

Statistica Sinica Preprint No: SS-2015-0293R2

Title	Singular prior distributions in Bayesian D-optimal design for nonlinear models
Manuscript ID	SS-2015-0293R2
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202015.0293
Complete List of Authors	Timothy Waite
Corresponding Author	Timothy Waite
E-mail	timothy.waite@manchester.ac.uk
Notice: Accepted version subject to English editing.	

SINGULAR PRIOR DISTRIBUTIONS AND ILL-CONDITIONING IN BAYESIAN D -OPTIMAL DESIGN FOR SEVERAL NONLINEAR MODELS

Timothy W. Waite

University of Manchester

Abstract: For Bayesian D -optimal design, we define a *singular prior distribution* for the model parameters as a prior distribution such that the determinant of the Fisher information matrix has a prior geometric mean of zero for all designs. For such a prior distribution, the Bayesian D -optimality criterion fails to select a design. For the exponential decay model, we characterize singularity of the prior distribution in terms of the expectations of a few elementary transformations of the parameter. For a compartmental model and several multi-parameter generalized linear models, we establish sufficient conditions for singularity of a prior distribution. For the generalized linear models we also obtain sufficient conditions for non-singularity. In the existing literature, weakly informative prior distributions are commonly recommended as a default choice for inference in logistic regression. Here it is shown that some of the recommended prior distributions are singular, and hence should not be used for Bayesian D -optimal design. Additionally, methods are developed to derive and assess Bayesian D -efficient designs when numerical evaluation of the objective function fails due to ill-conditioning, as often occurs for heavy-tailed prior distributions. These numerical methods are illustrated for logistic regression.

Key words and phrases: Compartmental model, exponential decay model, generalized linear model, ill-conditioning, logistic regression.

1. Introduction

In recent years, much effort has been devoted to the development of D -optimal design methods for nonlinear problems; for example, nonlinear models (e.g. Yang and Stufken (2009, 2012); Yang (2010)), generalized linear models (Khuri, Mukherjee, Sinha and Ghosh (2006); Woods, Lewis, Eccleston and Russell (2006); Yang, Zhang and Huang (2011); Yang and Mandal (2015)), and linear models with mixed effects (Jones and Goos (2009)). In each of these areas, the

choice of a D -optimal design depends on the unknown vector of model parameters, $\theta \in \Theta \subseteq \mathbb{R}^p$.

One approach to choosing a design is to make a ‘best guess’ of the parameter values, and calculate a corresponding locally D -optimal design (Chernoff, 1953), i.e. $\xi_\theta^* \in \arg \max_{\xi \in \Xi} |M(\xi; \theta)|$, where $M(\xi; \theta)$ is the Fisher information matrix for design $\xi \in \Xi$, and Ξ is the set of all competing designs. However, the performance of a locally optimal design may be highly sensitive to misspecification of the value of θ . Then a Bayesian approach is often used to derive designs that are efficient for a variety of plausible values for θ . This approach requires the adoption of a prior distribution, \mathcal{P} , on the parameters, and maximization of the value of an objective function that quantifies the expected information contained in the experiment. Throughout, we assume that \mathcal{P} is a probability measure on the measure space (Θ, Σ) , with Σ the Borel σ -algebra over Θ . A widely used objective function is

$$\phi(\xi; \mathcal{P}) = \int_{\Theta} \log |M(\xi; \theta)| d\mathcal{P}(\theta), \quad (1.1)$$

see, for example, Chaloner and Larntz (1989) and Gotwalt, Jones and Steinberg (2009). We adopt the measure-theoretic formulation of integration, under which the notation $\int_{\Theta} g(\theta) d\mathcal{P}(\theta) = \infty$ is standard when $g : \Theta \rightarrow \mathbb{R}$ is a non-negative Σ -measurable function (Capinski and Kopp (2004), pp. 77-8). When $g : \Theta \rightarrow \mathbb{R}$ is a general Σ -measurable function, it is said that $\int_{\Theta} g(\theta) d\mathcal{P}(\theta) = -\infty$ if and only if $\int_{\Theta} g^+(\theta) d\mathcal{P}(\theta) < \infty$ and $\int_{\Theta} g^-(\theta) d\mathcal{P}(\theta) = \infty$, where $g^+(\theta) = \max\{0, g(\theta)\}$ and $g^-(\theta) = \max\{0, -g(\theta)\}$.

A design that maximizes (1.1) is said to be *(pseudo-)Bayesian D -optimal*, and may be used whether or not a Bayesian analysis will be performed (e.g. Woods, Lewis, Eccleston and Russell (2006)). Maximization of (1.1) is equivalent to maximization of an asymptotic approximation to the Shannon information gain from prior to posterior (Chaloner and Verdinelli (1995)).

In nonlinear problems, a *singular parameter vector* is a θ such that $M(\xi; \theta)$ has determinant zero for *any* design $\xi \in \Xi$. For such θ , it is difficult to estimate the parameters no matter which design is used, often because of a lack of model identifiability (see Section 2.3). In this situation, the local D -optimality criterion fails to select a design. The analogue of a singular parameter vector for Bayesian D -optimality is defined through:

- (a) Given $\xi \in \Xi$ and a prior distribution, \mathcal{P} , we say that ξ is a *Bayesian singular design with respect to \mathcal{P}* if $\phi(\xi; \mathcal{P}) = -\infty$.
- (b) Given a prior distribution, \mathcal{P} , we say that \mathcal{P} is a *singular prior distribution* if *all* $\xi \in \Xi$ are Bayesian singular with respect to \mathcal{P} , or equivalently if the geometric mean of $|M(\xi; \theta)|$ under \mathcal{P} is zero for all $\xi \in \Xi$. Above, the geometric mean of a non-negative random variable X is defined as $E^{\mathcal{G}}(X) = \exp[E\{\log(X)\}]$, with $E^{\mathcal{G}}(X) = 0$ if $E \log(X) = -\infty$ (Feng et al. (2017)).

For a singular prior distribution \mathcal{P} , Bayesian D -optimality cannot be used to select a design, since all designs have the same objective function value $\phi(\xi; \mathcal{P}) = -\infty$. In many models, such as the exponential decay model and logistic regression, it is straightforward to detect singular parameter vectors, θ , by inspection of the information matrix. However, as shown below, it is more difficult to detect whether \mathcal{P} is a singular prior distribution, except in the case of point priors.

A different, but related, problem is the presence of ill-conditioned information matrices in a quadrature approximation to (1.1). For several models, this is likely to occur for a heavy-tailed prior distribution \mathcal{P} , even if \mathcal{P} is theoretically non-singular. Such ill-conditioning causes failure of numerical selection of Bayesian D -optimal designs.

In this paper, we clarify and extend the set of prior distributions for which Bayesian D -optimal design is feasible for three important classes of models. In Sections 2.1, 2.2, and 2.3, respectively, we give examples of singular prior distributions for the one-factor exponential decay model, a three-parameter compartmental model, and several multi-factor generalized linear models. In Section 2.3 the default weakly informative prior proposed for logistic regression by Gelman, Jakulin, Pittau and Su (2008) is shown to be singular. For the exponential and generalized linear models, sufficient conditions for a prior distribution to be non-singular are established. These conditions are easily checked to determine if the Bayesian D -optimality criterion can be used to select designs under \mathcal{P} . In Section 3, novel methods are developed that enable the selection of highly Bayesian D -efficient designs for logistic regression when the quadrature approximation to (1.1) is ill-conditioned, thereby facilitating design for heavy-tailed prior distributions. Finally, in Section 4 we discuss alternative approaches to finding efficient

designs when \mathcal{P} is a singular prior distribution.

2. Singularity of prior distributions for some standard models

2.1. Exponential decay model

We derive necessary and sufficient conditions for a prior distribution to be singular for the exponential decay model which is used, for example, to model the concentration of a chemical compound over time. This model is commonly used as a simple illustrative example of a nonlinear model in the optimal design of experiments literature, e.g. Dette and Neugebauer (1997); Atkinson (2003); the results here help to develop our intuition. Here, two parameterizations are considered: by rate, $\beta > 0$, and by ‘lifetime’, $\theta = 1/\beta > 0$. For the former, the model for the response, y , in terms of explanatory variable, $x > 0$, is

$$y_i = e^{-\beta x_i} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2),$$

where $i = 1, \dots, n$, $x_i \geq 0$, and $\sigma > 0$.

Assume that $\Xi = \mathcal{X}^n$, where $\mathcal{X} = [0, \infty)$. Then design $\xi = (x_1, \dots, x_n) \in \Xi$ has information matrix

$$M_\beta(\xi; \beta) = \sum_{i=1}^n x_i^2 e^{-2\beta x_i}.$$

Suppose that at least one $x_i > 0$ and let $S_{xx} = \sum_{i=1}^n x_i^2$. Then

$$-2\beta \max_{i=1, \dots, n} \{x_i\} \leq \log |M_\beta(\xi; \beta)| - \log S_{xx} \leq -2\beta \min_{i: x_i > 0} \{x_i\}. \quad (2.1)$$

By taking expectations, the following result is obtained.

Proposition 1. *Suppose that at least one $x_i > 0$. Then, for the β -parameterization, $\phi(\xi; \mathcal{P}) > -\infty$ if and only if $E_{\mathcal{P}}(\beta) < \infty$.*

Here the prior, \mathcal{P} , is non-singular provided the rate parameter has finite expectation, but \mathcal{P} can be singular if the distribution of β is heavy-tailed with infinite mean, e.g. if β is half-Cauchy (cf. Polson and Scott (2012)).

For the θ -parameterization, a change-of-variable argument shows that

$$\log |M_\theta(\xi; \theta)| = \log |M_\beta(\xi; \beta)| - 4 \log \theta. \quad (2.2)$$

This enables derivation of the following result; for proof see the supplementary material.

Proposition 2. *For the θ -parameterization, the prior distribution \mathcal{P} is singular if and only if either $E_{\mathcal{P}}(1/\theta) = \infty$ or $E_{\mathcal{P}}(\log \theta) = \infty$.*

In the context of designs that maximize $\phi(\xi; \mathcal{P})$ for nonlinear models, Chaloner and Verdinelli (1995) refer to potential ‘technical problems using prior distributions with unbounded support where [...] $M(\xi; \theta)$ may be arbitrarily close to being singular’. Corollary 1 below shows that, even with bounded support, seemingly innocuous prior distributions can cause Bayesian D -optimality to fail as a design selection criterion.

Corollary 1. *For the θ -parameterization, the prior distribution $\mathcal{P} = U(0, a)$, $a > 0$, is singular.*

Note that under the prior for θ in Corollary 1, the corresponding implied prior distribution for β has a proper density, $p(\beta) = 1/(a\beta^2)$ for $\beta \geq 1/a$. However, this implied distribution for β has unbounded support and is heavy tailed, such that $E(\beta) = \infty$. In other words, the implied a priori expectation is that the decay is very rapid.

2.2. Compartmental model

In this section, we derive sufficient conditions for a prior distribution to be singular for the following three-parameter compartmental model:

$$y_i = \theta_3 \{e^{-\theta_1 x_i} - e^{-\theta_2 x_i}\} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2), \quad (2.3)$$

where $x_i \geq 0$, $i = 1, \dots, n$, $\theta_2 > \theta_1 > 0$, $\theta_3 > 0$ and $\sigma > 0$. Here, $\Xi = [0, \infty)^n$. As with the exponential model, often the response y_i is a concentration of a compound in a system, and the x_i are the observation times. For example, Atkinson et al. (1993) consider a theophylline kinetics experiment on horses, finding optimal sampling times for model (2.3) under several different (pseudo-)Bayesian criteria.

The information matrix for the i th time point is

$$M(x_i; \theta) = \begin{pmatrix} x_i^2 \theta_3^2 e^{-2\theta_1 x_i} & -x_i^2 \theta_3^2 e^{-(\theta_1 + \theta_2)x_i} & -f_i x_i e^{-\theta_1 x_i} \\ -x_i^2 \theta_3^2 e^{-(\theta_1 + \theta_2)x_i} & x_i^2 \theta_3^2 e^{-2\theta_2 x_i} & f_i x_i e^{-\theta_2 x_i} \\ -f_i x_i e^{-\theta_1 x_i} & f_i x_i e^{-\theta_2 x_i} & f_i^2 / \theta_3^2 \end{pmatrix},$$

where $f_i = \theta_3\{e^{-\theta_1 x_i} - e^{-\theta_2 x_i}\}$. We have $|M(\xi; \theta)| = 0$ when (i) $\theta_1 = \theta_2$ or (ii) $\theta_3 = 0$, and $|M(\xi; \theta)| \rightarrow 0$ when (iii) $\theta_1 \rightarrow \infty$. Physically, conditions (i) and (iii) correspond to situations where the flow rates in and out of the compartment are either exactly balanced, or both very rapid. Each of the potentially very different parameter scenarios in (i)–(iii) results in a similar response profile, in which the concentration is close to zero throughout the duration of the experiment. Thus, if such a profile is observed, it is difficult to ascertain which values of the parameters generated the data.

From the above, it is clear that for \mathcal{P} to be a non-singular prior distribution its probability density must not be too highly concentrated near regions where $\theta_2 = \theta_1$ or $\theta_3 = 0$, nor can the prior for θ_1 be too heavy-tailed. This is formalized by Proposition 3, for which the following two lemmas are required; proofs are given in the supplementary material. Let $\delta = \theta_2 - \theta_1 > 0$.

Lemma 1. *We have the following bounds on $\log |M(\xi; \theta)|$,*

$$-6\theta_1 x_{\max} \leq \log |M(\xi; \theta)| - 4 \log \theta_3 - \log |\tilde{M}_{\delta,1}| \leq -6\theta_1 x_{\min},$$

where $\tilde{M}_{\delta,1}$ is $\tilde{M}_{\delta,\theta_3}$ evaluated at $\theta_3 = 1$, and

$$\tilde{M}_{\delta,\theta_3} = \sum_{i=1}^n \tilde{M}_{\delta,\theta_3}^{(i)}, \quad x_{\min} = \min_{i: x_i > 0} \{x_i\}, \quad x_{\max} = \max_{i=1, \dots, n} \{x_i\},$$

$$\tilde{M}_{\delta,\theta_3}^{(i)} = \begin{pmatrix} x_i^2 \theta_3^2 & -x_i^2 \theta_3^2 e^{-\delta x_i} & -x_i \theta_3 (1 - e^{-\delta x_i}) \\ -x_i^2 \theta_3^2 e^{-\delta x_i} & x_i^2 \theta_3^2 e^{-2\delta x_i} & x_i \theta_3 e^{-\delta x_i} (1 - e^{-\delta x_i}) \\ -x_i \theta_3 (1 - e^{-\delta x_i}) & x_i \theta_3 e^{-\delta x_i} (1 - e^{-\delta x_i}) & (1 - e^{-\delta x_i})^2 \end{pmatrix}.$$

Lemma 2. *If $\int_{\delta < 1} \log \delta d\mathcal{P}(\theta) = -\infty$, then $E_{\mathcal{P}}(\log |\tilde{M}_{\delta,1}|) = -\infty$.*

Proposition 3. *Suppose $\int_{\theta_3 > 1} \log \theta_3 d\mathcal{P}(\theta) < \infty$. For model (2.3), the prior parameter distribution \mathcal{P} is singular if $E_{\mathcal{P}}(\theta_1) = \infty$, $\int_{\theta_3 < 1} \log \theta_3 d\mathcal{P}(\theta) = -\infty$, or $\int_{\delta < 1} \log \delta d\mathcal{P}(\theta) = -\infty$.*

Heavy-tailed priors such as the half-Cauchy are increasingly recommended as weakly informative priors in various models (Gelman et al. (2008); Polson and Scott (2012)). For model (2.3), \mathcal{P} is singular if θ_1 is half-Cauchy distributed, although for physiological compartmental models more specific prior information is often used (Gelman et al. (1996)).

2.3. Generalized linear models

Suppose there are n design points $x_i = (x_{i1}, \dots, x_{iq})^T \in \mathcal{X}$, with responses y_i , $i = 1, \dots, n$. We assume a generalized linear model (GLM; McCullagh and Nelder (1989)), thus y_i has an exponential family distribution with mean $\mu_i = \mu(x_i; \beta)$ and variance $\gamma v(\mu_i)$, where μ satisfies

$$h[\mu(x; \beta)] = \eta(x; \beta) = f^T(x)\beta, \quad (2.4)$$

with h the link function, γ a dispersion parameter, v the variance function, and $\eta_i = \eta(x_i; \beta)$ the linear predictor. For binomial and Poisson responses, $\gamma = 1$ with variance function $v(\mu) = \mu(1-\mu)$ and $v(\mu) = \mu$, respectively. Above, $f(x) = (f_0(x), \dots, f_{p-1}(x))^T$ contains regression functions $f_j : \mathcal{X} \rightarrow \mathbb{R}$, $j = 0, \dots, p-1$, and $\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})^T \in \Theta$ is a vector of p regression parameters. We let $\mathcal{X} = [-1, 1]^q$ and $\Xi = \mathcal{X}^n$.

For design $\xi = (x_1, \dots, x_n)$ and model (2.4)

$$M(\xi; \beta) = \sum_{i=1}^n w_i f(x_i) f^T(x_i)$$

$$w(\eta) = \frac{1}{\gamma v(\mu)} \left(\frac{\partial \mu}{\partial \eta} \right)^2, \quad (2.5)$$

with $w_i = w(\eta_i)$, $i = 1, \dots, n$ (e.g. Khuri et al. (2006), Atkinson and Woods (2015), Yang and Mandal (2015)).

The following lemmas are important first steps towards the derivation of results on singular prior distributions. Lemma 3 also facilitates the development of numerical methods to overcome ill-conditioning in Section 3. The proofs are straightforward; the details are omitted. Let F be the model matrix with rows $f^T(x_i)$, noting that $\sum_{i=1}^n f(x_i) f^T(x_i) = F^T F$ is the information matrix of ξ under a linear model with regressors specified by f . The inequality below is with respect to the Loewner partial ordering on real symmetric matrices, in which $M_1 \preceq M_2$ if and only if $M_2 - M_1$ is non-negative definite (for example, Pukelsheim (1993, p.11)).

Lemma 3. *For a generalized linear model, the information matrix satisfies*

$$\min_{i=1, \dots, n} \{w_i\} F^T F \preceq M(\xi; \beta) \preceq \max_{i=1, \dots, n} \{w_i\} F^T F.$$

Thus, since the log-determinant respects the Loewner ordering,

$$p \log \min_i \{w_i\} + \log |F^T F| \leq \log |M(\xi; \beta)| \leq p \log \max_i \{w_i\} + \log |F^T F|.$$

Lemma 4. *Suppose that ξ is non-singular for the linear model with regressors given by f , that is $|F^T F| > 0$. Then we have the following:*

- (i) *If $E_{\mathcal{P}}\{\log \min_i w_i\} > -\infty$, then $\phi(\xi; \mathcal{P}) > -\infty$, i.e. ξ is Bayesian non-singular with respect to \mathcal{P} under the GLM.*
- (ii) *If $E_{\mathcal{P}}\{\log \max_i w_i\} = -\infty$, then $\phi(\xi; \mathcal{P}) = -\infty$, i.e. ξ is Bayesian singular with respect to \mathcal{P} under the GLM.*

Lemma 4 can often be used to identify clear conditions on the prior distribution that lead to singularity (or non-singularity). However, to do so it is necessary to analyse the tail behaviour of the GLM weight function, $w(\eta)$, as $|\eta| \rightarrow \infty$ in order to establish whether (i) or (ii) above holds. Thus, the results depend upon which link function is chosen. In the remainder of Section 2.3, results are given for logistic, probit and Poisson regression.

2.3.1. Logistic regression

For logistic regression, $y_i | \beta \sim \text{Bernoulli}(\pi_i)$, where $\pi_i = \Pr(y_i = 1 | \beta) = \mu(x_i; \beta)$. The link function is the logit, $h(\pi) = \log\{\pi/(1 - \pi)\}$, and

$$\begin{aligned} w(\eta) &= \exp(-|\eta|) \text{expit}(|\eta|)^2 \\ &\sim \exp(-|\eta|) \text{ as } |\eta| \rightarrow \infty. \end{aligned} \tag{2.6}$$

Above, $\text{expit}(\eta) = 1/\{1 + e^{-\eta}\}$. Lemma 4 is now used to establish sufficient conditions for the prior distribution to be non-singular for logistic regression.

Theorem 1. *Suppose that \mathcal{P} is such that $E_{\mathcal{P}}(|\beta_j|) < \infty$, for $j = 0, \dots, p-1$. If ξ is non-singular for the linear model with regressors given by f , that is $|F^T F| > 0$, then $\phi(\xi; \mathcal{P}) > -\infty$, i.e. ξ is also Bayesian non-singular with respect to \mathcal{P} for the logistic model.*

Note that there is no requirement for \mathcal{P} to have bounded support. In particular, this result provides theoretical reassurance that Bayesian D -optimality can

be used to select a design with a normal or log-normal prior on the parameters. There is also no requirement for the parameters to be independent a priori. For example, the result applies to a normal-mixture hierarchical variable selection prior distribution (Chipman et al. (1997)).

Other important prior distributions do not satisfy the conditions of Theorem 1; for example that proposed by Gelman, Jakulin, Pittau and Su (2008) which we denote by \mathcal{P}_G . These authors recommended rescaling before fitting the model. For observational studies, each explanatory variable is transformed to have mean zero and a standard deviation of $1/2$. This ensures that the method reflects the widely-held default prior belief that higher order interactions are likely to make a smaller contribution to the linear predictor. The combination of \mathcal{P}_G and this scaling was shown to give improved predictive performance relative to both maximum likelihood and penalized logistic regression. An analogue of the above method for designed experiments is to combine \mathcal{P}_G with a standardization of the design variables to have range $[-1/2, 1/2]$. This achieves a similar penalization of higher order interactions.

It is possible to obtain a partial inverse result to Theorem 1.

Proposition 4. *Given $j \in \{0, \dots, p-1\}$, suppose that:*

- (i) \mathcal{P} is such that $\Pr(\beta_j > 1) > 0$
- (ii) \mathcal{P} is such that, for all $\delta > 0$,

$$\Pr(|\beta_k| < \delta \text{ for all } k \neq j \mid \beta_j > 1) > 0$$

- (iii) \mathcal{P} is such that, for all $\delta > 0$,

$$E_{\mathcal{P}}[\beta_j \mid \beta_j > 1, |\beta_k| < \delta, \text{ for all } k \neq j] = \infty$$

- (iv) ξ is such that $\min_{i=1, \dots, n} |f_j(x_i)| > 0$.

Then ξ is Bayesian singular with respect to \mathcal{P} , i.e. $\phi(\xi; \mathcal{P}) = -\infty$.

A more intuitive understanding of the reason that the above conditions lead to a singular prior distribution can be obtained by considering locally optimal design. There, we have that $|M(\xi; \beta)| \approx 0$ if the responses are close to deterministic, i.e. if for all design points the success probability $\Pr(y_i = 1 \mid \beta)$ is close to either 0 or 1. In that case, there is also a high probability of separation (Albert

and Anderson (1984)) and thus non-existence of maximum likelihood estimates. For Bayesian design, a heavy-tailed prior satisfying the conditions of Proposition 4 leads to similarly extreme values of the success probability, which is now a random variable owing to dependence on β . Specifically, the implied distribution on $\Pr(y_i = 1 | \beta)$ has high concentration near either 0 or 1, in the following sense:

Proposition 5. *Under the conditions in Proposition 4, there exists an event $\mathcal{E} \subseteq \Theta$, with $\Pr(\mathcal{E}) > 0$, conditional upon which either $\Pr(y_i = 1 | \beta)$ or $1 - \Pr(y_i = 1 | \beta)$ has prior geometric mean zero, according to whether $f_j(x_i) < 0$ or $f_j(x_i) > 0$ respectively.*

The proofs of Propositions 4 and 5 both rest on the identification of a region, \mathcal{E} , of parameter space where the linear predictor η_i can be approximated by the contribution, $\beta_j f_j(x_i)$, from the j th predictor.

The Gelman prior distribution, \mathcal{P}_G , places independent standard Cauchy distributions on $(1/10)\beta_0, (2/5)\beta_1, \dots, (2/5)\beta_{p-1}$. Thus, the prior distributions for the regression coefficients are heavy-tailed, with undefined prior mean. The parameters are expected a priori to have large magnitude, i.e. $E|\beta_k| = \infty$, $k = 0, \dots, p - 1$. For a model with an intercept term, $f_0(x) = 1$, and Proposition 4 may be applied with $j = 0$; conditions (ii) and (iii) follow since β_0 is both heavy-tailed and independent of the other parameters, hence:

Corollary 2. *For a logistic model with an intercept term, the prior distribution \mathcal{P}_G is singular.*

Often prior independence of parameters is not a reasonable assumption. For example, Chipman et al. (1997) define a hierarchical variable selection prior in which the probability of an interaction term being active is dependent on whether the parent terms are active, thereby satisfying the weak heredity principle. Proposition 4 can be used to show that, for logistic regression, a prior with this hierarchical structure is singular if the prior distribution of the intercept parameter is a mixture of two scaled zero-mode Cauchy distributions rather than a mixture of two scaled zero-mean normal distributions. In this case, the intercept is again both heavy-tailed and (typically) independent of the other parameters.

For logistic models with a single controllable variable, scalar $x \in \mathcal{X}$, $\mathcal{X} = \mathbb{R}$, Bayesian D -optimal design has also been studied for a different parameterization

(for example, Chaloner and Larntz (1989)):

$$h(\pi_i) = \beta_1(x_i - \mu), \quad (2.7)$$

which can be obtained from (2.4) via $\beta_0 = -\beta_1\mu$. When $\beta_1 = 0$, μ is not identifiable and $|M_\theta(\xi; \theta)| = 0$ for all $\xi \in \Xi$, with $\theta = (\mu, \beta_1)^\top$, $\Xi = \mathcal{X}^n$. The following result, which is straightforward to prove using Theorem 1, provides sufficient conditions for a prior distribution to be non-singular for this form of the model.

Proposition 6. *For the (μ, β_1) -parameterization in (2.7), if (i) $E_{\mathcal{P}}(|\mu\beta_1|) < \infty$, (ii) $E_{\mathcal{P}}(|\beta_1|) < \infty$ and (iii) $E_{\mathcal{P}}(\log|\beta_1|) > -\infty$, then any design with two or more support points is Bayesian non-singular with respect to \mathcal{P} . Hence (i)–(iii) are sufficient for \mathcal{P} to be non-singular. In this case, ξ is Bayesian D-optimal for (β_0, β_1) if and only if it is Bayesian D-optimal for (μ, β_1) .*

2.3.2. Probit regression

For probit regression, $y_i | \beta \sim \text{Bernoulli}(\pi_i)$, $\pi_i = \mu(x_i; \beta)$, with link $h(\pi) = \Phi^{-1}(\pi)$, where Φ is the standard normal c.d.f.. Here,

$$w(\eta) = \frac{\varphi(\eta)^2}{\Phi(\eta)(1 - \Phi(\eta))},$$

where $\varphi(\eta) = \frac{1}{\sqrt{2\pi}}e^{-\eta^2/2}$ is the standard normal p.d.f.. The following asymptotic approximation holds (e.g. Abramowitz and Stegun (1964), p.298)

$$1 - \Phi(\eta) \sim \frac{1}{\eta\sqrt{2\pi}}e^{-\eta^2/2} \quad \text{as } \eta \rightarrow \infty.$$

Also, as $\eta \rightarrow \infty$, $\Phi(\eta) \rightarrow 1$, and so by symmetry of $w(\eta)$

$$w(\eta) \sim \frac{1}{\sqrt{2\pi}}|\eta|e^{-\eta^2/2} \quad \text{as } |\eta| \rightarrow \infty. \quad (2.8)$$

This asymptotic approximation can be used with Lemma 4 to obtain analogues of the results for logistic regression, with different conditions on the prior distribution.

Theorem 2. *If $E_{\mathcal{P}}|\beta_k\beta_l| < \infty$, for $k, l = 0, 1, \dots, p - 1$, then \mathcal{P} is non-singular for the probit regression model.*

Proposition 7. *Given $j \in \{0, \dots, p - 1\}$, suppose that:*

(i) \mathcal{P} is such that $\Pr(\beta_j > 1) > 0$ and, for all $\delta > 0$,

$$\Pr(|\beta_k| < \delta \text{ for all } k \neq j \mid \beta_j > 1) > 0$$

(ii) \mathcal{P} is such that, for all $\delta > 0$,

$$E_{\mathcal{P}}[|\beta_j|^2 \mid \beta_j > 1, |\beta_k| < \delta \text{ for all } k \neq j] = \infty$$

(iii) ξ is such that $\min_i |f_j(x_i)| > 0$.

Then, for the probit link the design ξ is Bayesian singular with respect to \mathcal{P} , i.e. $E_{\mathcal{P}} \log |M(\xi; \beta)| = -\infty$.

Corollary 3. For a probit model with an intercept term, the prior distribution \mathcal{P}_G is singular.

Again, a heavy-tailed prior on the intercept parameter results in the prior being singular for Bayesian D -optimality. The intuitive interpretation is similar to that for the logistic model. Note that \mathcal{P}_G would remain singular even if it were made somewhat less heavy-tailed, for example by replacing the Cauchy prior on β_0 with a $t(2)$ prior. In this case, condition (ii) above will still hold because β_0 has infinite variance.

2.3.3. Poisson regression

Consider the model $y_i \mid \beta \sim \text{Poisson}(\lambda_i)$, with $\mu_i = \lambda_i$ and $h(\mu) = \log \mu$. Optimal designs for this model were considered by Russell et al. (2009) and McGree and Eccleston (2012). Here, $w(\eta) = \exp(\eta)$ and we have the following results.

Theorem 3. For the Poisson regression model with log link, if $E_{\mathcal{P}}|\beta_k| < \infty$, $k = 0, \dots, p-1$, and $|F^T F| > 0$ then $E_{\mathcal{P}} \log |M(\xi; \beta)| > -\infty$. Hence, if the first moments for β_k are finite then \mathcal{P} is non-singular.

Proposition 8. Given $j \in \{0, \dots, p-1\}$, suppose that:

(i) \mathcal{P} is such that β_j is supported on $(-\infty, 0)$

(ii) \mathcal{P} is such that $\Pr(\beta_j < -1) > 0$ and, for all $\delta > 0$,

$$\Pr(|\beta_k| < \delta \text{ for all } k \neq j \mid \beta_j < -1) > 0$$

(iii) \mathcal{P} is such that for all $\delta > 0$,

$$E_{\mathcal{P}}[\beta_j | \beta_j < -1, |\beta_k| < \delta \text{ for all } k \neq j] = -\infty$$

(iv) \mathcal{P} is such that $E_{\mathcal{P}}|\beta_k| < \infty$, $k \neq j$

(v) ξ is such that $f_j(x_i) > 0$ for $i = 1, \dots, n$.

Then $E_{\mathcal{P}} \log |M(\xi; \beta)| = -\infty$, i.e. ξ is singular for the Poisson model with log link under Bayesian D -optimality.

Corollary 4. For a Poisson model with log link containing an intercept, i.e. $f_0(x) = 1$, if β_0 has a negative half-Cauchy prior independently of β_k , $k = 1, \dots, p-1$, with $E_{\mathcal{P}}|\beta_k| < \infty$, then \mathcal{P} is singular.

Here, a heavy-tailed negative intercept parameter can result in a singular prior. Intuitively, it is clear that large negative values of β_0 will lead to experiments where most of the responses are zero, leading to difficulties obtaining precise estimates of β_0 and the other parameters.

3. Numerical methods to overcome ill-conditioning

3.1. Objective function approximation

In performing a numerical search for Bayesian D -optimal designs it is necessary to approximate the objective function, usually via a weighted sum,

$$\phi(\xi; \mathcal{P}) \approx \phi(\xi; \mathcal{Q}) = \sum_{l=1}^{N_{\mathcal{Q}}} v_l \log |M(\xi; \beta^{(l)})|, \quad (3.1)$$

over a weighted sample,

$$\mathcal{Q} = \left\{ \begin{array}{ccc} \beta^{(1)} & \dots & \beta^{(N_{\mathcal{Q}})} \\ v_1 & \dots & v_{N_{\mathcal{Q}}} \end{array} \right\},$$

of parameter vectors, $\beta^{(l)}$, $l = 1, \dots, N_{\mathcal{Q}}$, with corresponding integration weights v_l , satisfying $\sum_{l=1}^{N_{\mathcal{Q}}} v_l = 1$.

The sample \mathcal{Q} may be obtained, for example, by space-filling criteria, as used by Woods, Lewis, Eccleston and Russell (2006), Latin hypercube sampling, or a quadrature scheme, such as that applied by Gotwalt, Jones and Steinberg (2009). Quadrature methods, and in particular the Gotwalt method, can often yield highly accurate approximations.

A problem with approximation (3.1) is that for multi-parameter models numerical evaluation of $\phi(\xi; \mathcal{Q})$ can fail due to the presence of ill-conditioned matrices $M(\xi; \beta^{(l)})$, whose determinant will be estimated numerically as zero. Note this can occur even for non-singular \mathcal{P} ; for singular \mathcal{P} there is little point in evaluating $\phi(\xi; \mathcal{Q})$ since $\phi(\xi; \mathcal{P}) = -\infty$. When numerical evaluation of $\phi(\xi; \mathcal{Q})$ fails for all $\xi \in \Xi$, we say that \mathcal{Q} is an *ill-conditioned quadrature scheme*. In principle, \mathcal{Q} can be ill-conditioned for any prior distribution. However, for the models considered here, such as logistic regression, ill-conditioning of \mathcal{Q} is more likely if the underlying prior distribution is heavy-tailed. In that case, there is high probability of large β , and so also of $M(\xi; \beta)$ being ill-conditioned. For any prior, even without heavy tails, other circumstances that may lead to ill-conditioning of \mathcal{Q} include: (i) use of a quadrature scheme, such as the Gotwalt method, which oversamples the tails of \mathcal{P} ; and (ii) use of a large number of quadrature points. In most integration problems, an increased number of quadrature points leads to improved approximation of the integral; paradoxically, in Bayesian D -optimal design this may cause numerical evaluation to fail due to ill-conditioning.

3.2. Objective function bounds for logistic regression

For some important models, it is possible to obtain bounds that allow approximation of $\phi(\xi; \mathcal{Q})$ when \mathcal{Q} is ill-conditioned, as often occurs for heavy-tailed priors. These bounds may be applied to enable straightforward selection of Bayesian D -efficient designs for such priors (see Section 3.3). Here we focus on the case of logistic regression, but a similar approach can be used for the compartmental model (using Lemma 1), and other GLMs. From Lemma 3 and (2.6), we see that $\phi(\xi; \beta) = \log |M(\xi; \beta)|$ lies in $[\phi_L(\xi; \beta), \phi_U(\xi; \beta)]$, where

$$\begin{aligned}\phi_L(\xi; \beta) &= \log |F^T F| + p \min_{i=1, \dots, n} \{-|\eta_i| + 2 \log \expit |\eta_i|\} \\ \phi_U(\xi; \beta) &= \log |F^T F| + p \max_{i=1, \dots, n} \{-|\eta_i| + 2 \log \expit |\eta_i|\}.\end{aligned}$$

Let \mathcal{S} be the set of $l \in \{1, \dots, N_{\mathcal{Q}}\}$ for which $M(\xi; \beta^{(l)})$ is ill-conditioned, then:

$$\phi_L(\xi; \mathcal{Q}) \leq \phi(\xi; \mathcal{Q}) \leq \phi_U(\xi; \mathcal{Q}), \quad (3.2)$$

where

$$\begin{aligned}\phi_L(\xi; \mathcal{Q}) &= \sum_{l \in \{1, \dots, N_{\mathcal{Q}}\} \setminus \mathcal{S}} v_l \log |M(\xi; \beta^{(l)})| + \sum_{l \in \mathcal{S}} v_l \log |F^T F| \\ &\quad + \sum_{l \in \mathcal{S}} v_l p \min_{i=1, \dots, n} \{-|f^T(x_i)\beta^{(l)}| + 2 \log \text{expit } |f^T(x_i)\beta^{(l)}|\} \\ \phi_U(\xi; \mathcal{Q}) &= \sum_{l \in \{1, \dots, N_{\mathcal{Q}}\} \setminus \mathcal{S}} v_l \log |M(\xi; \beta^{(l)})| + \sum_{l \in \mathcal{S}} v_l \log |F^T F| \\ &\quad + \sum_{l \in \mathcal{S}} v_l p \max_{i=1, \dots, n} \{-|f^T(x_i)\beta^{(l)}| + 2 \log \text{expit } |f^T(x_i)\beta^{(l)}|\}.\end{aligned}$$

The bounds $\phi_L(\xi; \mathcal{Q})$, $\phi_U(\xi; \mathcal{Q})$ are much better conditioned than $\phi(\xi; \mathcal{Q})$. The bounds for $\log |M(\xi; \beta^{(l)})|$, $l \in \mathcal{S}$, are often wide. However, as the corresponding v_l is often very small, we may nonetheless obtain from (3.2) a relatively narrow interval for $\phi(\xi; \mathcal{Q})$. Note that (3.2) specifies an interval that contains the approximation $\phi(\xi; \mathcal{Q})$, and not necessarily the value of $\phi(\xi; \mathcal{P})$.

In the remainder of Section 3, we use the following example to show how the bounds enable an extension of the set of prior distributions for which Bayesian D -efficient designs can be obtained in practice. We begin by illustrating the use of bounds for the objective function.

Example 1. Potato-packing experiment (Woods, Lewis, Eccleston and Russell (2006)). We use one of the authors' models, defined by

$$\begin{aligned}f(x) &= (1, x_1, x_2, x_3, x_1x_2, x_1x_3, x_2x_3)^T \\ \beta &= (\beta_0, \beta_1, \beta_2, \beta_3, \beta_{12}, \beta_{13}, \beta_{23})^T,\end{aligned}$$

where $q = 3$, $x = (x_1, x_2, x_3)^T$. We adopt a different prior distribution, namely $\log \beta_0 \sim N(-1, 2)$, $\beta_1 \sim N(2, 2)$, $\beta_2 \sim N(1, 2)$, $\beta_3 \sim N(-1, 2)$, and $\beta_{12}, \beta_{13}, \beta_{23} \sim N(0.5, 2)$ independently. Note that the log-normal prior for the intercept parameter is heavy-tailed. However, from Theorem 1, the above joint prior distribution is non-singular.

For a double-replicate of the 2^3 full factorial design, the value of $\phi(\xi; \mathcal{P})$ was approximated using the Gotwalt quadrature scheme, with 5 radial points and 4 random rotations. Direct numerical evaluation of $\phi(\xi; \mathcal{Q})$ failed, since \mathcal{S} was non-empty: it contained 39 parameter vectors. However, from (3.2), $\phi(\xi; \mathcal{Q}) \in [-6.85, -6.78]$.

3.3. Use of bounds in design optimization and assessment

The bounds from (3.2) may also be used within an optimization algorithm to help find Bayesian D -efficient designs. The *Bayesian D -efficiency* of ξ is

$$\text{Bayes-eff}(\xi; \mathcal{P}) = \exp\{[\phi(\xi; \mathcal{P}) - \phi(\xi_{\mathcal{P}}^*; \mathcal{P})]/p\} \times 100\%,$$

where $\xi_{\mathcal{P}}^* \in \arg \max_{\xi \in \Xi} \phi(\xi; \mathcal{P})$ is a Bayesian D -optimal design. Bayesian D -efficiencies near 100% indicate that ξ achieves a near-optimal trade-off in performance across the support of the prior distribution for β .

When \mathcal{Q} is well-conditioned, the Bayesian D -efficiency may be approximated by numerical search for a $\xi_{\mathcal{Q}}^* \in \arg \max_{\xi \in \Xi} \phi(\xi; \mathcal{Q})$ that maximizes the quadrature approximated objective function, and substitution of the design found into

$$\text{Bayes-eff}(\xi; \mathcal{Q}) = \exp\{[\phi(\xi; \mathcal{Q}) - \phi(\xi_{\mathcal{Q}}^*; \mathcal{Q})]/p\} \times 100\%.$$

However, if \mathcal{Q} is ill-conditioned, for example if \mathcal{P} is heavy-tailed, then this method fails since (i) $\phi(\xi; \mathcal{Q})$ cannot be evaluated directly, and (ii) $\xi_{\mathcal{Q}}^*$ cannot be found using a numerical search. We may nonetheless use numerical methods to find designs $\xi_{\mathcal{Q},L}^*$ and $\xi_{\mathcal{Q},U}^*$ that maximize the lower and upper bounds respectively, i.e. $\xi_{\mathcal{Q},L}^* \in \arg \max_{\xi \in \Xi} \phi_L(\xi; \mathcal{Q})$ and $\xi_{\mathcal{Q},U}^* \in \arg \max_{\xi \in \Xi} \phi_U(\xi; \mathcal{Q})$. Then a lower bound for the Bayesian efficiency of $\xi_{\mathcal{Q},L}^*$ can be approximated, via substitution of the designs found into

$$\text{Bayes-eff}(\xi_{\mathcal{Q},L}^*; \mathcal{Q}) \geq \exp\{[\phi_L(\xi_{\mathcal{Q},L}^*; \mathcal{Q}) - \phi_U(\xi_{\mathcal{Q},U}^*; \mathcal{Q})]/p\} \times 100\%. \quad (3.3)$$

To find exact designs that maximize the bounds, we use a continuous co-ordinate exchange algorithm similar to that of Gotwalt, Jones and Steinberg (2009).

Example 1 (continued). A co-ordinate exchange algorithm was used, with 100 random starts, to search for $\xi_{\mathcal{Q},L}^*$, $\xi_{\mathcal{Q},U}^*$ among exact designs with $n = 16$ runs. The quadrature scheme \mathcal{Q} was generated using the Gotwalt method, with 3 radial points and one random rotation, yielding a total of 217 support points for \mathcal{Q} . The design $\xi_{\mathcal{Q},L}^*$, given in Table 3.1, is very similar to $\xi_{\mathcal{Q},U}^*$: to 2 d.p. the two are identical. For this \mathcal{Q} , the objective function $\phi(\xi; \mathcal{Q})$ cannot be computed exactly due to ill-conditioning. Thus, given an alternative design ξ' , e.g. a 16-run combination of $\xi_{\mathcal{Q},L}^*$ and $\xi_{\mathcal{Q},U}^*$, it is not possible to evaluate whether ξ' has higher

Run	x_1	x_2	x_3	Run	x_1	x_2	x_3
1	0.456	1.000	1.000	9	-1.000	-1.000	1.000
2	-1.000	-1.000	-1.000	10	-0.269	1.000	1.000
3	-1.000	0.512	-1.000	11	1.000	-1.000	-1.000
4	-0.137	-1.000	-1.000	12	1.000	-1.000	0.045
5	1.000	-1.000	1.000	13	-1.000	-1.000	-0.124
6	1.000	1.000	-1.000	14	0.085	-1.000	1.000
7	1.000	-0.038	1.000	15	-1.000	1.000	-0.213
8	-1.000	1.000	1.000	16	-0.149	1.000	-1.000

Table 3.1: Example 1, Bayesian design, $\xi_{\mathcal{Q},L}^*$, that maximizes the lower bound $\phi_L(\xi; \mathcal{Q})$.

Bayesian D -efficiency than $\xi_{\mathcal{Q},L}^*$. However, the lower bound on the Bayesian D -efficiency is $\text{Bayes-eff}(\xi_{\mathcal{Q},L}^*; \mathcal{Q}) \gtrsim 99.4\%$, so any improvement to be gained by using a different design will be very small.

Note that the computation of the numerical value of the lower bound in (3.3) is approximate since we cannot be certain to have found the global optimum $\xi_{\mathcal{Q},U}^*$, although in the above example an assessment of the objective function values from the different random initializations of the algorithm suggests that the number of starts is adequate.

To assess the performance of a given design, ξ , for different β , we use the local D -efficiency,

$$\text{eff}(\xi; \beta) = \{|M(\xi; \beta)|/|M(\xi_{\beta}^*; \beta)|\}^{1/p}, \quad (3.4)$$

where $\xi_{\beta}^* \in \arg \max_{\xi \in \Xi} |M(\xi; \beta)|$ is a locally D -optimal design. For some β , $M(\xi; \beta)$ is well-conditioned for most $\xi \in \Xi$. In this case, the local D -efficiency can be approximated by searching numerically for ξ_{β}^* , and substituting the design found into (3.4). For other β , $M(\xi; \beta)$ is ill-conditioned for all $\xi \in \Xi$. Then approximate bounds on the efficiency can be derived by numerical search for the designs $\xi_{U,\beta}^* \in \arg \max_{\xi \in \Xi} \phi_U(\xi; \beta)$ and $\xi_{L,\beta}^* \in \arg \max_{\xi \in \Xi} \phi_L(\xi; \beta)$, and from the fact that

$$\exp \frac{1}{p} [\phi_L(\xi; \beta) - \phi_U(\xi_{U,\beta}^*; \beta)] \leq \text{eff}(\xi; \beta) \leq \exp \frac{1}{p} [\phi_U(\xi; \beta) - \phi_L(\xi_{L,\beta}^*; \beta)]. \quad (3.5)$$

To visualize the dependence of the local efficiency on the individual parameters, for each regression coefficient β_j we plot the approximate mean and 10% and

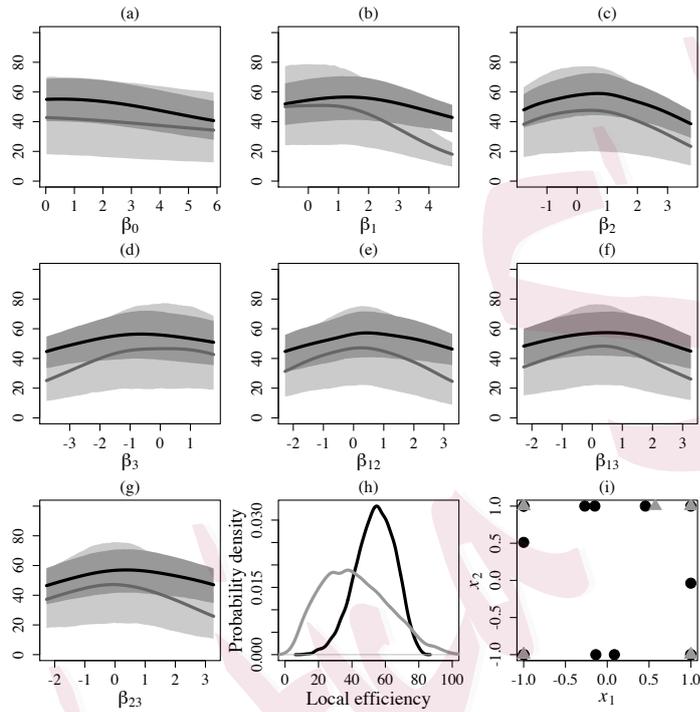


Figure 3.1: Robustness comparison of the Bayesian D -efficient design $\xi_{Q,L}^*$ (black lines/points) versus the EW-optimal design (grey lines/points). *Panels (a)–(g)*: conditional distribution, given β_k , of the local efficiency, $\text{eff}(\xi; \beta)$, induced by the prior on β (solid line, conditional mean; shaded region, 10% and 90% quantiles). *Panel (h)*: marginal distribution of the local efficiency. *Panel (i)*: 2-dimensional projection of the design points.

90% quantiles of the conditional distribution of $\text{eff}(\xi; \beta)$ given β_j . Owing to the need to search for a locally D -optimal design, evaluation of $\text{eff}(\xi; \beta)$ is computationally intensive. Thus, before computing the conditional mean and quantiles it is advantageous to first build a statistical emulator of $\text{eff}(\xi; \beta)$ as a function of β , using Gaussian process interpolation. This is analogous to the approach followed in the computer experiments literature when the main effects of a computationally expensive simulator are visualized (e.g. Santner, Williams and Notz (2003, Ch.7)). A similar method was used by Waite and Woods (2015) to visualize the efficiency profile of Bayesian designs for logistic models with random effects.

Example 1 (continued). We consider further the performance of the design, $\xi_{\mathcal{Q},L}^*$, that maximizes the lower bound for $\phi(\xi; \mathcal{Q})$. The support points of the quadrature scheme are used to train the emulator of $\text{eff}(\xi_{\mathcal{Q},L}^*; \beta)$. In the example, only three out of the 217 β vectors in \mathcal{Q} led to $M(\xi; \beta)$ being ill-conditioned for all $\xi \in \Xi$. For these vectors, the efficiency bounds in (3.5) gave no additional information beyond $\text{eff}(\xi_{\mathcal{Q},L}^*; \beta) \in [0\%, 100\%]$. Thus we decided to omit these β vectors from the training set, as including the bounds $[0\%, 100\%]$ would not substantively reduce uncertainty about the efficiency at these β . Figure 3.1 shows approximations to the conditional mean and conditional quantiles (given β_j) of the local efficiency, obtained using the emulator and Monte Carlo sampling. Also shown is a kernel density estimate for the marginal distribution of local efficiencies of $\xi_{\mathcal{Q},L}^*$ induced by the prior distribution on β . This is derived by computing the Kriging-based estimates of $\text{eff}(\xi_{\mathcal{Q},L}^*; \beta)$ for a Monte Carlo sample of 10,000 β vectors from the prior distribution. From Figure 3.1, it appears that the modal local efficiency of $\xi_{\mathcal{Q},L}^*$ is in the range 55-60%. The lower and upper quartiles of the local efficiency distribution are approximately 46% and 62%. Although at first glance the typical local efficiencies may appear fairly low, it is important to remember that due to the large amount of prior uncertainty here, there will exist no design whose local efficiency is significantly higher than $\xi_{\mathcal{Q},L}^*$ uniformly across the entire parameter space. The design $\xi_{\mathcal{Q},L}^*$ achieves a near-optimal trade-off in performance, as quantified by the high estimated Bayesian D -efficiency obtained earlier, across the very different parameter scenarios that are possible under the prior for β . The design is thus relatively robust. There appear to be no significant areas of the parameter space where the design performance is very poor and, for

example, the prior probability that $\text{eff}(\xi_{Q,L}^*; \beta) < 0.2$ appears negligible.

For comparison, results are included for the EW-optimal design, ξ_{EW}^* , advocated by Yang et al. (2016), which maximizes

$$\psi_{\text{EW}}(\xi) = \log |EM_{\beta}(\xi; \beta)| = \log \left| \sum_{i=1}^n E[w(\eta_i)] f(x_i) f^T(x_i) \right|.$$

In factorial experiments for logistic regression, Yang et al. (2016) found EW-optimal designs to be of comparable statistical efficiency to Bayesian D -optimal designs, while requiring less computational effort to obtain. Here, as shown in Figure 3.1(a)–(h), the EW-optimal design is much less robust than the Bayesian D -efficient design, at least in this case; its local efficiency $\text{eff}(\xi_{\text{EW}}^*; \theta)$ has generally lower mean, both conditionally on β_j and marginally, and the local efficiency also exhibits higher variability. The smaller difference in robustness between Bayesian D -optimal and EW-optimal designs observed by Yang et al. (2016) may be due to their restriction to a factorial design space: here, the greater performance of the Bayesian D -efficient design appears to be due to the inclusion of a greater number of factor settings in the interior of $(-1, 1)$ (see Figure 3.1(i)).

In addition to the worse statistical performance of the EW-optimal design in this case, here the computational convenience of EW-optimal designs over Bayesian D -optimal designs is much reduced. For factorial problems, a reduction in computational cost is achieved by precomputing $E[w(\eta)]$ for every point in the finite design space, enabling faster evaluation of $\psi_{\text{EW}}(\xi)$. Here, precomputation is not possible since we have continuous factors and therefore an uncountably infinite design space. One computational benefit of the EW criterion is that it successfully avoids problems with ill-conditioning, but this offers only a minor advantage over Bayesian D -optimality, for which ill-conditioning problems can now be overcome using the bounds developed in Section 3.2.

4. Discussion

The central tenet of this paper is that it is not permissible to use a singular prior distribution in conjunction with Bayesian D -optimality as a design selection criterion. Our new theoretical results, summarized below, can help to ascertain whether a prior is singular or non-singular. This is useful, since if it can be demonstrated that the prior is non-singular, then we may proceed to find

Bayesian D -optimal designs either using standard methods, or using the new numerical techniques developed in Section 3 if problems are encountered with ill-conditioning. If it is instead demonstrated that the current prior is singular, then $\phi(\xi; \mathcal{P}) = -\infty$ for all $\xi \in \Xi$, meaning that design selection fails. One rough intuitive interpretation of this is that the parameter uncertainty under this prior is so great that any design will have low (local) efficiency across a significant portion of the parameter space. In this case, there are two possibilities: (i) consider a different prior distribution, or (ii) adopt a different design selection criterion. These alternatives are considered in Sections 4.1 and 4.2 respectively.

To summarize our theoretical results, for three generalized linear models we have given conditions that can easily be checked to establish non-singularity of \mathcal{P} and, importantly, identified that a prominent class of default prior distributions for logistic regression should not be used for Bayesian D -optimal design. For the compartmental model in Section 2.2, sufficient conditions were established only for singularity of \mathcal{P} , thus highlighting only prior distributions that should not be used. Though desirable, the proof of an inverse result guaranteeing non-singularity seems highly involved and is beyond the scope of this paper. Future work could seek to develop results on singular prior distributions for population pharmacokinetic models, for which optimal sampling times are more commonly sought (e.g. Mentré et al. (1997)). Such models extend (2.3) by allowing subject-specific kinetic parameters.

4.1. Alternative prior distributions

Often there are multiple plausible candidates for a suitable prior distribution. For example, in the subjectivist framework, informative priors are elicited from expert knowledge by obtaining summaries to which a probability distribution may be fitted (e.g. Garthwaite et al. (2005), Oakley and O'Hagan (2007)). Typically there will be multiple distributions that fit the observed summaries. Away from this approach, if using uninformative or weakly informative priors there are still often multiple possible candidate priors. Thus, if design selection fails because \mathcal{P} is singular but there exists an alternative candidate prior \mathcal{P}' that is non-singular, then Bayesian D -optimality may be used with \mathcal{P}' instead. Nonetheless, if adopting a subjectivist viewpoint we should be careful to avoid selecting prior distributions purely for analytical convenience if they do not accurately represent

the available expert belief or knowledge.

As an example, in Section 2.1 we found that $\mathcal{P} : \theta \sim U(0, a)$ is singular for the exponential regression model. A natural question is whether it is sufficient to find designs for the non-singular prior $\mathcal{P}_\epsilon : \theta \sim U(\epsilon, a)$ for some small value of ϵ (e.g. 10^{-3} or 10^{-6}). The adequacy of \mathcal{P}_ϵ as a representation of the expert's beliefs will depend substantially on the specifics of the application. For small ϵ , the quartiles of \mathcal{P} and \mathcal{P}_ϵ are similar, thus for example it is possible for both distributions to fit expert statements obtained by the bisection method (Garthwaite et al. (2005)). However, the implication of \mathcal{P}_ϵ that there is zero probability that $\theta < \epsilon$ is too strong unless the expert is certain that $\theta \geq \epsilon$. The fidelity of the representation \mathcal{P}_ϵ would be less important if the resulting design decision were insensitive to the choice of ϵ . Unfortunately this is not the case, as shown by the proposition below and its proof (in the supplementary material). Intuitively, as $\epsilon \rightarrow 0$, some points in the Bayesian D -optimal design for \mathcal{P}_ϵ will converge to zero (while never being equal to zero).

Proposition 9. *For the exponential model, if ξ does not vary with ϵ then*

$$\text{Bayes-eff}(\xi; \mathcal{P}_\epsilon) \rightarrow 0 \quad \text{as } \epsilon \rightarrow 0.$$

Thus, even if one were to compute the Bayesian D -optimal design for $\mathcal{P}_{\epsilon'}$, with say $\epsilon' = 10^{-6}$, the resulting design would be highly inefficient when evaluated under \mathcal{P}_ϵ for $\epsilon \ll \epsilon'$.

The situation above is somewhat similar to problems in the objective Bayesian approach with improper uninformative priors (e.g. Berger (1985, Ch.3); Berger (2006)), which one may need to modify in order to obtain a proper posterior. For example, if an improper prior, say $U(10, \infty)$, does not give a proper posterior, one might attempt to replace it with $U(10, M)$, with M large, e.g. 10^5 or 10^6 . However, the results would often be highly sensitive to the value chosen for M , which is arbitrary and typically has no objective justification. For further discussion on the role of prior information in design of experiments, see Woods et al. (2016).

4.2. Alternative selection criteria

If all candidate prior distributions that agree with the elicited prior knowledge or beliefs are singular, then a Bayesian D -optimal design cannot be found

and it is necessary to use an alternative design selection criterion that suffers from fewer problems with singularities. One such criterion that has already been mentioned is EW D -optimality (Yang et al. (2016)). Note that the numerical results in Section 3.3 suggest that if the problem is with an ill-conditioned \mathcal{Q} rather than a singular \mathcal{P} , then the EW D -optimal design may be less robust than a Bayesian D -efficient design found using the numerical methods developed here. Another alternative is to select ξ to maximize the *mean local efficiency*,

$$\Psi(\xi; \mathcal{P}) = E_{\mathcal{P}}\{\text{eff}(\xi; \theta)\},$$

which is fairly insensitive to the presence of θ with $|M(\xi; \theta)| \approx 0$. This is a special case of the objective function discussed by Dette and Wong (1996) (Φ_1 in their notation). Unlike Bayesian D -optimality, neither of the above alternative criteria has the interpretation of approximate equivalence to the maximization of Shannon information gain. As an example of the use of the mean local efficiency criterion, consider again the exponential decay model from Section 2.1. From Corollary 1, when $\mathcal{P} = U(0, a)$, $a > 0$, all designs are Bayesian (D -)singular with respect to \mathcal{P} for the θ -parameterization. By contrast, it is shown easily that the design with a single support point $x = a/2$ is Ψ -optimal with a mean efficiency of approximately 67%. This design is locally D -optimal when θ is equal to its prior mean, but highly inefficient when θ is small. Thus, Ψ -optimal designs are much less strongly driven by their worst-case behaviour. As with EW D -optimality, if the problem is in fact with an ill-conditioned \mathcal{Q} rather than a singular \mathcal{P} , then it is possible that Ψ -optimal designs may be less robust than Bayesian D -optimal or Bayesian D -efficient designs.

A further alternative approach to design selection under parameter uncertainty is to consider maximin designs. In the case of greatest interest in this paper, Θ is such that $\inf_{\theta \in \Theta} |M(\xi; \theta)| = 0$ for all $\xi \in \Xi$, thus design selection clearly fails using the unstandardized maximin D -criterion (Imhof (2001)). Often design selection will also fail when using the standardized maximin D -criterion. It is clear that the Bayesian approach, under suitable prior distributions, benefits from greater robustness to the presence of singular θ than the use of maximin criteria. For related results see Braess and Dette (2007), where conditions are established under which the number of support points in a standardized maximin or Bayesian D -optimal approximate design grows arbitrarily large as Θ is

expanded. In contrast, in the work presented here it is supposed that Θ is fixed and the focus is instead on examining the adequacy of the set of competing exact designs under different prior distributions. Here, results have also been developed for additional multiparameter models and numerical methods proposed.

Supplementary material

The online supplementary material for this paper contains proofs of the analytical results described in the text.

Acknowledgement

I am grateful to two anonymous referees whose insightful and helpful comments prompted several improvements to the paper, and to Professors David Woods and Susan Lewis for several invaluable discussions and comments. This work was supported by the UK Engineering and Physical Sciences Research Council, via a PhD studentship, Doctoral Prize, and a project grant. The work made use of the Iridis computational cluster at the University of Southampton.

References

- Abramowitz, M. and Stegun, I. A. (1964), *Handbook of Mathematical Functions: with Formulas, Graphs, and Mathematical Tables*, Dover, New York. 10th printing.
- Albert, A. and Anderson, J. (1984), ‘On the existence of maximum likelihood estimates in logistic regression models’, *Biometrika* **71**, 1–10.
- Atkinson, A. C. (2003), ‘Horwitz’s rule, transforming both sides and the design of experiments for mechanistic models’, *J. Roy. Stat. Soc. C* **52**, 261–278.
- Atkinson, A. C., Chaloner, K., Herzberg, A. M. and Juritz, J. (1993), ‘Optimum experimental designs for properties of a compartmental model’, *Biometrics* **49**, 325–337.
- Atkinson, A. C. and Woods, D. C. (2015), Designs for generalized linear models, in A. Dean, M. Morris, J. Stufken and D. Bingham, eds, ‘Handbook of Design and Analysis of Experiments’, Chapman and Hall/CRC, chapter 13, pp. 471–548.
- Berger, J. O. (1985), *Statistical decision theory and Bayesian analysis*, Springer, New York.
- Berger, J. O. (2006), ‘The case for objective Bayesian analysis’, *Bayesian Anal.* **1**, 385–402.
- Braess, D. and Dette, H. (2007), ‘On the number of support points of maximin and Bayesian optimal designs’, *Ann. Statist.* **35**, 772–792.

- Capinski, M. and Kopp, P. E. (2004), *Measure, Integral and Probability*, 2nd edn, Springer, New York.
- Chaloner, K. and Larntz, K. (1989), ‘Optimal Bayesian design applied to logistic regression experiments’, *J. Statist. Plann. Infer.* **21**, 191–208.
- Chaloner, K. and Verdinelli, I. (1995), ‘Bayesian experimental design: a review’, *Statist. Sci.* **10**, 273–304.
- Chernoff, H. (1953), ‘Locally optimal designs for estimating parameters’, *Ann. Math. Statist.* **24**, 586–602.
- Chipman, H., Hamada, M. and Wu, C. (1997), ‘A Bayesian variable-selection approach for analyzing designed experiments with complex aliasing’, *Technometrics* **39**, 372–381.
- Dette, H. and Neugebauer, H.-M. (1997), ‘Bayesian D -optimal designs for exponential regression models’, *J. Statist. Plann. Infer.* **60**, 331–349.
- Dette, H. and Wong, W. K. (1996), ‘Optimal Bayesian designs for models with partially specified heteroscedastic structure’, *Ann. Statist.* **24**, 2108–2127.
- Feng, C., Wang, H., Zhang, Y., Han, Y., Liang, Y. and Tu, X. M. (2017), ‘Generalized definition of the geometric mean of a nonnegative random variable’, *Comm. Statist. Th. Meth.* **46**, 3614–3620.
- Garthwaite, P. H., Kadane, J. B. and O’Hagan, A. (2005), ‘Statistical methods for eliciting probability distributions’, *J. Amer. Statist. Assoc.* **100**, 680–701.
- Gelman, A., Bois, F. and Jiang, J. (1996), ‘Physiological pharmacokinetic analysis using population modeling and informative prior distributions’, *J. Amer. Statist. Assoc.* **91**, 1400–1412.
- Gelman, A., Jakulin, A., Pittau, M. G. and Su, Y.-S. (2008), ‘A weakly informative default prior distribution for logistic and other regression models’, *Ann. Appl. Statist.* **2**, 1360–1383.
- Gotwalt, C. M., Jones, B. A. and Steinberg, D. M. (2009), ‘Fast computation of designs robust to parameter uncertainty for nonlinear settings’, *Technometrics* **51**, 88–95.
- Imhof, L. A. (2001), ‘Maximin designs for exponential growth models and heteroscedastic polynomial models’, *Ann. Statist.* **29**, 561–576.
- Jones, B. and Goos, P. (2009), ‘ D -optimal design of split-split-plot experiments’, *Biometrika* **96**, 67–82.
- Khuri, A. I., Mukherjee, B., Sinha, B. K. and Ghosh, M. (2006), ‘Design issues for generalized linear models: a review’, *Statist. Sci.* **21**, 376–399.
- McCullagh, P. and Nelder, J. A. (1989), *Generalized Linear Models*, 2nd edn, Chapman and Hall, Boca Raton, FL.

- McGree, J. M. and Eccleston, J. A. (2012), ‘Robust designs for Poisson regression models’, *Technometrics* **54**, 64–72.
- Mentré, F., Mallet, A. and Baccar, D. (1997), ‘Optimal design in random-effects regression models’, *Biometrika* **84**, 429–442.
- Oakley, J. E. and O’Hagan, A. (2007), ‘Uncertainty in prior elicitation: a nonparametric approach’, *Biometrika* **94**, 427–441.
- Polson, N. G. and Scott, J. G. (2012), ‘On the half-Cauchy prior for a global scale parameter’, *Bayesian Anal.* **7**, 887–902.
- Pukelsheim, F. (1993), *Optimal design of experiments*, SIAM, Philadelphia.
- Russell, K. G., Woods, D. C., Lewis, S. M. and Eccleston, J. A. (2009), ‘D-optimal designs for Poisson regression models’, *Statist. Sinica* **19**, 721–730.
- Santner, T. J., Williams, B. J. and Notz, W. (2003), *The Design and Analysis of Computer Experiments*, Springer, New York.
- Waite, T. W. and Woods, D. C. (2015), ‘Designs for generalized linear models with random block effects via information matrix approximations’, *Biometrika* **102**, 677–693.
- Woods, D. C., Lewis, S. M., Eccleston, J. A. and Russell, K. G. (2006), ‘Designs for generalized linear models with several variables and model uncertainty’, *Technometrics* **48**, 284–292.
- Woods, D. C., Overstall, A. M., Adamou, M. and Waite, T. W. (2016), ‘Bayesian design of experiments for generalised linear models and dimensional analysis with industrial and scientific application (with discussion)’, *Quality Engineering* **29**, 91–118.
- Yang, J. and Mandal, A. (2015), ‘D-optimal factorial designs under generalized linear models’, *Comm. Statist. Simul. Comput.* **44**, 2264–2277.
- Yang, J., Mandal, A. and Majumdar, D. (2016), ‘Optimal designs for 2^k factorial experiments with binary response’, *Statist. Sinica* **26**, 385–411.
- Yang, M. (2010), ‘On the de la Garza Phenomenon’, *Ann. Statist.* **38**, 2499–2524.
- Yang, M. and Stufken, J. (2009), ‘Support points of locally optimal designs for nonlinear models with two parameters’, *Ann. Statist.* **37**, 518–541.
- Yang, M. and Stufken, J. (2012), ‘Identifying locally optimal designs for nonlinear models: a simple extension with profound consequences’, *Ann. Statist.* **40**, 1665–1681.
- Yang, M., Zhang, B. and Huang, S. (2011), ‘Optimal designs for generalized linear models with multiple design variables’, *Statist. Sinica* **21**, 1415–1430.

School of Mathematics, University of Manchester, Manchester, M13 9PL, U.K.

E-mail: timothy.waite@manchester.ac.uk