

# Experimental design in nonlinear models: small-sample properties

LUC PRONZATO

CNRS, Université Côte d'Azur

# 1 Introduction

## Regression model

$$\underbrace{y_i = y(x_i)}_{\text{observation at } x_i} = \underbrace{\eta(x_i, \bar{\theta})}_{\text{model response at } x_i} + \underbrace{\varepsilon_i}_{\text{error}}$$

where the  $\varepsilon_i$  are i.i.d., with  $E\{\varepsilon_i\} = 0$  and  $E\{\varepsilon_i^2\} = \sigma^2$

# 1 Introduction

## Regression model

$$\underbrace{y_i = y(x_i)}_{\text{observation at } x_i} = \underbrace{\eta(x_i, \bar{\theta})}_{\text{model response at } x_i} + \underbrace{\varepsilon_i}_{\text{error}}$$

where the  $\varepsilon_i$  are i.i.d., with  $E\{\varepsilon_i\} = 0$  and  $E\{\varepsilon_i^2\} = \sigma^2$

$\mathbf{X}_n = (x_1, \dots, x_n)$  the design

$\mathbf{y} = (y_1, \dots, y_n)^\top$  the vector of observations

$\boldsymbol{\eta}(\theta) = (\eta(x_1, \theta), \dots, \eta(x_n, \theta))^\top$  the vector of model responses

$\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$  the errors ( $\rightarrow E\{\boldsymbol{\varepsilon}\} = \mathbf{0}$  and  $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$ )

$\bar{\theta}$  = true value of the model parameters  $\theta \in \mathbb{R}^p$

Least Squares (LS) estimator:  $\hat{\theta}^n = \arg \min_{\theta} \|\mathbf{y} - \boldsymbol{\eta}(\theta)\|^2$

**Information matrix** (at  $\theta^0$ , normalised — per observation)

$$\mathbf{M}(\mathbf{X}_n, \theta^0) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \eta(x_i, \theta)}{\partial \theta} \Big|_{\theta^0} \frac{\partial \eta(x_i, \theta)}{\partial \theta^\top} \Big|_{\theta^0} = \frac{1}{n} \frac{\partial \boldsymbol{\eta}^\top(\theta)}{\partial \theta} \Big|_{\theta^0} \frac{\partial \boldsymbol{\eta}(\theta)}{\partial \theta^\top} \Big|_{\theta^0}$$

(a  $p \times p$  matrix, with  $p = \dim(\theta)$ )

## A. Linear regression

$$\eta(x, \theta) = \mathbf{f}^\top(x)\theta \rightarrow \frac{\partial \eta^\top(\theta)}{\partial \theta} = \mathbf{F}^\top = (\mathbf{f}(x_1), \dots, \mathbf{f}(x_n)) \text{ and}$$

$$\hat{\theta}^n = (\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{y}$$

normalised information matrix:  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \mathbf{F}^\top \mathbf{F}$

→ choose the  $x_i$  such that  $\mathbf{M}_n$  has full rank

## A. Linear regression

$$\eta(x, \theta) = \mathbf{f}^\top(x)\theta \rightarrow \frac{\partial \eta^\top(\theta)}{\partial \theta} = \mathbf{F}^\top = (\mathbf{f}(x_1), \dots, \mathbf{f}(x_n)) \text{ and}$$

$$\hat{\theta}^n = (\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{y}$$

normalised information matrix:  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \mathbf{F}^\top \mathbf{F}$

→ choose the  $x_i$  such that  $\mathbf{M}_n$  has full rank

$$y_i = \mathbf{f}^\top(x_i)\bar{\theta} + \varepsilon_i \text{ for all } i, \text{ with } E\{\varepsilon_i\} = 0 \text{ and } E\{\varepsilon_i^2\} = \sigma^2$$

$$\Rightarrow E\{\hat{\theta}^n\} = \bar{\theta}$$

$$\Rightarrow \text{Var}(\hat{\theta}^n) = E\{(\hat{\theta}^n - \bar{\theta})(\hat{\theta}^n - \bar{\theta})^\top\} = \frac{\sigma^2}{n} \mathbf{M}_n^{-1}$$

→ choose the  $x_i$  to minimise a scalar function of  $\mathbf{M}_n^{-1}$   
or maximise a function  $\Phi(\mathbf{M}_n)$  (information function (Pukelsheim, 1993))

## A. Linear regression

$$\eta(x, \theta) = \mathbf{f}^\top(x)\theta \rightarrow \frac{\partial \eta^\top(\theta)}{\partial \theta} = \mathbf{F}^\top = (\mathbf{f}(x_1), \dots, \mathbf{f}(x_n)) \text{ and}$$

$$\hat{\theta}^n = (\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{y}$$

normalised information matrix:  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \mathbf{F}^\top \mathbf{F}$

→ choose the  $x_i$  such that  $\mathbf{M}_n$  has full rank

$$y_i = \mathbf{f}^\top(x_i)\bar{\theta} + \varepsilon_i \text{ for all } i, \text{ with } \mathbb{E}\{\varepsilon_i\} = 0 \text{ and } \mathbb{E}\{\varepsilon_i^2\} = \sigma^2$$

$$\Rightarrow \mathbb{E}\{\hat{\theta}^n\} = \bar{\theta}$$

$$\Rightarrow \text{Var}(\hat{\theta}^n) = \mathbb{E}\{(\hat{\theta}^n - \bar{\theta})(\hat{\theta}^n - \bar{\theta})^\top\} = \frac{\sigma^2}{n} \mathbf{M}_n^{-1}$$

→ choose the  $x_i$  to minimise a scalar function of  $\mathbf{M}_n^{-1}$   
or maximise a function  $\Phi(\mathbf{M}_n)$  (information function (Pukelsheim, 1993))

$$\text{Normal errors } \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n) \implies \hat{\theta}^n \sim \mathcal{N}(\bar{\theta}, \frac{\sigma^2}{n} \mathbf{M}^{-1}(\mathbf{X}_n))$$

→ no particular problem with *small data*

## B. Nonlinear regression

$\eta(x, \theta)$  nonlinear in  $\theta$

Under «standard» assumptions ( $\theta \in \Theta$  compact,  $\eta(x, \theta)$  continuous in  $\theta$  for any  $x \dots$ ), for a suitable sequence  $(x_i)$ ,

$\hat{\theta}^n \xrightarrow{\text{a.s.}} \bar{\theta}$  as  $n \rightarrow \infty$  (strong consistency) [but  $E\{\hat{\theta}^n\} \neq \bar{\theta}$  ( $\hat{\theta}^n$  is biased)]

## B. Nonlinear regression

$\eta(x, \theta)$  nonlinear in  $\theta$

Under «standard» assumptions ( $\theta \in \Theta$  compact,  $\eta(x, \theta)$  continuous in  $\theta$  for any  $x \dots$ ), for a suitable sequence  $(x_i)$ ,

$$\boxed{\hat{\theta}^n \xrightarrow{\text{a.s.}} \bar{\theta} \text{ as } n \rightarrow \infty} \text{ (strong consistency) [but } E\{\hat{\theta}^n\} \neq \bar{\theta} \text{ (}\hat{\theta}^n \text{ is biased)]}$$

normalised information matrix at  $\theta$ :  $\mathbf{M}_n(\theta) = \mathbf{M}(\mathbf{X}_n, \theta) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \eta(x_i, \theta)}{\partial \theta} \frac{\partial \eta(x_i, \theta)}{\partial \theta^\top}$

Under «standard» regularity assumptions ( $\eta(x, \theta)$  twice continuously differentiable w.r.t.  $\theta$  for any  $x \dots$ ), for a suitable sequence  $(x_i)$ ,

$$\boxed{\sqrt{n}(\hat{\theta}^n - \bar{\theta}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{M}^{-1}(\bar{\theta})) \text{ as } n \rightarrow \infty} \text{ (asymptotic normality)}$$

with  $\mathbf{M}(\theta) = \lim_{n \rightarrow \infty} \mathbf{M}_n(\theta)$

## B. Nonlinear regression

$\eta(x, \theta)$  nonlinear in  $\theta$

Under «standard» assumptions ( $\theta \in \Theta$  compact,  $\eta(x, \theta)$  continuous in  $\theta$  for any  $x \dots$ ), for a suitable sequence  $(x_i)$ ,

$$\hat{\theta}^n \xrightarrow{\text{a.s.}} \bar{\theta} \text{ as } n \rightarrow \infty \quad (\text{strong consistency}) \quad [\text{but } E\{\hat{\theta}^n\} \neq \bar{\theta} \text{ (}\hat{\theta}^n \text{ is biased)}]$$

normalised information matrix at  $\theta$ :  $\mathbf{M}_n(\theta) = \mathbf{M}(\mathbf{X}_n, \theta) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \eta(x_i, \theta)}{\partial \theta} \frac{\partial \eta(x_i, \theta)}{\partial \theta^\top}$

Under «standard» regularity assumptions ( $\eta(x, \theta)$  twice continuously differentiable w.r.t.  $\theta$  for any  $x \dots$ ), for a suitable sequence  $(x_i)$ ,

$$\sqrt{n}(\hat{\theta}^n - \bar{\theta}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{M}^{-1}(\bar{\theta})) \text{ as } n \rightarrow \infty \quad (\text{asymptotic normality})$$

with  $\mathbf{M}(\theta) = \lim_{n \rightarrow \infty} \mathbf{M}_n(\theta)$

→ choose the  $x_i$  to minimise a scalar function of  $\mathbf{M}_n^{-1}(\theta^0)$ ,  
or maximise a function  $\Phi(\mathbf{M}_n(\theta^0))$ , for a prior guess  $\theta^0$  (local design)

= **classical approach for DoE in nonlinear models**  
(based on asymptotic normality)

- 1) DoE for linear models (local design for nonlinear models, for a given  $\theta^0$ )
  - Which information function  $\Phi$ ?
  - How to construct an optimal design for  $\Phi$ ?

- 1) DoE for linear models (local design for nonlinear models, for a given  $\theta^0$ )
  - Which information function  $\Phi$ ?
  - How to construct an optimal design for  $\Phi$ ?
- 3,4,5,6) Small-sample issues

- 1) DoE for linear models (local design for nonlinear models, for a given  $\theta^0$ )
  - Which information function  $\Phi$ ?
  - How to construct an optimal design for  $\Phi$ ?
- 3,4,5,6) [Small-sample issues](#)
- 7) nonlocal DoE for nonlinear models (based on asymptotic normality)

## 2 DoE for linear models

### Design criterion $\Phi$

- A-optimality: minimise  $\text{trace}[\mathbf{M}^{-1}] \Leftrightarrow$  maximise  $\Phi(\mathbf{M}) = 1/\text{trace}[\mathbf{M}^{-1}]$   
 $\Leftrightarrow$  minimise sum of lengths<sup>2</sup> of axes of (asymptotic) confidence ellipsoids

## 2 DoE for linear models

### Design criterion $\Phi$

- **A-optimality:** minimise  $\text{trace}[\mathbf{M}^{-1}] \Leftrightarrow$  maximise  $\Phi(\mathbf{M}) = 1/\text{trace}[\mathbf{M}^{-1}]$   
 $\Leftrightarrow$  minimise sum of lengths<sup>2</sup> of axes of (asymptotic) confidence ellipsoids
- **E-optimality:** maximise  $\Phi(\mathbf{M}) = \lambda_{\min}(\mathbf{M})$   
 $\Leftrightarrow$  minimise longest axis of (asymptotic) confidence ellipsoids

## 2 DoE for linear models

### Design criterion $\Phi$

- **A-optimality:** minimise  $\text{trace}[\mathbf{M}^{-1}] \Leftrightarrow$  maximise  $\Phi(\mathbf{M}) = 1/\text{trace}[\mathbf{M}^{-1}]$   
 $\Leftrightarrow$  minimise sum of lengths<sup>2</sup> of axes of (asymptotic) confidence ellipsoids
- **E-optimality:** maximise  $\Phi(\mathbf{M}) = \lambda_{\min}(\mathbf{M})$   
 $\Leftrightarrow$  minimise longest axis of (asymptotic) confidence ellipsoids
- **D-optimality:** maximise  $\Phi(\mathbf{M}) = \det^{1/p}(\mathbf{M})$  [ $p = \dim(\theta)$ ]  
 $\Leftrightarrow$  minimise volume of (asymptotic) confidence ellipsoids  
 (proportional to  $1/\sqrt{\det(\mathbf{M})}$ )

Very much used:

- a *D*-optimum design is invariant by reparameterisation

$$\det \mathbf{M}'(\beta(\theta)) = \det \mathbf{M}(\theta) \det^{-2} \left( \frac{\partial \beta}{\partial \theta^T} \right)$$

- often leads to repeat the same experimental conditions (replications)

## Construction of an optimal design

### A/ Exact design

$n$  observations at  $\mathbf{X}_n = (x_1, \dots, x_n)$  in a regression model (for simplicity)  
Each design point  $x_i$  can be anything, e.g. a point in a subset  $\mathcal{X}$  of  $\mathbb{R}^d$

Maximise  $\Phi(\mathbf{M}_n)$  w.r.t.  $\mathbf{X}_n$  with  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$

## Construction of an optimal design

### A/ Exact design

$n$  observations at  $\mathbf{X}_n = (x_1, \dots, x_n)$  in a regression model (for simplicity)  
Each design point  $x_i$  can be anything, e.g. a point in a subset  $\mathcal{X}$  of  $\mathbb{R}^d$

Maximise  $\Phi(\mathbf{M}_n)$  w.r.t.  $\mathbf{X}_n$  with  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$

- If  $n \times d$  not too large  $\rightarrow$  standard algorithm  
[but there exist constraints ( $x_i \in \mathcal{X}$  for all  $i$ ), local optimas...]

## Construction of an optimal design

### A/ Exact design

$n$  observations at  $\mathbf{X}_n = (x_1, \dots, x_n)$  in a regression model (for simplicity)  
 Each design point  $x_i$  can be anything, e.g. a point in a subset  $\mathcal{X}$  of  $\mathbb{R}^d$

Maximise  $\Phi(\mathbf{M}_n)$  w.r.t.  $\mathbf{X}_n$  with  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$

- If  $n \times d$  not too large  $\rightarrow$  standard algorithm  
 [but there exist constraints ( $x_i \in \mathcal{X}$  for all  $i$ ), local optimas...]
- Otherwise  $\rightarrow$  take the particular form of the problem into account

Exchange methods: (Fedorov, 1972; Mitchell, 1974)

At iteration  $k$ , exchange **one support point**  $x_j$  by a better one  $x^*$  in  $\mathcal{X}$  in the sense of  $\Phi(\cdot)$

$$\mathbf{X}_n^k = (x_1, \dots, \boxed{\begin{array}{c} x_j \\ \updownarrow \\ x^* \end{array}}, \dots, x_n)$$

## Construction of an optimal design

### A/ Exact design

$n$  observations at  $\mathbf{X}_n = (x_1, \dots, x_n)$  in a regression model (for simplicity)  
 Each design point  $x_i$  can be anything, e.g. a point in a subset  $\mathcal{X}$  of  $\mathbb{R}^d$

Maximise  $\Phi(\mathbf{M}_n)$  w.r.t.  $\mathbf{X}_n$  with  $\mathbf{M}_n = \mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$

- If  $n \times d$  not too large  $\rightarrow$  standard algorithm  
 [but there exist constraints ( $x_i \in \mathcal{X}$  for all  $i$ ), local optimas...]
- Otherwise  $\rightarrow$  take the particular form of the problem into account

Exchange methods: (Fedorov, 1972; Mitchell, 1974)

At iteration  $k$ , exchange **one support point**  $x_j$  by a better one  $x^*$  in  $\mathcal{X}$  in the sense of  $\Phi(\cdot)$

$$\mathbf{X}_n^k = (x_1, \dots, \boxed{\begin{array}{c} x_j \\ \updownarrow \\ x^* \end{array}}, \dots, x_n)$$

- Branch and bound (Welch, 1982), rounding an optimal design measure (Pukelsheim and Reider, 1992)

## B/ Design measures: approximate design theory

(Chernoff, 1953; Kiefer and Wolfowitz, 1960; Fedorov, 1972; Silvey, 1980; Pázman, 1986; Pukelsheim, 1993; Fedorov and Leonov, 2014)

$$\mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

[with  $\mathbf{M}(\mathbf{X}_n) = \mathbf{M}(\mathbf{X}_n, \theta^0)$  and  $\mathbf{f}(x_i) = \frac{\partial \eta(x_i, \theta)}{\partial \theta} \Big|_{\theta^0}$  in a nonlinear model]

The additive form is essential (comes from the independence of observations)

## B/ Design measures: approximate design theory

(Chernoff, 1953; Kiefer and Wolfowitz, 1960; Fedorov, 1972; Silvey, 1980; Pázman, 1986; Pukelsheim, 1993; Fedorov and Leonov, 2014)

$$\mathbf{M}(\mathbf{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

[with  $\mathbf{M}(\mathbf{X}_n) = \mathbf{M}(\mathbf{X}_n, \theta^0)$  and  $\mathbf{f}(x_i) = \frac{\partial \eta(x_i, \theta)}{\partial \theta} \Big|_{\theta^0}$  in a nonlinear model]

The additive form is essential (comes from the independence of observations)

Repeat  $r_i$  observations at the same  $x_i \in \mathcal{X}$  ( $r_i$  replications):

→ only  $m \leq n$  different  $x_i$

$$\mathbf{M}(\mathbf{X}_n) = \sum_{i=1}^m \frac{r_i}{n} \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

- $\frac{r_i}{n}$  = proportion of observations collected at  $x_i$
- = «percentage of experimental effort» at  $x_i$
- = weight  $w_i$  of support point  $x_i$

$$\mathbf{M}(\mathbf{X}_n) = \sum_{i=1}^m w_i \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

→ design  $\mathbf{X}_n \Leftrightarrow \left\{ \begin{array}{ccc} x_1 & \cdots & x_m \\ w_1 & \cdots & w_m \end{array} \right\}$  with  $\sum_{i=1}^m w_i = 1$

→ normalised discrete distribution on  $\mathcal{X}$ , with constraints  $r_i/n = w_i$

$$\mathbf{M}(\mathbf{X}_n) = \sum_{i=1}^m w_i \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

→ design  $\mathbf{X}_n \Leftrightarrow \left\{ \begin{array}{ccc} x_1 & \cdots & x_m \\ w_1 & \cdots & w_m \end{array} \right\}$  with  $\sum_{i=1}^m w_i = 1$

→ normalised discrete distribution on  $\mathcal{X}$ , with constraints  $r_i/n = w_i$

Release the constraints →  $w_i \geq 0$  with  $\sum_{i=1}^m w_i = 1$

→  $\xi =$  discrete probability measure on  $\mathcal{X}$  (= design space)  
 with support points  $x_i$  and associated weights  $w_i$   
 = «approximate design»

$$\mathbf{M}(\mathbf{X}_n) = \sum_{i=1}^m w_i \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

→ design  $\mathbf{X}_n \Leftrightarrow \left\{ \begin{array}{ccc} x_1 & \cdots & x_m \\ w_1 & \cdots & w_m \end{array} \right\}$  with  $\sum_{i=1}^m w_i = 1$

→ normalised discrete distribution on  $\mathcal{X}$ , with constraints  $r_i/n = w_i$

Release the constraints →  $w_i \geq 0$  with  $\sum_{i=1}^m w_i = 1$

→  $\xi =$  discrete probability measure on  $\mathcal{X}$  (= design space)  
with support points  $x_i$  and associated weights  $w_i$   
= «approximate design»

More general expression:  $\xi =$  any probability measure on  $\mathcal{X}$  ( $\int_{\mathcal{X}} \xi(dx) = 1$ )

$$\mathbf{M}(\xi) = \int_{\mathcal{X}} \mathbf{f}(x) \mathbf{f}^\top(x) \xi(dx)$$

$$\mathbf{M}(\mathbf{X}_n) = \sum_{i=1}^m w_i \mathbf{f}(x_i) \mathbf{f}^\top(x_i)$$

→ design  $\mathbf{X}_n \Leftrightarrow \left\{ \begin{array}{ccc} x_1 & \cdots & x_m \\ w_1 & \cdots & w_m \end{array} \right\}$  with  $\sum_{i=1}^m w_i = 1$

→ normalised discrete distribution on  $\mathcal{X}$ , with constraints  $r_i/n = w_i$

Release the constraints →  $w_i \geq 0$  with  $\sum_{i=1}^m w_i = 1$

→  $\xi =$  discrete probability measure on  $\mathcal{X}$  (= design space)  
with support points  $x_i$  and associated weights  $w_i$   
= «approximate design»

More general expression:  $\xi =$  any probability measure on  $\mathcal{X}$  ( $\int_{\mathcal{X}} \xi(dx) = 1$ )

$$\mathbf{M}(\xi) = \int_{\mathcal{X}} \mathbf{f}(x) \mathbf{f}^\top(x) \xi(dx)$$

$\mathbf{M}(\xi) \in$  convex closure of the set of rank 1 matrices  $\mathbf{f}(x) \mathbf{f}^\top(x)$

$\mathbf{M}(\xi)$  is symmetric  $p \times p$ , belongs to a  $\frac{p(p+1)}{2}$ -dimensional space

Caratheodory Theorem → for any  $\xi$ , there exists a discrete probability measure  $\xi_d$   
with  $\frac{p(p+1)}{2} + 1$  support points at most, such that  $\mathbf{M}(\xi_d) = \mathbf{M}(\xi)$   
(true in particular for the optimum design)

Maximise  $\Phi[\mathbf{M}(\xi)]$ ,  $\Phi(\cdot)$  concave (e.g.,  $A$ ,  $E$ ,  $D$ -optimality) and  $\mathbf{M}(\xi)$  linear in  $\xi$   
→ convex programming

Usually,  $\mathcal{X}$  is first discretised  
→ optimise a vector of weights  
(possibly high dimensional, but the solution is sparse)

Typical algorithm when  $\Phi$  is differentiable ( $A$ ,  $D$ -optimality):  
Frank-Wolfe conditional gradient (called vertex-direction algorithm in DoE), with predefined (Wynn, 1970) or optimal (Fedorov, 1972) step-size  
[but there exist more efficient methods]

More difficult if  $\Phi$  not differentiable ( $E$ -optimality), but feasible

## Application to models with complete product-type interactions

Single factor models:  $\eta_k(x, \theta^{(k)}) \triangleq [\mathbf{f}^{(k)}(x)]^\top \theta^{(k)}$

global model for  $d$  factors  $\mathbf{x} = (\{\mathbf{x}\}_1, \{\mathbf{x}\}_2, \dots, \{\mathbf{x}\}_d)^\top$ :

$$\eta(\mathbf{x}, \gamma) = [\mathbf{f}_1(\{\mathbf{x}\}_1) \otimes \dots \otimes \mathbf{f}_d(\{\mathbf{x}\}_d)]^\top \gamma$$

## Application to models with complete product-type interactions

Single factor models:  $\eta_k(\mathbf{x}, \theta^{(k)}) \triangleq [\mathbf{f}^{(k)}(\mathbf{x})]^\top \theta^{(k)}$

global model for  $d$  factors  $\mathbf{x} = (\{\mathbf{x}\}_1, \{\mathbf{x}\}_2, \dots, \{\mathbf{x}\}_d)^\top$ :

$$\eta(\mathbf{x}, \gamma) = [\mathbf{f}_1(\{\mathbf{x}\}_1) \otimes \dots \otimes \mathbf{f}_d(\{\mathbf{x}\}_d)]^\top \gamma$$

In particular, if  $\eta_k =$  polynomial of degree  $d_k$  ( $\dim(\theta^{(k)}) = p_k = 1 + d_k$ ),

$$\eta = \text{polynomial with total degree } \sum_{k=1}^d d_k \quad (\dim(\gamma) = \prod_{k=1}^d p_k)$$

Example:

$$\begin{aligned} \mathbf{f}^\top(\mathbf{x})\gamma &= (\theta_0^{(1)} + \theta_1^{(1)}\{\mathbf{x}\}_1 + \theta_2^{(1)}\{\mathbf{x}\}_1^2) \times (\theta_0^{(2)} + \theta_1^{(2)}\{\mathbf{x}\}_2 + \theta_2^{(2)}\{\mathbf{x}\}_2^2) \\ &= \gamma_0 + \gamma_1\{\mathbf{x}\}_1 + \gamma_2\{\mathbf{x}\}_2 + \gamma_{12}\{\mathbf{x}\}_1\{\mathbf{x}\}_2 + \gamma_{11}\{\mathbf{x}\}_1^2 + \gamma_{22}\{\mathbf{x}\}_2^2 \\ &\quad + \gamma_{112}\{\mathbf{x}\}_1^2\{\mathbf{x}\}_2 + \gamma_{122}\{\mathbf{x}\}_1\{\mathbf{x}\}_2^2 + \gamma_{1122}\{\mathbf{x}\}_1^2\{\mathbf{x}\}_2^2 \end{aligned}$$

## Application to models with complete product-type interactions

Single factor models:  $\eta_k(\mathbf{x}, \theta^{(k)}) \triangleq [\mathbf{f}^{(k)}(\mathbf{x})]^\top \theta^{(k)}$

global model for  $d$  factors  $\mathbf{x} = (\{\mathbf{x}\}_1, \{\mathbf{x}\}_2, \dots, \{\mathbf{x}\}_d)^\top$ :

$$\eta(\mathbf{x}, \gamma) = [\mathbf{f}_1(\{\mathbf{x}\}_1) \otimes \dots \otimes \mathbf{f}_d(\{\mathbf{x}\}_d)]^\top \gamma$$

In particular, if  $\eta_k = \text{polynomial of degree } d_k$  ( $\dim(\theta^{(k)}) = p_k = 1 + d_k$ ),

$$\eta = \text{polynomial with total degree } \sum_{k=1}^d d_k \quad (\dim(\gamma) = \prod_{k=1}^d p_k)$$

Example:

$$\begin{aligned} \mathbf{f}^\top(\mathbf{x})\gamma &= (\theta_0^{(1)} + \theta_1^{(1)}\{\mathbf{x}\}_1 + \theta_2^{(1)}\{\mathbf{x}\}_1^2) \times (\theta_0^{(2)} + \theta_1^{(2)}\{\mathbf{x}\}_2 + \theta_2^{(2)}\{\mathbf{x}\}_2^2) \\ &= \gamma_0 + \gamma_1\{\mathbf{x}\}_1 + \gamma_2\{\mathbf{x}\}_2 + \gamma_{12}\{\mathbf{x}\}_1\{\mathbf{x}\}_2 + \gamma_{11}\{\mathbf{x}\}_1^2 + \gamma_{22}\{\mathbf{x}\}_2^2 \\ &\quad + \gamma_{112}\{\mathbf{x}\}_1^2\{\mathbf{x}\}_2 + \gamma_{122}\{\mathbf{x}\}_1\{\mathbf{x}\}_2^2 + \gamma_{1122}\{\mathbf{x}\}_1^2\{\mathbf{x}\}_2^2 \end{aligned}$$

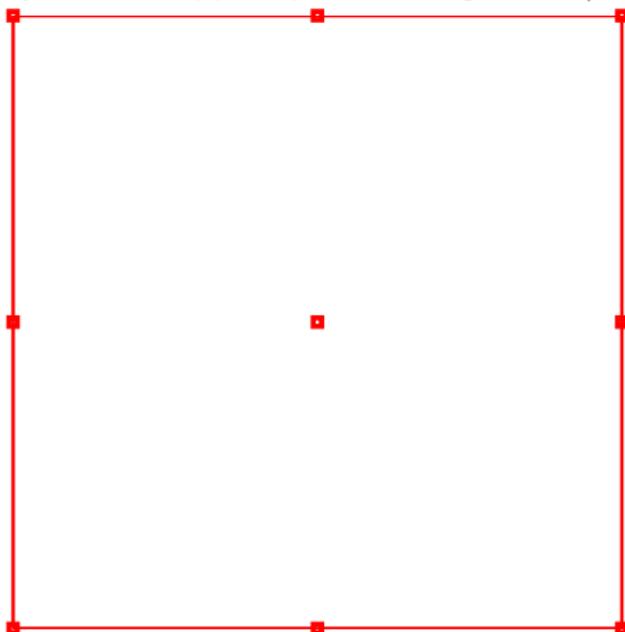
$D$ ,  $A$  and  $E$ -optimal design measure = tensor product of the  $d$  optimal design measures (Schwabe, 1996)

(true for any complete product-type interaction model — not only for polynomials)

Polynomial with degree  $k$ :  $D$ -optimal design supported on  $k + 1$  points  
 (on  $[-1, 1]$ : roots of  $(1 - t^2)P'_k(t)$ , with  $P_k(t) \triangleq k$ -th de Legendre polynomial),  
 all with the same weight  $1/(k + 1)$

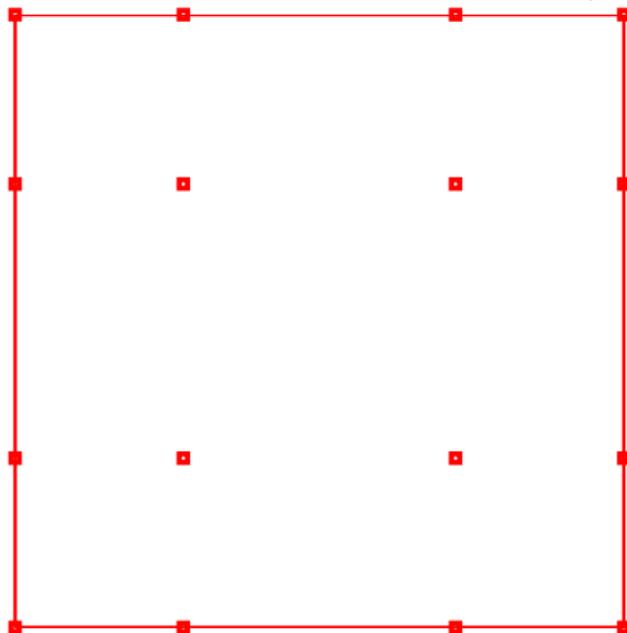
dimension 2,  $d_1 = d_2 = 2$

$\xi^*$  has 9 support points, weights =  $1/9$



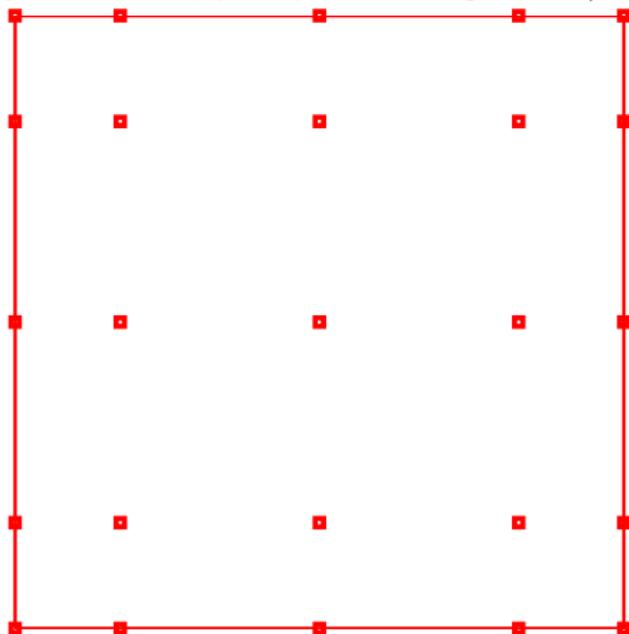
Polynomial with degree  $k$ :  $D$ -optimal design supported on  $k + 1$  points  
 (on  $[-1, 1]$ : roots of  $(1 - t^2)P'_k(t)$ , with  $P_k(t) \triangleq k$ -th de Legendre polynomial),  
 all with the same weight  $1/(k + 1)$

dimension 2,  $d_1 = d_2 = 3$   
 $\xi^*$  has 16 support points, weights  $=1/16$



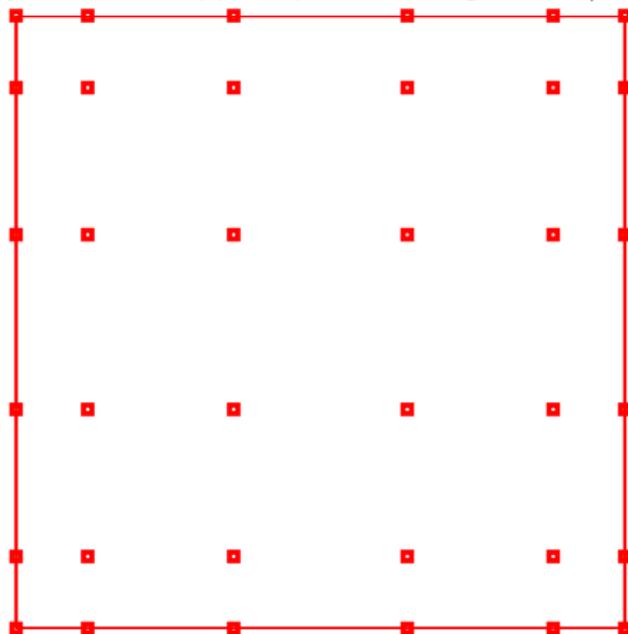
Polynomial with degree  $k$ :  $D$ -optimal design supported on  $k + 1$  points  
 (on  $[-1, 1]$ : roots of  $(1 - t^2)P'_k(t)$ , with  $P_k(t) \triangleq k$ -th de Legendre polynomial),  
 all with the same weight  $1/(k + 1)$

dimension 2,  $d_1 = d_2 = 4$   
 $\xi^*$  has 25 support points, weights  $=1/25$



Polynomial with degree  $k$ :  $D$ -optimal design supported on  $k + 1$  points  
 (on  $[-1, 1]$ : roots of  $(1 - t^2)P'_k(t)$ , with  $P_k(t) \triangleq k$ -th de Legendre polynomial),  
 all with the same weight  $1/(k + 1)$

dimension 2,  $d_1 = d_2 = 5$   
 $\xi^*$  has 36 support points, weights  $=1/36$



## Application to models with intercept, no interaction

Single factor models:  $\eta_k(\mathbf{x}, \theta^{(k)}) \triangleq \theta_0^{(k)} + \sum_{i=1}^{d_k} \theta_i^{(k)} f_i^{(k)}(\mathbf{x})$

global model for  $d$  factors:  $\eta(\mathbf{x}, \gamma) = \theta_0 + \sum_{k=1}^d \sum_{i=1}^{d_k} \theta_i^{(k)} f_i^{(k)}(\{\mathbf{x}\}_k)$

## Application to models with intercept, no interaction

Single factor models:  $\eta_k(\mathbf{x}, \theta^{(k)}) \triangleq \theta_0^{(k)} + \sum_{i=1}^{d_k} \theta_i^{(k)} f_i^{(k)}(\mathbf{x})$

global model for  $d$  factors:  $\eta(\mathbf{x}, \gamma) = \theta_0 + \sum_{k=1}^d \sum_{i=1}^{d_k} \theta_i^{(k)} f_i^{(k)}(\{\mathbf{x}\}_k)$

In particular, if  $\eta_k =$  polynomial of degree  $d_k$  ( $\dim(\theta^{(k)}) = p_k = 1 + d_k$ ),

$\eta =$  polynomial with total degree  $\max_k^d d_k$  ( $\dim(\gamma) = 1 + \sum_{k=1}^d d_k$ )

Example:

$$\begin{aligned} \mathbf{f}^\top(\mathbf{x})\gamma &= (\theta_0^{(1)} + \theta_1^{(1)}\{\mathbf{x}\}_1 + \theta_2^{(1)}\{\mathbf{x}\}_1^2) + (\theta_0^{(2)} + \theta_1^{(2)}\{\mathbf{x}\}_2 + \theta_2^{(2)}\{\mathbf{x}\}_2^2) \\ &= \gamma_0 + \gamma_1\{\mathbf{x}\}_1 + \gamma_2\{\mathbf{x}\}_2 + \gamma_{11}\{\mathbf{x}\}_1^2 + \gamma_{22}\{\mathbf{x}\}_2^2 \end{aligned}$$

$D$ -optimal design measure = tensor product of  $d$   $D$ -optimal measures (Schwabe, 1996)

## Application to models with intercept, no interaction

Single factor models:  $\eta_k(\mathbf{x}, \theta^{(k)}) \triangleq \theta_0^{(k)} + \sum_{i=1}^{d_k} \theta_i^{(k)} f_i^{(k)}(\mathbf{x})$

global model for  $d$  factors:  $\eta(\mathbf{x}, \gamma) = \theta_0 + \sum_{k=1}^d \sum_{i=1}^{d_k} \theta_i^{(k)} f_i^{(k)}(\{\mathbf{x}\}_k)$

In particular, if  $\eta_k =$  polynomial of degree  $d_k$  ( $\dim(\theta^{(k)}) = p_k = 1 + d_k$ ),

$\eta =$  polynomial with total degree  $\max_k d_k$  ( $\dim(\gamma) = 1 + \sum_{k=1}^d d_k$ )

Example:

$$\begin{aligned} \mathbf{f}^\top(\mathbf{x})\gamma &= (\theta_0^{(1)} + \theta_1^{(1)}\{\mathbf{x}\}_1 + \theta_2^{(1)}\{\mathbf{x}\}_1^2) + (\theta_0^{(2)} + \theta_1^{(2)}\{\mathbf{x}\}_2 + \theta_2^{(2)}\{\mathbf{x}\}_2^2) \\ &= \gamma_0 + \gamma_1\{\mathbf{x}\}_1 + \gamma_2\{\mathbf{x}\}_2 + \gamma_{11}\{\mathbf{x}\}_1^2 + \gamma_{22}\{\mathbf{x}\}_2^2 \end{aligned}$$

$D$ -optimal design measure = tensor product of  $d$   $D$ -optimal measures (Schwabe, 1996)

Hardly manageable in high dimension

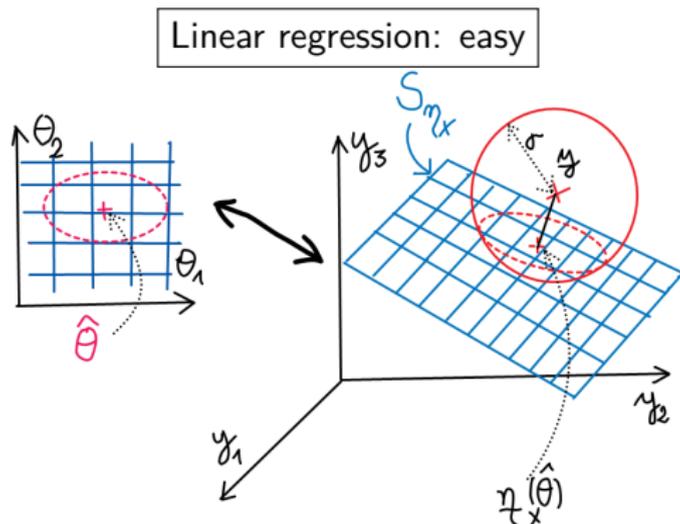
( $d$  polynomials of degree  $k \rightsquigarrow (k+1)^d$  support points),

but maybe a useful message for Gaussian Process models and kriging:

→ put more points along the boundaries than deeply inside

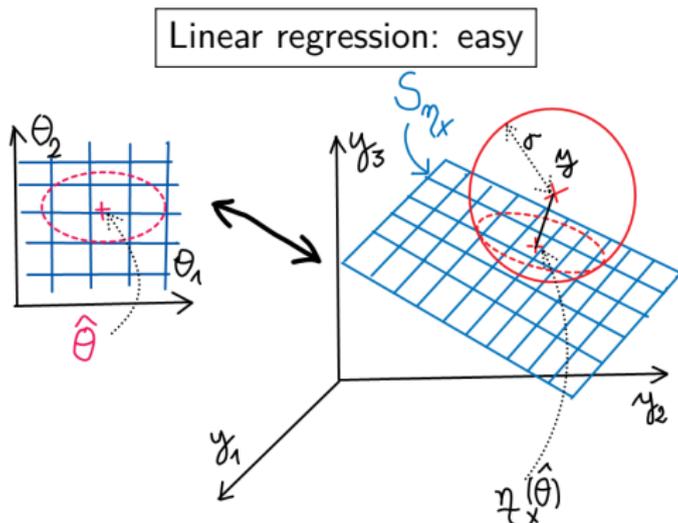
(Dette and Pepelyshev, 2010)

## 3 Linear and nonlinear models



The expectation surface  $S_{\eta} = \{\eta(\theta) = (\eta(x_1, \theta), \dots, \eta(x_n, \theta))^T : \theta \in \mathbb{R}^p\}$  is flat and linearly parameterised

# 3 Linear and nonlinear models

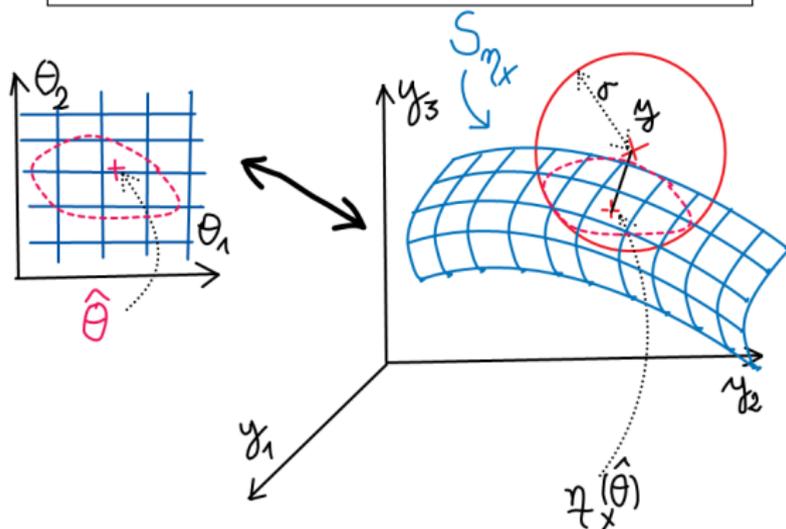


The expectation surface  $\mathcal{S}_\eta = \{\boldsymbol{\eta}(\theta) = (\eta(x_1, \theta), \dots, \eta(x_n, \theta))^T : \theta \in \mathbb{R}^p\}$  is flat and linearly parameterised

$\mathbf{M}(\mathbf{X}_n, \theta)$  does not depend on  $\theta$

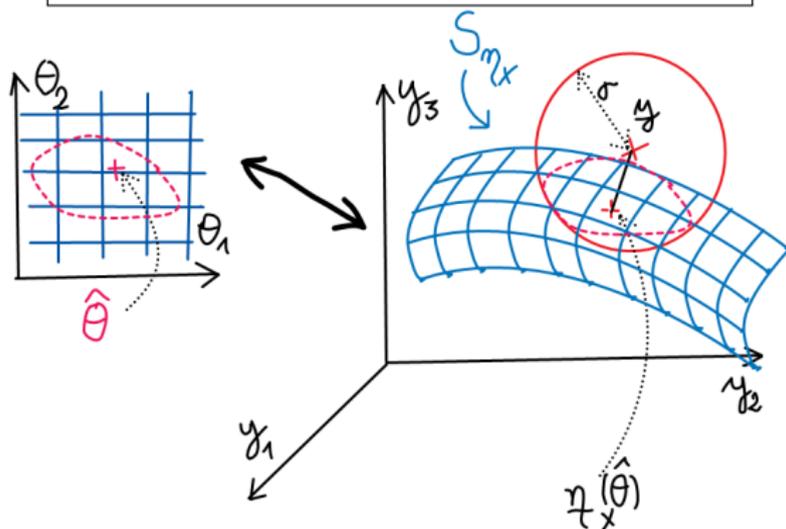
Normal errors  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n) \implies \hat{\boldsymbol{\theta}}^n \sim \mathcal{N}(\bar{\boldsymbol{\theta}}, \frac{\sigma^2}{n} \mathbf{M}^{-1}(\mathbf{X}_n))$

Nonlinear regression: maybe a bit tricky...



$S_{\eta}$  is curved (intrinsic curvature) and nonlinearly parameterised (parametric curvature) (Bates and Watts, 1980)

Nonlinear regression: maybe a bit tricky...



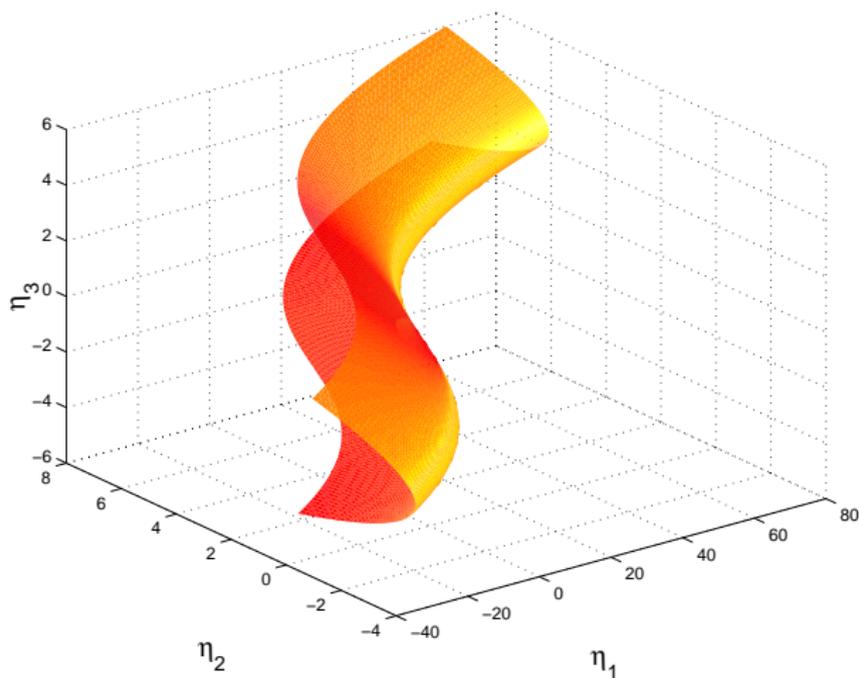
$S_n$  is curved (intrinsic curvature) and nonlinearly parameterised (parametric curvature) (Bates and Watts, 1980)

$M(\mathbf{X}_n, \theta)$  does depend on  $\theta$

Normal errors  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n) \implies \hat{\theta}^n \sim ?$

$$\text{Ex: } \eta(\mathbf{x}, \theta) = \theta_1 \{\mathbf{x}\}_1 + \theta_1^3 (1 - \{\mathbf{x}\}_1) + \theta_2 \{\mathbf{x}\}_2 + \theta_2^2 (1 - \{\mathbf{x}\}_2)$$

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3), \quad \mathbf{x}_1 = (0 \ 1), \quad \mathbf{x}_2 = (1 \ 0), \quad \mathbf{x}_3 = (1 \ 1), \quad \theta \in [-3, 4] \times [-2, 2]$$



## Two major difficulties with nonlinear models:

- ① Asymptotically ( $n \rightarrow \infty$ ) — or if  $\sigma^2$  small enough — all seems fine:  
use linear approximation

But the distribution of  $\hat{\theta}^n$  may be far from normal for small  $n$  (or for  $\sigma^2$  large)  
    ▶ small-sample properties

## Two major difficulties with nonlinear models:

- ❶ Asymptotically ( $n \rightarrow \infty$ ) — or if  $\sigma^2$  small enough — all seems fine:  
use linear approximation

But the distribution of  $\hat{\theta}^n$  may be far from normal for small  $n$  (or for  $\sigma^2$  large)  
    ➡ small-sample properties

- ❷ Everything is local (depends on  $\theta$ ): if we linearise, **where do we linearise?**  
(choice of a nominal value  $\theta^0$ )  
    ➡ nonlocal optimum design

## 4 Small-sample properties

Asymptotically ( $n \rightarrow \infty$ )  $\implies$   $\boxed{\sqrt{n}(\hat{\theta}^n - \bar{\theta}) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta}))}$

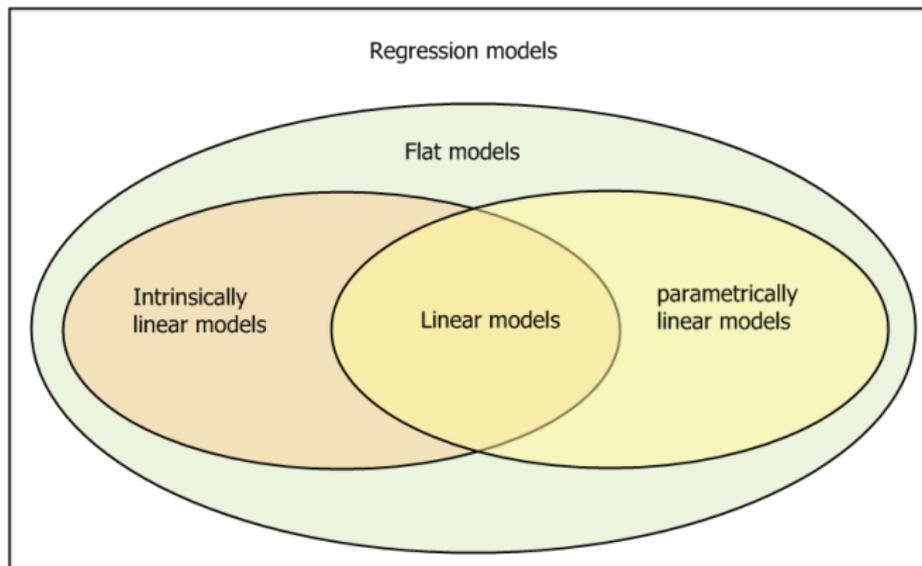
but what is the small sample precision?

## 4 Small-sample properties

Asymptotically ( $n \rightarrow \infty$ )  $\implies \sqrt{n}(\hat{\theta}^n - \bar{\theta}) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta}))$

but what is the small sample precision?

A classification of regression models (Pázman, 1993)



→ Consider projection on the expectation surface  $\mathbb{S}_\eta$ :

►  $\mathbf{P}_{\theta^0}$  = orthogonal projector onto the tangent space to  $\mathbb{S}_\eta$  at  $\boldsymbol{\eta}(\theta^0)$ :

$$\mathbf{P}_{\theta^0} = \frac{1}{n} \frac{\partial \boldsymbol{\eta}(\theta)}{\partial \boldsymbol{\theta}^\top} \Big|_{\theta^0} \mathbf{M}^{-1}(\mathbf{X}_n, \theta^0) \frac{\partial \boldsymbol{\eta}^\top(\theta)}{\partial \boldsymbol{\theta}} \Big|_{\theta^0}$$

(an  $n \times n$  matrix, depends on  $\mathbf{X}_n$ )

Bates and Watts (1980) intrinsic and parametric-effect measures of nonlinearity:

$$C_{int}(\mathbf{X}_n, \theta; \mathbf{u}) = \frac{\|[\mathbf{I}_n - \mathbf{P}_\theta] \sum_{i,j=1}^p u_i \mathbf{H}_{ij}(\theta) u_j\|}{n \mathbf{u}^\top \mathbf{M}(\mathbf{X}_n, \theta) \mathbf{u}}$$

$$C_{par}(\mathbf{X}_n, \theta; \mathbf{u}) = \frac{\|\mathbf{P}_\theta \sum_{i,j=1}^p u_i \mathbf{H}_{ij}(\theta) u_j\|}{n \mathbf{u}^\top \mathbf{M}(\mathbf{X}_n, \theta) \mathbf{u}}$$

with  $\mathbf{u} \in \mathbb{R}^p$  and  $\mathbf{H}_{ij}(\theta) = \frac{\partial^2 \boldsymbol{\eta}(\theta)}{\partial \theta_i \partial \theta_j}$

Intrinsic curvature:  $C_{int}(\mathbf{X}_n, \theta) = \sup_{\mathbf{u} \in \mathbb{R}^p \setminus \{\mathbf{0}\}} C_{int}(\mathbf{X}_n, \theta; \mathbf{u})$

Parametric curvature:  $C_{par}(\mathbf{X}_n, \theta) = \sup_{\mathbf{u} \in \mathbb{R}^p \setminus \{\mathbf{0}\}} C_{par}(\mathbf{X}_n, \theta; \mathbf{u})$

## Intrinsically linear models

- ▶ The expectation surface  $\mathbb{S}_\eta = \{\boldsymbol{\eta}(\boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathbb{R}^p\}$  is flat (plane)
  - intrinsic curvature  $\equiv 0$
- ▶ There exists a reparameterisation (continuously differentiable) that makes the model linear
- ▶  $\mathbf{P}_\theta \mathbf{H}_{ij}(\boldsymbol{\theta}) = \mathbf{H}_{ij}(\boldsymbol{\theta})$ , where  $\mathbf{H}_{ij}(\boldsymbol{\theta}) = \frac{\partial^2 \boldsymbol{\eta}(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}$

## Intrinsically linear models

- ▶ The expectation surface  $\mathbb{S}_\eta = \{\boldsymbol{\eta}(\boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathbb{R}^p\}$  is flat (plane)
  - intrinsic curvature  $\equiv 0$
- ▶ There exists a reparameterisation (continuously differentiable) that makes the model linear
- ▶  $\mathbf{P}_\theta \mathbf{H}_{ij}(\boldsymbol{\theta}) = \mathbf{H}_{ij}(\boldsymbol{\theta})$ , where  $\mathbf{H}_{ij}(\boldsymbol{\theta}) = \frac{\partial^2 \boldsymbol{\eta}(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}$

Observing at  $p$  different  $x_i$  only (replications) makes the model intrinsically linear  
 $[p = \dim(\boldsymbol{\theta})]$

## Parametrically linear models

- ▶  $\mathbf{M}(\mathbf{X}_n, \theta) = \text{constant}$
- ▶  $\mathbf{P}_\theta \mathbf{H}_{ij}(\theta) = \mathbf{0}$  — parametric curvature  $\equiv 0$

## Parametrically linear models

- ▶  $\mathbf{M}(\mathbf{X}_n, \theta) = \text{constant}$
- ▶  $\mathbf{P}_\theta \mathbf{H}_{ij}(\theta) = \mathbf{0}$  — parametric curvature  $\equiv 0$

## Linear models

- ▶  $\eta(x, \theta) = \mathbf{f}^\top(x)\theta + c(x)$
- ▶ the model is intrinsically and parametrically linear

## Parametrically linear models

- ▶  $\mathbf{M}(\mathbf{X}_n, \theta) = \text{constant}$
- ▶  $\mathbf{P}_\theta \mathbf{H}_{ij}(\theta) = \mathbf{0}$  — parametric curvature  $\equiv 0$

## Linear models

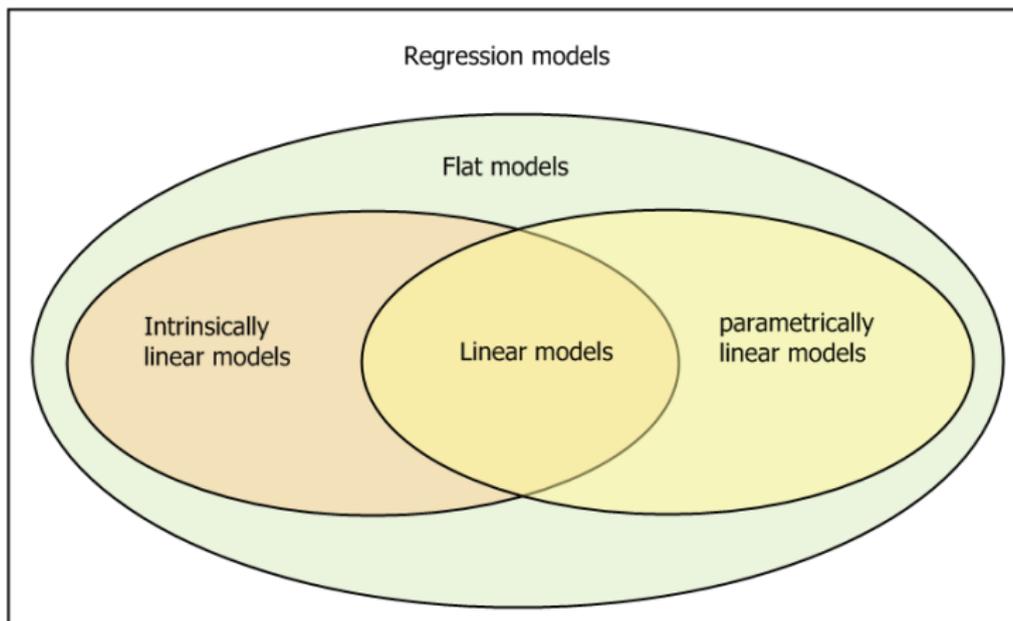
- ▶  $\eta(x, \theta) = \mathbf{f}^\top(x)\theta + c(x)$
- ▶ the model is intrinsically and parametrically linear

## Flat models

- ▶ A reparameterisation exists that makes the information matrix constant
- ▶ Riemannian curvature tensor  $R_{hijk}(\theta) = T_{hjik}(\theta) - T_{hkij}(\theta) \equiv 0$   
with  $T_{hjik}(\theta) = [\mathbf{H}_{hj}(\theta)]^\top [\mathbf{I}_n - \mathbf{P}_\theta] \mathbf{H}_{ik}(\theta)$

If all parameters but one appear linearly, then the model is flat

## A classification of regression models (Pázman, 1993)



## Density of the LS estimator (we suppose $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ )

Intrinsically linear models (in particular, repetitions at  $p$  points):

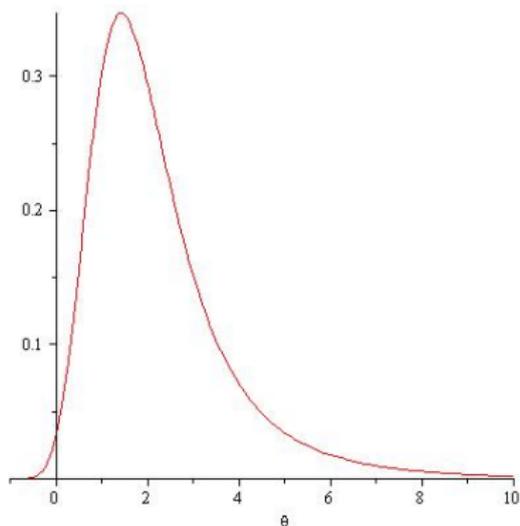
→ exact distribution  $\hat{\theta}^n \sim q(\theta|\bar{\theta}) = \frac{n^{p/2} \det^{1/2} \mathbf{M}(\mathbf{X}_n, \theta)}{(2\pi)^{p/2} \sigma^p} \exp \left\{ -\frac{1}{2\sigma^2} \|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|^2 \right\}$

## Density of the LS estimator (we suppose $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ )

Intrinsically linear models (in particular, repetitions at  $p$  points):

→ exact distribution  $\hat{\theta}^n \sim q(\theta|\bar{\theta}) = \frac{n^{p/2} \det^{1/2} \mathbf{M}(\mathbf{X}_n, \theta)}{(2\pi)^{p/2} \sigma^p} \exp \left\{ -\frac{1}{2\sigma^2} \|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|^2 \right\}$

**Ex:**  $\eta(x, \theta) = \exp(-\theta x)$ ,  $\bar{\theta} = 2$ , 15 observations at the same  $x = 1/2$  ( $\sigma^2 = 1$ )

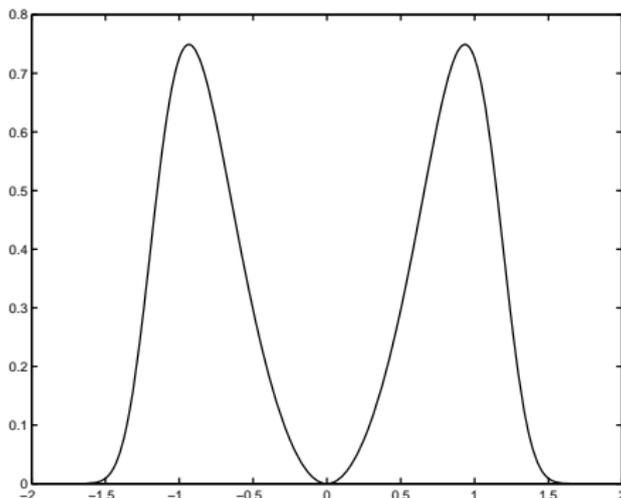


## Density of the LS estimator (we suppose $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ )

Intrinsically linear models (in particular, repetitions at  $p$  points):

→ exact distribution  $\hat{\theta}^n \sim q(\theta|\bar{\theta}) = \frac{n^{p/2} \det^{1/2} \mathbf{M}(\mathbf{X}_n, \theta)}{(2\pi)^{p/2} \sigma^p} \exp \left\{ -\frac{1}{2\sigma^2} \|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|^2 \right\}$

**Ex:**  $\eta(x, \theta) = x\theta^3$ ,  $\bar{\theta} = 0$ , all observations at the same  $x \neq 0$



**Flat models:** approximate density of  $\hat{\theta}^n$

$$q(\theta|\bar{\theta}) = \frac{\det[\mathbf{Q}(\theta, \bar{\theta})]}{(2\pi)^{p/2} \sigma^p n^{p/2} \det^{1/2} \mathbf{M}(\mathbf{X}_n, \theta)} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{P}_\theta[\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})]\|^2 \right\}$$

where  $\{\mathbf{Q}(\theta, \bar{\theta})\}_{ij} = \{n \mathbf{M}(\mathbf{X}_n, \theta)\}_{ij} + [\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})]^\top [\mathbf{I}_n - \mathbf{P}_\theta] \mathbf{H}_{ij}(\theta)$

There exists other approximations (more complicated) for models with  $R_{hijk}(\theta) \neq 0$  (non-flat)

**Flat models:** approximate density of  $\hat{\theta}^n$

$$q(\theta|\bar{\theta}) = \frac{\det[\mathbf{Q}(\theta, \bar{\theta})]}{(2\pi)^{p/2} \sigma^p n^{p/2} \det^{1/2} \mathbf{M}(\mathbf{X}_n, \theta)} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{P}_\theta [\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})]\|^2 \right\}$$

where  $\{\mathbf{Q}(\theta, \bar{\theta})\}_{ij} = \{n \mathbf{M}(\mathbf{X}_n, \theta)\}_{ij} + [\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})]^\top [\mathbf{I}_n - \mathbf{P}_\theta] \mathbf{H}_{ij}(\theta)$

There exists other approximations (more complicated) for models with  $R_{hijk}(\theta) \neq 0$  (non-flat)

⇒ Design of experiments? (since  $q(\theta|\bar{\theta})$  depends on  $\mathbf{X}_n$ )

# 5 DOE based on small sample precision

(P & Pázman, 2013, Chap. 6)

## 1) Minimise the MSE $E\{\|\hat{\theta}^n(\mathbf{y}) - \bar{\theta}\|^2\}$

The approximation of Clarke (1980) requires the 4th-order derivatives of  $\eta(\theta)$

# 5 DOE based on small sample precision

(P & Pázman, 2013, Chap. 6)

## 1) Minimise the MSE $E\{\|\hat{\theta}^n(\mathbf{y}) - \bar{\theta}\|^2\}$

The approximation of Clarke (1980) requires the 4th-order derivatives of  $\eta(\theta)$

→ Use the (approximate) density  $q(\theta|\bar{\theta})$

→ Minimise  $\int_{\Theta} \|\theta - \bar{\theta}\|^2 q(\theta|\bar{\theta}) d\theta$  w.r.t.  $\mathbf{X}_n$  using stochastic approximation

## 5 DOE based on small sample precision

(P & Pázman, 2013, Chap. 6)

### 1) Minimise the MSE $E\{\|\hat{\theta}^n(\mathbf{y}) - \bar{\theta}\|^2\}$

The approximation of Clarke (1980) requires the 4th-order derivatives of  $\eta(\theta)$

→ Use the (approximate) density  $q(\theta|\bar{\theta})$

→ Minimise  $\int_{\Theta} \|\theta - \bar{\theta}\|^2 q(\theta|\bar{\theta}) d\theta$  w.r.t.  $\mathbf{X}_n$  using stochastic approximation

Problem: we need to force  $\theta$  to remain in  $\Theta$

→ the integral can be made equal to 0

# 5 DOE based on small sample precision

(P & Pázman, 2013, Chap. 6)

## 1) Minimise the MSE $E\{\|\hat{\theta}^n(\mathbf{y}) - \bar{\theta}\|^2\}$

The approximation of Clarke (1980) requires the 4th-order derivatives of  $\eta(\theta)$

→ Use the (approximate) density  $q(\theta|\bar{\theta})$

→ Minimise  $\int_{\Theta} \|\theta - \bar{\theta}\|^2 q(\theta|\bar{\theta}) d\theta$  w.r.t.  $\mathbf{X}_n$  using stochastic approximation

Problem: we need to force  $\theta$  to remain in  $\Theta$

→ the integral can be made equal to 0

Solution: approximate the density  $\tilde{q}_w(\theta|\bar{\theta})$  of a penalised LS estimator  $\tilde{\theta}^n$

$$\tilde{\theta}^n = \arg \min_{\theta} \{\|\mathbf{y} - \boldsymbol{\eta}(\theta)\|^2 + 2w(\theta)\}$$

where  $w(\theta)$  forces  $\theta$  to remain in  $\Theta$  [ $w(\theta) = +\infty$  outside  $\Theta$ ]

→ Minimise  $\int_{\Theta} \|\theta - \bar{\theta}\|^2 \tilde{q}_w(\theta|\bar{\theta}) d\theta$  w.r.t.  $\mathbf{X}_n$

[also covers the case of max. *a posteriori* estimation (relate  $w(\theta)$  to the prior on  $\theta$ )]

(P & Pázman, 1992; Pázman and Gauchi, 2006)

## 2) Use a small-sample variant of $D$ -optimal design

A  $D$ -optimal design minimises

- (i) the volume of asymptotic (ellipsoidal) confidence regions
- (ii) the (Shannon) entropy of the asymptotic distribution of  $\hat{\theta}^n$

## 2) Use a small-sample variant of $D$ -optimal design

A  $D$ -optimal design minimises

- (i) the volume of asymptotic (ellipsoidal) confidence regions
- (ii) the (Shannon) entropy of the asymptotic distribution of  $\hat{\theta}^n$

Hamilton and Watts (1985) minimize the (approximate) volume  $V(\mathbf{X}_n, \theta^0)$  of (approximate) confidence regions ( $V(\mathbf{X}_n, \theta^0)$  has an explicit form and a geometrical interpretation)

Vila (1990); Vila and Gauchi (2007) minimize the expected volume of exact confidence regions (not ellipsoidal, not necessarily of minimum volume), using stochastic approximation

→ Choose  $\mathbf{X}_n$  that minimises the approximate entropy of the approximate distribution of  $\hat{\theta}^n$  (P & Pázman, 1994b)

$$\text{Minimise Ent}[q(\cdot|\bar{\theta})] = - \int_{\mathbb{R}^n} \log[q(\hat{\theta}^n(\mathbf{y})|\bar{\theta})] \varphi(\mathbf{y}|\mathbf{X}_n, \bar{\theta}) d\mathbf{y} \text{ w.r.t. } \mathbf{X}_n$$

where  $\varphi(\mathbf{y}|\mathbf{X}_n, \bar{\theta})$  corresponds to  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\eta}(\bar{\theta}), \sigma^2 \mathbf{I}_n)$

→ Choose  $\mathbf{X}_n$  that minimises the approximate entropy of the approximate distribution of  $\hat{\theta}^n$  (P & Pázman, 1994b)

$$\text{Minimise Ent}[q(\cdot|\bar{\theta})] = - \int_{\mathbb{R}^n} \log[q(\hat{\theta}^n(\mathbf{y})|\bar{\theta})] \varphi(\mathbf{y}|\mathbf{X}_n, \bar{\theta}) d\mathbf{y} \text{ w.r.t. } \mathbf{X}_n$$

where  $\varphi(\mathbf{y}|\mathbf{X}_n, \bar{\theta})$  corresponds to  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\eta}(\bar{\theta}), \sigma^2 \mathbf{I}_n)$

Use a 2nd order Taylor development of  $\log[q(\hat{\theta}^n(\mathbf{y})|\bar{\theta})]$  around  $\mathbf{y} = \boldsymbol{\eta}(\bar{\theta})$ :

$$\text{Ent}[q(\cdot|\bar{\theta})] = - \log q(\bar{\theta}|\bar{\theta}) - \frac{\sigma^2}{2} \sum_{i=1}^N \frac{\partial^2 \log q[\hat{\theta}^n(\mathbf{y})|\bar{\theta}]}{\partial y_i^2} \Big|_{\boldsymbol{\eta}(\bar{\theta})} + \mathcal{O}(\sigma^4)$$

After some (lengthy) calculations...

$$\begin{aligned}
 \text{Ent}[q(\cdot|\bar{\theta})] &= \overbrace{\frac{p}{2}[1 + \log(2\pi\sigma^2)] - \frac{1}{2} \log \det[n\mathbf{M}(\mathbf{X}_n, \bar{\theta})]}^{\text{entropy of asymptotic normal distribution}} \\
 &- \frac{\sigma^2}{2n} \sum_{h,i,j,k=1}^p \left( \{\mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta})\}_{ij} \left[ \frac{1}{n} \{\mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta})\}_{kh} [R_{kjhi}(\bar{\theta}) + U_{kij}^h(\bar{\theta})] \right. \right. \\
 &\left. \left. - G_{ki}^h(\bar{\theta}) G_{hj}^k(\bar{\theta}) - G_{kh}^k(\bar{\theta}) G_{ij}^h(\bar{\theta}) \right] \right) + \mathcal{O}(\sigma^4)
 \end{aligned}$$

$$\text{where } U_{kij}^h(\theta) = \frac{\partial^3 \boldsymbol{\eta}^\top(\theta)}{\partial \theta_k \partial \theta_i \partial \theta_j} \frac{\partial \boldsymbol{\eta}(\theta)}{\partial \theta_h}$$

$$G_{ij}^k(\theta) = \frac{1}{n} \sum_{h=1}^p \frac{\partial \boldsymbol{\eta}^\top(\theta)}{\partial \theta_h} \mathbf{H}_{ij}^\cdot \{\mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta})\}_{hk}$$

with  $R_{hijk}(\theta) = T_{hjik}(\theta) - T_{hkij}(\theta)$ ,  $T_{hjik}(\theta) = [\mathbf{H}_{ij}^\cdot(\theta)]^\top [\mathbf{I}_n - \mathbf{P}_\theta] \mathbf{H}_{ik}^\cdot(\theta)$  and  $\mathbf{H}_{ij}^\cdot(\theta) = \frac{\partial^2 \boldsymbol{\eta}(\theta)}{\partial \theta_i \partial \theta_j}$

### 3) Related work using the approximate density $q(\theta|\bar{\theta})$

3a) (approximate) marginal densities of  $\hat{\theta}^n$  (Pázman & P, 1996)

Denote  $\gamma = h(\theta)$  [with  $\gamma = \theta_i$  for some  $i \in \{1, \dots, p = \dim(\theta)\}$ ] as particular case]

$$q(\gamma|\bar{\theta}) = \frac{1}{\sqrt{2\pi\sigma}\|\mathbf{b}_\gamma\|} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{P}_\gamma[\boldsymbol{\eta}(\theta_\gamma) - \boldsymbol{\eta}(\bar{\theta})]\|^2 \right\}$$

where

$$\theta_\gamma = \arg \min_{\theta: h(\theta)=\gamma} \|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|^2$$

$$\mathbf{b}_\gamma = \frac{1}{n} \frac{\partial \boldsymbol{\eta}(\theta)}{\partial \theta^\top} \Big|_{\theta_\gamma} \mathbf{M}^{-1}(\mathbf{X}_n, \theta_\gamma) \frac{\partial h(\theta)}{\partial \theta} \Big|_{\theta_\gamma}$$

$$\mathbf{P}_\gamma = \frac{\mathbf{b}_\gamma \mathbf{b}_\gamma^\top}{\|\mathbf{b}_\gamma\|^2}$$

[There also exist more precise approximations, more complicated; the difficulty compared to (Tierney et al., 1989) is that  $\hat{\theta}^n(\mathbf{y})$  is not known explicitly]

→ Can be used to compare experiments

**Ex:** a two-compartment model in pharmacokinetics (P & Pázman, 2001)

Observe  $y(t) = x_C(t)/V + \varepsilon(t)$  where  $x_C(t)$  evolves according to

$$\begin{cases} \frac{dx_C(t)}{dt} = (-K_{EL} - K_{CP})x_C(t) + K_{PC}x_P(t) + u(t) \\ \frac{dx_P(t)}{dt} = K_{CP}x_C(t) - K_{PC}x_P(t) \end{cases}$$

errors  $\varepsilon(t_i)$  i.i.d.  $\mathcal{N}(0, \sigma^2)$

→ Can be used to compare experiments

**Ex:** a two-compartment model in pharmacokinetics (P & Pázman, 2001)

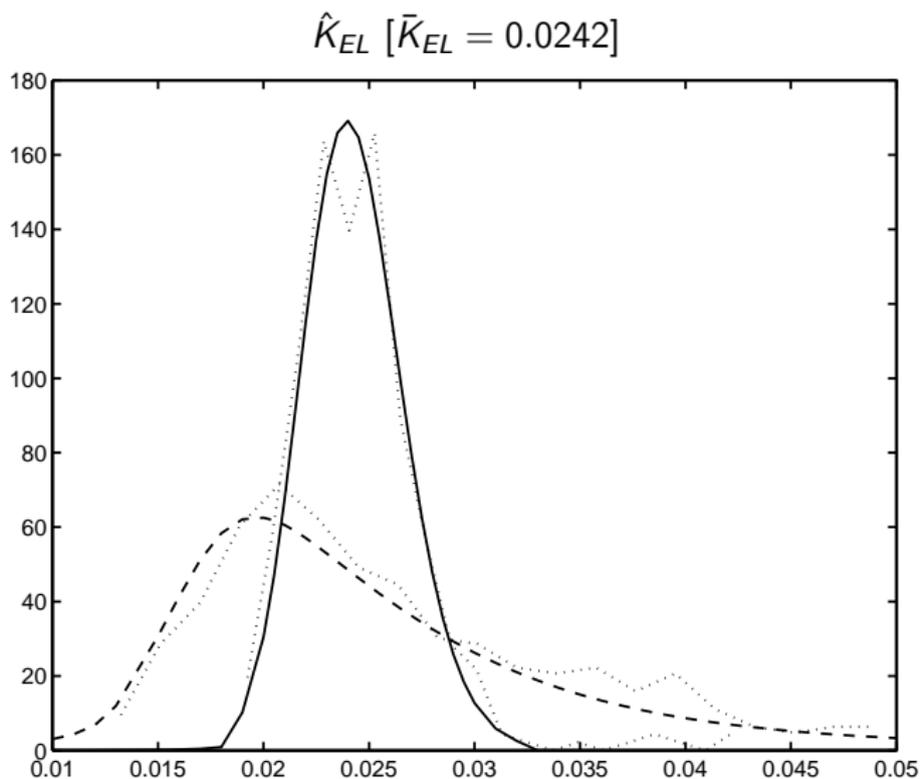
Observe  $y(t) = x_C(t)/V + \varepsilon(t)$  where  $x_C(t)$  evolves according to

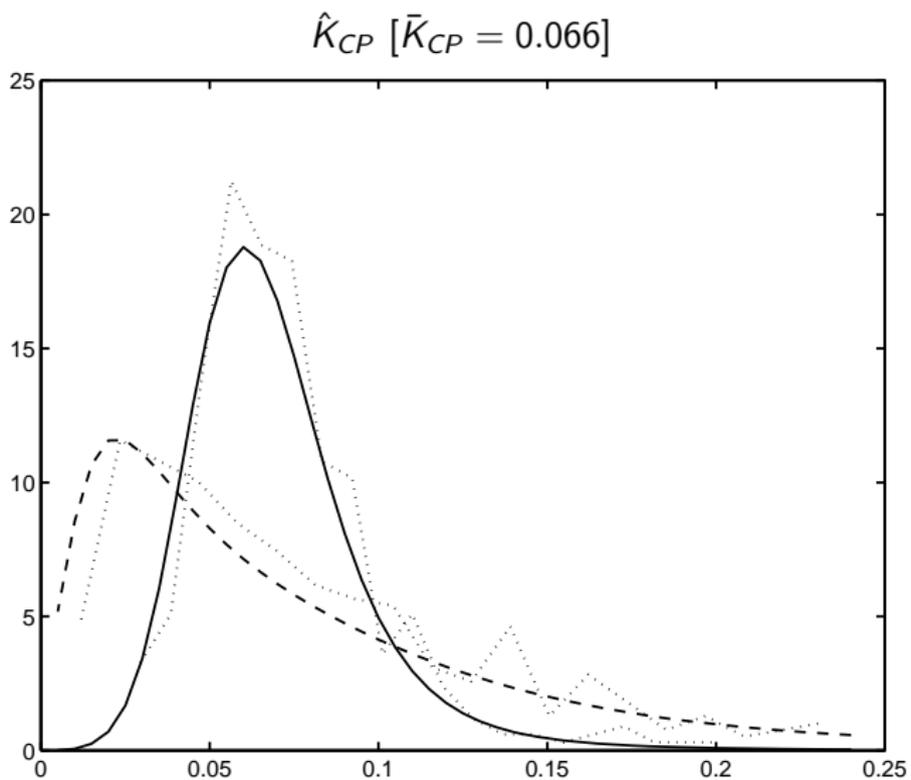
$$\begin{cases} \frac{dx_C(t)}{dt} = (-K_{EL} - K_{CP})x_C(t) + K_{PC}x_P(t) + u(t) \\ \frac{dx_P(t)}{dt} = K_{CP}x_C(t) - K_{PC}x_P(t) \end{cases}$$

errors  $\varepsilon(t_i)$  i.i.d.  $\mathcal{N}(0, \sigma^2)$

→ 4 unknown parameters  $\theta = (K_{CP}, K_{PC}, K_{EL}, V)^T$

Compare 2 designs (8 observation times each) using simulated experiments with a given true  $\bar{\theta}$





3b) Bias correction for LS estimation in nonlinear regression

$$\begin{aligned}
 \mathbf{b}(\bar{\theta}) &= \text{bias of } \hat{\theta}^n = E_{\mathbf{X}_n, \bar{\theta}}\{\hat{\theta}^n(\mathbf{y})\} - \bar{\theta} \\
 &= \underbrace{-\frac{\sigma^2}{2n^2} \mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta}) \frac{\partial \boldsymbol{\eta}^\top(\theta)}{\partial \theta} \Big|_{\bar{\theta}} \sum_{i,j=1}^p \mathbf{H}_{ij}(\bar{\theta}) \{\mathbf{M}^{-1}(\mathbf{X}_n, \bar{\theta})\}_{ij}}_{= \tilde{\mathbf{b}}(\bar{\theta}) \text{ (Box, 1971)}} + \mathcal{O}(\sigma^4)
 \end{aligned}$$

We can write  $\hat{\theta}^n = \mathbf{b}(\bar{\theta}) + \bar{\theta} + \omega$ , with  $E_{\mathbf{X}_n, \bar{\theta}}\{\omega\} = \mathbf{0}$

Two-stage LS: solve  $\hat{\theta}^n = \mathbf{b}(\theta) + \theta$  for  $\theta \rightarrow \hat{\theta}^{n,*}$

$[\hat{\theta}^{n,*}$  unbiased when  $\mathbf{b}(\theta) = \mathbf{A}\theta + \mathbf{c}$  for all  $\theta$  with  $\mathbf{I}_p + \mathbf{A}$  nonsingular]

1st method:  $\hat{\theta}^{n,0} = \hat{\theta}^n$  given, then

$$\hat{\theta}^{n,1} = \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,0})$$

[... sometimes more biased than  $\hat{\theta}^n$  (Picard and Prum, 1992)]

1st method:  $\hat{\theta}^{n,0} = \hat{\theta}^n$  given, then

$$\begin{aligned}
 \hat{\theta}^{n,1} &= \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,0}) \\
 & \quad [\dots \text{sometimes more biased than } \hat{\theta}^n \text{ (Picard and Prum, 1992)}] \\
 \hat{\theta}^{n,2} &= \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,1}) \\
 \vdots &= \vdots \\
 \hat{\theta}^{n,*} = \hat{\theta}^{n,\infty} &= \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,\infty})
 \end{aligned}$$

that is,  $\hat{\theta}^{n,*} + \mathbf{b}(\hat{\theta}^{n,*}) = \hat{\theta}^n$ , or

$$\mathbb{E}_{\mathbf{X}_n, \hat{\theta}^{n,*}} \{ \hat{\theta}^n(\mathbf{y}) \} = \int_{\mathbb{R}^n} \hat{\theta}^n(\mathbf{y}) \varphi(\mathbf{y} | \mathbf{X}_n, \hat{\theta}^{n,*}) d\mathbf{y} = \hat{\theta}^n$$

1st method:  $\hat{\theta}^{n,0} = \hat{\theta}^n$  given, then

$$\begin{aligned}
 \hat{\theta}^{n,1} &= \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,0}) \\
 &\quad [\dots \text{sometimes more biased than } \hat{\theta}^n \text{ (Picard and Prum, 1992)}] \\
 \hat{\theta}^{n,2} &= \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,1}) \\
 \vdots &= \vdots \\
 \hat{\theta}^{n,*} = \hat{\theta}^{n,\infty} &= \hat{\theta}^n - \mathbf{b}(\hat{\theta}^{n,\infty})
 \end{aligned}$$

that is,  $\hat{\theta}^{n,*} + \mathbf{b}(\hat{\theta}^{n,*}) = \hat{\theta}^n$ , or

$$\mathbb{E}_{\mathbf{X}_n, \hat{\theta}^{n,*}} \{ \hat{\theta}^n(\mathbf{y}) \} = \int_{\mathbb{R}^n} \hat{\theta}^n(\mathbf{y}) \varphi(\mathbf{y} | \mathbf{X}_n, \hat{\theta}^{n,*}) d\mathbf{y} = \hat{\theta}^n$$

Solve for  $\hat{\theta}^{n,*}$  using stochastic approximation (P & Pázman, 1994a)

2nd method (approximate): use  $\tilde{\mathbf{b}}$  instead of  $\mathbf{b}$

2nd method (approximate): use  $\tilde{\mathbf{b}}$  instead of  $\mathbf{b}$

$$\text{Solve for } \tilde{\theta}^{n,*}: \quad \tilde{\theta}^{n,*} + \tilde{\mathbf{b}}(\tilde{\theta}^{n,*}) = \hat{\theta}^n$$

that is

$$\tilde{\theta}^{n,*} - \frac{\sigma^2}{2n^2} \mathbf{M}^{-1}(\mathbf{X}_n, \tilde{\theta}^{n,*}) \left. \frac{\partial \boldsymbol{\eta}^\top(\theta)}{\partial \theta} \right|_{\tilde{\theta}^{n,*}} \sum_{i,j=1}^p \mathbf{H}_{ij}(\tilde{\theta}^{n,*}) \{ \mathbf{M}^{-1}(\mathbf{X}_n, \tilde{\theta}^{n,*}) \}_{ij} = \hat{\theta}^n$$

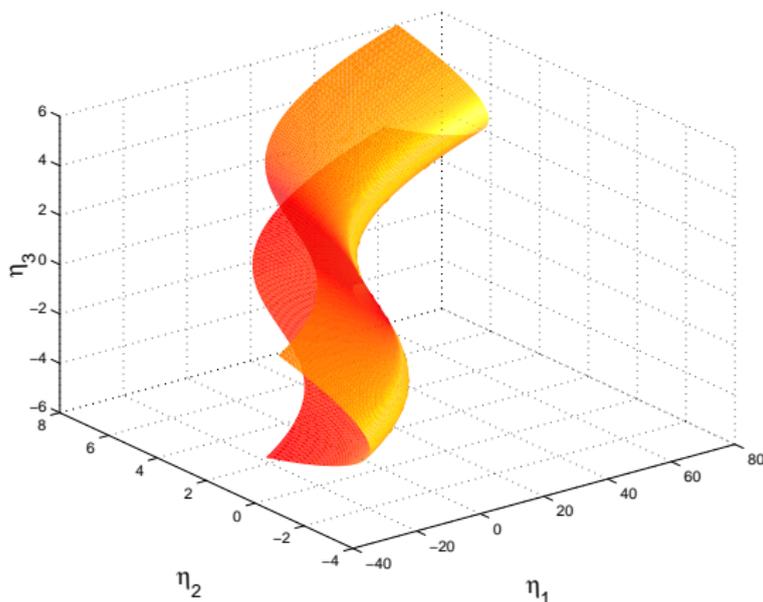
[Different from the score-corrected estimator  $\hat{\theta}_{sc}^n$  of (Firth, 1993):

$$\rightarrow \text{solve } \left[ \frac{\partial \boldsymbol{\eta}^\top(\theta)}{\partial \theta} [\mathbf{y} - \boldsymbol{\eta}(\theta)] - \mathbf{M}(\mathbf{X}_n, \theta) \tilde{\mathbf{b}}(\theta) = \mathbf{0} \right] \text{ for } \theta$$

(Pázman & P, 1998) gives the (approximate) joint and marginal densities of  $\tilde{\theta}^{n,*}$  and  $\hat{\theta}_{sc}^n$

## 6 Extended optimality criteria

(P & Pázman, 2013, Chap. 7)



$\mathbb{S}_\eta$  may overlap, there may be local minimisers for the LS problem. . .

Important and difficult problem, often neglected

## What can we do at the design stage?

▣ extensions of usual optimality criteria

→ Avoid situations where  $\|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|$  can be small when  $\|\theta - \bar{\theta}\|$  is large:

$$\text{maximise } \phi_{eE}(\mathbf{X}_n, \theta^0) = \min_{\theta} \frac{\|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\theta^0)\|^2}{\|\theta - \theta^0\|^2}$$

corresponds to  $E$ -optimal design ( $\Leftrightarrow$  maximise  $\lambda_{\min}[\mathbf{M}(\mathbf{X}_n)]$ ) when  $\boldsymbol{\eta}$  is linear

## What can we do at the design stage?

▣ extensions of usual optimality criteria

→ Avoid situations where  $\|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|$  can be small when  $\|\theta - \bar{\theta}\|$  is large:

$$\text{maximise } \phi_{eE}(\mathbf{X}_n, \theta^0) = \min_{\theta} \frac{\|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\theta^0)\|^2}{\|\theta - \theta^0\|^2}$$

corresponds to  $E$ -optimal design ( $\Leftrightarrow$  maximise  $\lambda_{\min}[\mathbf{M}(\mathbf{X}_n)]$ ) when  $\boldsymbol{\eta}$  is linear

Extensions of  $E$ -,  $G$ - and  $c$ -optimal design in (Pázman & P, 2014)

## What can we do at the design stage?

▣ extensions of usual optimality criteria

→ Avoid situations where  $\|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\bar{\theta})\|$  can be small when  $\|\theta - \bar{\theta}\|$  is large:

$$\text{maximise } \phi_{eE}(\mathbf{X}_n, \theta^0) = \min_{\theta} \frac{\|\boldsymbol{\eta}(\theta) - \boldsymbol{\eta}(\theta^0)\|^2}{\|\theta - \theta^0\|^2}$$

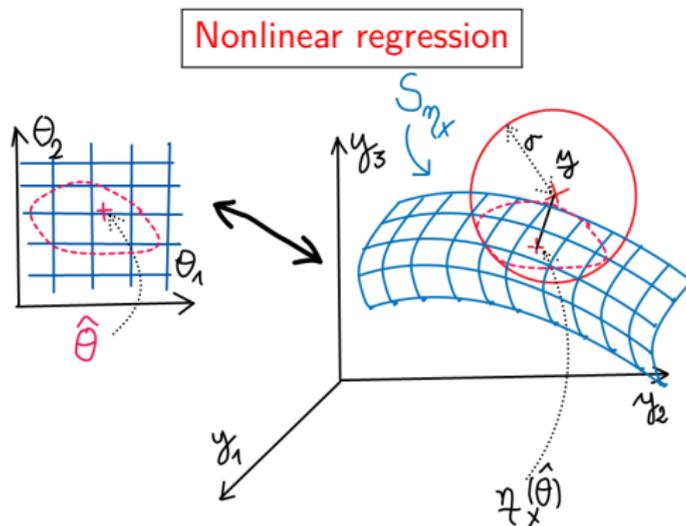
corresponds to  $E$ -optimal design ( $\Leftrightarrow$  maximise  $\lambda_{\min}[\mathbf{M}(\mathbf{X}_n)]$ ) when  $\boldsymbol{\eta}$  is linear

Extensions of  $E$ -,  $G$ - and  $c$ -optimal design in (Pázman & P, 2014)

Extensions to generalised regression models and other design criteria in the Ph.D. thesis (Sternmüllerová, 2019)

# 7 Nonlocal DoE for nonlinear models

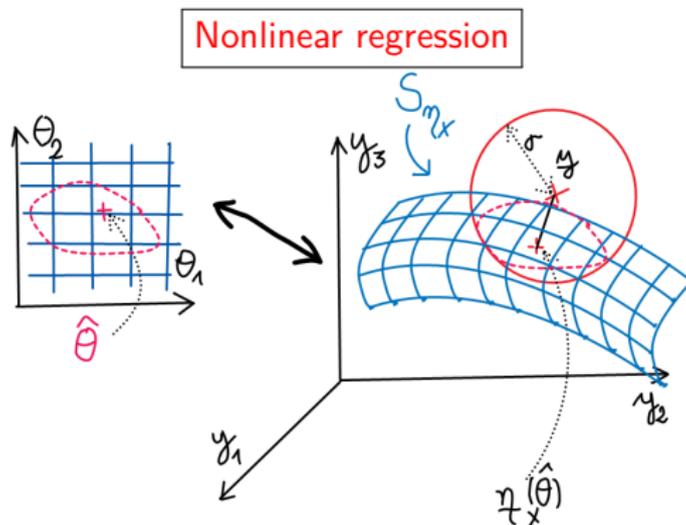
(P & Pázman, 2013, Chap. 8)



Nonlinear model  $\implies$  **everything is local**

# 7 Nonlocal DoE for nonlinear models

(P & Pázman, 2013, Chap. 8)



Nonlinear model  $\implies$  **everything is local**

$\phi(\cdot)$  an information criterion, to be maximised with respect to the design  $\mathbf{X}_n$ :

$\phi(\mathbf{X}_n) = \phi(\mathbf{X}_n, \theta)$ , but **which  $\theta$** ?

Local optimum design: based on a nominal value  $\theta^0 \rightarrow$  maximize  $\phi(\mathbf{X}_n, \theta^0)$   
[concerns all methods considered so far,  
based on asymptotic normality (AN) or small-sample properties]

Local optimum design: based on a nominal value  $\theta^0 \rightarrow$  maximize  $\phi(\mathbf{X}_n, \theta^0)$   
[concerns all methods considered so far,  
based on asymptotic normality (AN) or small-sample properties]

**Objective of nonlocal DoE: remove the dependence in  $\theta^0$**

3 main classes, essentially for  $\phi(\xi, \theta) = \Phi[\mathbf{M}(\mathbf{X}_n, \theta)]$  (based on AN)

Local optimum design: based on a nominal value  $\theta^0 \rightarrow$  maximize  $\phi(\mathbf{X}_n, \theta^0)$   
 [concerns all methods considered so far,  
 based on asymptotic normality (AN) or small-sample properties]

**Objective of nonlocal DoE: remove the dependence in  $\theta^0$**

3 main classes, essentially for  $\phi(\xi, \theta) = \Phi[\mathbf{M}(\mathbf{X}_n, \theta)]$  (based on AN)

① Average optimum design: maximise  $E_{\theta}\{\phi(\mathbf{X}_n, \theta)\}$  (or  $E_{\theta}\{\phi(\xi, \theta)\}$ )

Local optimum design: based on a nominal value  $\theta^0 \rightarrow$  maximize  $\phi(\mathbf{X}_n, \theta^0)$   
 [concerns all methods considered so far,  
 based on asymptotic normality (AN) or small-sample properties]

**Objective of nonlocal DoE: remove the dependence in  $\theta^0$**

3 main classes, essentially for  $\phi(\xi, \theta) = \Phi[\mathbf{M}(\mathbf{X}_n, \theta)]$  (based on AN)

- ❶ Average optimum design: maximise  $E_{\theta}\{\phi(\mathbf{X}_n, \theta)\}$  (or  $E_{\theta}\{\phi(\xi, \theta)\}$ )
- ❷ Maximin optimum design: maximise  $\min_{\theta}\{\phi(\mathbf{X}_n, \theta)\}$  (or  $\min_{\theta}\{\phi(\xi, \theta)\}$ )
- ➡ Between ❶ and ❷: regularised maximin criteria, quantiles and probability level criteria

Local optimum design: based on a nominal value  $\theta^0 \rightarrow$  maximize  $\phi(\mathbf{X}_n, \theta^0)$   
 [concerns all methods considered so far,  
 based on asymptotic normality (AN) or small-sample properties]

### Objective of nonlocal DoE: remove the dependence in $\theta^0$

3 main classes, essentially for  $\phi(\xi, \theta) = \Phi[\mathbf{M}(\mathbf{X}_n, \theta)]$  (based on AN)

- ❶ Average optimum design: maximise  $E_\theta\{\phi(\mathbf{X}_n, \theta)\}$  (or  $E_\theta\{\phi(\xi, \theta)\}$ )
- ❷ Maximin optimum design: maximise  $\min_\theta\{\phi(\mathbf{X}_n, \theta)\}$  (or  $\min_\theta\{\phi(\xi, \theta)\}$ )
- ➡ Between ❶ and ❷: regularised maximin criteria, quantiles and probability level criteria
- ❸ Sequential design

## 1 Average Optimum design

Probability measure  $\mu(d\theta)$  on  $\Theta \subseteq \mathbb{R}^p$  ( $\neq$  Bayesian estimation)

$$\phi(\cdot, \theta^0) \rightarrow \phi_A(\cdot) = \int_{\Theta} \phi(\cdot, \theta) \mu(d\theta)$$

[No difficulty if  $\Theta$  is finite and  $\mu = \sum_{i=1}^M \alpha_i \delta_{\theta}^{(i)}$  (integral  $\rightarrow$  finite sum); otherwise, use stochastic approximation to avoid evaluations of integrals]

## ① Average Optimum design

Probability measure  $\mu(d\theta)$  on  $\Theta \subseteq \mathbb{R}^p$  ( $\neq$  Bayesian estimation)

$$\phi(\cdot, \theta^0) \rightarrow \phi_A(\cdot) = \int_{\Theta} \phi(\cdot, \theta) \mu(d\theta)$$

[No difficulty if  $\Theta$  is finite and  $\mu = \sum_{i=1}^M \alpha_i \delta_{\theta}^{(i)}$  (integral  $\rightarrow$  finite sum); otherwise, use stochastic approximation to avoid evaluations of integrals]

$\phi_A(\xi)$  is concave in  $\xi$  when each  $\phi(\xi, \theta)$  is concave

▣ same properties and same algorithms as for local design

## 1 Average Optimum design

Probability measure  $\mu(d\theta)$  on  $\Theta \subseteq \mathbb{R}^p$  ( $\neq$  Bayesian estimation)

$$\phi(\cdot, \theta^0) \rightarrow \phi_A(\cdot) = \int_{\Theta} \phi(\cdot, \theta) \mu(d\theta)$$

[No difficulty if  $\Theta$  is finite and  $\mu = \sum_{i=1}^M \alpha_i \delta_{\theta}^{(i)}$  (integral  $\rightarrow$  finite sum); otherwise, use stochastic approximation to avoid evaluations of integrals]

$\phi_A(\xi)$  is concave in  $\xi$  when each  $\phi(\xi, \theta)$  is concave

▀ same properties and same algorithms as for local design

## 2 Maximin Optimum design

$$\phi(\cdot, \theta^0) \rightarrow \phi_M(\cdot) = \min_{\theta \in \Theta} \phi(\cdot, \theta)$$

$\phi_M(\xi)$  is concave in  $\xi$  when each  $\phi(\xi, \theta)$  is concave, but  $\phi_M(\cdot)$  is non-differentiable

## ① Average Optimum design

Probability measure  $\mu(d\theta)$  on  $\Theta \subseteq \mathbb{R}^p$  ( $\neq$  Bayesian estimation)

$$\phi(\cdot, \theta^0) \rightarrow \phi_A(\cdot) = \int_{\Theta} \phi(\cdot, \theta) \mu(d\theta)$$

[No difficulty if  $\Theta$  is finite and  $\mu = \sum_{i=1}^M \alpha_i \delta_{\theta}^{(i)}$  (integral  $\rightarrow$  finite sum); otherwise, use stochastic approximation to avoid evaluations of integrals]

$\phi_A(\xi)$  is concave in  $\xi$  when each  $\phi(\xi, \theta)$  is concave

▀ same properties and same algorithms as for local design

## ② Maximin Optimum design

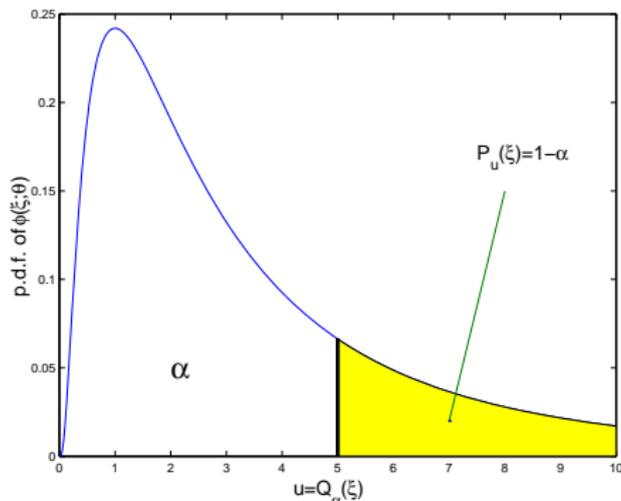
$$\phi(\cdot, \theta^0) \rightarrow \phi_M(\cdot) = \min_{\theta \in \Theta} \phi(\cdot, \theta)$$

$\phi_M(\xi)$  is concave in  $\xi$  when each  $\phi(\xi, \theta)$  is concave, but  $\phi_M(\cdot)$  is non-differentiable

## Problems

- ① Optimal design for  $\phi_A(\cdot)$  not invariant by a monotone transformation of  $\phi(\cdot, \theta)$
- ② Optimal design for  $\phi_M(\cdot)$  very sensitive to the choice of the boundary of  $\Theta$

## Between ① and ②: quantiles and probability level criteria



→ maximise  $P_u$  for a given  $u$ , or maximise  $Q_\alpha$  for a given  $\alpha$   
 (when  $\alpha \rightarrow 0$ , tends to maximin optimality)

Directional derivatives, algorithms . . . but the criteria are not concave:

→ no guarantee of successful maximisation

Directional derivatives, algorithms ... but the criteria are not concave:

→ no guarantee of successful maximisation

A related very promising approach: maximise the conditional value at risk (or superquantile) as proposed by Valenzuela et al. (2015)

$$\phi_\alpha(\mathbf{X}_n) = \max_{t \in \mathbb{R}} \left\{ t + \frac{1}{\alpha} \int_{\Theta} \min [0, \phi(\mathbf{X}_n; \theta) - t] \mu(d\theta) \right\}$$

When  $\mu$  has a density (w.r.t. Lebesgue measure on  $\Theta$ ) then

$$\phi_\alpha(\mathbf{X}_n) = \frac{1}{\alpha} \int_{\{\theta: \phi(\mathbf{X}_n; \theta) < Q_\alpha(\mathbf{X}_n)\}} \phi(\mathbf{X}_n; \theta) \mu(d\theta)$$

$\phi(\xi, \theta)$  concave in  $\xi \Rightarrow \phi_\alpha(\xi)$  concave

$\phi_1(\mathbf{X}_n) = \phi_A(\mathbf{X}_n)$  and  $\phi_\alpha(\mathbf{X}_n) \rightarrow \phi_M(\mathbf{X}_n)$  as  $\alpha \rightarrow 0$

[part of the Ph.D. thesis (Sternmüllerová, 2019)]

### ③ Sequential design

$\theta^0 \rightarrow$  design:  $\mathbf{X}^1 = \arg \max_{\mathbf{X}} \phi(\mathbf{X}, \theta^0)$

$\rightarrow$  observe:  $\mathbf{y}^1 = \mathbf{y}^1(\mathbf{X}^1)$

$\rightarrow$  estimate:  $\hat{\theta}^1 = \arg \min_{\theta} LS(\theta; \mathbf{y}^1, \mathbf{X}^1)$

$\rightarrow$  design:  $\mathbf{X}^2 = \arg \max_{\mathbf{X}} \phi(\{\mathbf{X}^1, \mathbf{X}\}, \hat{\theta}^1)$

$\rightarrow$  observe:  $\mathbf{y}^2 = \mathbf{y}^2(\mathbf{X}^2)$

$\rightarrow$  estimate:  $\hat{\theta}^2 = \arg \min_{\theta} LS(\theta; \underbrace{\{\mathbf{y}^1, \mathbf{y}^2\}}_{\text{growing}}, \underbrace{\{\mathbf{X}^1, \mathbf{X}^2\}}_{\text{growing}})$

$\rightarrow$  design:  $\mathbf{X}^3 = \arg \max_{\mathbf{X}} \phi(\{\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}\}, \hat{\theta}^2)$

... etc.

### ③ Sequential design

$\theta^0 \rightarrow$  design:  $\mathbf{X}^1 = \arg \max_{\mathbf{X}} \phi(\mathbf{X}, \theta^0)$

$\rightarrow$  observe:  $\mathbf{y}^1 = \mathbf{y}^1(\mathbf{X}^1)$

$\rightarrow$  estimate:  $\hat{\theta}^1 = \arg \min_{\theta} LS(\theta; \mathbf{y}^1, \mathbf{X}^1)$

$\rightarrow$  design:  $\mathbf{X}^2 = \arg \max_{\mathbf{X}} \phi(\{\mathbf{X}^1, \mathbf{X}\}, \hat{\theta}^1)$

$\rightarrow$  observe:  $\mathbf{y}^2 = \mathbf{y}^2(\mathbf{X}^2)$

$\rightarrow$  estimate:  $\hat{\theta}^2 = \arg \min_{\theta} LS(\theta; \underbrace{\{\mathbf{y}^1, \mathbf{y}^2\}}_{\text{growing}}, \underbrace{\{\mathbf{X}^1, \mathbf{X}^2\}}_{\text{growing}})$

$\rightarrow$  design:  $\mathbf{X}^3 = \arg \max_{\mathbf{X}} \phi(\{\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}\}, \hat{\theta}^2)$

... etc.

$\rightarrow$  Replace unknown  $\theta$  by best current guess  $\hat{\theta}^k$

(there exist variants with Bayesian estimation and average optimality)

### ③ Sequential design

$\theta^0 \rightarrow$  design:  $\mathbf{X}^1 = \arg \max_{\mathbf{X}} \phi(\mathbf{X}, \theta^0)$

$\rightarrow$  observe:  $\mathbf{y}^1 = \mathbf{y}^1(\mathbf{X}^1)$

$\rightarrow$  estimate:  $\hat{\theta}^1 = \arg \min_{\theta} LS(\theta; \mathbf{y}^1, \mathbf{X}^1)$

$\rightarrow$  design:  $\mathbf{X}^2 = \arg \max_{\mathbf{X}} \phi(\{\mathbf{X}^1, \mathbf{X}\}, \hat{\theta}^1)$

$\rightarrow$  observe:  $\mathbf{y}^2 = \mathbf{y}^2(\mathbf{X}^2)$

$\rightarrow$  estimate:  $\hat{\theta}^2 = \arg \min_{\theta} LS(\theta; \underbrace{\{\mathbf{y}^1, \mathbf{y}^2\}}_{\text{growing}}, \underbrace{\{\mathbf{X}^1, \mathbf{X}^2\}}_{\text{growing}})$

$\rightarrow$  design:  $\mathbf{X}^3 = \arg \max_{\mathbf{X}} \phi(\{\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}\}, \hat{\theta}^2)$

... etc.

$\rightarrow$  Replace unknown  $\theta$  by best current guess  $\hat{\theta}^k$

(there exist variants with Bayesian estimation and average optimality)

Consistency of  $\hat{\theta}^n$ ? Asymptotic normality (for designs based on  $\mathbf{M}$ )?

(difficulty:  $\mathbf{X}^k$  depends on  $\mathbf{y}^1, \dots, \mathbf{y}^{k-1} \implies$  independence is lost)

Each  $\mathbf{X}^i$  has size  $q$

⇒ No big difficulty if  $q \geq p = \dim(\theta)$  (batch sequential design)

Each  $\mathbf{X}^i$  has size  $q$

⇒ No big difficulty if  $q \geq p = \dim(\theta)$  (batch sequential design)

If  $n$  observation in total, two stages only: size of first batch?

→ should be proportional to  $\sqrt{n}$  (*which does not say much ...*)

Each  $\mathbf{X}^i$  has size  $q$

⇒ No big difficulty if  $q \geq p = \dim(\theta)$  (batch sequential design)

If  $n$  observation in total, two stages only: size of first batch?

→ should be proportional to  $\sqrt{n}$  (*which does not say much ...*)

⇒ Full sequential design:  $\mathbf{X}^k = \{x_k\}$  ( $q = 1$ )

→ convergence properties are difficult to investigate...

Each  $\mathbf{X}^i$  has size  $q$

⇒ No big difficulty if  $q \geq p = \dim(\theta)$  (batch sequential design)

If  $n$  observation in total, two stages only: size of first batch?

→ should be proportional to  $\sqrt{n}$  (which does not say much ...)

⇒ Full sequential design:  $\mathbf{X}^k = \{x_k\}$  ( $q = 1$ )

→ convergence properties are difficult to investigate...

When

$$\mathbf{M}(\mathbf{X}_{k+1}, \hat{\theta}^k) = \frac{k}{k+1} \mathbf{M}(\mathbf{X}_k, \hat{\theta}^k) + \frac{1}{k+1} \frac{\partial \eta(x_{k+1}, \theta)}{\partial \theta} \Big|_{\hat{\theta}^k} \frac{\partial \eta(x_{k+1}, \theta)}{\partial \theta^\top} \Big|_{\hat{\theta}^k}$$

with  $x_{k+1} = \arg \max_{\mathcal{X}} \underbrace{F_\phi(\xi^k; \delta_x | \hat{\theta}^k)}_{\text{directional derivative}} \Leftrightarrow \text{conditional gradient algorithm}$   
with step-size  $\frac{1}{k+1}$  (Wynn, 1970)

➤ some CV results for Bayesian estimation (Hu, 1998)

➤ no general CV results for LS and ML estimation,  
[unless  $\mathcal{X} = \{x^{(1)}, \dots, x^{(\ell)}\}$  finite (P 2009, 2010)]

# 8 Conclusions

DoE for nonlinear models with small data:

- ✓ Using the small-sample properties of the estimator can be a bit tricky

## 8 Conclusions

DoE for nonlinear models with small data:

- ✓ Using the small-sample properties of the estimator can be a bit tricky
- ✓ Numerical simulations are useful: for instance,
  - construct a locally optimum design (at  $\theta^0$ ),
  - simulate data (for another  $\theta^1$ ),
  - estimate  $\theta$  (correct the bias),
  - check closeness to  $\theta^1$  (plot marginals),
  - repeat (for other  $\theta^0$  and  $\theta^1$ ), etc.

## 8 Conclusions

DoE for nonlinear models with small data:

- ✓ Using the small-sample properties of the estimator can be a bit tricky
- ✓ Numerical simulations are useful: for instance,
  - construct a locally optimum design (at  $\theta^0$ ),
  - simulate data (for another  $\theta^1$ ),
  - estimate  $\theta$  (correct the bias),
  - check closeness to  $\theta^1$  (plot marginals),
  - repeat (for other  $\theta^0$  and  $\theta^1$ ), etc.
- ✓ When applicable, sequential (adaptive) design is a good remedy to the dependence of the optimal design in  $\theta$

## 8 Conclusions

DoE for nonlinear models with small data:

- ✓ Using the small-sample properties of the estimator can be a bit tricky
- ✓ Numerical simulations are useful: for instance,
  - construct a locally optimum design (at  $\theta^0$ ),
  - simulate data (for another  $\theta^1$ ),
  - estimate  $\theta$  (correct the bias),
  - check closeness to  $\theta^1$  (plot marginals),
  - repeat (for other  $\theta^0$  and  $\theta^1$ ), etc.
- ✓ When applicable, sequential (adaptive) design is a good remedy to the dependence of the optimal design in  $\theta$

Thank you for your attention !

# References I

- Bates, D., Watts, D., 1980. Relative curvature measures of nonlinearity. *Journal of Royal Statistical Society* B42, 1–25.
- Box, M., 1971. Bias in nonlinear estimation. *Journal of Royal Statistical Society* B33, 171–201.
- Chernoff, H., 1953. Locally optimal designs for estimating parameters. *Annals of Math. Stat.* 24, 586–602.
- Clarke, G., 1980. Moments of the least-squares estimators in a non-linear regression model. *Journal of Royal Statistical Society* B42, 227–237.
- Dette, H., Pepelyshev, A., 2010. Generalized latin hypercube design for computer experiments. *Technometrics* 52 (4), 421–429.
- Fedorov, V., 1972. *Theory of Optimal Experiments*. Academic Press, New York.
- Fedorov, V., Leonov, S., 2014. *Optimal Design for Nonlinear Response Models*. CRC Press, Boca Raton.
- Firth, D., 1993. Bias reduction of maximum likelihood estimates. *Biometrika* 80 (1), 27–38.
- Gauchi, J.-P., Pázman, A., 2006. Designs in nonlinear regression by stochastic minimization of functionals of the mean square error matrix. *Journal of Statistical Planning and Inference* 136, 1135–1152.
- Hamilton, D., Watts, D., 1985. A quadratic design criterion for precise estimation in nonlinear regression models. *Technometrics* 27, 241–250.
- Hu, I., 1998. On sequential designs in nonlinear problems. *Biometrika* 85 (2), 496–503.
- Kiefer, J., Wolfowitz, J., 1960. The equivalence of two extremum problems. *Canadian Journal of Mathematics* 12, 363–366.
- Mitchell, T., 1974. An algorithm for the construction of “*D*-optimal” experimental designs. *Technometrics* 16, 203–210.

# References II

- Pázman, A., 1986. Foundations of Optimum Experimental Design. Reidel (Kluwer group), Dordrecht (co-pub. VEDA, Bratislava).
- Pázman, A., 1993. Nonlinear Statistical Models. Kluwer, Dordrecht.
- Pázman, A., Pronzato, L., 1992. Nonlinear experimental design based on the distribution of estimators. *Journal of Statistical Planning and Inference* 33, 385–402.
- Pázman, A., Pronzato, L., 1996. A Dirac function method for densities of nonlinear statistics and for marginal densities in nonlinear regression. *Statistics & Probability Letters* 26, 159–167.
- Pázman, A., Pronzato, L., 1998. Approximate densities of two bias-corrected nonlinear LS estimators. In: Atkinson, A., Pronzato, L., Wynn, H. (Eds.), *MODA'5 – Advances in Model-Oriented Data Analysis and Experimental Design, Proceedings of the 5th Int. Workshop, Marseille, 22–26 juin 1998*. Physica Verlag, Heidelberg, pp. 145–152.
- Pázman, A., Pronzato, L., 2014. Optimum design accounting for the global nonlinear behavior of the model. *Annals of Statistics* 42 (4), 1426–1451.
- Picard, D., Prum, B., 1992. The bias of the MLE, an example of the behaviour of different corrections in genetic models. *Statistics* 23, 159–169.
- Pronzato, L., 2009. Asymptotic properties of nonlinear estimates in stochastic models with finite design space. *Statistics & Probability Letters* 79, 2307–2313.
- Pronzato, L., 2010. One-step ahead adaptive  $D$ -optimal design on a finite design space is asymptotically optimal. *Metrika* 71 (2), 219–238, ( DOI: 10.1007/s00184-008-0227-y).
- Pronzato, L., Pázman, A., July 1994a. Bias correction in nonlinear regression via two-stages least-squares estimation. In: Blanke, M., Söderström, T. (Eds.), *Prep. 10th IFAC/IFORS Symposium on Identification and System Parameter Estimation. Vol. 1*. Danish Automation Society, Copenhagen, pp. 137–142.

# References III

- Pronzato, L., Pázman, A., 1994b. Second-order approximation of the entropy in nonlinear least-squares estimation. *Kybernetika* 30 (2), 187–198, *Erratum* 32(1):104, 1996.
- Pronzato, L., Pázman, A., 2001. Using densities of estimators to compare pharmacokinetic experiments. *Computers in Biology and Medicine* 31 (3), 179–195.
- Pronzato, L., Pázman, A., 2013. *Design of Experiments in Nonlinear Models. Asymptotic Normality, Optimality Criteria and Small-Sample Properties.* Springer, LNS 212, New York.
- Pukelsheim, F., 1993. *Optimal Experimental Design.* Wiley, New York.
- Pukelsheim, F., Reider, S., 1992. Efficient rounding of approximate designs. *Biometrika* 79 (4), 763–770.
- Schwabe, R., 1996. *Optimum Designs for Multi-Factor Models.* Springer, New York.
- Silvey, S., 1980. *Optimal Design.* Chapman & Hall, London.
- Sternmüllerová, K., 2019. *Optimum design in nonlinear models.* Ph.D. Thesis, Comenius University, Bratislava, Slovakia.
- Tierney, L., Kass, R., Kadane, J., 1989. Approximate marginal densities of nonlinear functions. *Biometrika* 76 (3), 425–433.
- Valenzuela, P., Rojas, C., Hjalmarsson, H., 2015. Uncertainty in system identification: learning from the theory of risk. *IFAC-PapersOnLine* 48 (28), 1053–1058.
- Vila, J.-P., 1990. Exact experimental designs via stochastic optimization for nonlinear regression models. In: *Proc. Compstat, Int. Assoc. for Statistical Computing.* Physica Verlag, Heidelberg, pp. 291–296.
- Vila, J.-P., Gauchi, J.-P., 2007. Optimal designs based on exact confidence regions for parameter estimation of a nonlinear regression model. *Journal of Statistical Planning and Inference* 137, 2935–2953.
- Welch, W., 1982. Branch-and-bound search for experimental designs based on  $D$ -optimality and other criteria. *Technometrics* 24 (1), 41–28.
- Wynn, H., 1970. The sequential generation of  $D$ -optimum experimental designs. *Annals of Math. Stat.* 41, 1655–1664.