

Extending the use of statistical emulators in Bayesian experimental design

James McGree

Associate Professor of Statistics

Queensland University of Technology

james.mcgree@qut.edu.au, jamesmcgree.com, [@j_mcgree](https://twitter.com/j_mcgree)

Joint work with

Antony Overstall

Associate Professor of Statistics

University of Southampton, UK.



April 18, 2018



Introduction

Aim

- Efficiently derive high dimensional Bayesian designs
- Large scale model and parameter uncertainty

Methodology

- Extend the use of emulators in Bayesian design
- Propose the new k -dimensional approximate coordinate exchange algorithm

Motivation

- Current methods can be inefficient
- Extensions are needed in general but particularly for high dimensional designs such as screening experiments



Bayesian inference

- Predominately focused the **posterior distribution**:

$$p(\theta_m | \mathbf{y}, \mathbf{d}, m) = p(\theta_m | m) p(\mathbf{y} | \theta_m, \mathbf{d}, m) / Z_m$$

where \mathbf{y} represent the data, \mathbf{d} the design, θ_m the model parameters, $p(\theta_m | m)$ and $p(\mathbf{y} | \theta_m, \mathbf{d}, m)$ are the prior and the likelihood for model m and model m has prior model probability $p(m)$, for $m = 1, \dots, K$.

- Z_m is the **model evidence**, defined as

$$Z_m = p(\mathbf{y} | \mathbf{d}, m) = \int_{\theta_m} p(\mathbf{y} | \theta_m, \mathbf{d}, m) p(\theta_m | m) d\theta_m$$

- Z_m used for **model choice** (proportional to posterior model probabilities)



Bayesian design

- Find a design \mathbf{d} to address particular experimental aims
- Aim is encapsulated in a utility function $u(\mathbf{d}, \mathbf{y}, m)$
- Could include parameter estimation, model selection and/or prediction
- Maximise expected utility $\mathbf{d}^* = \arg \max_{\mathbf{d}} u(\mathbf{d})$, where

$$u(\mathbf{d}) = \sum_{m=1}^K p(m) \int_{\mathbf{y}} u(\mathbf{d}, \mathbf{y}, m) p(\mathbf{y}|\mathbf{d}, m) d\mathbf{y}.$$

- $u(\mathbf{d}, \mathbf{y}, m)$ is some measure of information gained from \mathbf{d} given model m and observed data \mathbf{y} .

Estimation utilities

- **Shannon information gain** on θ (Shannon, 1948)

$$u(\mathbf{d}, \mathbf{y}, m) = \int_{\theta_m} p(\theta_m|m, \mathbf{y}, \mathbf{d}) \log p(\mathbf{y}|\theta_m, m, \mathbf{d}) d\theta_m - \log p(\mathbf{y}|m, \mathbf{d})$$

- **Negative squared loss** (Overstall and Woods, 2017)

$$u(\mathbf{d}, \theta_m, \mathbf{y}, m) = -(\theta_m - E[\theta_m])'(\theta_m - E[\theta_m])$$

where $E[\theta_m] = \sum_{m=1}^K p(m) \int_{\theta_m} \theta_m p(\theta_m|\mathbf{y}, \mathbf{d}, m) d\theta_m$

- Need $p(\theta_m|m, \mathbf{y}, \mathbf{d})$ to get $u(\mathbf{d}, \mathbf{y}, m)$
- Difficult to approximate $\log p(\mathbf{y}|m, \mathbf{d})$
- Computationally difficult to efficiently approximate $p(\theta_m|m, \mathbf{y}, \mathbf{d})$ (more details later)



Model discrimination utilities

- **Mutual information** (Box and Hill, 1967, Drovandi, McGree, Pettitt, 2014)

$$u(\mathbf{d}, \mathbf{y}, m) = \log p(m|\mathbf{y}, \mathbf{d})$$

- **01 utility** (Overstall, McGree, Drovandi, 2018)

$$u(\mathbf{d}, \mathbf{y}, m) = I(m = \arg \max_w p(w|\mathbf{y}\mathbf{d})), w = 1, \dots, K$$

- Need $p(m|\mathbf{y}, \mathbf{d})$ to get $u(\mathbf{d}, \mathbf{y}, m)$
- Computationally difficult

Approximating expected utility

- $u(\mathbf{d})$ typically **cannot be solved analytically**
- Can be approximated using Monte Carlo integration

$$u(\mathbf{d}) \approx \sum_{m=1}^K p(m) \frac{1}{B} \sum_{b=1}^B u(\mathbf{d}, \mathbf{y}_{mb}, m),$$

where $\mathbf{y}_{mb} \sim p(\mathbf{y}|\theta_{mb}, m, \mathbf{d})$ and $\theta_{mb} \sim p(\theta_m|m)$.

- *BK* evaluation of $u(\mathbf{d}, \mathbf{y}_{mb}, m)$ needed
- Hence, ***BK* posterior distributions need to be approximated or sampled from** to approximate $u(\mathbf{d})$.
- **Computationally challenging task.** How can this be achieved efficiently?



Approximating utility

- Long et al. (2013) and Overstall, McGree and Drovandi (2018) used the **Laplace approximation** for efficiently estimating $u(\mathbf{d}, m, \mathbf{y})$;
- The main result is that the approximation to the posterior distribution of θ_m has the following **multivariate Normal** form:

$$\hat{p}(\theta_m | \mathbf{y}, \mathbf{d}, m) = (2\pi)^{-\frac{q_m}{2}} |\Sigma_{my}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\theta_m - \hat{\theta}_{my})^t \Sigma_{my}^{-1} (\theta_m - \hat{\theta}_{my})\right)$$

where q_m denotes the number of parameters in model m , $\hat{\theta}_{my}$ and Σ_{my} denote the posterior mode and posterior variance-covariance matrix, respectively, for model m upon the observation of \mathbf{y} from design \mathbf{d} , for $m = 1, 2, \dots, K$.



Approximating utility

- For posterior inference on m , the posterior model probability ($p(m|\mathbf{y}, \mathbf{d})$) can be considered;
- Proportional to the model evidence;
- Based on Laplace approximation, the **model evidence can be approximated as follows**:

$$\hat{p}(\mathbf{y}|m, \mathbf{d}) = (2\pi)^{\frac{q_m}{2}} |\boldsymbol{\Sigma}_{my}|^{\frac{1}{2}} p(\mathbf{y}|\hat{\boldsymbol{\theta}}_{my}, \mathbf{d})p(\hat{\boldsymbol{\theta}}_{my}|m). \quad (1)$$

- Thus, posterior summaries such as $u(\mathbf{d}, m, \mathbf{y})$ can be evaluated based on the above Laplace approximation facilitating a relatively efficient approximation to $u(\mathbf{d})$;
- However, non-smooth utility function leading to a **difficult optimisation problem**;
- Try using **statistical emulators**?



The use of emulators in locating Bayesian designs

- Muller and Parmigiani (1995) - Curve fitting methods to smooth Monte Carlo draws (2D)
- Weaver et al. (2016) - Gaussian process with EQI (3D)
- Jones et al. (2016) - Bayes linear analysis (9D)
- **Approximate coordinate exchange** (Overstall and Woods, 2017) (\approx 200D);
 - Extension of the coordinate exchange algorithm;
 - Emulator used to interpolate $\hat{u}(\mathbf{d})$ in one dimension at a time;
 - Can be inefficient if one-dimensional solution changes depending on the values of other design elements.

The ACE algorithm

- Optimises each design element one-at-a-time
- Instead of exchanging a discrete set of points for each design element, **fit an emulator and optimise**
- With k explanatory variables, n design points,

$$\mathbf{d} = \begin{bmatrix} d_{11} & d_{21} & \dots & d_{k1} \\ d_{12} & d_{22} & \dots & d_{k2} \\ \vdots & & & \\ d_{1n} & d_{2n} & \dots & d_{kn} \end{bmatrix}$$

- Emulate $u(\mathbf{d})$ as a function of \mathbf{d}_{ip} given \mathbf{d}_{-ij} is fixed. Denote as $u(\mathbf{d}_{ij}|\mathbf{d}_{-ij})$, for $i = 1, \dots, n$ and $j = 1, \dots, k$
- Optimise emulator (brute force)
- Cycle through all nk design points R times



Emulation

- Chosen emulator is the **Gaussian Process** (GP, Rasmussen and Williams, 2006)
- Defined by a mean and covariance function:

$$f(\mathbf{z}) \sim \mathcal{GP}(m(\mathbf{z}), k(\mathbf{z}, \gamma)),$$

where \mathbf{z} are predictor variables, mean function $m(\mathbf{z})$ is the expected value at \mathbf{z} and $k(\mathbf{z}, \gamma)$ is the covariance function which models the dependence between function values at z_k and z_l , for $k, l = 1, 2, \dots, T$.

- Each element of the **covariance matrix**:

$$k(z_k, z_l; \gamma) = \begin{cases} \gamma_0 + \gamma_1 & \text{if } x = 0 \\ \gamma_1 k(z_k, z_l, \gamma_2) & \text{if } x > 0, \end{cases}$$

where x is the Euclidean distance between z_k and z_l , γ_0 is nugget, γ_1 is partial sill and γ_2 is an additional parameter.



The ACE algorithm

- Some popular choices for $k(\mathbf{z}_k, \mathbf{z}_l, \gamma_2)$ include:

$$k(\mathbf{z}_k, \mathbf{z}_l, \gamma_2) = \exp(-0.5(x/\gamma_2)^2)$$

Gaussian

$$k(\mathbf{z}_k, \mathbf{z}_l, \gamma_2) = \exp(-x/\gamma_2)$$

Exponential

$$k(\mathbf{z}_k, \mathbf{z}_l, \gamma_2) = (1 + \sqrt{5}x\gamma_2 + (5/3)(x/\gamma_2)^2) \exp(-\sqrt{5}x\gamma_2)$$

Matérn.

- Efficient** when $f(\mathbf{z})$ is expensive as can interpolate $f(\tilde{\mathbf{z}})$ (based on the fitted zero mean GP):

$$f(\tilde{\mathbf{z}}) = k(\tilde{\mathbf{z}}, \mathbf{Z}, \boldsymbol{\gamma})(k(\mathbf{Z}, \boldsymbol{\gamma}) + \gamma_0 \mathbf{I})^{-1} \mathbf{z},$$

where \mathbf{Z} denotes training set, $k(\tilde{\mathbf{z}}, \mathbf{Z}, \boldsymbol{\gamma})$ evaluates the covariances between $\tilde{\mathbf{z}}$ and \mathbf{Z} , and \mathbf{I} denotes the identity matrix.



The ACE algorithm

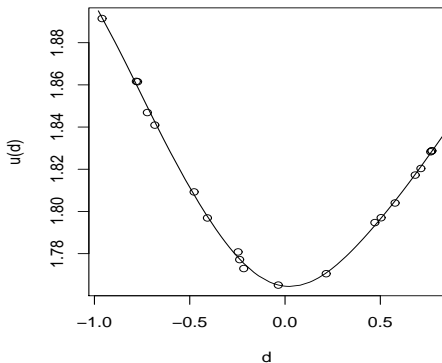


Figure: Emulated utility surface and training points

The k -dimensional ACE algorithm

- Recall

$$\mathbf{d} = \begin{bmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \\ \vdots \\ \mathbf{d}_n \end{bmatrix} = \begin{bmatrix} d_{11} & d_{21} & \dots & d_{k1} \\ d_{12} & d_{22} & \dots & d_{k2} \\ \vdots & & & \\ d_{1n} & d_{2n} & \dots & d_{kn} \end{bmatrix}$$

- Propose optimising \mathbf{d}_i simultaneously
- Extend the GP to be k -dimensional
- Emulate $u(\mathbf{d})$ as a function of \mathbf{d}_i given \mathbf{d}_{-i} is fixed. Denote as $u(\mathbf{d}_i|\mathbf{d}_{-i})$

$$u(\mathbf{d}_i|\mathbf{d}_{-i}) \sim \mathcal{GP}(\mathbf{0}, k(\mathbf{D}, \gamma)), \quad \mathbf{D} \text{ set of training points}$$

- Optimise k -dimensional GP
- Cycle through all n design points R times



The k -dimensional ACE algorithm

Training the emulator

- For each \mathbf{d}_j , need a training set \mathbf{D} to fit GP (for $\hat{u}(\mathbf{d}_j|\mathbf{d}_{-j})$)
- Consider \mathbf{d}_j
 - Propose \mathbf{D} through random draws from $q(\cdot)$;
 - For each proposal, evaluate $\hat{u}(\mathbf{d}_j|\mathbf{d}_{-j})$. Denote as \mathbf{u}
 - Fit GP for $\hat{u}(\mathbf{d}_j|\mathbf{d}_{-j})$ (via maximum likelihood);
- After training, need to **optimise** the k -dimensional emulator

The k -dimensional ACE algorithm

Optimising the emulator

- GP is a smooth function, use gradient-based methods;
- Efficient (as above):

$$\tilde{\mathbf{u}} = k(\tilde{\mathbf{d}}, \mathbf{D}, \boldsymbol{\gamma})(k(\mathbf{D}, \boldsymbol{\gamma}) + \gamma_0 \mathbf{I})^{-1} \mathbf{u}. \quad (2)$$

- Once training is complete, loop through the following
 - Optimise the GP based on many random starts
 - Evaluate the 'actual' expected utility $\hat{u}(\mathbf{d}_i | \mathbf{d}_{-i})$ of maximum
 - Re-fit the GP
- The 'best' design from the above loop is either accepted or rejected
- Then move to next row of the design



The k -dimensional ACE algorithm

Training and optimisation

- Two components of the proposed algorithm; training and optimisation
- Need to consider trade-off in computational resources
- How many utility evaluations for training? How many for optimisation?
- Explored through examples (results omitted, 75% training, 25% optimising).

Choice of covariance function

- k -dimensional emulator, (in general) $k > 1$
- Choice of covariance function is potentially important
- Does this have implications for locating Bayesian designs?
- Explored through examples (some results shown)

Examples

Details of examples

- Two examples to be considered; test problem and motivating example
- Will compare the performance of ACE and k -dimensional ACE
- In both examples, both algorithms will be run in a similar way (equal number of utility evaluations)



Test problem

Logistic regression

- GLMs which are of general interest to the design community (Woods et al., 2006, Dror and Steinberg, 2008)
- Often used to benchmark new computational algorithms in Bayesian inference (Minka, 2001, Cabras et al., 2015)
- Emerging in Bayesian design (Overstall and Woods, 2017, McGree, 2017)
- We consider independent and dependent data settings

Test problem

Logistic regression

- Four factor logistic regression model
- Two scenarios; all data are independent and data are collected in blocks of $n_G = 10$, such that $n = n_G G$, G is total number of blocks
- Model defined as

$$\left(\frac{\pi_{ig}}{1 - \pi_{ig}} \right) = \theta_0 + \beta_{0g} + \sum_{j=1}^{k=4} (\theta_j + \beta_{jg}) d_{ijg}$$

where θ are regression parameters, β are block specific parameters, $d_{ij} \in [-1, 1]$, $g = 1, \dots, G$

- Assumed unknown sign of effect ($\theta = \mathbf{0}$) with variances of 2, $\beta_{jg} \sim N(0, \sigma_{jg}^2)$, $\sigma_{jg}^2 \sim G(2, 2)$
- For blocked data setting, likelihood approximated (details omitted)



Example 1

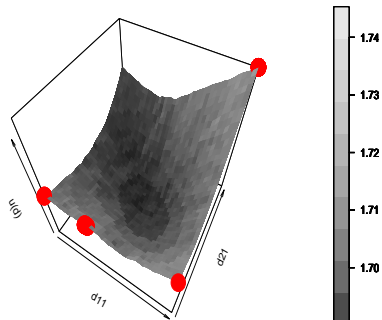


Figure: Utility surface and proposed optimal design points (only in 2D)

Example 1

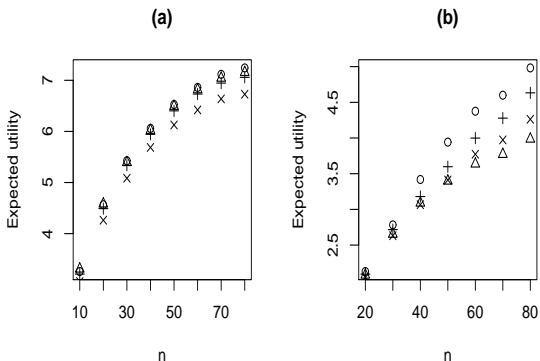


Figure: ACE with exp (Δ), k -ACE with Matérn (\circ), exp (\times) and quantile improvement ($+$) for Shannon information gain on θ for standard (a) and hierarchical (b) logistic regression

Motivating example

Screening experiments

- Eight factor logistic regression model
- Unknown if any or all factors are important, **256 competing models**
- Model defined as

$$\left(\frac{\pi_i}{1 - \pi_i} \right) = \theta_0 + \sum_{j=1}^{k=8} \delta_j \theta_j d_{ij}$$

where δ_j is a binary indicator for whether factor j is active or not, and $d_{ij} \in [-1, 1]$

- $p(m)$ correct for Bayesian multiplicity (Scott and Berger, 2010)
- Assume uniform priors for θ , lower and upper bounds $(-3, 4, 5, -6, -2.5, -2, -4, -5, -6)$ and $(3, 10, 11, 0, 3.5, 4, 2, 1, 0)$, respectively.

Example 2

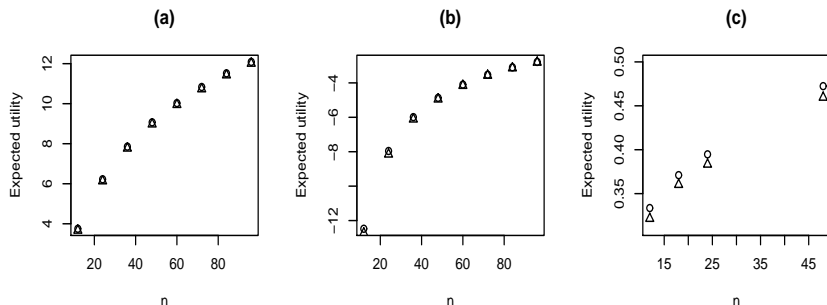


Figure: ACE (Δ) vs k -ACE (\circ) for (a) Shannon information gain on θ (b) Negative squared loss and (c) 01-utility

Discussion

- Proposed an extension to the ACE algorithm
- k -dimensional ACE worked better or at least as well as ACE in all examples considered
- Bayesian designs differed based on choice of covariance function - propose using CV to choose
- Also tried quantile improvement for all examples.

Future research

- Potential to consider **non-parametric emulators**.
- Should avoid some concerns about goodness-of-fit of the emulator
- Could think about **more efficient or adaptive proposal distribution** for training the emulator
- This is a design problem within itself
- Potential benefits in considering **space filling** approaches

Selected references

- Box and Hill (1967). *Technometrics*, 9, 57-71.
- Drovandi, McGree and Pettitt (2014). *Journal of Computational and Graphical Statistics*, 23, 3-24.
- McGree (2017). *Computational Statistics & Data Analysis*, 113, 207-225.
- Muller and Parmigiani (1995). *Journal of the American Statistical Association*, 90, 1322-1330.
- Overstall and Woods (2017). *Technometrics*, 59, 458-470.
- Overstall, McGree and Drovandi (2018). *Statistics and Computing*, 28, 343-358.
- Woods, Lewis, Eccleston and Russell (2006) *Technometrics*, 48, 284-292.

