

Design of Experiments when Gradient Information Is Available

Yiou Li

Department of Applied Mathematics

Summary Message

- For some computer experiments, evaluating a single data point can be computationally expensive, limiting the number of data that one can afford.
- Evaluating the gradient of the function using adjoint techniques, and involving gradient information in the model can substantially improve the accuracy of the prediction.
- The gradient of a d -variate function provides d more scalar pieces of information, at a cost of perhaps only several times the cost of a function value alone.
- Choosing the experimental design from two perspectives, *robustness* and *efficiency*.

Statistical Model

- The experimental region, Ω , is some measurable subset of $\Omega_1 \times \dots \times \Omega_d$.
- \mathcal{H} be some vector space of real-valued functions defined on Ω , and assumed to be a separable Hilbert space with a reproducing kernel $K : \Omega \times \Omega \rightarrow \mathbb{R}$.
- Define an operator $\mathbf{L}_x : \mathcal{H} \rightarrow \mathbb{R}^{d+1}$, which when applied to a d -variate function $f \in \mathcal{H}$, returns

$$\mathbf{L}_x f = \left(f(\mathbf{x}), \frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_d}(\mathbf{x}) \right)^T.$$

- For a vector function $\mathbf{f} = (f_1, \dots, f_\ell)^T : \Omega \rightarrow \mathbb{R}^\ell$, the definition of this operator is extended:

$$\mathbf{L}_x \mathbf{f}^T = (\mathbf{L}_x f_1, \dots, \mathbf{L}_x f_\ell).$$

- A linear regression model with gradient information:

$$\tilde{\mathbf{y}}_i = (\mathbf{L}_x \mathbf{g}^T) \beta + \tilde{\epsilon}_i, \quad i = 1, \dots, n,$$

- $\tilde{\mathbf{y}}_i$: observed vector response at the design point \mathbf{x}_i
- $\mathbf{g} = (g_1, \dots, g_k)^T$: the vector of basis functions
- β : regression coefficient to be estimated
- $\tilde{\epsilon}_i$: the error in estimating the response by the linear combination of k basis functions. It is assumed that $\tilde{\epsilon}_1, \dots, \tilde{\epsilon}_n$ are i.i.d with zero mean and covariance matrix $\sigma^2 \tilde{\Lambda}$.

- Vector-matrix notation:

$$\mathbf{y} = \mathbf{G}\beta + \boldsymbol{\epsilon},$$

where

$$\mathbf{G} = \begin{pmatrix} \mathbf{L}_x \mathbf{g}^T \\ \vdots \\ \mathbf{L}_x \mathbf{g}^T \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \tilde{\mathbf{y}}_1 \\ \vdots \\ \tilde{\mathbf{y}}_n \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \tilde{\epsilon}_1 \\ \vdots \\ \tilde{\epsilon}_n \end{pmatrix},$$

and $\boldsymbol{\epsilon}$ has zero mean and covariance matrix $\sigma^2 \tilde{\Lambda}$, where $\tilde{\Lambda} = \text{diag}(\tilde{\Lambda}_1, \dots, \tilde{\Lambda}_n)$.

- The weighted least squares estimate of the regression coefficient β :

$$\hat{\beta} = \mathbf{B}\mathbf{y} = \beta + \mathbf{B}\boldsymbol{\epsilon}, \quad \mathbf{B} = (\mathbf{G}^T \tilde{\Lambda}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \tilde{\Lambda}^{-1}.$$

Scaled Integrated Mean Squared Error

- Experimental design:

$$\xi = \left\{ \begin{matrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_N \\ w_1 & w_2 & \dots & w_N \end{matrix} \right\}, \text{ where } N \leq n.$$

- The goal of the linear regression is assumed to be the estimation of $\mathbf{T}(\mathbf{g}^T \beta)$, where $\mathbf{T} : \mathcal{H} \rightarrow \mathcal{L}_2^s$.

- For any vectors of functions, $\mathbf{u} = (u_1, \dots, u_s)^T$ and $\mathbf{v} = (v_1, \dots, v_s)^T$, let $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}_2^s} = \int_{\Omega} \mathbf{u}^T(\mathbf{x}) \mathbf{v}(\mathbf{x}) dF_{\text{IMSE}}(\mathbf{x})$.

- $\text{IMSE}(\xi, \mathbf{g}) = \frac{n}{\sigma^2} E \left\| \mathbf{T}(\mathbf{g}^T \beta) - \mathbf{T}(\mathbf{g}^T \hat{\beta}) \right\|_{\mathcal{L}_2^s}^2$.

Proposition

$$\text{IMSE}(\xi, \mathbf{g}) = \text{tr}(\mathbf{M}^{-1} \mathbf{A}),$$

with

$$\mathbf{M} = \mathbf{M}_\xi = \frac{1}{n} \mathbf{G}^T \tilde{\Lambda}^{-1} \mathbf{G}, \quad \text{and } \mathbf{A} = \left(\langle \mathbf{T} g_i, \mathbf{T} g_j \rangle_{\mathcal{L}_2^s} \right)_{i,j=1}^k$$

Low Discrepancy Design Bounds IMSE

Theorem

Suppose that F_T is a probability distribution function defined on Ω , which may be different from F_{IMSE} , and that \mathcal{H} is a reproducing kernel Hilbert space of functions defined on Ω with reproducing kernel K . Consider the information matrix for F_T ,

$$\mathbf{M}_{F_T} = \int_{\Omega} (\mathbf{L}_x \mathbf{g}^T)^T \tilde{\Lambda}^{-1} (\mathbf{L}_x \mathbf{g}^T) dF_T(\mathbf{x}),$$

and suppose that the function $h_\alpha : \mathbf{x} \mapsto \alpha^T (\mathbf{M}_{F_T})^{-\frac{1}{2}} \mathbf{M}_x (\mathbf{M}_{F_T})^{-\frac{1}{2}} \alpha$ lies in \mathcal{H} for any $\alpha \in \mathbb{R}^k$. Define a variation over the basis \mathbf{g} as

$$V_{\mathbf{g}, F_T} = \sup_{\|\alpha\|_2 \leq 1} V(h_\alpha),$$

where V is the variation. Then it follows that the integrated mean square error is bounded above by

$$\text{IMSE}(\xi, \mathbf{g}) \leq \frac{\text{tr}(\mathbf{M}_{F_T}^{-1} \mathbf{A})}{1 - D_{F_T}(\xi) V_{\mathbf{g}, F_T}},$$

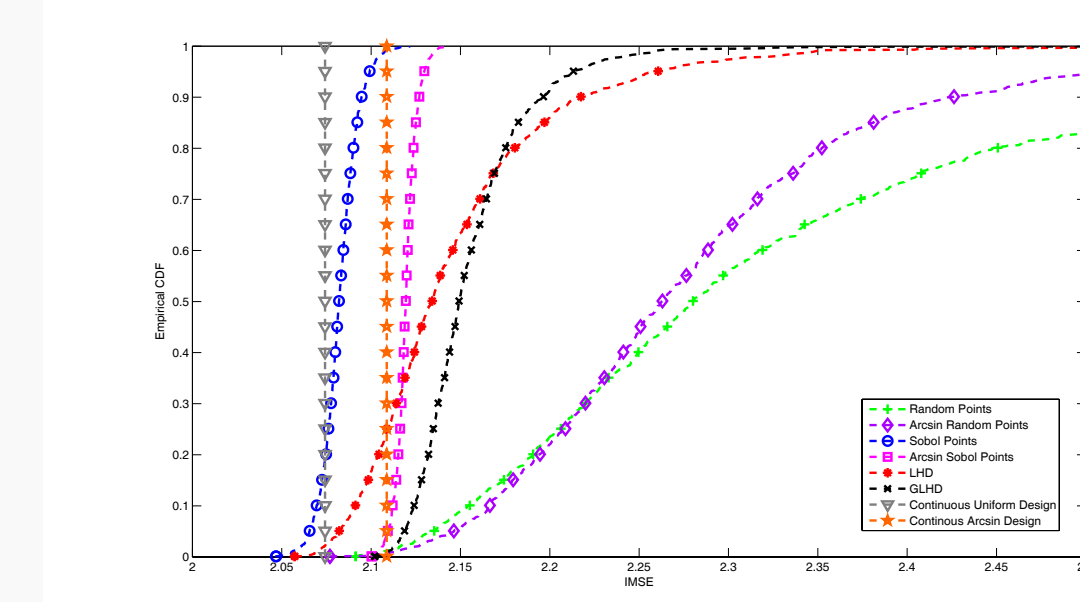
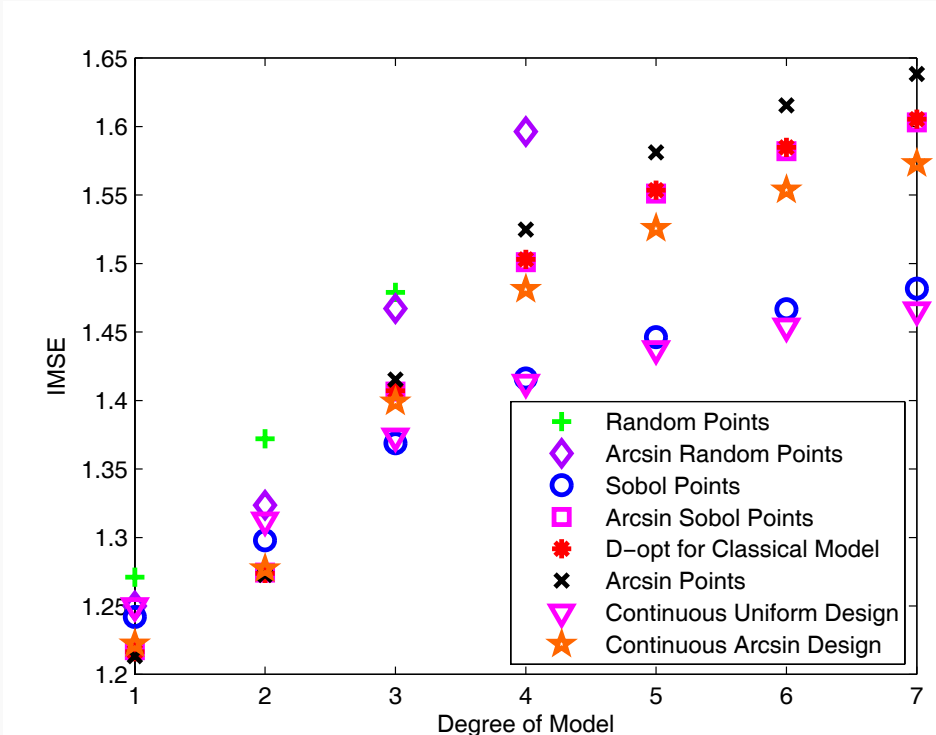
provided that $D_{F_T}(\xi) V_{\mathbf{g}, F_T} < 1$.

Remark

Note that, in some cases, there exists F_T , such that $\text{tr}(\mathbf{M}_{F_T}^{-1} \mathbf{A}) < \text{tr}(\mathbf{M}_{F_{\text{IMSE}}}^{-1} \mathbf{A})$. This means that choosing the design to match the distribution F_T will yield to smaller upper bound compared to choosing the design to match distribution F_{IMSE} .

Numerical Experiments on Low Discrepancy Design

- $\Omega = [-1, 1]$, polynomial basis, 16 sample points, repeated 1000 times, and $\Omega = [-1, 1] \times [-1, 1]$, orthogonal basis up to degree 4, 32 sample points, repeated 1000 times.



Semi-definite Programming with Gradient Information

- Define $\mathbf{G}^T(\mathbf{x}) = \mathbf{L}_x \mathbf{g}^T$, $\mathbf{F}(\mathbf{x}_j) = \mathbf{A}^{-\frac{1}{2}} \mathbf{G}(\mathbf{x}_j)$.

- Equivalent SDP model for I-optimal design:

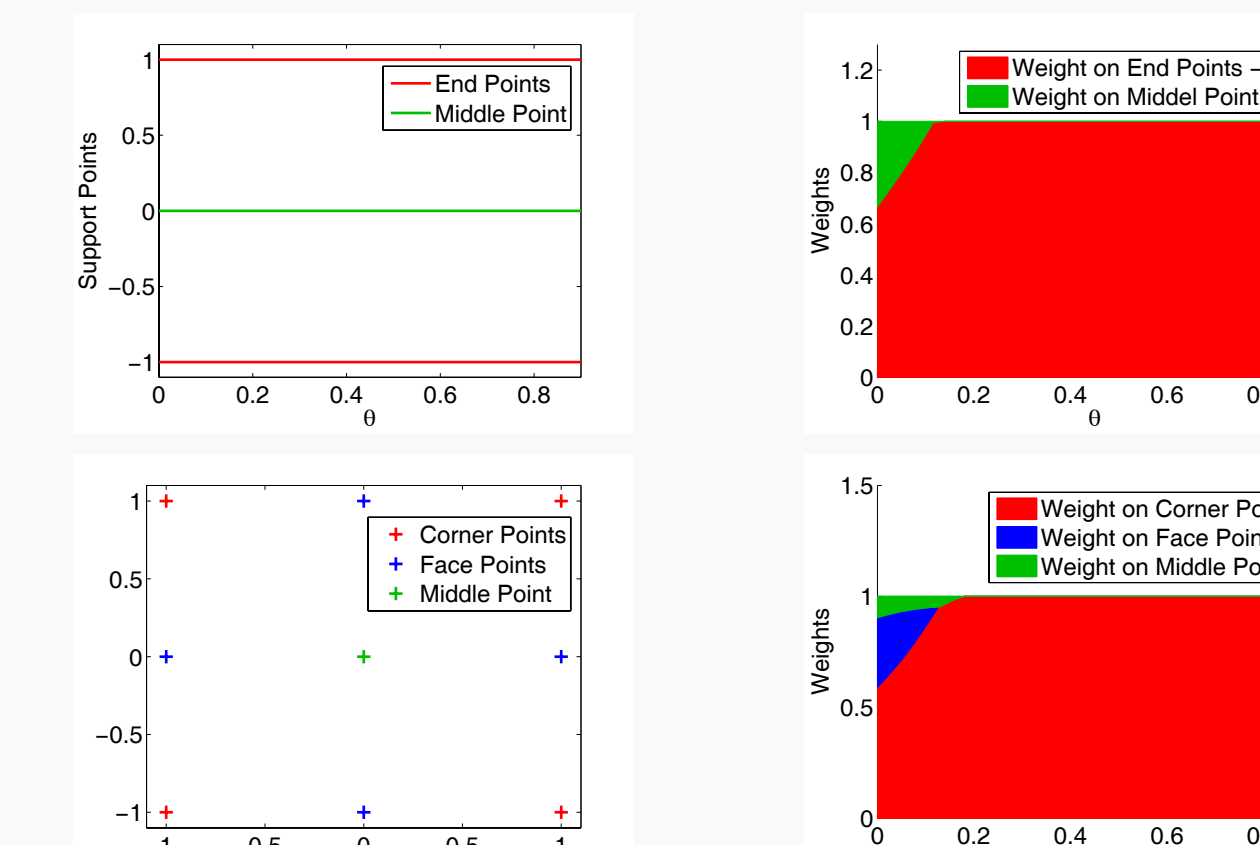
$$\begin{aligned} & \text{Minimize}_{w_j, \gamma} \mathbf{e}^T \boldsymbol{\gamma} \\ & \text{Subject to: } \begin{bmatrix} \sum_{j=1}^N w_j \mathbf{F}(\mathbf{v}_j) \tilde{\Lambda}^{-1} \mathbf{F}^T(\mathbf{v}_j) & \mathbf{I} \\ \mathbf{I} & \text{diag}(\boldsymbol{\gamma}) \end{bmatrix} \succeq 0, \\ & \mathbf{e}^T \mathbf{w} = 1, \\ & \mathbf{w} \geq 0. \end{aligned}$$

- SDP model for D-optimal design:

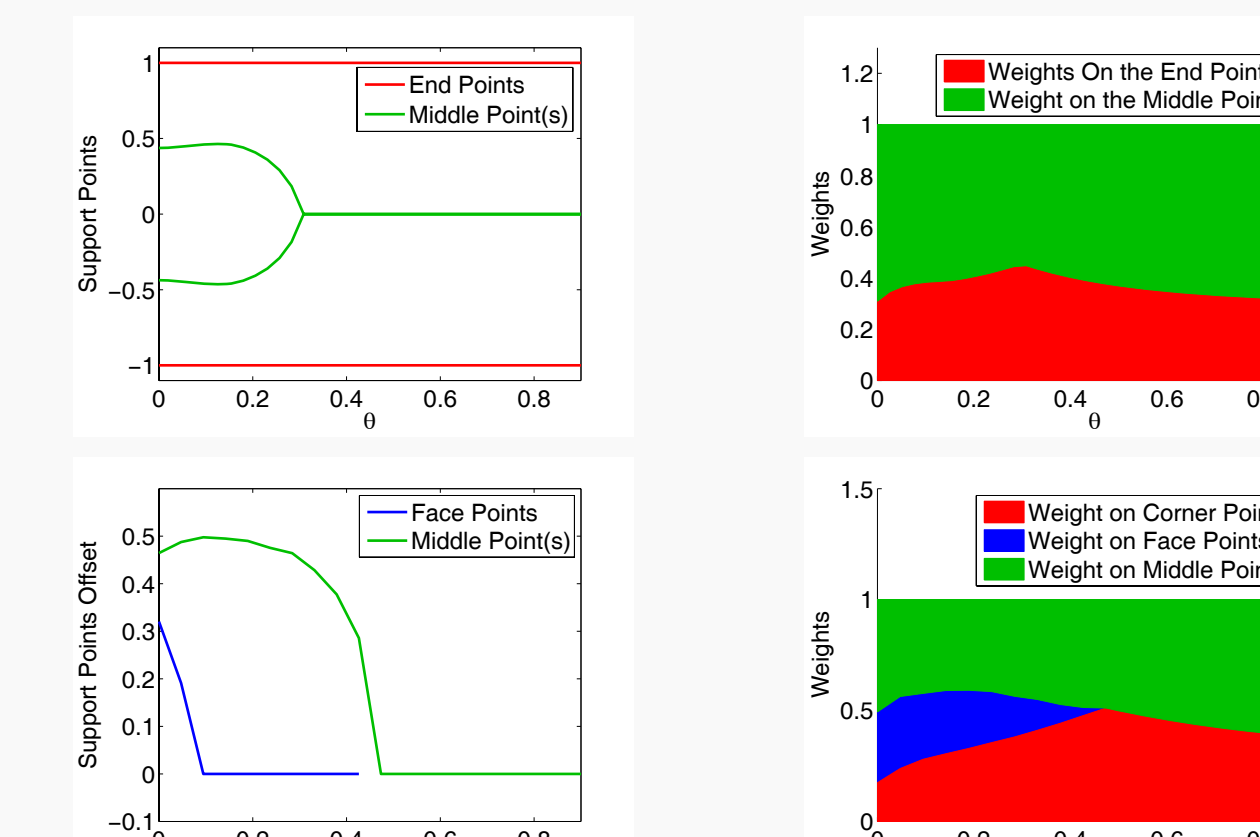
$$\begin{aligned} & \text{Minimize}_{w_j} -\log \det \left[\sum_{j=1}^K w_j \mathbf{G}(\mathbf{v}_j) \tilde{\Lambda}^{-1} \mathbf{G}^T(\mathbf{v}_j) \right] \\ & \text{Subject to: } \mathbf{e}^T \mathbf{w} = 1, \\ & \mathbf{w} \geq 0. \end{aligned}$$

Numerical Results for Semi-definite Programming

- D optimal design for quadratic model with variable in 1-d and 2-d:

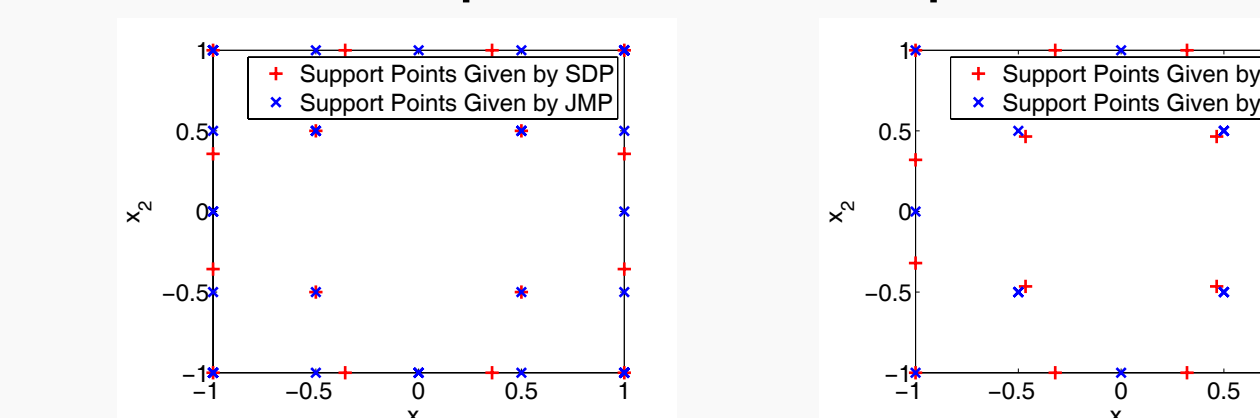


- I optimal design for cubic model with variable in 1-d and 2-d:



- JMP vs SDP (cubic 2-d model)

- Support Points for D-optimal and I-optimal design



- Efficiency ratio:

	I Efficiency Ratio				D Efficiency Ratio				
	1 - d Cubic		2 - d Cubic		1 - d Cubic		2 - d Cubic		
	Cont	Exact	Cont	Exact	Cont	Exact	Cont	Exact	
Number of Sample Points	12	1.0029	1.0004	1.1544	NA	1.0000	1.0000	1.0482	NA
	16	1.0185	0.9994	1.0359	1.0359	1.0000	1.0000	1.0496	1.0168
	20	1.0007	1.0003	1.0574	1.0471	1.0000	1.0000	1.0403	1.0359
	24	1.0030	1.0005	1.0349	1.0302	1.0000	1.0000	1.0171	1.0171
	28	1.0031	1.0003	1.0280	1.0280	1.0000	1.0000	1.0256	1.0256
	32	1.0004	1.0004	1.0295	1.0295	1.0000	1.0000	1.0249	1.0211
	36	1.0031	1.0006	1.0327	1.0327	1.0000	1.0000	1.0185	1.0137
	40	1.0009	1.0004	1.0277	1.0277	1.0000	1.0000	1.0201	1.0159

Acknowledgements

I would like to thank my advisor Fred J. Hickernell, and our collaborator Dr. Mihai Anitescu, Dr. Jieqiu Chen from Argonne National Lab for the assistance in my research.