{0i}| = c_i$ and $n_0 = 0$. In the case on the right we are not adding additional (pseudo) observations, but we are considering a scenario in which an amount c_i is summed/subtracted to the sum $\sum_j y_{ji}$ but without changing the sample size. For instance, when y_{ij} are frequency data from an election poll (this example will be discussed in Section 8.2), we can interpret the lower and upper expectations as the result of a swing scenario in which we test what happens if c_i votes are subtracted (lower) to a candidate and assigned to another one (upper). The interpretation of c_i as a pseudo-observation is the main avenue to be followed for the choice of a value for c_i .

The choice of c_i can also be based on measures of posterior robustness such as [6]: (a) the convergence rate of the lower and/or upper expectations to suitable limits; (b) the convergence rate of the posterior imprecision, i.e., the difference between upper and lower expectations. Here the expectations are computed w.r.t. some function of interest g and the convergence is defined w.r.t. the number of samples n . We will give an example of this approach in Section 8.1. Another possible requirement for the choice of c_i is that the family of priors \mathcal{M} should be large enough to encompass frequentist or objective Bayesian inferences, but not too large to avoid obtaining too weak inferences. These are the approaches that we have followed in [13] for the univariate one-parameter exponential family. In the multivariate case the only difference is that the value of k parameters have to be fixed. However, as discussed previously, for symmetry reasons in many cases it can be assumed $c_i = c$ for (all) a subset of parameters $i = 1, \dots, k$. Finally, we

can choose c to have some desirable frequentist properties (e.g., robust credible intervals to be calibrated frequentist intervals, hypothesis tests to be calibrated for the Type I error, etc.). We may also place an upper bound on the value of c based on considerations regarding the power of the hypothesis test.

5. Prior Near-Ignorance: Some Further Useful Properties

Notice that in all the above examples the set of priors (16) guarantees prior ignorance w.r.t. ∇b . This choice has also been motivated by the meaning of ∇b for exponential families. Remember in fact from Section 3 that ∇b (equal to b' in the univariate case) is the mean of Y and, thus, is the quantity about which we will have prior beliefs before seeing the data and posterior beliefs after observing the data. Furthermore, to state prior ignorance, invariance, learning and converge in the natural space \mathcal{W} allows us to treat jointly all the members of the exponential families and, thus, to derive general properties. The natural parametrization is in fact a preferred parametrization for exponential families. In this section, it will be shown that the set of priors (16) satisfies (A2) also for other functions of interest in statistical analysis (i.e., one- and two-sided hypotheses testing and credible intervals). Therefore, the choice of imposing prior ignorance only on $\mathcal{G}_0 = \{\nabla b\}$, is not very restrictive for exponential families.

Corollary 5.1. *The family of priors \mathcal{M} in Theorem 4.6 satisfies the following properties:*

- (B1) *For each ball $B_\zeta(\mathcal{W}) = \{w : \|w - w_0\| \leq \zeta\} \subset \mathcal{W}$, $w_0 \in \mathcal{W}$ and $\zeta > 0$, it holds that $\underline{E}[I_{B_\zeta(\mathcal{W})}] = 0$ and $\bar{E}[I_{B_\zeta(\mathcal{W})}] = 0$ but $\underline{E}[\mathbb{R}^d] = 1$.*
- (B2) *For each subset $\mathcal{W}' \subset \mathcal{W}$ such that $\mathcal{W}' = \times_i \Pi_i(\mathcal{W}')$, where $\Pi_i(\mathcal{W}') = [a_i, \sup \Pi_i(\mathcal{W}'))$ or $\Pi_i(\mathcal{W}') = (\inf \Pi_i(\mathcal{W}'), a_i]$ and $a_i \in \mathbb{R}$, with $\Pi_i(\mathcal{W}')$ denoting the orthogonal projection of \mathcal{W}' w.r.t. the i -th component, then $\underline{E}[I_{\{\mathcal{W}'\}}] = 0$ and $\bar{E}[I_{\{\mathcal{W}'\}}] = 1$.*

Proof. The first part of (B1) follows straightforwardly from the fact the set of priors in Theorem 4.6 concentrate their mass on the border of \mathcal{W} . In fact, the set of priors is composed by the following truncated densities

$$\exp(\ell^T w) \prod_{i=1}^k I_{\mathcal{W}_{r_i}}(w_i) \frac{\ell_i}{\exp(\ell_i r_i)}.$$

Then, for any ball in \mathcal{W} , it follows that

$$\bar{E}[I_{B_\zeta(\mathcal{W})}] = \limsup_{r_1 \rightarrow \infty} \dots \limsup_{r_k \rightarrow \infty} \int \dots \int_{B_\zeta(\mathcal{W})} \exp(\ell^T w) \prod_{i=1}^k I_{\mathcal{W}_{r_i}}(w_i) \frac{\ell_i}{\exp(\ell_i r_i)} dw_i = 0,$$

since the integral of $I_{B_\zeta(\mathcal{W})}$ is bounded and $1/\exp(\ell_i r_i)$ goes to zero for $|r_i| \rightarrow \infty$. Concerning (B2), \mathcal{W}' is the Cartesian product of $\Pi_i(\mathcal{W}')$. Thus, if we take $\ell_i > 0$ for $\Pi_i(\mathcal{W}') = [a_i, \sup \Pi_i(\mathcal{W}'))$ or $\ell_i < 0$ for $\Pi_i(\mathcal{W}') = (\inf \Pi_i(\mathcal{W}'), a_i]$ for all $i = 1, \dots, k$, then $\bar{E}[I_{\{\mathcal{W}'\}}] = 1$. While if we take $\ell_i > 0$ for $\Pi_i(\mathcal{W}') = (\inf \Pi_i(\mathcal{W}'), a_i]$ or and $\ell_i < 0$ for $\Pi_i(\mathcal{W}') = [a_i, \sup \Pi_i(\mathcal{W}'))$ for some i , then $\underline{E}[I_{\{\mathcal{W}'\}}] = 0$. The second part of (B1) follows from the fact that any truncated prior integrates on its support to 1 and so it is the limit when its support tends to \mathcal{W} . \square

Corollary 5.1 has several important implications in statistical inference. Since it is common to represent credible regions with balls $B_\zeta(\mathcal{W}) \subset \mathcal{W}$, from (B1) it follows that a priori for any $\zeta < \infty$, $\underline{E}[I_{B_\zeta(\mathcal{W})}] = 0$. That is, the lower probability of the true w_0 to be in $B_\zeta(\mathcal{W})$ is zero, and is only greater than zero for $\zeta \rightarrow \infty$. The only convex set that has prior lower probability greater than zero is \mathcal{W} , which means that a priori we only know that $w_0 \in \mathcal{W}$. Concerning (B2), it

implies that the lower and upper distribution of w are vacuous and also that the marginal lower and upper distribution of any component w_i are vacuous.

6. The Case $\mathcal{W} \neq \mathbb{R}^k$

In Theorems 4.2–4.6 we have defined a model of prior ignorance which satisfies (A1)–(A4) in case $\mathcal{W} = \mathbb{R}^k$. For instance in case $k = 1$, we have seen that translation invariance and prior ignorance can be satisfied because the set of priors \mathcal{M} includes limits of truncated increasing/decreasing exponentials which, at the limit, accumulate to $\pm\infty$. However, for many exponential families of interest it holds that $\mathcal{W} \neq \mathbb{R}^k$, e.g., the centred-normal, multivariate-normal with unknown variance and the exponential distribution, etc. Since \mathcal{W} is semi-bounded in these three cases, translation invariance is not well defined. For example, in case $\mathcal{W} = \mathbb{R}^-$ then $w + a \notin \mathcal{W}$ for $a > 0$ and $|w| < a$. Also prior ignorance does not hold in these cases for the set of priors (16). In fact, for $\mathcal{W} = \mathbb{R}^-$ the prior $\exp(\ell w)$ accumulates to minus infinity for $\ell < 0$, but it is integrable for $\ell > 0$, since $\int_{-\infty}^0 \exp(\ell w) = 1/\ell$ for $\ell > 0$. Thus, in these cases, properties (A1)–(A2) cannot be satisfied by the set of priors (16). Observe in fact that for the centred normal $b' = \lambda^{-1}$ and $\int_{-\infty}^0 \lambda^{-1} \exp(-\ell\lambda) d\lambda = \infty$ for both $\ell > 0$ and $\ell < 0$. Thus, prior ignorance is not satisfied in this case, being $\inf b' = 0$ and $\sup b' = \infty$. Translation invariance is also not satisfied for $\ell > 0$, since $0 \neq \int_A \exp(-\ell\lambda) d\lambda \neq \int_{A+a} \exp(-\ell\lambda) d\lambda \neq 0$ for any bounded set $A \subset \mathbb{R}^+$ and $a > 0$. How can we impose prior near-ignorance in case $\mathcal{W} \neq \mathbb{R}^k$?

In Proposition 3.2, it is stated that the prior kernel $\exp(n_0 y_0 w - n_0 b(w))$ is integrable if $y_0 \in \mathcal{Y}_0$ and $n_0 > 0$. Diaconis and Ylvisaker in [20] give more conditions on the integrability of $\exp(n_0 y_0 w - n_0 b(w))$. In particular, they prove the following result: (i) if $n_0 > 0$ and $y_0 \in \mathcal{Y}_0$ then $\exp(n_0 y_0 w - n_0 b(w))$ is integrable; (ii) if $\exp(n_0 y_0 w - n_0 b(w))$ is integrable and $\mathcal{W} = \mathbb{R}^d$ then $n_0 > 0$ and (iii) if $\exp(n_0 y_0 w - n_0 b(w))$ is integrable and $n_0 > 0$, then $y_0 \in \mathcal{Y}_0$.

Thus, even in the case $y_0 \in \mathcal{Y}_0$, $\exp(n_0 y_0 w - n_0 b(w))$ may not be integrable if $n_0 < 0$ (for $\mathcal{W} = \mathbb{R}^d$ the kernel is certainly not integrable if $n_0 < 0$, see condition (ii) of Diaconis and Ylvisaker, while for $\mathcal{W} \neq \mathbb{R}^d$ there is not a general result, i.e., the integrability of the kernel in case $n_0 < 0$ depends on the value of n_0 and on the particular exponential family). This means that for some value of $n_0 < 0$, it can happen that $\int_{\mathcal{W}} \exp(n_0 y_0 w - n_0 b(w)) dw = \infty$. Therefore, in this case, if we can find a sequence of subsets $\mathcal{W}_1, \mathcal{W}_2, \dots, \mathcal{W}_r$ of \mathcal{W} converging to \mathcal{W} from below and such that the truncated prior satisfies

$$\int_{\mathcal{W}_i} I_{\mathcal{W}_i} \exp(n_0 y_0 w - n_0 b(w)) dw < \infty,$$

for any \mathcal{W}_i , then we can define a sequence of countably additive truncated densities on \mathcal{W}_r whose lower limit for $r \rightarrow \infty$ defines a lower expectation model as in (9). Since $\bar{E}[A] = 0$ for any bounded set A (as it is shown in Appendix A), the set of priors \mathcal{M} associated to this lower expectation model includes only finitely additive probabilities.

Therefore, for $\mathcal{W} \neq \mathbb{R}^d$ we may be able to satisfy (A1)–(A4) with a family of finitely additive probabilities obtained as limit of truncated conjugate exponential priors $\exp(n_0 y_0 w - n_0 b(w))$ obtained by choosing $\ell = n_0 y_0 \in \mathbb{L}$ and $n_0 < 0$. In the following sections, we show that this is possible for the centred-normal and multivariate-normal with unknown variance.¹ Our conjecture is that this result may be extended to other distributions in the exponential families by following the same approach. However, we cannot develop a general theory because the non-integrability of $\exp(n_0 y_0 w - n_0 b(w))$ for $\mathcal{W} \neq \mathbb{R}^d$ and $n_0 < 0$ depends on the value of n_0 and on the particular exponential family. Thus, a general treatment does not seem to be possible.

¹This also holds for the exponential distribution for $n_0 < -1$.

6.1 Centred-Normal Distribution

Assume that the likelihood is a centred-normal distribution and consider the following set of priors

$$\mathcal{M} \propto \left\{ \lambda^{-\frac{\nu}{2}} \exp(-\ell \frac{\lambda}{2}), \ell \in [-c, c] \right\}, \quad (26)$$

defined by the parameters $c > 0$, $\nu \in \mathbb{R}^+$ and $n_0 = -\nu < 0$. It can be verified that for any value of ℓ and $\nu \geq 2$ all the kernels in (26) are not integrable. The kernels are integrable in $\mathcal{W}_r = [1/r, r]$ which, for $r \rightarrow \infty$, defines an infinite sequence of subsets of \mathbb{R}^+ converging from below to \mathbb{R}^+ . By truncating these kernels in $[1/r, r]$ and taking the limit of the expectations for $r \rightarrow \infty$, one can verify that for $\nu > 2$ and $\ell > 0$

$$\overline{E}[\lambda^{-1}] = \limsup_{r \rightarrow \infty} \frac{\int_{1/r}^r \lambda^{-\frac{\nu+2}{2}} \exp(-\ell \frac{\lambda}{2}) d\lambda}{\int_{1/r}^r \lambda^{-\frac{\nu}{2}} \exp(-\ell \frac{\lambda}{2}) d\lambda} = \infty,$$

while for $\ell < 0$

$$\underline{E}[\lambda^{-1}] = \liminf_{r \rightarrow \infty} \frac{\int_{1/r}^r \lambda^{-\frac{\nu+2}{2}} \exp(-\ell \frac{\lambda}{2}) d\lambda}{\int_{1/r}^r \lambda^{-\frac{\nu}{2}} \exp(-\ell \frac{\lambda}{2}) d\lambda} = 0.$$

Notice that we are interested in the expectation of λ^{-1} because $1/\lambda = \sigma^2$. For any bounded subset A of \mathbb{R}^+ , it holds that $\overline{E}[I_{\{A+a\}}] = \overline{E}[I_{\{A\}}] = 0$ for any $a > 0$ (which is a weaker version of translation invariance). Finally, a-posteriori by combining the centred-normal likelihood with the priors in (26), one obtains

$$\mathcal{M} \propto \left\{ \lambda^{\frac{n-\nu+2}{2}-1} \exp\left(-n \frac{\ell + n\hat{y}_n}{n} \frac{\lambda}{2}\right), \ell \in [-c, c] \right\}. \quad (27)$$

It can be noticed that for $n > \nu$ and $\frac{\ell + n\hat{y}_n}{n} > 0$ for all $\ell \in [-c, c]$, all the kernels in (27) are integrable and the expected value of $1/\lambda = \sigma^2$ is equal to $(\ell + n\hat{y}_n)/(n - \nu)$ in this case. Thus, it follows that for $n > \nu$:

$$\underline{E}[\lambda^{-1} | n, \hat{y}_n] = \max\left(0, \frac{-c + n\hat{y}_n}{n - \nu}\right), \quad \overline{E}[\lambda^{-1} | n, \hat{y}_n] = \frac{c + n\hat{y}_n}{n - \nu}, \quad (28)$$

while if $n < \nu$, the posterior expectation of λ^{-1} is vacuous. Thus, there exists a $\delta > 0$ such that for any $n > \delta$ it holds that $n > \nu$, which implies that properties (A3)–(A4) hold for inferences derived from (26).

6.2 Normal with Unknown Mean and Variance

In the case the variance is unknown, $\prod_{i=1}^n p(y_i | x) = \prod_{i=1}^n N(y_i; x, \sigma^2)$, with $\sigma^2 > 0$ and $x, y_i \in \mathbb{R}$ (the following results can be easily generalized to the case $x, y_i \in \mathbb{R}^d$). The likelihood model can

be expressed in the canonical form by setting $w = [\lambda x, -\frac{1}{2}\lambda]^T$,

$$b(w) = \frac{\lambda x^2 - \ln(\lambda)}{2}, \quad \nabla b(w) = [x, x^2 + \lambda^{-1}]^T,$$

where $\lambda = 1/\sigma^2$ and $\mathcal{W} = [\mathbb{R}, \mathbb{R}^-]^T$. Unfortunately in this case $\mathcal{W} \neq \mathbb{R}^2$ and, thus, Theorem 4.6 cannot be applied directly. However, this issue can be solved by instead considering the kernel

$$\mathcal{M} = \left\{ \exp(x\lambda \ell_{.1}) \lambda^{\frac{1}{2}} \exp(-\frac{1}{2}\lambda \ell_{.2}) \lambda^{-\frac{\nu}{2}}, \ell_{.i} \in [-c_{.i}, c_{.i}] \right\},$$

with $\nu \geq 2$. Notice in fact that $\exp(-\frac{1}{2}\lambda \ell_{.2}) \lambda^{-\frac{\nu}{2}}$ is the same model discussed in (26). The likelihood can be rewritten as:

$$p(\text{data}|x, \lambda) \propto \exp(n\hat{y}_{n1}\lambda x - n\frac{1}{2}\lambda x^2) \exp(-n\hat{y}_{n2}\frac{1}{2}\lambda) \lambda^{\frac{n}{2}},$$

where $n\hat{y}_{n1} = \sum_{i=1}^n y_i$ and $n\hat{y}_{n2} = \sum_{i=1}^n y_i^2$. By multiplication of likelihood and prior, one gets:

$$p(\text{data}|x, \lambda) p(x, \lambda) \propto \lambda^{\frac{1}{2}} \exp((n\hat{y}_{n1} + \ell_{.1})\lambda x - n\frac{1}{2}\lambda x^2) \exp(-(n\hat{y}_{n2} + \ell_{.2})\frac{1}{2}\lambda) \lambda^{\frac{n-\nu+2}{2}-1},$$

which gives the posterior

$$N(x; \hat{y}_{p1}, (n\lambda)^{-1}) G(\lambda; \frac{n-\nu+2}{2}, \frac{n}{2}(\hat{y}_{p2} - \hat{y}_{p1}^2)), \quad (29)$$

where $\hat{y}_{p1} = \frac{n\hat{y}_{n1} + \ell_{.1}}{n}$ and $\hat{y}_{p2} = \frac{n\hat{y}_{n2} + \ell_{.2}}{n}$. Observe that (29) makes sense only if $\hat{y}_{p2} - \hat{y}_{p1}^2 > 0$ and $n - \nu + 2 > 0$, i.e., the posterior is a proper Normal-Gamma density. For $n - \nu + 2 > 0$, by marginalizing out λ in (29), one can compute the posterior distribution of X

$$p(x|\text{data}) \propto \left[1 + \frac{(x - \hat{y}_{p1})^2}{\hat{y}_{p2} - \hat{y}_{p1}^2} \right]^{-(n-\nu+3)/2}, \quad (30)$$

which is a Student t distribution. Thus, for $n - \nu + 2 > 0$, it follows that

$$\underline{E}[X|n, \hat{y}_n] = \frac{n\hat{y}_{n1} - c_{.1}}{n}, \quad \bar{E}[X|n, \hat{y}_n] = \frac{n\hat{y}_{n1} + c_{.1}}{n},$$

$$\begin{aligned} \underline{E}[\lambda^{-1}|n, \hat{y}_n] &= \max \left(0, \frac{n}{n-\nu} \left[\left(\frac{\sum_{i=1}^n y_i^2 - c_{.2}}{n} \right) - \max_{|c| < c_{.1}} \left(\frac{\sum_{i=1}^n y_i + c}{n} \right)^2 \right] \right), \\ \bar{E}[\lambda^{-1}|n, \hat{y}_n] &= \frac{n}{n-\nu} \left(\frac{\sum_{i=1}^n y_i^2 + c_{.2}}{n} - \min_{|c| < c_{.1}} \left(\frac{\sum_{i=1}^n y_i + c}{n} \right)^2 \right). \end{aligned}$$

Observe that for $\ell_{.1} = \ell_{.2} = 0$ and $\nu = 3$, one obtains $E[x|n, \hat{y}_n] = \hat{y}_{n1}$ and $E[\lambda^{-1}|n, \hat{y}_n] = \frac{1}{n-3} \sum_{i=1}^n (y_i - \hat{y}_{n1})^2$ which are the estimates of X and σ^2 one would obtain from the improper prior $p(x, \sigma^2) = (\sigma^2)^{-1}$, while for $\nu = 2$ one obtains the estimates from the improper prior $p(x, \sigma^2) = (\sigma)^{-1}$, which gives the same inferences as the frequentist ones. The advantage of this model over the ones based on improper priors is that it satisfies prior near-ignorance.

7. Comparison with Other Prior Near-Ignorance Models

The aim of this section is to compare the model in Theorems 4.2–4.6 with other near-ignorance models available in the literature.

7.1 Near-Ignorance with a Set of Proper Priors in the One-Parameter Exponential Families

In [13], we have defined a model of prior ignorance for one-parameter exponential families based on a set of proper conjugate priors:

$$\mathcal{M} = \left\{ k(n_0, y_0) \exp(n_0(y_0 w - b(w))) : y_0 \in \mathcal{Y}_0, 0 < n_0 \leq \min(\bar{n}_0, c/|y_0|) \right\}. \quad (31)$$

We have proven that this set satisfies prior ignorance (A2) w.r.t. $g = b'$, learning (A3) and convergence (A4).¹ More precisely, we have proven that prior ignorance and learning/convergence can be both satisfied provided that y_0 is free to vary in \mathcal{Y}_0 and $|n_0 y_0| \leq c$, i.e., n_0 must depend on y_0 . It can be observed that the set of priors in Theorem 4.2 can be obtained from (31) by relaxing the constraint $y_0 \in \mathcal{Y}_0$ to $y_0 \in \mathbb{R}$ and by taking the limit of the inferences for $\bar{n}_0 \rightarrow 0$. In fact, $y_0 \notin \mathcal{Y}_0$ allows us to include at the limit also improper priors in (31), while the limit $\bar{n}_0 \rightarrow 0$ guarantees the fulfillment of (A1) as it is shown in Appendix A. Translation invariance is important because it guarantees that the posterior imprecision does not depend on the sample mean. This in general does not hold for (31).

Consider for instance the Poisson case. The lower and upper posterior means obtained from (31) and the Poisson likelihood are [13]:

$$\frac{n\hat{y}_n}{n + \bar{n}_0} \leq E[\lambda | n, \hat{y}_n] \leq \frac{n\hat{y}_n + c}{n}. \quad (32)$$

It can be noticed that the difference between the upper and lower means depends on \hat{y}_n for any $\bar{n}_0 > 0$. In fact, it can be verified that the set of priors (31) does not satisfy translation invariance, i.e., $\overline{E}[I_{\{A\}}] \neq \overline{E}[I_{\{A+a\}}]$. Translation invariance is only satisfied if we take the limit $\bar{n}_0 \rightarrow 0$ of the posterior inferences. However, in this case, the lower and upper posterior means reduce to:

$$\hat{y}_n \leq E[\lambda | n, \hat{y}_n] \leq \frac{n\hat{y}_n + c}{n}, \quad (33)$$

and, thus, the lower and upper means are not symmetric w.r.t. \hat{y}_n : the left bound is equal to the sample mean \hat{y}_n ; this is not desirable because it implies weakening the robustness of the posterior inferences. Conversely, by relaxing $y_0 \notin \mathcal{Y}_0$, the lower and upper posterior means at the limit $\bar{n}_0 \rightarrow 0$ become:

$$\max\left(0, \frac{n\hat{y}_n - c}{n}\right) \leq E[\lambda | n, \hat{y}_n] \leq \frac{n\hat{y}_n + c}{n}. \quad (34)$$

Now the lower and upper posterior means coincide with \hat{y}_n only for $n \rightarrow \infty$ and the posterior imprecision depends only on n and c for a sufficiently large n . In other words, posterior robustness and independence of posterior imprecision to the sample mean can both be guaranteed with conjugate prior models only by allowing finitely additive priors in the set \mathcal{M} . This means that they are incompatible with countable additivity.

¹Since the priors in \mathcal{M} are all countably additive, it also satisfies strong coherence as defined in [6, Ch. 7].

7.2 Imprecise Dirichlet Model

The *imprecise Dirichlet model* (IDM) [21], [22] is a model of prior near-ignorance for multinomial observations. The IDM considers the following set of priors:

$$\mathcal{M} = \left\{ \frac{\Gamma(s)}{\prod_{i=1}^d \Gamma(st_i)} \prod_{i=1}^d \theta_i^{st_i-1} : t_i > 0, \sum_{i=1}^d t_i = 1, s > 0 \right\}, \quad (35)$$

where $\Gamma(\cdot)$ is the Gamma function. Therefore \mathcal{M} includes Dirichlet densities parametrized by t_i and s . It follows that a priori $\underline{E}[\theta_i] = 0$ and $\bar{E}[\theta_i] = 1$, while a-posteriori:

$$\underline{E}[\theta_i | n, \hat{y}_n] = \frac{n\hat{y}_{ni}}{n+s}, \quad \bar{E}[\theta_i | n, \hat{y}_n] = \frac{n\hat{y}_{ni} + s}{n+s}. \quad (36)$$

The IDM is a model of prior ignorance that guarantees conjugacy between likelihood and priors and satisfies (A2)–(A4) but it does not satisfy translation invariance (A1) in the transformed domain \mathcal{W} . Observe that although the IDM does not satisfy translation invariance, the posterior imprecision does not depend on the sample mean. Furthermore, the IDM is invariant to permutations of the components of θ and satisfies RIP. Thus the IDM satisfies permutation invariance and independence of the posterior imprecision to the sample mean without the need of finitely additive priors.

The question is what happens if we impose translation invariance to the IDM in \mathcal{W} . Fixed $s > 0$, the only way to satisfy also (A1) is not constraining t to lie in $\mathcal{B}_0 = \{t : t_i > 0, \sum_{i=1}^d t_i = 1\}$. Assuming that $t_i \in \mathbb{R}$ and imposing the constraint $0 < |st_i| \leq c_i = c$ for $i = 1, \dots, d-1$, the set of priors (35) reduces to

$$\mathcal{M} = \left\{ \prod_{i=1}^{d-1} \theta_i^{\ell_i-1} \left(1 - \sum_{j=1}^{d-1} \theta_j \right)^{-\sum_{j=1}^{d-1} \ell_j-1} : \ell_i \in [-c, c] \right\}, \quad (37)$$

which is the model in (20). If we impose one-step ahead RIP to (20), we obtain (23).

7.3 Nonparametric Predictive Inference for Multinomial Data and Known Number of Categories

In the case of multinomial data, it is worth pointing out also the relationship between the model in Theorem 4.6 and another model for robust inferences, the so-called *nonparametric predictive inference* (NPI) model [15]. The NPI is a nonparametric model based on a post-data assumption about the uncertainty associated to a future observation and can only yield predictive inferences (inferences on the space of the observations). Consider $k \geq 3$ possible categories denoted by C_1, \dots, C_k . Without loss of generality, we assume that the first o of these have already been observed and the last $k - o$ have not yet been observed. There are n observations $\{y_1, \dots, y_n\}$ such that $y_{ij} \in \{0, 1\}$ for $i = 1, \dots, n$ and $j = 1, \dots, k$ and $\sum_{i,j} y_{ij} = n$. The event of interest can generally be denoted by

$$Y_{n+1} \in \bigcup_{j \in J} C_j,$$

i.e., the next observation belongs to the subset of categories $J \subseteq \{1, \dots, k\}$. Let $OJ = J \cap \{C_1, \dots, C_o\}$ denote the index-set for the categories in the event of interest that have already been observed, and $UJ = J \cap \{C_{o+1}, \dots, C_k\}$ the corresponding index-set for the categories in the event of interest that have not yet been observed. Let r be the number of elements of OJ and

l the number of elements of UJ , so $0 \leq r \leq o$ and $0 \leq l \leq k - o$. The NPI-based lower and upper probabilities for the event of interest, based on the n observations, are:

$$\begin{aligned} \underline{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \frac{1}{n} (\sum_{i=1}^n \sum_{j \in J} y_{ij} + \max(r+l-k, -r)), \\ \bar{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \frac{1}{n} (\sum_{i=1}^n \sum_{j \in J} y_{ij} + \min(r+l, o-r)). \end{aligned} \quad (38)$$

It can be noticed that in case $l = 0$, i.e., all the categories have been observed, then $r = |J|$ and $\sum_{i=1}^n \sum_{j \in J} y_{ij} \geq |J|$ and, thus, one has:

$$\begin{aligned} \underline{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \frac{1}{n} (\sum_{i=1}^n \sum_{j \in J} y_{ij} + \max(|J| - k, -|J|)), \\ \bar{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \frac{1}{n} (\sum_{i=1}^n \sum_{j \in J} y_{ij} + \min(|J|, k - |J|)), \end{aligned} \quad (39)$$

which coincides exactly with (22) in the case $c = 1$. It is thus extremely interesting that in this case the same inferences of the NPI can be obtained by our set of conjugate parametric priors which includes, at the limit, improper densities. In the case $l \neq 0$, the two models still coincide provided that $\sum_{i=1}^n \sum_{j \in J} y_{ij} = r$, i.e., each category in OJ has been observed at most once. In the other cases, the NPI is in general more precise. This can be easily proven by comparing (38) and (39) and by noticing that $-r \geq -|J|$, $r+l-k = |J| - k$ and $r+l = |J|$, $o-r \leq k - |J| = o-r+k-o-l$. For instance, consider the case $k = 4$, $J = \{1, 2\}$, $n = 2$ and only the category 1 has been observed. Then, since $l = 1$, the lower and upper probability of the event of interest obtained with the NPI are:

$$\begin{aligned} \underline{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \frac{1}{2} (2 + \max(1 + 1 - 4, -1)) = \frac{1}{2}, \\ \bar{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \frac{1}{2} (2 + \min(1 + 1, 0)) = 1, \end{aligned}$$

while the inferences obtained with the model in (22) are:

$$\begin{aligned} \underline{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \max(0, \frac{1}{2} (2 + \max(2 - 4, -2))) = 0, \\ \bar{E}[Y_{n+1} \in \cup_{j \in J} C_j | n, \hat{y}_n] &= \min(1, \frac{1}{2} (2 + \min(2, 2))) = 1, \end{aligned}$$

and it remains vacuous up to 3 observations. In [15] the authors compare the NPI with the IDM highlighting the following differences between them. First, the IDM lower probability for the second observation to be equal to the first, is $1/(1+s)$. Thus, small values of s , e.g., $s = 1$ or $s = 2$, lead to surprisingly high values for this lower probability. In the NPI, this lower probability is 0 and the same holds for our model (22).

Second, the IDM predictive lower and upper probabilities depend only on the observed frequency of that category and the total number of observations (RIP). This is not the case for NPI-based lower and upper probabilities and also for our model (21). However, in Section 4.2, we have pointed out that one-step ahead RIP can be satisfied by our model by adding an additional constraint w.r.t. the NPI, see (23).

Third, the IDM upper probabilities for events that the next observation is in an as yet unseen category do not depend on the number of categories seen so far. This is not the case for the NPI, while this is also the case for our model. The difference is the term l in the computation of the lower and upper probabilities, which makes NPI inferences depend on the unobserved categories. This behavior could be included in our model only by letting the set of priors depend on the data.

The main advantage of the NPI w.r.t. the IDM and our models is that it is a nonparametric model. In fact, the NPI can also be used if one does not know the total number of possible categories, and wishes to distinguish in the event of interest between fully defined categories that have not yet been observed, and any new category occurring at the next observation [15]. In the IDM and in our model, the lower and upper probabilities of these two cases are the same.

7.4 Bounded Derivative Model

Another comparison is with the so-called *bounded derivative model* (BDM) [14]. The BDM is a prior near-ignorance model for a scalar variable $w \in \mathbb{R}$ in which the set of priors \mathcal{M} includes all continuous proper probability density functions for which the derivative of the log-density is bounded by a positive constant:

$$\left| \frac{d}{dw} \ln p(w) \right| \leq c. \quad (40)$$

Observe that the BDM satisfies (A1)–(A4) as stated in Theorem 4.2. It is clear that the set of priors in Theorem 4.2 also satisfies (40), since

$$\left| \frac{d}{dw} \ln \exp(\ell w) \right| = |\ell| \leq c.$$

However, the BDM set of priors is larger, since it does not restrict the set of priors to be $\{\exp(\ell w), |\ell| \leq c\}$. We have imposed this set of priors in order to exploit conjugacy. This is not the case of the BDM in which the extreme priors that obtain the lower and upper expectation of a RVBF cannot always be expressed as simple limits of conjugate priors. Walley in fact proves in [14] that in the BDM the extreme priors which obtain the upper and lower expectations of a generic RVBF g are piecewise exponential, i.e., the real line can be divided into intervals on which $p(w) = \exp(cw)$ or $p(w) = \exp(-cw)$. These piecewise priors do not belong to $\{\exp(\ell w), |\ell| \leq c\}$ in general. It is however worth to observe that in the case $g = b^l$ or g is an indicator over a subset of $\mathcal{W} = \mathbb{R}$ (in general if g is monotone), both the BDM and our model produce the same lower and upper expectations. This holds since the lower and upper expectations are obtained in correspondence of either the kernel $p(w) = \exp(cw)$ or $p(w) = \exp(-cw)$. For a generic RVBF g , we expect that the inferences obtained with the BDM are more conservative, i.e., it has a larger posterior imprecision. The price to be paid for this gain of robustness and generality is the increase of the computational cost for the inferences.

8. Some Practical Examples

In this section we give some insight of the possible applications of the proposed models in statistical inference.

8.1 One-Sample Location Test

Consider i.i.d. observations y_1, \dots, y_n from a real variable Y with unknown mean x and unknown precision λ ($\lambda = 1/\sigma^2$). Our goal is to test the hypotheses:

$$H_0 : x \leq 0, \quad H_1 : x > 0.$$

A way to tackle this problem is by means of a frequentist t-test. The t-test considers the statistic:

$$t = \frac{\hat{y}_{p1} - \Delta_0}{\sqrt{\frac{\hat{y}_{p2} - \hat{y}_{p1}^2}{n-1}}},$$

where Δ_0 is the value of the mean under the null hypothesis, $\hat{y}_{p1} = \frac{1}{n} \sum_i y_i$, $\hat{y}_{p2} = \frac{1}{n} \sum_i y_i^2$ and, thus, $\hat{y}_{p2} - \hat{y}_{p1}^2 = \frac{1}{n} \sum_i (y_i - \hat{y}_{p1})^2$. Then it computes the p -value of t under the null hypothesis,

i.e., t is Student-t distributed with mean $\Delta_0 = 0$ and $n - 1$ degrees of freedom. Hypothesis H_0 is rejected when the p -value is less than the given threshold $\alpha = 0.05$. An alternative approach is to perform a Bayesian test, using a Normal-Gamma model based on the improper prior $p(x, \lambda) = \lambda^{\frac{-\nu+1}{2}}$, with parameter $\nu > 0$. In this case we accept hypothesis H_1 if the posterior probability $P(x > 0 | y_1, \dots, y_n) > 1 - \alpha$, with

$$P(x > 0 | y_1, \dots, y_n) = 1 - tcdf \left(\frac{0 - \hat{y}_{p1}}{\sqrt{\frac{\hat{y}_{p2} - \hat{y}_{p1}^2}{n - \nu + 1}}}, n - \nu + 1 \right),$$

where $tcdf$ is the cumulative distribution function of Student distribution with $n - \nu + 1$ degrees of freedom computed in 0. Note that for $\nu = 3$ we obtain the left Haar invariant prior, while for $\nu = 2$ we obtain the right Haar invariant prior [2, Sec. 3.3.3]. In the experiments we have selected $\nu = 2$ because the Bayesian test and the frequentist test coincide for this value. We can also devise a Bayesian test based on a hierarchical Normal model [23, Sec. 11.7]:

$$\text{likelihood: } N(y; x, \sigma^2), \quad \text{prior: } N(x; \mu, \tau^2), \quad \text{hyper-prior: } p(\mu, \log(\sigma), \log(\tau)) \propto \tau.$$

Note that the hyperprior has been defined on $\log(\sigma), \log(\tau)$ to avoid that the posterior distribution is improper. However several other choices are possible. The posterior probability $P(x > 0 | y_1, \dots, y_n)$ (necessary for the hypothesis test) is computed numerically by Gibbs sampling.

Finally we can perform a test based on the *imprecise exponential family model* (IEM) presented in Section 6.2. In this case we compute the lower and upper probability:¹

$$\begin{aligned} \underline{P}(x > 0 | y_1, \dots, y_n) &= \min_{\ell_{.1} \in [-c_{.1}, c_{.1}], \ell_{.2} \in [-c_{.2}, c_{.2}]} 1 - tcdf \left(\frac{0 - \hat{y}_{p1} + \frac{\ell_{.1}}{n}}{\sqrt{\frac{\hat{y}_{p2} + \frac{\ell_{.2}}{n} - (\hat{y}_{p1} + \frac{\ell_{.1}}{n})^2}{n - \nu + 1}}}, n - \nu + 1 \right), \\ \bar{P}(x > 0 | y_1, \dots, y_n) &= \max_{\ell_{.1} \in [-c_{.1}, c_{.1}], \ell_{.2} \in [-c_{.2}, c_{.2}]} 1 - tcdf \left(\frac{0 - \hat{y}_{p1} + \frac{\ell_{.1}}{n}}{\sqrt{\frac{\hat{y}_{p2} + \frac{\ell_{.2}}{n} - (\hat{y}_{p1} + \frac{\ell_{.1}}{n})^2}{n - \nu + 1}}}, n - \nu + 1 \right), \end{aligned} \quad (41)$$

and we perform the hypothesis test as follows:

- return H_1 if $\underline{P}(x > 0 | y_1, \dots, y_n) > 1 - \alpha$;
- return H_0 if $\bar{P}(x > 0 | y_1, \dots, y_n) < 1 - \alpha$;
- otherwise issue an indeterminate answer, i.e., $\{H_0, H_1\}$.

In the latter case we cannot decide: the result of the hypothesis test depends on the choice of the prior. To compare the frequentist/Bayesian and IEM-based test we have considered a Monte Carlo experiment in which n observations Y are generated based on $Y \sim N(\Delta, 1)$, with Δ ranging from 0 to 1.5; for $\Delta = 0$, we are under the null hypothesis of the frequentist/Bayesian tests. In particular, for each value of Δ we have performed 20000 Monte Carlo runs by generating in each run $n = 15$ and $n = 30$ observations for Y and computed the percentage of cases in which the alternative hypothesis is returned, using $\alpha = 0.05$ for both the tests. This means that for $\Delta = 0$

¹In the formulation we need to impose the additional constraint that the argument of the square root is positive. This implies that the parameters $\ell_{.1}$ and $\ell_{.2}$ cannot vary independently in $[-c_{.1}, c_{.1}]$ and, respectively, $[-c_{.2}, c_{.2}]$. However, for suitably large n and non-degenerate distributions, the argument of the square root is usually positive for any $\ell_{.1}$ and $\ell_{.2}$ in the intervals.

we are computing the Type-I error of the test and for $\Delta > 0$ its power. The IEM test has been implemented by choosing $\nu = 2$ (equal to the frequentist/Bayesian tests). About the choice of the c parameters, we have selected c_1 equal to 0.75 so that the posterior imprecision for the estimate of x is reduced to 1.5 after one observation, while $c_2 \approx 0.28$, so that $c_2 = (c_1)^2/n$ for $n = 2$. The expression $(c_1)^2/2$ gives the sample second non-central moment of the pseudo-observations $\pm c_1$; remember in fact that ℓ_2 can be interpreted as the square of a pseudo-observation and we need two observations to estimate both x and λ . This is the meaning of $c_2 = (c_1)^2/2$. Note that with this choice our model has only one free parameter c_1 . We will empirically show that this choice is a good compromise between robustness and conservativeness (indeterminacy).

The average power¹ is shown in Figure 1 as a function of Δ for the two cases. In particular,

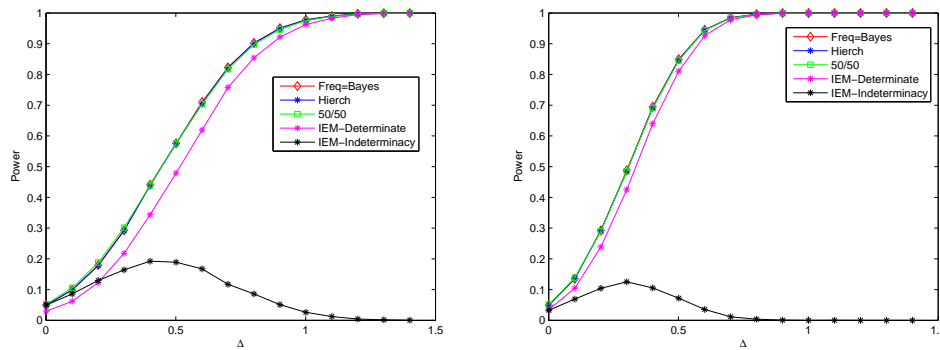


Figure 1.: Power as a function of Δ for the case $n = 15$ (left) and $n = 30$ (right).

Figure 1 reports (i) the power of the frequentist/Bayesian tests; (ii) the power of the hierarchical Bayesian test; (iii) the power of the IEM test when it is determinate; (iv) the indeterminacy of the IEM test, i.e., the number of times it returns an indeterminate response divided by the total number of Monte Carlo runs; (v) the power of a new test (called “50/50”) which returns the same response as IEM when IEM is determinate, and issues a random answer (with 50/50 chance) when IEM is indeterminate. We have introduced this new artificial test to facilitate the comparison of the frequentist/Bayesian, hierarchical Bayesian and IEM tests. From Figure 1 (left and right), it is evident that the performance of the frequentist/Bayesian, hierarchical Bayesian and 50/50 tests practically coincide. Furthermore, since in all cases in which IEM is determinate, the frequentist/Bayesian and hierarchical Bayesian return the same response as IEM (i.e., in the determinate cases the accuracies of the frequentist/Bayesian and hierarchical Bayesian tests are also equal to the (magenta) star curve in Figure 1), the difference between the three tests is only in the runs where IEM is indeterminate. In these runs, the frequentist/Bayesian and hierarchical Bayesian tests have essentially the same accuracy as the 50/50 test. Therefore, the IEM is able to isolate some instances in which the frequentist/Bayesian and hierarchical Bayesian tests are virtually guessing at random, in the sense that they have the same accuracy of the 50/50 test. Assume for instance that we are trying to evaluate the effects of a medical treatments (“X is greater than zero”) and that, given the available data, the IEM is indeterminate. In such a situation the frequentist/Bayesian and hierarchical tests always issue a determinate response (I can tell if “X is greater than zero”), but it turns out that their response is completely random (like if we were tossing a coin). On the other side, the IEM acknowledges the impossibility of making a decision (I do not know whether “X is greater than zero”) and thus, although the frequentist/Bayesian and hierarchical tests and the IEM (more precisely the 50/50 test derived by IEM) have the same accuracy, the IEM provides more information.

¹For $\Delta = 0$ the plot actually reports the Type I error.

Note that in Figure 1 (left) for $\Delta \approx 0.4$, the percentage of runs in which IEM is indeterminate is about 18% ; this means that the frequentist/Bayesian and hierarchical tests are issuing a random answer in 18% of the cases, which is a large percentage. Looking at the curve of the IEM indeterminacy, it can be observed that it reaches its maximum value at $\Delta = 0.4$ and then goes to zero. For $|\Delta| > 1.2$ the power of the frequentist/Bayesian, hierarchical Bayesian and IEM-determinate is practically one, which means that when the hypothesis test is easy (large Δ), there are not indeterminate instances and both tests have power one. The conclusions for Figure 1 (right) are similar.

For the frequentist/Bayesian test, we have also evaluated the distribution of the p -values in the instances where the IEM is indeterminate. The histogram of the p -values is shown in Figures 2 (a) and (b) for $\Delta = 0$ and, respectively, $\Delta = 0.4$. It can be observed that the IEM test is indeterminate when the p -value is close to the threshold 0.05, but also in instances where it is quite far (e.g., $p = 0.02$ or $p = 0.1$). Moreover, it can be noticed that the distribution of the p -values in the IEM indeterminate instances is roughly symmetric around 0.05 (more precisely, the area of the distribution of the p -values to the left of 0.05 is (approximately) equal to the area to the right of 0.05), no matter the value of Δ . This explains why the “50/50” test and the frequentist/Bayesian test have the same power. This behavior cannot be reproduced by simply adding a “no decision zone” to the frequentist/Bayesian test, i.e., by suspending the decision when $a < p < b$. For instance, Figures 2 (c) and (d) show the distribution of the p -values of the frequentist/Bayesian test in the interval with boundary $a = 0.02$ and $b = 0.095$.¹ By comparing (c) and (d), it can be observed that the distribution of the p -values in the “no decision zone” $0.02 < p < 0.095$ changes with Δ . It is uniform in the case $\Delta = 0$ (it is in fact well known that the distribution of the p -values is uniform under the null hypothesis), but is strongly skewed in the case $\Delta = 0.4$. This is not true for the distribution of the p -values in the IEM indeterminate instances, which instead does not change much with Δ , see Figures 2 (a) and (b). Therefore, since the distribution of the p -values in the “no decision zone” is not symmetric around 0.05 (the case $\Delta = 0.4$ in Figure 2 (d)), a “50/50” test designed over the “no decision zone” $0.02 < p < 0.095$ would not have the same behavior of IEM “50/50” test, e.g., it would not have the same power as the frequentist/Bayesian test. In other words, the “no decision zone” $0.02 < p < 0.095$ is not able to isolate the instances that are critical for the frequentist/Bayesian test, i.e., it is both classifying as indeterminate the instances that are easy and those that are difficult. This is evident comparing Figures 2 (b) and (d) for $p = 0.02$: the IEM is classifying as indeterminate only few instances with p -value close to $p = 0.02$, while the “no decision zone” $0.02 < p < 0.095$ is classifying as indeterminate all the instances with p -value close to $p = 0.02$. The only way a frequentist test with a “no decision zone” can mimic the behavior of the IEM test is with $a(data) < p < b(data)$. This means that the boundary must depend on the observations. Mathematically, this is evident from (41). The indeterminacy of the IEM test cannot be reproduced by $a < P(x > 0|y_1, \dots, y_n) < b$, since the values ℓ_1 and ℓ_2 that give the maximum and the minimum depend on the data (through the nonlinear function tcd).

Finally, in case $n = 15$ and $\Delta = 0$, we have computed the error for the IEM-based test as a function of c_1 (with $c_2 = (c_1)^2/2$). The error of the frequentist test is in this case equal to 0.05 (we are under the null hypothesis). From Figure 3, it can be observed that the error of the IEM test when determinate decreases with c_1 , because of the increase of the indeterminacy (note that for $c_1 = 0$, the error of IEM when determinate and the error of the 50/50 test coincide, since there are not indeterminate instances in this case). Looking at the plot relative to the 50/50 test, it can be noticed that it is almost constant for $c_1 < 0.75$ and then it increases. This (together with the other results of this section) may be seen as an empirical confirmation that the choice of $c_1 = 0.75$ is appropriate. It guarantees the maximum robustness at the minimum of conservativeness (indeterminacy). We have also evaluated frequentist properties of this choice.

¹We have chosen this interval in analogy with that of the IEM test. The distribution of the p -values for a different boundary of the “no decision zone” can easily be deduced from Figures 2 (c) and (d).

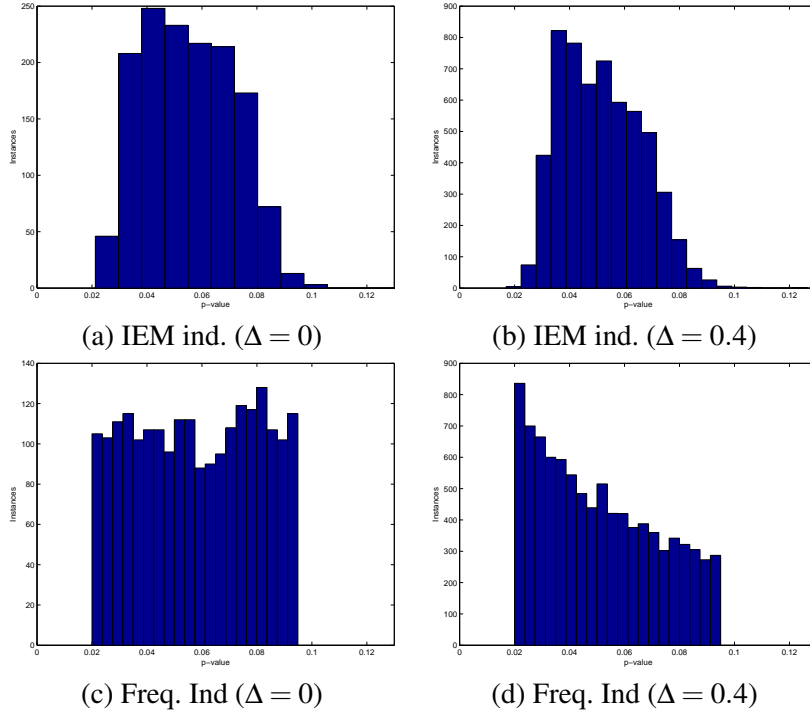


Figure 2.: Histogram of the p -values ($n = 15$).

	Calibration
$n = 15$	0.9686
$n = 30$	0.9593
$n = 150$	0.9552

Table 1.: Calibration of the IEM robust (symmetric) 95% high posterior density interval.

In particular, we have evaluated the calibration of the robust (symmetric) high posterior density interval (HDI) of the IEM, i.e., the minimum length interval that has lower probability equal to 0.95. The results are shown in Table 1 as a function of n . It can be observed that the robust HDI includes the true mean with a probability always greater than 0.95, which means that it is calibrated. The probability does not excessively exceed the prescribed value 0.95, which is another empirical confirmation that the choice $c_{.1} = 0.75$ is a good compromise between robustness and conservativeness.

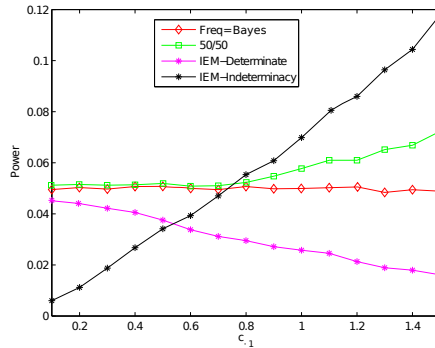


Figure 3.: Power as a function of $c_{.1}$ for the case $n = 15$.

It should be stressed again that the test “50/50” has been introduced only for the sake of comparison. We are not suggesting that when IEM is indeterminate we should toss a coin to take the decision. On the contrary we claim that the indeterminacy of IEM is an additional useful information that our approach gives to the analyst. In these cases she/he knows that (i) her/his posterior decisions would depend on the choice of the prior; (ii) the hypothesis test is difficult as shown by the comparison with the frequentist test. Based on this additional information, the analyst can for example decide to collect additional measurements to eliminate the indeterminacy (in fact we have seen that when the number of observations goes to infinity the indeterminacy goes to zero).

8.2 Election Polls

A total of n adults are polled to indicate their preference for two candidates A and B . Let \hat{y}_{n1} denote the proportion of the population that supports A , \hat{y}_{n2} denote the proportion of the population that supports B and $1 - \hat{y}_{n1} - \hat{y}_{n2}$ denote the proportion of the population that is either undecided or vote for someone else. The counts $n\hat{y}_{n1}$, $n\hat{y}_{n2}$ and $n(1 - \hat{y}_{n1} - \hat{y}_{n2})$ are assumed to have a multinomial distribution with sample size n and respectively parameters θ_1 , θ_2 and θ_3 . Thus, the likelihood model is:

$$p(n, \hat{y}_n | \theta) = \theta_1^{n\hat{y}_{n1}} \theta_2^{n\hat{y}_{n2}} \theta_3^{n\hat{y}_{n3}},$$

where $\theta_1 + \theta_2 + \theta_3 = 1$ are the unknown non-negative chances to be estimated. The focus is to compare the proportions of voters for A and B by considering the difference $\theta_1 - \theta_2$. A Dirichlet conjugate prior can be assumed on θ_1 , θ_2 and θ_3 :

$$p(\theta) \propto \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \theta_3^{\alpha_3-1},$$

where in the case of lack of prior information the prior parameters are commonly selected as follows: Haldane’s prior $\alpha_1 = \alpha_2 = \alpha_3 = 0$; Jeffreys’ prior $\alpha_1 = \alpha_2 = \alpha_3 = \frac{1}{2}$; uniform prior $\alpha_1 = \alpha_2 = \alpha_3 = 1$. The expected value of $E[\theta_1 - \theta_2]$ is equal to 0 and the prior probability $P(\theta_1 > \theta_2) = 0.5$ for both Jeffreys and uniform priors (for Haldane’s prior they are not defined). These priors express indifference between candidate A and B , but not prior ignorance. To see that, consider $P[\theta_1 + 0.5\theta_3 > \theta_2 + 0.4\theta_3]$, this is the probability that the proportion of votes of A exceeds the votes for B assuming a “swing” scenario in which 50% of the undecideds vote for A and 40% of the undecideds for B . This probability is equal to 0.76 in the case of the uniform prior and 0.66 in the case of Jeffreys’ prior. It depends on the choice of the prior and this shows that the uniform and Jeffrey’s priors are not really uninformative for this kind of poll.

Combining likelihood and prior, the resulting posterior is

$$p(x|n, \hat{y}_n) \propto \theta_1^{n\hat{y}_{n1} + \alpha_1 - 1} \theta_2^{n\hat{y}_{n2} + \alpha_2 - 1} \theta_3^{n\hat{y}_{n3} + \alpha_3 - 1},$$

which is always proper in the case of Jeffreys’ and uniform prior and in the case of Haldane’s prior provided that $\hat{y}_{n1}, \hat{y}_{n2}, \hat{y}_{n3} > 0$. The posterior expected value of $\theta_1 - \theta_2$ is:

$$E[\theta_1 - \theta_2 | n, \hat{y}_n] = \frac{n\hat{y}_{n1} + \alpha_1}{n + \alpha_1 + \alpha_2 + \alpha_3} - \frac{n\hat{y}_{n2} + \alpha_2}{n + \alpha_1 + \alpha_2 + \alpha_3},$$

while the posterior probability of the event $\theta_1 - \theta_2 > 0$ is

$$P[\theta_1 > \theta_2 | n, \hat{y}_n] = \frac{\int_{\{\theta_1 > \theta_2\}} \theta_{.1}^{n\hat{y}_{n1} + \alpha_1 - 1} \theta_{.2}^{n\hat{y}_{n2} + \alpha_2 - 1} \theta_{.3}^{n\hat{y}_{n3} + \alpha_3 - 1} d\theta}{\int \theta_{.1}^{n\hat{y}_{n1} + \alpha_1 - 1} \theta_{.2}^{n\hat{y}_{n2} + \alpha_2 - 1} \theta_{.3}^{n\hat{y}_{n3} + \alpha_3 - 1} d\theta},$$

which can be computed numerically by sampling from the Dirichlet distribution.

Consider now the model of near-ignorance \mathcal{M} defined in Theorem 4.6. From these results it follows that a priori $\underline{E}[\theta_1 - \theta_2] = -1$, $\bar{E}[\theta_1 - \theta_2] = 1$. From Corollary 5.1, one has that $\underline{P}(\theta_1 > \theta_2) = 0$, $\bar{P}(\theta_1 > \theta_2) = 1$ and that $\underline{P}[\theta_1 + 0.5\theta_3 > \theta_2 + 0.4\theta_3] = 0$, and $\bar{P}[\theta_1 + 0.5\theta_3 > \theta_2 + 0.4\theta_3] = 1$.¹ This is a more correct expression of the lack of prior information on the election result. The set of posteriors produced by \mathcal{M} is:

$$\mathcal{M}_p = \left\{ p(x | n, \hat{y}_n) \propto \theta_{.1}^{n\hat{y}_{n1} + \ell_1 - 1} \theta_{.2}^{n\hat{y}_{n2} + \ell_2 - 1} \theta_{.3}^{n\hat{y}_{n3} + \ell_3 - 1} \right\},$$

with $\ell \in \mathbb{L}$. Assuming that \mathbb{L} is the set defined in (23), two of the extreme densities belonging to this set are:

$$\theta_{.1}^{n\hat{y}_{n1} - c - 1} \theta_{.2}^{n\hat{y}_{n2} + c - 1} \theta_{.3}^{n\hat{y}_{n3} - 1}, \quad (42)$$

$$\theta_{.1}^{n\hat{y}_{n1} + c - 1} \theta_{.2}^{n\hat{y}_{n2} - c - 1} \theta_{.3}^{n\hat{y}_{n3} - 1}. \quad (43)$$

These extreme densities give respectively the lower and upper posterior expected values of $\theta_1 - \theta_2$:

$$\underline{E}[\theta_1 - \theta_2 | n, \hat{y}_n] = \max \left(-1, \frac{n\hat{y}_{n1}}{n} - \frac{n\hat{y}_{n2}}{n} - \frac{2c}{n} \right), \quad (44)$$

$$\bar{E}[\theta_1 - \theta_2 | n, \hat{y}_n] = \min \left(1, \frac{n\hat{y}_{n1}}{n} - \frac{n\hat{y}_{n2}}{n} + \frac{2c}{n} \right), \quad (45)$$

where it has been assumed that $n\hat{y}_{n1}, n\hat{y}_{n2} > 0$. The lower expectation considers the case in which c votes move from candidate A to B . The upper expectation considers the case in which c votes move from candidate B to A . Clearly this exchange of votes gives also the the lower and upper posterior probabilities of the event $\theta_1 - \theta_2 > 0$:

$$\underline{P}[\theta_1 > \theta_2 | n, \hat{y}_n] = \frac{\int_{\{\theta_1 > \theta_2\}} \theta_{.1}^{n\hat{y}_{n1} - c - 1} \theta_{.2}^{n\hat{y}_{n2} + c - 1} \theta_{.3}^{n\hat{y}_{n3} - 1} d\theta}{\int \theta_{.1}^{n\hat{y}_{n1} - c - 1} \theta_{.2}^{n\hat{y}_{n2} + c - 1} \theta_{.3}^{n\hat{y}_{n3} - 1} d\theta},$$

$$\bar{P}[\theta_1 > \theta_2 | n, \hat{y}_n] = \frac{\int_{\{\theta_1 > \theta_2\}} \theta_{.1}^{n\hat{y}_{n1} + c - 1} \theta_{.2}^{n\hat{y}_{n2} - c - 1} \theta_{.3}^{n\hat{y}_{n3} - 1} d\theta}{\int \theta_{.1}^{n\hat{y}_{n1} + c - 1} \theta_{.2}^{n\hat{y}_{n2} - c - 1} \theta_{.3}^{n\hat{y}_{n3} - 1} d\theta}.$$

Above it has been assumed that the resulting posteriors are both proper. Similar calculations can be used to determine $\underline{P}[\theta_1 + 0.5\theta_3 > \theta_2 + 0.4\theta_3 | n, \hat{y}_n]$ and $\bar{P}[\theta_1 + 0.5\theta_3 > \theta_2 + 0.4\theta_3 | n, \hat{y}_n]$.

¹These lower and upper probabilities are obtained by densities in \mathcal{M} approaching the extreme priors in (42)–(43).

Compare (44)–(45) with the posterior means obtained by the IDM:

$$\underline{E}[\theta_1 - \theta_2 | n, \hat{y}_n] = \frac{n\hat{y}_{n1}}{n+s} - \frac{n\hat{y}_{n2} + s}{n+s}, \quad (46)$$

$$\bar{E}[\theta_1 - \theta_2 | n, \hat{y}_n] = \frac{n\hat{y}_{n1} + s}{n+s} - \frac{n\hat{y}_{n2}}{n+s}. \quad (47)$$

In this case the lower and upper expectations consider the case in which s more adults are polled and assume that these adults vote either for candidate B (lower case) or for candidate A (upper case). Although by choosing an appropriate value of s we can suitably enlarge $\bar{E}[\theta_1 - \theta_2 | n, \hat{y}_n] - \underline{E}[\theta_1 - \theta_2 | n, \hat{y}_n]$, the increase of s decreases the variance of the extreme densities and, thus, makes the inferences sharper.

As numerical example we consider the problem of estimating the number of electoral votes in the 2004 USA Presidential Election using the polling data from `realclearpolitics.com` (last-day poll). We compute the winner prediction of a Bayesian estimator that: (i) using Haldane's prior computes the win probabilities $P[\theta_1 > \theta_2 | n, \hat{y}_n]$ for each state; (ii) we use these probabilities to compute the (sampling) distribution of the total electoral votes for the two candidates.

Then we compute the winner prediction for the two near-ignorance prior models: IDM ($s = 2$) and the set of priors \mathcal{M} with $c = 1$ hereafter denoted as IEM (imprecise exponential model), using the same procedure as in the Bayesian estimator but considering the two cases $\underline{P}[\theta_1 > \theta_2 | n, \hat{y}_n]$ and $\bar{P}[\theta_1 > \theta_2 | n, \hat{y}_n]$, which give the lower and, respectively, upper probability that Bush wins in each state.

A number of 200000 simulated elections are generated. In each simulation, for each State we have first computed the posterior probability of θ_1 and θ_2 (using the Bayesian prior or the lower and upper priors). For each State, we have used this posterior to sample θ_1 and θ_2 and assigned the electoral vote of the State to Bush whenever $\theta_1 > \theta_2$ (to Kerry whenever $\theta_2 < \theta_1$). Finally, we have computed Bush's total electoral votes. We have repeated this process for each of the 200000 runs. Figure 4 reports the histograms of Bush's total electoral votes that we have obtained on the 200000 runs for the three estimator together with the break-even equal to 269 electoral votes. It is evident from the histogram of the Bayesian estimator that Kerry is slightly favorite (for Bush the posterior mean of electoral votes is 267.1 and the posterior median 266) but also that there is a high uncertainty. Because of this uncertainty the contribution of the prior on the final result is crucial. In fact, notice that if we just add two (pseudo-)samples ($s = 2$) in favor of Bush in each state, then Bush becomes slightly favorite (posterior upper mean equal to 271.2 electoral votes and posterior upper median equal to 270 electoral votes for Bush), see IDM light red histogram. This corresponds to the case $\bar{P}[\theta_1 > \theta_2 | n, \hat{y}_n]$ in the IDM and the best scenario for Bush. The lower probability $\underline{P}[\theta_1 > \theta_2 | n, \hat{y}_n]$ (IDM light blue histogram) gives instead the worst scenario for Bush. Similar results hold for the IEM. In this case, we assume that in each state one elector ($c = 1$) in the sample changes mind and instead of voting for Kerry he/she votes for Bush. This is enough for Bush to become slightly favorite. Notice that the IDM and IEM give the same result in this case. In fact, being the sample size large (around 600 in each state), $n + s \approx n$ for $s = 2$ and, thus, the inferences of the IDM and IEM almost coincide.

9. Conclusions

In this paper, we have proposed a model of prior ignorance about a multivariate variable based on a set of distributions \mathcal{M} . In particular, we have discussed four properties that a prior model should satisfy to model lack of prior information: invariance, near-ignorance, learning and convergence. Near-ignorance and invariance ensure that our prior model behaves as a vacuous model with respect to some statistical inferences (e.g., mean, credible intervals, etc.) and some transfor-

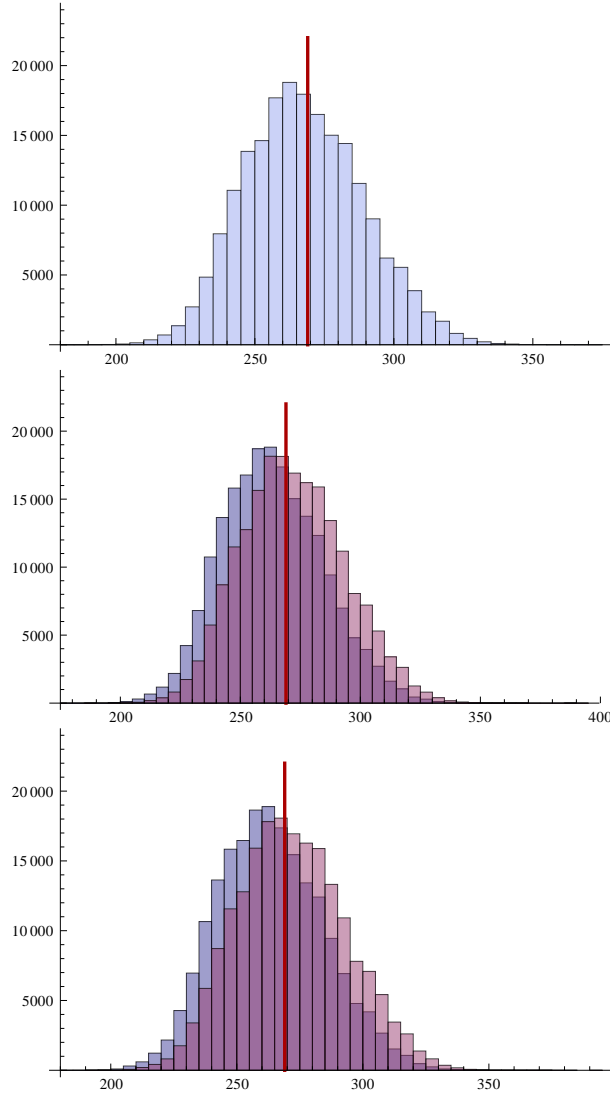


Figure 4.: Histogram of Bush's total electoral votes for the Bayesian estimator (top), IDM (center) and IEM (bottom). The vertical line denotes the break-even 269 electoral votes.

mation of the parameter space. Learning and convergence ensure that our prior model can learn from data and, in particular, that the influence of \mathcal{M} on the posterior inferences vanishes with increasing number of observations. Furthermore, for exponential families, we have shown that translation invariance, near-ignorance, learning and convergence can all be satisfied by a set of conjugate priors if the set of priors \mathcal{M} includes finitely additive probabilities obtained as limits of truncated increasing/decreasing exponential functions.

In future work it would be useful to address the following issues: (1) the extension of this model to more cases for which $\mathcal{W} \neq \mathbb{R}^k$; (2) employing the model to represent prior ignorance in *generalized linear models*.

The extension to the case $\mathcal{W} \neq \mathbb{R}^k$ is important because, for many exponential families of interest $\mathcal{W} \neq \mathbb{R}^k$. Furthermore, with this extension we may employ our approach to model prior ignorance in *generalized linear models*. Generalized linear models comprise the traditional analyses such as *t*-tests, analysis of variance, linear regression and logistic regression. Consider for instance a normal linear regression model. In this case, the probabilistic relationship between observations and variables of interests are expressed via a multivariate normal density. The Bayesian approach to normal linear regression assumes that the prior information on the

variables of interest is expressed with a Gaussian-inverse Gamma distribution (or Gaussian-inverse Wishart distribution) which belongs to k -parameters exponential families. In the case of lack of prior information, prior ignorance is commonly modelled by selecting the parameters of the prior to make it noninformative. As pointed out in this paper, it is not possible to express prior ignorance with noninformative priors. In our opinion, a better approach is that of using a set of priors satisfying the property of near-ignorance. Notice that the Gaussian-inverse Gamma distribution belongs to k -parameters exponential families. Therefore, we can employ the multivariate models of near-ignorance for k -parameters exponential families to develop new robust regressors based on near-ignorance priors.

Acknowledgements

This work was partly supported by the Swiss NSF grants nos. 200021_146606 / 1 and 200020_137680 / 1. The authors would like to thank the anonymous referees for comments and constructive criticism that helped us to improve the presentation of the paper.

Appendix A. Invariance and Conjugate Improper Priors

In Bayesian statistics, in case of lack of prior information, it is often common to select the prior in order to satisfy some invariance property, see for instance [3, 24, 25]. Consider for instance the group of transformations $\mathcal{F} = \{f_a : a \in \mathbb{R}\}$ with $f_a(w) = w + a$, i.e., a shift of the parameter. From Definition 2.2 with $\underline{E} = \bar{E} = E$, it results that a prior $p(w)$ with support \mathcal{W} is \mathcal{F} -invariant if:

$$\int g(w+a)p(w)dw = \int g(w)p(w)dw \quad \forall a. \quad (\text{A1})$$

If $g = I_{\{A\}}$, with $A \subseteq \mathbb{R}^k$ bounded and measurable, the above inequality means that the set A and the shifted set $A + a$ should have the same probability for any value of a . By a change of variables and considering the case $g = I_{\{A\}}$, (A1) can be rewritten as:

$$\int_{A-a} p(w-a)dw = \int_A p(w)dw \quad \forall a. \quad (\text{A2})$$

For a bounded space \mathcal{W} , the only prior that satisfies (A2) is the uniform distribution. Conversely, for instance in case $\mathcal{W} = \mathbb{R}$, it is known that there is no countably additive probability measure that is translation invariant, even on the intervals. To see that, note that \mathbb{R} is the countable union of the intervals $[n, n+1)$ with $n \in \mathbb{Z}$ and these intervals must have the same probability, because of (A1), which must be zero by countable additivity.

However, we can define translation-invariant lower and upper expectations as limits of uniform distributions on finite intervals [6, Sec. 2.9.7], [26, Sec. 3]. Let \mathcal{K} be the linear space of Borel-measurable RVBFs on \mathbb{R} , and define:

$$\underline{E}[g] = \liminf_{r \rightarrow \infty} \frac{1}{2r} \int g(w)I_{[-r,r]}(w)dw, \quad \forall g \in \mathcal{K}, \quad (\text{A3})$$

where $r > 0$ and \underline{E} denotes the lower expectation of a RVBF g belonging to \mathcal{K} . This is a lower limit of uniform distributions on intervals centered at zero as their length $2r \rightarrow \infty$. By replacing infimum with supremum in the definition of \underline{E} , we can obtain the upper expectation \bar{E} . Observe that, since $\underline{E}[g] = -\bar{E}[-g]$, the upper expectation is completely determined by the

lower expectation and vice versa. In case $g = I_{\{A\}}$ for some subset A of \mathcal{W} , $\underline{E}[I_{\{A\}}]$ represents the lower probability of A and in this case the upper probability of A , i.e., $\overline{E}[I_{\{A\}}]$, is defined as $1 - \underline{E}[I_{\{A^c\}}]$, where A^c is the complementary set of A .

Lemma A.1. *Define the translations f_a by $f_a(w) = w + a$ for any $a \in \mathbb{R}$ and $g(f_a(w)) = g(w + a)$ for any $g \in \mathcal{K}$. The lower expectation model in (A3) satisfies $\underline{E}[g(f_a)] = \underline{E}[g]$ (translation invariance) for any $g \in \mathcal{K}$ and $a \in \mathbb{R}$.*

The proof is given in the first part of the proof of Lemma 4.1. Observe that \underline{E} defined in (A3) satisfies the following properties: (i) for $g = I_{\mathbb{R}}$, $\underline{E}[I_{\mathbb{R}}] = 1$; (ii) for $A = [0, \infty)$, $\underline{E}[I_A] = \frac{1}{2}$; (iii) for $g = I_A$, where A is a bounded subset of \mathbb{R} , it holds that $\underline{E}[I_{\{A+a\}}] = \underline{E}[I_A] = 0$ for any $a \in \mathbb{R}$ and, furthermore, by denoting with A^c the complementary set of A it holds that $\underline{E}[I_{A^c}] = \underline{E}[I_{\{A^c+a\}}] = 1$. Hence, since $\overline{E}[I_A] = 1 - \underline{E}[I_{A^c}]$, it follows that $\overline{E}[I_{\{A+a\}}] = \overline{E}[I_A] = 0$. As discussed in the introduction, to each functional \underline{E} (satisfying some regularity properties [6, Ch. 2]), it is possible to associate a closed and convex set \mathcal{M} of probabilities that generates the lower expectation $\underline{E}[g]$ for any g . This is the set of probabilities that dominate \underline{E} (equivalently, that are dominated by \overline{E}), i.e., the lower (upper) envelope of the expectation of g with respect to the closed and convex set \mathcal{M} of probabilities is equal to $\underline{E}[g]$ ($\overline{E}[g]$) for any g . In the case of \underline{E} defined in (A3), since $\overline{E}[I_A] = 0$ for any bounded subset A of \mathbb{R} , none of the probabilities dominated by \overline{E} can be countably additive. In other words the set \mathcal{M} in this case includes finitely additive probabilities.

In order to be able to interchange the operations of limit and integration even in case of unbounded functions g from \mathcal{W} to \mathbb{R} , in the following we take \mathcal{K} to be the linear space of Borel-measurable functions on \mathbb{R} satisfying $\int |g(w)| \exp(n(\hat{y}_n w - b(w))) dw < \infty$ for any $n > 0$ and $\hat{y}_n \in Cl(\mathcal{B}_0)$ ($Cl(\cdot)$ denotes the closure of the set). This is again a linear space.¹ Now consider the likelihood $p(n, \hat{y}_n | w) \propto \exp(n(\hat{y}_n w - b(w)))$; the lower posterior expectation of g can be obtained from (A3):

$$\underline{E}[g|n, \hat{y}_n] = \liminf_{r \rightarrow \infty} \frac{\int g(w) \exp(n(\hat{y}_n w - b(w))) I_{[-r, r]}(w) dw}{\int \exp(n(\hat{y}_n w - b(w))) I_{[-r, r]}(w) dw}. \quad (\text{A4})$$

Assuming $\mathcal{W} = \mathbb{R}$, since $\int |g(w)| \exp(n(\hat{y}_n w - b(w))) dw < \infty$ for any $n > 0$ and $\hat{y}_n \in Cl(\mathcal{B}_0)$, applying Lebesgue's dominated convergence theorem, one has:

$$\underline{E}[g|n, \hat{y}_n] = \frac{\int g(w) \exp(n(\hat{y}_n w - b(w))) dw}{\int \exp(n(\hat{y}_n w - b(w))) dw}, \quad (\text{A5})$$

which is equivalent to the posterior inference obtained from the limit kernel (improper uniform prior) $p(w) = 1 = \lim_{r \rightarrow \infty} I_{[-r, r]}(w)$. Observe that $\underline{E}[g|n, \hat{y}_n] = \overline{E}[g|n, \hat{y}_n]$. It is worth to notice that (A5) holds for any likelihood in the exponential families such that $\mathcal{W} = \mathbb{R}$. For instance, in the one-parameter Normal, Beta and Gamma case, $\mathcal{W} = \mathbb{R}$ and thus the limit kernel $p(w) = 1$ guarantees translation invariance. Note that $p(w) = 1$, transformed back to the original parameter space, becomes $p(x) = 1$ in the Normal case, $p(\theta) = \theta^{-1}(1 - \theta)^{-1}$ in the Beta case and $p(\lambda) = \lambda^{-1}$ in the Gamma case. These are the kernels that, when multiplied with the Normal, Binomial and Poisson likelihoods, give a Normal, Beta and, respectively, Gamma posterior density.

Note also that $p(w) = 1$ is not the only limit prior that guarantees translation invariance and preserves conjugacy, i.e., the posterior is the conjugate of the likelihood. Assume again $\mathcal{W} = \mathbb{R}$ and that g satisfies $\int |g(w)| \exp(n(\hat{y}_n w - b(w))) \exp(\ell w) dw < \infty$ for some $\ell > 0$ and any

¹Notice in fact that if g belongs to \mathcal{K} also λg does for any real λ . If g_1, g_2 belong to \mathcal{K} , also $g_1 + g_2$ belongs to \mathcal{K} , since $\int |g_1(w) + g_2(w)| \exp(n(\hat{y}_n w - b(w))) dw \leq \int |g_1(w)| \exp(n(\hat{y}_n w - b(w))) dw + \int |g_2(w)| \exp(n(\hat{y}_n w - b(w))) dw < \infty$.

$n > 0$, $\hat{y}_n \in Cl(\mathcal{Y}_0)$. Consider then the posterior limit:

$$\underline{E}[g|n, \hat{y}_n] = \liminf_{r \rightarrow \infty} \frac{\int g(w) \exp(n(\hat{y}_n w - b(w))) \exp(\ell w) I_{(-\infty, r]}(w) dw}{\int \exp(n(\hat{y}_n w - b(w))) \exp(\ell w) I_{(-\infty, r]}(w) dw}, \quad (\text{A6})$$

as before, since $\int |g(w)| \exp(n(\hat{y}_n w - b(w))) \exp(\ell w) dw < \infty$, by Lebesgue's dominated convergence theorem, the above limit is equal to:

$$\underline{E}[g|n, \hat{y}_n] = \frac{\int g(w) \exp\left(n\left(\frac{\ell + n\hat{y}_n}{n} w - b(w)\right)\right) dw}{\int \exp\left(n\left(\frac{\ell + n\hat{y}_n}{n} w - b(w)\right)\right) dw}. \quad (\text{A7})$$

Thus, (A6) is equivalent to considering the posterior expectation obtained from the kernel $p(w) = \exp(\ell w) = \lim_{r \rightarrow \infty} \exp(\ell w) I_{(-\infty, r]}(w)$. In analogy with (A3), we can define translation invariant lower and upper expectations as limits of exponential priors truncated on finite intervals:

$$\underline{E}[g] = \liminf_{r \rightarrow \infty} \frac{\ell}{\exp(\ell r)} \int g(w) \exp(\ell w) I_{(-\infty, r]}(w) dw, \quad (\text{A8})$$

for $r \in \mathbb{R}$, which satisfies translation invariance, $\underline{E}[g(f_a)] = \underline{E}[g]$, for any g in the set of Borel-measurable RVBFs \mathcal{K} .

Lemma A.2. *The lower expectation model in (A8) satisfies $\underline{E}[g(f_a)] = \underline{E}[g]$ (translation invariance) for any $g \in \mathcal{K}$ and $a \in \mathbb{R}$.*

The proof is given in the second part of the proof of Lemma 4.1. Observe that \underline{E} defined in (A8) satisfies the following properties: (i) for $g = I_{\mathbb{R}}$, $\underline{E}[I_{\mathbb{R}}] = 1$; (ii) for $A = [0, \infty)$, $\underline{E}[I_A] = 1$; (iii) for $g = I_A$, where A is a bounded subset of \mathbb{R} , it holds that $\overline{E}[I_{\{A+a\}}] = \overline{E}[I_A] = 0$ for any $a \in \mathbb{R}$. Thus, because of (iii), none of the probabilities dominated by \overline{E} can be countably additive and the set \mathcal{M} includes finitely additive probabilities. However comparing the lower probabilities in (ii) for (A3) and (A8), it can be noticed that the set of finitely additive probabilities induced by the two improper priors $p(w) = 1$ and, respectively, $p(w) = \exp(\ell w)$ are different. Observe that (A7) holds for any likelihood belonging to the exponential families provided that $\mathcal{W} = \mathbb{R}$. Note that $p(w) = \exp(\ell w)$, transformed back to the original parameter space, becomes $p(x) = \exp(\ell x)$ in the Normal case, $p(\theta) = \theta^{\ell-1} (1-\theta)^{-\ell-1}$ in the Beta case and $p(\lambda) = \lambda^{\ell-1}$ in the Gamma case. It can be noticed that for $\ell = 0$, $\exp(\ell w) = 1$, thus the improper uniform distribution can be seen as a particular case of $\exp(\ell w)$. We are interested in these exponential kernels $\exp(\ell w)$ since they preserve conjugacy with exponential families as shown in (A6). In summary, invariance properties can in general be satisfied by considering the lower expectation obtained as the limit of proper truncated priors.

References

- [1] H. Jeffreys. *Theory of Probability*. Clarendon Press, Oxford, 1983. 3rd edition.
- [2] J. O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer Series in Statistics, New York, 1985.
- [3] J. M. Bernardo and A. F. M. Smith. *Bayesian theory*. John Wiley & Sons, 1994.
- [4] Mervyn Stone. Review and analysis of some inconsistencies related to improper priors and finite additivity. *Studies in Logic and the Foundations of Mathematics*, 104:413–426, 1982.
- [5] Bruce M Hill and David Lane. Conglomerability and countable additivity. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 366–379, 1985.

- [6] P. Walley. *Statistical Reasoning with Imprecise Probabilities*. Chapman and Hall, New York, 1991.
- [7] J. O. Berger, E. Moreno, L. R. Pericchi, M. J. Bayarri, Bernardo, et al. An overview of robust Bayesian analysis. *Test*, 3(1):5–124, 1994.
- [8] P. J. Huber. The use of Choquet capacities in statistics. *Bull. Internat. Statist. Inst*, 45(4):181–191, 1973.
- [9] S. Sivaganesan and J.O. Berger. Ranges of posterior measures for priors with unimodal contaminations. *The Annals of Statistics*, 17(2):868–889, 1989.
- [10] L. DeRoberts and J.A. Hartigan. Bayesian inference using intervals of measures. *The Annals of Statistics*, 9(2):235–244, 1981.
- [11] L. R. Pericchi and P. Walley. Robust Bayesian credible intervals and prior ignorance. *International Statistical Review*, pages 1–23, 1991.
- [12] P. J. Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, pages 73–101, 1964.
- [13] A. Benavoli and M. Zaffalon. A model of prior ignorance for inferences in the one-parameter exponential family. *Journal of Statistical Planning and Inference*, 142(7):1960–1979, 2012.
- [14] P. Walley. A bounded derivative model for prior ignorance about a real-valued parameter. *Scandinavian Journal of Statistics*, 24(4):463–483, 1997.
- [15] F. P. A. Coolen and T. Augustin. A nonparametric predictive alternative to the imprecise dirichlet model: the case of a known number of categories. *International Journal of Approximate Reasoning*, 50(2):217–230, 2009.
- [16] Anirban DasGupta. *Asymptotic theory of statistics and probability*. Springer, 2008.
- [17] L. D. Brown. *Fundamentals of statistical exponential families: with applications in statistical decision theory*. IMS, Hayward, California, 1986.
- [18] O. Barndorff-Nielsen. *Information and exponential families in statistical inference*. Wiley, New York, 1978.
- [19] G. Letac. *Lectures on natural exponential families and their variance functions*. Number 50. Conselho Nacional de Desenvolvimento Científico e Tecnológico, Instituto de Matemática Pura e Aplicada, 1992.
- [20] P. Diaconis and D. Ylvisaker. Conjugate priors for exponential families. *The Annals of statistics*, 7(2):269–281, 1979.
- [21] P. Walley. Inferences from multinomial data: learning about a bag of marbles. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):3–57, 1996.
- [22] J.-M. Bernard. An introduction to the imprecise Dirichlet model for multinomial data. *International Journal of Approximate Reasoning*, 39(2–3):123–150, 2005.
- [23] A. Gelman, J.B. Carlin, H.S. Stern, and D.B. Rubin. *Bayesian data analysis*. CRC press, 2004.
- [24] J. Hartigan. Invariant prior distributions. *The Annals of Mathematical Statistics*, pages 836–845, 1964.
- [25] R. E. Kass and L. Wasserman. The selection of prior distributions by formal rules. *Journal of the American Statistical Association*, pages 1343–1370, 1996.
- [26] E. Regazzini. De finetti’s coherence and statistical inference. *The Annals of Statistics*, pages 845–864, 1987.