

Maximum Mean Discrepancy, Bayesian integration and kernel herding for space-filling design

Luc PRONZATO

(joint work with Anatoly ZHIGLJAVSKY, Cardiff Univ.)

Université Côte d'Azur, CNRS, France

Dec. 2, 2021

1 Space-filling design

Objective: approximate $f(\cdot)$ over \mathcal{X} (a compact subset of \mathbb{R}^d)
using pairs $(\mathbf{x}_i, f(\mathbf{x}_i))$, $i = 1, 2, \dots, n \rightarrow$ observe "everywhere"

Design $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

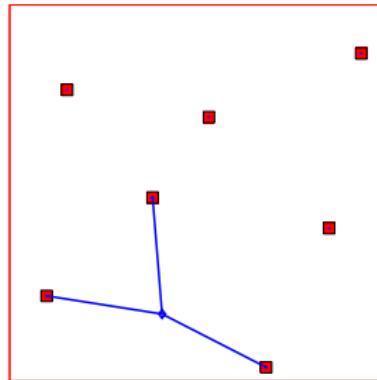
1 Space-filling design

Objective: approximate $f(\cdot)$ over \mathcal{X} (a compact subset of \mathbb{R}^d)
 using pairs $(\mathbf{x}_i, f(\mathbf{x}_i))$, $i = 1, 2, \dots, n \rightarrow$ observe "everywhere"

Design $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

Covering radius $\text{CR}(\mathbf{X}_n) \triangleq \max_{\mathbf{x} \in \mathcal{X}} \min_{\mathbf{x}_i} \|\mathbf{x} - \mathbf{x}_i\|$

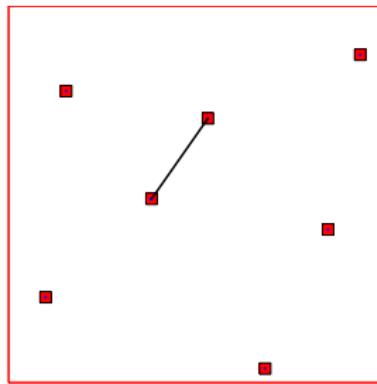
$\text{CR}(\mathbf{X}_n) = \text{fill distance} = \text{dispersion} = \text{miniMax distance criterion}$



→ we are never far from a design point

Packing radius $\text{PR}(\mathbf{X}_n) \triangleq \frac{1}{2} \min_{i \neq j} \|\mathbf{x}_i - \mathbf{x}_j\|$

$\text{PR}(\mathbf{X}_n)$ = separation radius = $\frac{1}{2}$ Maximin distance criterion

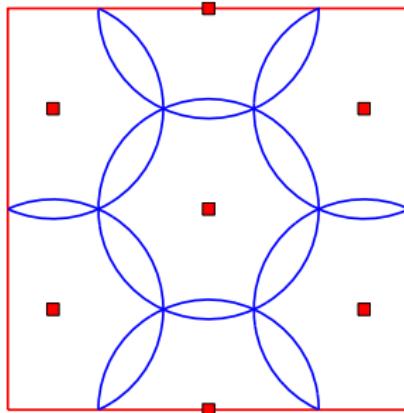


→ easier to compute, but pushes points to the boundary of \mathcal{X}

Examples:

① Covering, miniMax

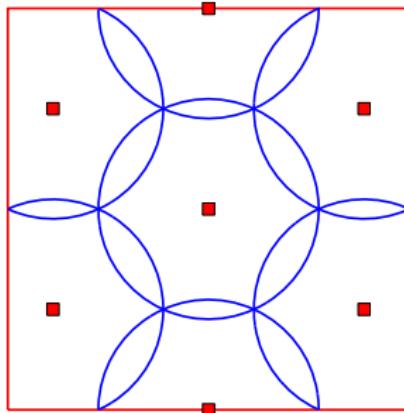
$d = 2, n = 7$ (radius= $\text{CR}(\mathbf{X}_n)$)



Examples:

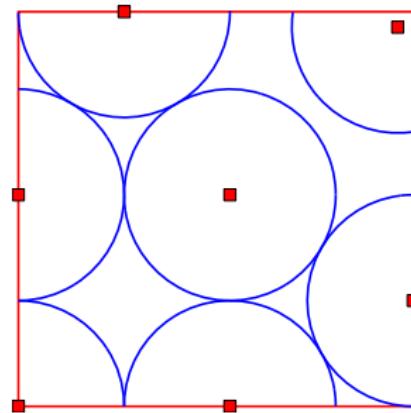
① Covering, miniMax

$$d = 2, n = 7 \text{ (radius} = \text{CR}(\mathbf{X}_n))$$



② Packing, Maximin

$$d = 2, n = 7 \text{ (radius} = \text{P}(\mathbf{X}_n))$$



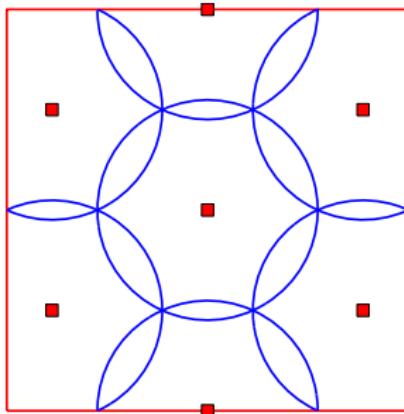
→ Minimising CR is preferable to maximising PR; both are difficult
 Bounds on approximation error are related to CR:

$$\|\text{error}\|_{L_\infty(\mathcal{X})} \leq \text{CR}(\mathbf{X}_n)^{\alpha-d/2} \text{ when } f \in \text{Sobolev space } W_2^\alpha(\mathcal{X})$$

Examples:

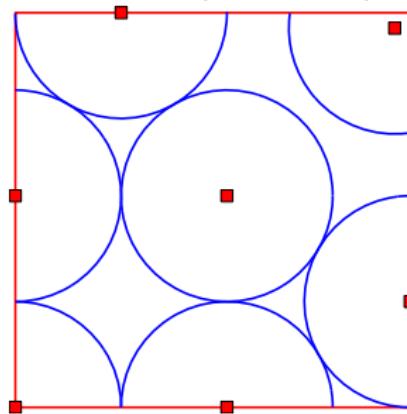
① Covering, miniMax

$$d = 2, n = 7 \text{ (radius} = \text{CR}(\mathbf{X}_n))$$



② Packing, Maximin

$$d = 2, n = 7 \text{ (radius} = \text{P}(\mathbf{X}_n))$$



→ Minimising CR is preferable to maximising PR; both are difficult
 Bounds on approximation error are related to CR:

$$\|\text{error}\|_{L_\infty(\mathcal{X})} \leq \text{CR}(\mathbf{X}_n)^{\alpha-d/2} \text{ when } f \in \text{Sobolev space } W_2^\alpha(\mathcal{X})$$

Any-time design: good-space-filling properties for $\mathbf{X}_k \subset \mathbf{X}_{k+1} \subset \mathbf{X}_{k+3} \subset \dots$?
 Can we apply a greedy gradient-type descent algorithm to a suitable criterion?

2 Bayesian integration

Notation:

- $\mathcal{M} = \mathcal{M}[\mathcal{X}]$ set of finite signed Borel measures on \mathcal{X}
- $\mathcal{M}(\alpha) = \{\mu \in \mathcal{M} : \mu(\mathcal{X}) = \alpha\}$ (signed measures with total mass $\alpha \in \mathbb{R}$)
- \mathcal{M}^+ positive measures in \mathcal{M}
- $\mathcal{M}^+(1)$ probability measures on \mathcal{X}

Objective: integrate f with respect to $\mu \in \mathcal{M}^+(1)$

→ estimate $I_\mu(f) = E_\mu\{f(\mathbf{X})\} \triangleq \int_{\mathcal{X}} f(\mathbf{x}) d\mu(\mathbf{x})$

2 Bayesian integration

Notation:

- $\mathcal{M} = \mathcal{M}[\mathcal{X}]$ set of finite signed Borel measures on \mathcal{X}
- $\mathcal{M}(\alpha) = \{\mu \in \mathcal{M} : \mu(\mathcal{X}) = \alpha\}$ (signed measures with total mass $\alpha \in \mathbb{R}$)
- \mathcal{M}^+ positive measures in \mathcal{M}
- $\mathcal{M}^+(1)$ probability measures on \mathcal{X}

Objective: integrate f with respect to $\mu \in \mathcal{M}^+(1)$

$$\rightarrow \text{estimate } I_\mu(f) = \mathbb{E}_\mu\{f(\mathbf{X})\} \triangleq \int_{\mathcal{X}} f(\mathbf{x}) d\mu(\mathbf{x})$$

Suppose $f(\mathbf{x}) = \beta + Z_x$, with Z_x a Gaussian RF,

$$\mathbb{E}\{Z_x\} = 0 \text{ and } \mathbb{E}\{Z_x Z_{x'}\} = \sigma^2 K(\mathbf{x}, \mathbf{x}')$$

and a vague prior on β ($\sim \mathcal{N}(0, \sigma^2 A)$ with $A \rightarrow +\infty$)

(K symmetric positive definite \rightarrow defines a RKHS \mathcal{H}_K)

\rightarrow observe f at $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

Bayesian integration of f :

Denote $\mathbf{1}_n = (1, \dots, 1)^T$, $\mathbf{y}_n = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^T$, $\{\mathbf{K}_n\}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$,

$$P_\mu(\mathbf{x}) \triangleq \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}') \quad (= \text{kernel imbedding of } \mu \text{ into } \mathcal{H}_K)$$

$$\mathcal{E}_K(\mu) \triangleq \int_{\mathcal{X}^2} K(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}) d\mu(\mathbf{x}') \quad (= \text{energy of } \mu)$$

Bayesian integration of f :

Denote $\mathbf{1}_n = (1, \dots, 1)^T$, $\mathbf{y}_n = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^T$, $\{\mathbf{K}_n\}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$,

$$P_\mu(\mathbf{x}) \triangleq \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}') \quad (= \text{kernel imbedding of } \mu \text{ into } \mathcal{H}_K)$$

$$\mathcal{E}_K(\mu) \triangleq \int_{\mathcal{X}^2} K(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}) d\mu(\mathbf{x}') \quad (= \text{energy of } \mu)$$

$I_\mu(f)$ has the normal posterior $\mathcal{N}(\hat{I}_n, \sigma^2 s_n^2)$ (O'Hagan, 1991)

$$\begin{aligned}\hat{I}_n &= \hat{\beta}^n + \mathbf{p}_n(\mu)^T \mathbf{K}_n^{-1} (\mathbf{y}_n - \hat{\beta}^n \mathbf{1}_n) \\ &= \hat{\mathbf{w}}_n^T \mathbf{y}_n \quad (\text{the estimation is linear})\end{aligned}$$

$$s_n^2 = \mathcal{E}_K(\mu) - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{p}_n(\mu) + \frac{(1 - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

where $\hat{\beta}^n = \frac{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{y}_n}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$ ($=$ BLUE) and $\mathbf{p}_n(\mu) = (P_\mu(\mathbf{x}_1), \dots, P_\mu(\mathbf{x}_n))^T$

The approximation of $f(\mathbf{x}_0)$ is a particular case!

Take $\mu = \delta_{\mathbf{x}_0} \rightarrow P_\mu(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}_0), \mathcal{E}_K(\mu) = K(\mathbf{x}_0, \mathbf{x}_0)$

and $f(\mathbf{x}_0)$ has the normal posterior $\mathcal{N}(\hat{\eta}_n(\mathbf{x}_0), \sigma^2 \rho_n^2(\mathbf{x}_0))$ with

$$\hat{\eta}_n(\mathbf{x}_0) = \hat{\beta}^n + \mathbf{k}_n^T(\mathbf{x}_0) \mathbf{K}_n^{-1} (\mathbf{y}_n - \hat{\beta}^n \mathbf{1}_n)$$

$$\rho_n^2(\mathbf{x}_0) = K(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}_n^T(\mathbf{x}_0) \mathbf{K}_n^{-1} \mathbf{k}_n(\mathbf{x}_0) + \frac{(1 - \mathbf{k}_n^T(\mathbf{x}_0) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

where $\{\mathbf{k}_n(\mathbf{x})\}_i = K(\mathbf{x}, \mathbf{x}_i)$

\equiv ordinary kriging

\rightarrow minimise $\boxed{\text{IMSPE}(\mathbf{X}_n) = \sigma^2 \int_{\mathcal{X}} \rho_n^2(\mathbf{x}) d\mu(\mathbf{x})}$

The approximation of $f(\mathbf{x}_0)$ is a particular case!

Take $\mu = \delta_{\mathbf{x}_0} \rightarrow P_\mu(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}_0), \mathcal{E}_K(\mu) = K(\mathbf{x}_0, \mathbf{x}_0)$

and $f(\mathbf{x}_0)$ has the normal posterior $\mathcal{N}(\hat{\eta}_n(\mathbf{x}_0), \sigma^2 \rho_n^2(\mathbf{x}_0))$ with

$$\hat{\eta}_n(\mathbf{x}_0) = \hat{\beta}^n + \mathbf{k}_n^T(\mathbf{x}_0) \mathbf{K}_n^{-1} (\mathbf{y}_n - \hat{\beta}^n \mathbf{1}_n)$$

$$\rho_n^2(\mathbf{x}_0) = K(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}_n^T(\mathbf{x}_0) \mathbf{K}_n^{-1} \mathbf{k}_n(\mathbf{x}_0) + \frac{(1 - \mathbf{k}_n^T(\mathbf{x}_0) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

where $\{\mathbf{k}_n(\mathbf{x})\}_i = K(\mathbf{x}, \mathbf{x}_i)$

\equiv ordinary kriging

\rightarrow minimise $\text{IMSPE}(\mathbf{X}_n) = \sigma^2 \int_{\mathcal{X}} \rho_n^2(\mathbf{x}) d\mu(\mathbf{x})$

- For μ uniform on \mathcal{X} , a design minimising s_n^2 should be well spread
- Minimise $s_n^2 = \mathcal{E}_K(\mu) - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{p}_n(\mu) + \frac{(1 - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$
instead of $\int_{\mathcal{X}} \rho_n^2(\mathbf{x}) d\mu(\mathbf{x})$ for space-filling design?

3 Kernel discrepancy, energy & potentials

Very much based on:

Damelin, S., Hickernell, F., Ragozin, D., Zeng, X., 2010. On energy, discrepancy and group invariant measures on measurable subsets of Euclidean space. *J. Fourier Anal. Appl.* 16, 813–839.

Sriperumbudur, B., Gretton, A., Fukumizu, K., Schölkopf, B., Lanckriet, G., 2010. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research* 11 (Apr), 1517–1561.

Sejdinovic, S., Sriperumbudur, B., Gretton, A., Fukumizu, K., 2013. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *The Annals of Statistics* 41 (5), 2263–2291.

Pronzato, L., Zhigljavsky, A., 2020. Bayesian quadrature, energy minimization and space-filling design. *SIAM/ASA J. Uncertainty Quantification* 8 (3), 959–1011.

3 Kernel discrepancy, energy & potentials

Very much based on:

Damelin, S., Hickernell, F., Ragozin, D., Zeng, X., 2010. On energy, discrepancy and group invariant measures on measurable subsets of Euclidean space. *J. Fourier Anal. Appl.* 16, 813–839.

Sriperumbudur, B., Gretton, A., Fukumizu, K., Schölkopf, B., Lanckriet, G., 2010. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research* 11 (Apr), 1517–1561.

Sejdinovic, S., Sriperumbudur, B., Gretton, A., Fukumizu, K., 2013. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *The Annals of Statistics* 41 (5), 2263–2291.

Pronzato, L., Zhigljavsky, A., 2020. Bayesian quadrature, energy minimization and space-filling design. *SIAM/ASA J. Uncertainty Quantification* 8 (3), 959–1011.

K = a symmetric positive definite kernel, defines a RKHS \mathcal{H}_K

For ν a signed measure,

$$\mathcal{E}_K(\nu) \triangleq \int_{\mathcal{X}^2} K(\mathbf{x}, \mathbf{x}') d\nu(\mathbf{x})d\nu(\mathbf{x}') = \text{energy of } \nu$$

$$P_\nu(\mathbf{x}) \triangleq \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{x}') d\nu(\mathbf{x}') = \text{potential of } \nu \text{ at } \mathbf{x}$$

$[P_\nu(\cdot) = \text{kernel imbedding of } \nu \text{ into } \mathcal{H}_K]$

Suppose $f \in \mathcal{H}_K$, μ, ν prob. measures with finite energy

$K_x(\cdot) = K(x, \cdot)$, reproducing property $f(x) = \langle f, K_x \rangle_{\mathcal{H}_K} \Rightarrow$

$$|I_\mu(f) - I_\nu(f)| = \left| \int_{\mathcal{X}} \langle f, K_x \rangle_{\mathcal{H}_K} d(\mu - \nu)(x) \right| = |\langle f, P_\mu - P_\nu \rangle_{\mathcal{H}_K}|$$

CS inequality \rightarrow a Kokosma-Hlawka type inequality:

$$\left| \int_{\mathcal{X}} f(x) d\nu(x) - \int_{\mathcal{X}} f(x) d\mu(x) \right| \leq \|f\|_{\mathcal{H}_K} \gamma_K(\mu, \nu)$$

where $\boxed{\gamma_K^2(\mu, \nu) \triangleq \|P_\mu - P_\nu\|_{\mathcal{H}_K}^2 = \mathcal{E}_K(\nu - \mu)}$

$$\gamma_K(\mu, \nu) = \sup_{\|f\|_{\mathcal{H}_K}=1} |\int_{\mathcal{X}} f(x) d\nu(x) - \int_{\mathcal{X}} f(x) d\mu(x)|$$

= **Maximum Mean Discrepancy** between (MMD) μ and ν
 (Sriperumbudur et al., 2010; Sejdinovic et al., 2013)

Suppose $f \in \mathcal{H}_K$, μ, ν prob. measures with finite energy

$K_x(\cdot) = K(x, \cdot)$, reproducing property $f(x) = \langle f, K_x \rangle_{\mathcal{H}_K} \Rightarrow$

$$|I_\mu(f) - I_\nu(f)| = \left| \int_{\mathcal{X}} \langle f, K_x \rangle_{\mathcal{H}_K} d(\mu - \nu)(x) \right| = |\langle f, P_\mu - P_\nu \rangle_{\mathcal{H}_K}|$$

CS inequality \rightarrow a Kokosma-Hlawka type inequality:

$$\left| \int_{\mathcal{X}} f(x) d\nu(x) - \int_{\mathcal{X}} f(x) d\mu(x) \right| \leq \|f\|_{\mathcal{H}_K} \gamma_K(\mu, \nu)$$

where $\boxed{\gamma_K^2(\mu, \nu) \triangleq \|P_\mu - P_\nu\|_{\mathcal{H}_K}^2 = \mathcal{E}_K(\nu - \mu)}$

$$\gamma_K(\mu, \nu) = \sup_{\|f\|_{\mathcal{H}_K}=1} |\int_{\mathcal{X}} f(x) d\nu(x) - \int_{\mathcal{X}} f(x) d\mu(x)|$$

= **Maximum Mean Discrepancy** between (MMD) μ and ν
 (Sriperumbudur et al., 2010; Sejdinovic et al., 2013)

Space-filling design: take μ uniform on \mathcal{X}

\rightarrow find ξ_n (with n support points) minimising $\mathcal{E}_K(\xi_n - \mu) = \text{MMD}^2(\xi_n, \mu)$

$\gamma_K(\cdot, \cdot)$ defines a pseudo-metric on $\mathcal{M}^+(1)$

Does it define a metric? ($\Leftrightarrow K$ is **characteristic**)

Definition

K is Integrally Strictly Positive Definite (ISPD) on \mathcal{M} when $\mathcal{E}_K(\nu) > 0$ for any nonzero $\nu \in \mathcal{M}$

Definition

K is Conditionally Integrally Strictly Positive Definite (CISPD) on \mathcal{M} when it is ISPD on $\mathcal{M}(0)$; that is, when $\mathcal{E}_K(\nu) > 0$ for all nonzero $\nu \in \mathcal{M}$ with $\nu(\mathcal{X}) = 0$

$\gamma_K(\cdot, \cdot)$ defines a pseudo-metric on $\mathcal{M}^+(1)$

Does it define a metric? ($\Leftrightarrow K$ is **characteristic**)

Definition

K is Integrally Strictly Positive Definite (ISPD) on \mathcal{M} when $\mathcal{E}_K(\nu) > 0$ for any nonzero $\nu \in \mathcal{M}$

Definition

K is Conditionally Integrally Strictly Positive Definite (CISPD) on \mathcal{M} when it is ISPD on $\mathcal{M}(0)$; that is, when $\mathcal{E}_K(\nu) > 0$ for all nonzero $\nu \in \mathcal{M}$ with $\nu(\mathcal{X}) = 0$

Sriperumbudur et al. (2010):

- K bounded & ISPD $\Rightarrow K$ is strictly positive definite
(\rightarrow defines a RKHS \mathcal{H}_K)
- if K uniformly bounded: characteristic \Leftrightarrow CISPD

$\gamma_K(\cdot, \cdot)$ defines a pseudo-metric on $\mathcal{M}^+(1)$

Does it define a metric? ($\Leftrightarrow K$ is **characteristic**)

Definition

K is *Integrally Strictly Positive Definite (ISPD)* on \mathcal{M} when $\mathcal{E}_K(\nu) > 0$ for any nonzero $\nu \in \mathcal{M}$

Definition

K is *Conditionally Integrally Strictly Positive Definite (CISPD)* on \mathcal{M} when it is ISPD on $\mathcal{M}(0)$; that is, when $\mathcal{E}_K(\nu) > 0$ for all nonzero $\nu \in \mathcal{M}$ with $\nu(\mathcal{X}) = 0$

Sriperumbudur et al. (2010):

- K bounded & ISPD $\Rightarrow K$ is strictly positive definite
(\rightarrow defines a RKHS \mathcal{H}_K)
- if $\underbrace{K \text{ uniformly bounded}}$: characteristic \Leftrightarrow CISPD
assumed in the following

- Spectral interpretation (Sriperumbudur et al., 2010):

if K translation invariant, $K(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x} - \mathbf{x}')$, with ψ bounded, continuous, pos. def., $\psi(\mathbf{x}) = \int_{\mathbb{R}^d} e^{-i\mathbf{x}^T \omega} d\Lambda(\omega)$ (Bochner)

$\gamma_K(\mu, \nu) = L_2$ distance between characteristic functions of μ and ν :

$$\gamma_K(\mu, \nu) = \left[\int_{\mathbb{R}^d} |\Phi_\mu(\omega) - \Phi_\nu(\omega)|^2 d\Lambda(\omega) \right]^{1/2}$$

K characteristic $\Leftrightarrow \Lambda$ supported on \mathbb{R}^d
 (for example, Matérn kernels)

For many kernels K :

- $\gamma_K(\cdot, \cdot)$ defines a metric for probability measures
- $\mathcal{E}_K(\cdot)$ is strictly convex

For many kernels K :

- $\gamma_K(\cdot, \cdot)$ defines a metric for probability measures
- $\mathcal{E}_K(\cdot)$ is strictly convex

Examples:

- Matérn kernels, e.g., $K_{1/2,\theta}(\mathbf{x}, \mathbf{x}') = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|)$,

$$K_{3/2,\theta}(\mathbf{x}, \mathbf{x}') = (1 + \sqrt{3}\theta \|\mathbf{x} - \mathbf{x}'\|) \exp(-\sqrt{3}\theta \|\mathbf{x} - \mathbf{x}'\|)$$

$$K_{5/2,\theta}(\mathbf{x}, \mathbf{x}') = [1 + \sqrt{5}\theta \|\mathbf{x} - \mathbf{x}'\| + 5\theta^2 \|\mathbf{x} - \mathbf{x}'\|^2/3] \exp(-\sqrt{5}\theta \|\mathbf{x} - \mathbf{x}'\|)$$

...

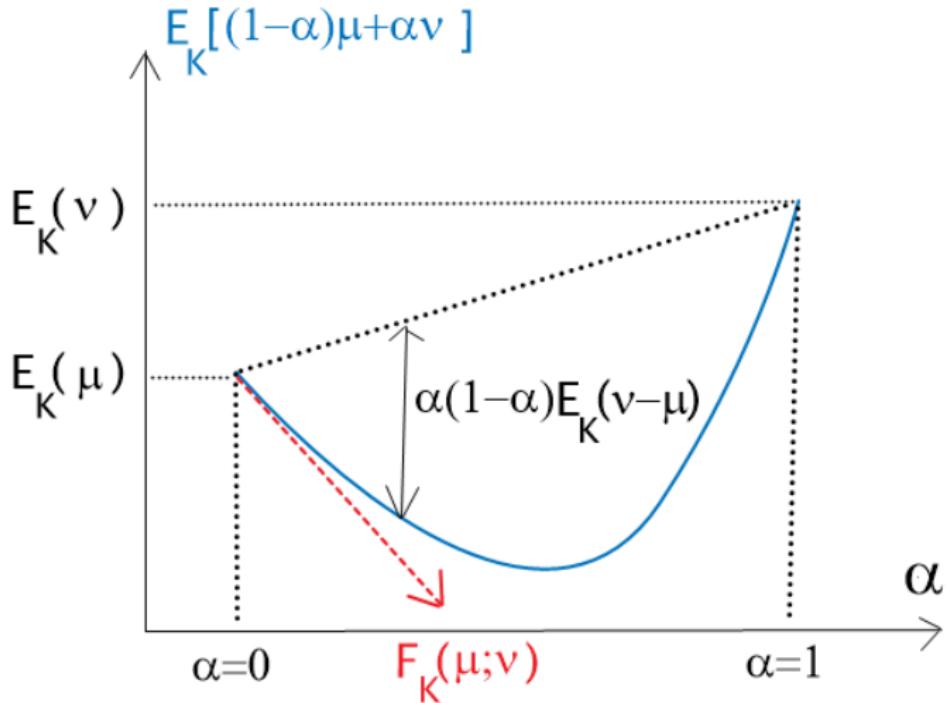
- distance-induced kernels (\rightarrow energy distance of Székely and Rizzo (2013)):

$$K'_{(s)}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x}\|^s + \|\mathbf{x}'\|^s - \|\mathbf{x} - \mathbf{x}'\|^s, \quad 0 < s < 2$$

- Riesz kernels (unbounded):

$$K_{(s)}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^{-s}, \quad s \in (0, d), \quad K_{(0)}(\mathbf{x}, \mathbf{x}') = -\log \|\mathbf{x} - \mathbf{x}'\|$$

- ...



For $\alpha \in (0, 1)$, $\mu, \nu \in \mathcal{M}$:

$$(1 - \alpha) \mathcal{E}_K(\mu) + \alpha \mathcal{E}_K(\nu) - \mathcal{E}_K[(1 - \alpha)\mu + \alpha\nu] = \alpha(1 - \alpha) \mathcal{E}_K(\nu - \mu)$$

- K ISPD $\Leftrightarrow \mathcal{E}_K(\cdot)$ strictly convex on \mathcal{M}
- K CISPD $\Leftrightarrow \mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}(1)$ (or $\mathcal{M}^+(1)$)

For $\alpha \in (0, 1)$, $\mu, \nu \in \mathcal{M}$:

$$(1 - \alpha) \mathcal{E}_K(\mu) + \alpha \mathcal{E}_K(\nu) - \mathcal{E}_K[(1 - \alpha)\mu + \alpha\nu] = \alpha(1 - \alpha) \mathcal{E}_K(\nu - \mu)$$

- K ISPD $\Leftrightarrow \mathcal{E}_K(\cdot)$ strictly convex on \mathcal{M}
- K CISPD $\Leftrightarrow \mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}(1)$ (or $\mathcal{M}^+(1)$)

Directional derivative of $\mathcal{E}_K(\cdot)$ at μ in the direction ν :

$$\begin{aligned} F_K(\mu; \nu) &= \lim_{\alpha \rightarrow 0^+} \frac{\mathcal{E}_K[(1 - \alpha)\mu + \alpha\nu] - \mathcal{E}_K(\mu)}{\alpha} \\ &= 2 \left[\int_{\mathcal{X}^2} K(\mathbf{x}, \mathbf{x}') d\nu(\mathbf{x})d\mu(\mathbf{x}') - \mathcal{E}_K(\mu) \right] \end{aligned}$$

For $\alpha \in (0, 1)$, $\mu, \nu \in \mathcal{M}$:

$$(1 - \alpha) \mathcal{E}_K(\mu) + \alpha \mathcal{E}_K(\nu) - \mathcal{E}_K[(1 - \alpha)\mu + \alpha\nu] = \alpha(1 - \alpha) \mathcal{E}_K(\nu - \mu)$$

- K ISPD $\Leftrightarrow \mathcal{E}_K(\cdot)$ strictly convex on \mathcal{M}
- K CISPD $\Leftrightarrow \mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}(1)$ (or $\mathcal{M}^+(1)$)

Directional derivative of $\mathcal{E}_K(\cdot)$ at μ in the direction ν :

$$\begin{aligned} F_K(\mu; \nu) &= \lim_{\alpha \rightarrow 0^+} \frac{\mathcal{E}_K[(1 - \alpha)\mu + \alpha\nu] - \mathcal{E}_K(\mu)}{\alpha} \\ &= 2 \left[\int_{\mathcal{X}^2} K(\mathbf{x}, \mathbf{x}') d\nu(\mathbf{x}) d\mu(\mathbf{x}') - \mathcal{E}_K(\mu) \right] \end{aligned}$$

$\Rightarrow F_K(\mu; \delta_{\mathbf{x}})$ is related to the potential $P_\mu(\mathbf{x})$:

$$P_\mu(\mathbf{x}) = \frac{1}{2} F_K(\mu; \delta_{\mathbf{x}}) + \mathcal{E}_K(\mu)$$

Suppose \mathcal{X} compact

Theorem (Equivalence theorem in DOE)

If $\mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}^+(1)$, $\mu_K^* \in \mathcal{M}^+(1)$ is the minimum-energy probability measure on \mathcal{X} iff $\forall \mathbf{x} \in \mathcal{X}$, $P_{\mu_K^*}(\mathbf{x}) \geq \mathcal{E}_K(\mu_K^*)$; moreover, $P_{\mu_K^*}(\mathbf{x}) = \mathcal{E}_K(\mu_K^*)$ on the support of μ_K^*

Suppose \mathcal{X} compact

Theorem (Equivalence theorem in DOE)

If $\mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}^+(1)$, $\mu_K^* \in \mathcal{M}^+(1)$ is the minimum-energy probability measure on \mathcal{X} iff $\forall \mathbf{x} \in \mathcal{X}$, $P_{\mu_K^*}(\mathbf{x}) \geq \mathcal{E}_K(\mu_K^*)$; moreover, $P_{\mu_K^*}(\mathbf{x}) = \mathcal{E}_K(\mu_K^*)$ on the support of μ_K^*

Potential theory (K singular, ISPD):

$\mu_K^* \in \mathcal{M}^+(1) =$ equilibrium measure of \mathcal{X}

$C_K = [\inf_{\mu \in \mathcal{M}^+(1)} \mathcal{E}_K(\mu)]^{-1} \geq 0 =$ capacity of \mathcal{X}

Consider now (signed) measures in $\mathcal{M}(1)$

Theorem

If $\mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}(1)$, $\tilde{\mu}_K^* \in \mathcal{M}(1)$ is the minimum-energy signed measure on \mathcal{X} with total charge 1 iff $\forall \mathbf{x} \in \mathcal{X}$, $P_{\tilde{\mu}_K^*}(\mathbf{x}) = \mathcal{E}_K(\tilde{\mu}_K^*)$

Consider now (signed) measures in $\mathcal{M}(1)$

Theorem

If $\mathcal{E}_K(\cdot)$ strictly convex on $\mathcal{M}(1)$, $\tilde{\mu}_K^* \in \mathcal{M}(1)$ is the minimum-energy signed measure on \mathcal{X} with total charge 1 iff $\forall \mathbf{x} \in \mathcal{X}$, $P_{\tilde{\mu}_K^*}(\mathbf{x}) = \mathcal{E}_K(\tilde{\mu}_K^*)$

... but $\tilde{\mu}_K^*$ does not always exist!

When minimum-energy signed measures are probability measures?

Theorem

Assume that K is ISPD and translation invariant, with $K(\mathbf{x}, \mathbf{x}') = \Psi(\mathbf{x} - \mathbf{x}')$ and Ψ continuous, twice differentiable except at the origin, with Laplacian $\Delta_\Psi(\mathbf{x}) = \sum_{i=1}^d \partial^2 \Psi(\mathbf{x}) / \partial x_i^2 \geq 0$, $\mathbf{x} \neq \mathbf{0}$. Then there exists a unique $\tilde{\mu}_K^*$ in $\mathcal{M}(1)$, and $\tilde{\mu}_K^* \in \mathcal{M}^+(1)$

(Generalises (Hájek, 1956) to $d > 1$, but $\Psi(\cdot)$ must have a singularity at zero.)

► minimise $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$

Relation with (continuous) BLUE of β in $f(\mathbf{x}) = \beta + Z_x$

with Z_x a Gaussian RF ($\mathbb{E}\{Z_x\} = 0$, $\mathbb{E}\{Z_x Z_{x'}\} = \sigma^2 K(\mathbf{x}, \mathbf{x}')$)
 Näther (1985, Sect. 4.2):

- Any linear estimator of β : $\hat{\beta} = \hat{\beta}(\xi) = \int_{\mathcal{X}} f(\mathbf{x}) d\xi(\mathbf{x})$, $\xi \in \mathcal{M}$
- $\hat{\beta}(\xi)$ unbiased: $\xi \in \mathcal{M}(1)$
- Variance: $V_\xi = \mathbb{E}\{(\hat{\beta}(\xi) - \beta)^2\} = \sigma^2 \mathcal{E}_K(\xi)$

Relation with (continuous) BLUE of β in $f(\mathbf{x}) = \beta + Z_x$

with Z_x a Gaussian RF ($\mathbb{E}\{Z_x\} = 0$, $\mathbb{E}\{Z_x Z_{x'}\} = \sigma^2 K(\mathbf{x}, \mathbf{x}')$)
 Näther (1985, Sect. 4.2):

- Any linear estimator of β : $\hat{\beta} = \hat{\beta}(\xi) = \int_{\mathcal{X}} f(\mathbf{x}) d\xi(\mathbf{x})$, $\xi \in \mathcal{M}$
 - $\hat{\beta}(\xi)$ unbiased: $\xi \in \mathcal{M}(1)$
 - Variance: $V_\xi = \mathbb{E}\{(\hat{\beta}(\xi) - \beta)^2\} = \sigma^2 \mathcal{E}_K(\xi)$
- Existence of a minimum-energy signed measure $\tilde{\mu}_K^* \in \mathcal{M}(1)$
 \Leftrightarrow existence of the BLUE $\hat{\beta}^* = \hat{\beta}(\tilde{\mu}_K^*)$, with variance $\sigma^2 \mathcal{E}_K(\tilde{\mu}_K^*)$

(The property $[\forall \mathbf{x} \in \mathcal{X}, P_{\tilde{\mu}_K^*}(\mathbf{x}) = \mathcal{E}_K(\tilde{\mu}_K^*)]$ corresponds to Wiener-Hopf equation in Grenander's theorem (1950)).

We do not want to minimise $\mathcal{E}_K(\nu)$ but $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$
for a given μ

Kernel reduction

Damelin et al. (2010): for μ with total mass 1, consider the kernel

$$\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - P_\mu(\mathbf{x}) - P_\mu(\mathbf{x}') + \mathcal{E}_K(\mu)$$

- $\forall \xi$ with total mass 1 , $\underbrace{\mathcal{E}_K(\xi - \mu)}_{\gamma_K^2(\xi, \mu)} = \mathcal{E}_{\tilde{K}_\mu}(\xi - \mu) = \mathcal{E}_{\tilde{K}_\mu}(\xi)$

For μ, ξ having total mass 1,

minimising $\text{MMD } \gamma_K(\xi, \mu)$ w.r.t. $\xi \Leftrightarrow$ minimising $\text{MMD } \gamma_{\tilde{K}_\mu}(\xi, \mu)$
 \Leftrightarrow minimising $\mathcal{E}_{\tilde{K}_\mu}(\xi)$

We do not want to minimise $\mathcal{E}_K(\nu)$ but $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$
for a given μ

Kernel reduction

Damelin et al. (2010): for μ with total mass 1, consider the kernel

$$\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - P_\mu(\mathbf{x}) - P_\mu(\mathbf{x}') + \mathcal{E}_K(\mu)$$

- $\forall \xi$ with total mass 1 , $\underbrace{\mathcal{E}_K(\xi - \mu)}_{\gamma_K^2(\xi, \mu)} = \mathcal{E}_{\tilde{K}_\mu}(\xi - \mu) = \mathcal{E}_{\tilde{K}_\mu}(\xi)$

For μ, ξ having total mass 1,

$$\begin{aligned} \text{minimising MMD } \gamma_K(\xi, \mu) \text{ w.r.t. } \xi &\Leftrightarrow \text{minimising MMD } \gamma_{\tilde{K}_\mu}(\xi, \mu) \\ &\Leftrightarrow \text{minimising } \mathcal{E}_{\tilde{K}_\mu}(\xi) \end{aligned}$$

[For $\mu = \delta_{\mathbf{x}_0}$ (function approximation):

$$\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - K(\mathbf{x}, \mathbf{x}_0) - K(\mathbf{x}', \mathbf{x}_0) + K(\mathbf{x}_0, \mathbf{x}_0)]$$

We do not want to minimise $\mathcal{E}_K(\nu)$ but $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$
for a given μ

Kernel reduction

Damelin et al. (2010): for μ with total mass 1, consider the kernel

$$\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - P_\mu(\mathbf{x}) - P_\mu(\mathbf{x}') + \mathcal{E}_K(\mu)$$

- $\forall \xi$ with total mass 1 , $\underbrace{\mathcal{E}_K(\xi - \mu)}_{\gamma_K^2(\xi, \mu)} = \mathcal{E}_{\tilde{K}_\mu}(\xi - \mu) = \mathcal{E}_{\tilde{K}_\mu}(\xi)$

For μ, ξ having total mass 1,

minimising MMD $\gamma_K(\xi, \mu)$ w.r.t. $\xi \Leftrightarrow$ minimising MMD $\gamma_{\tilde{K}_\mu}(\xi, \mu)$
 \Leftrightarrow minimising $\mathcal{E}_{\tilde{K}_\mu}(\xi)$

[For $\mu = \delta_{\mathbf{x}_0}$ (function approximation):

$$\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - K(\mathbf{x}, \mathbf{x}_0) - K(\mathbf{x}', \mathbf{x}_0) + K(\mathbf{x}_0, \mathbf{x}_0)]$$

▶ compute $\tilde{K}_\mu(\mathbf{x}, \mathbf{x}')$

Back to the (continuous) BLUE of β

$$f(\mathbf{x}) = \beta + Z_x, \quad Z_x \text{ a Gaussian RF}, \quad \mathbb{E}\{Z_x\} = 0, \quad \mathbb{E}\{Z_x Z_{x'}\} = \sigma^2 K(\mathbf{x}, \mathbf{x}')$$

Denote $\mathfrak{p} \triangleq$ orthogonal projection of $L^2(\mathcal{X}, \mu)$ onto the linear space spanned by constant 1

$$\begin{aligned} a) \quad f(\mathbf{x}) &= \beta + \underbrace{Z_x}_{\mathfrak{p}Z_x + (\text{Id}_{L^2} - \mathfrak{p})Z_x} \rightarrow \sigma^2 K \end{aligned}$$

a) $\text{var}\{\text{BLUE}(\beta)\} = \sigma^2 \mathcal{E}_K(\tilde{\mu}_K^*)$ (if $\text{BLUE}(\beta)$ exists...)

Back to the (continuous) BLUE of β

$$f(\mathbf{x}) = \beta + Z_x, \quad Z_x \text{ a Gaussian RF}, \quad \mathbb{E}\{Z_x\} = 0, \quad \mathbb{E}\{Z_x Z_{x'}\} = \sigma^2 K(\mathbf{x}, \mathbf{x}')$$

Denote $\mathfrak{p} \triangleq$ orthogonal projection of $L^2(\mathcal{X}, \mu)$ onto the linear space spanned by constant 1

$$a) \quad f(\mathbf{x}) = \underbrace{\beta}_{\mathfrak{p}Z_x} + \underbrace{Z_x}_{(\text{Id}_{L^2} - \mathfrak{p})Z_x} \rightarrow \sigma^2 K$$

$$b) \quad f(\mathbf{x}) = \underbrace{\beta'}_{\beta + \mathfrak{p}Z_x} + \underbrace{\tilde{Z}_x}_{(\text{Id}_{L^2} - \mathfrak{p})Z_x} \rightarrow \boxed{\sigma^2 \tilde{K}_\mu}$$

a) $\text{var}\{\text{BLUE}(\beta)\} = \sigma^2 \mathcal{E}_K(\tilde{\mu}_K^*)$ (if $\text{BLUE}(\beta)$ exists...)

b) $\text{var}\{\text{BLUE}(\beta')\} = \sigma^2 \mathcal{E}_{\tilde{K}_\mu}(\mu) = 0$ (with $\text{BLUE}(\beta') = I_\mu(f)$)

How to compute $\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - P_\mu(\mathbf{x}) - P_\mu(\mathbf{x}') + \mathcal{E}_K(\mu)$?

How to compute $\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - P_\mu(\mathbf{x}) - P_\mu(\mathbf{x}') + \mathcal{E}_K(\mu)$?

Tensor product kernels:

Take K separable on $\mathcal{X} = \times_{i=1}^d \mathcal{X}_i$:

$$K(\mathbf{x}, \mathbf{x}') = K^{(t)}(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d K_i(x_i, x'_i)$$

Szabó and Sriperumbudur (2018): each K_i continuous, bounded, translation invariant and CISPD on $\mathcal{M}^{(i)} = \mathcal{M}[\mathcal{X}_i] \Rightarrow K^{(t)}$ is CISPD on $\mathcal{M}[\mathcal{X}]$

How to compute $\tilde{K}_\mu(\mathbf{x}, \mathbf{x}') \triangleq K(\mathbf{x}, \mathbf{x}') - P_\mu(\mathbf{x}) - P_\mu(\mathbf{x}') + \mathcal{E}_K(\mu)$?

Tensor product kernels:

Take K separable on $\mathcal{X} = \times_{i=1}^d \mathcal{X}_i$:

$$K(\mathbf{x}, \mathbf{x}') = K^{(t)}(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d K_i(x_i, x'_i)$$

Szabó and Sriperumbudur (2018): each K_i continuous, bounded, translation invariant and CISPD on $\mathcal{M}^{(i)} = \mathcal{M}[\mathcal{X}_i] \Rightarrow K^{(t)}$ is CISPD on $\mathcal{M}[\mathcal{X}]$

When $\mu = \otimes_{i=1}^d \mu^{(i)}$ is a product measure on $\mathcal{X} = \times_{i=1}^d \mathcal{X}_i$, then

$$\mathcal{E}_{K^{(t)}}(\mu) = \prod_{i=1}^d \mathcal{E}_{K_i}(\mu^{(i)}), \quad P_\mu(\mathbf{x}) = \prod_{i=1}^d P_{\mu^{(i)}}(x_i)$$

and $\mathcal{E}_{K_i}(\mu^{(i)})$ and $P_{\mu^{(i)}}(x)$ easily computed (analytical expressions available for some $\mu^{(i)}$ and K_i on $\mathcal{X}_i = [a, b]$ — Matérn, Riesz...)

4 Design of experiments

μ is uniform on \mathcal{X}

Minimise $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$ w.r.t. $\xi \rightarrow \xi = \mu$, not very useful...

4 Design of experiments

μ is uniform on \mathcal{X}

Minimise $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$ w.r.t. $\xi \rightarrow \xi = \mu$, not very useful...

We want a discrete measure !

- $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ a n -point design ($\mathbf{x}_i \in \mathcal{X}$ for all i)
- $\xi_n = \sum_{i=1}^n w_i \delta_{\mathbf{x}_i}$ a finite signed measure supported on \mathbf{X}_n ,
 $\mathbf{w}_n = (w_1, \dots, w_n)^\top$

$$\rightarrow \mathcal{E}_K(\xi_n - \mu) = \mathcal{E}_K(\mu) - 2\mathbf{w}_n^\top \mathbf{p}_n(\mu) + \mathbf{w}_n^\top \mathbf{K}_n \mathbf{w}_n$$

where $\mathbf{p}_n(\mu) = (P_\mu(\mathbf{x}_1), \dots, P_\mu(\mathbf{x}_n))^\top$ and $\{\mathbf{K}_n\}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$

4 Design of experiments

μ is uniform on \mathcal{X}

Minimise $\mathcal{E}_K(\xi - \mu) = \text{MMD}^2(\xi, \mu)$ w.r.t. $\xi \rightarrow \xi = \mu$, not very useful...

We want a discrete measure !

- $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ a n -point design ($\mathbf{x}_i \in \mathcal{X}$ for all i)
- $\xi_n = \sum_{i=1}^n w_i \delta_{\mathbf{x}_i}$ a finite signed measure supported on \mathbf{X}_n ,
 $\mathbf{w}_n = (w_1, \dots, w_n)^\top$

$$\rightarrow \mathcal{E}_K(\xi_n - \mu) = \mathcal{E}_K(\mu) - 2\mathbf{w}_n^\top \mathbf{p}_n(\mu) + \mathbf{w}_n^\top \mathbf{K}_n \mathbf{w}_n$$

where $\mathbf{p}_n(\mu) = (P_\mu(\mathbf{x}_1), \dots, P_\mu(\mathbf{x}_n))^\top$ and $\{\mathbf{K}_n\}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$

$$\rightarrow \mathcal{E}_K(\xi_n - \mu) = \mathcal{E}_{\tilde{\mathbf{K}}_\mu}(\xi_n) = \mathbf{w}_n^\top \tilde{\mathbf{K}}_\mu \mathbf{w}_n$$

\mathbf{X}_n fixed, optimal weights $\hat{\mathbf{w}}_n$ summing to 1: Bayesian integration

$$\hat{\mathbf{w}}_n = \frac{\tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^\top \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$$

\mathbf{X}_n fixed, optimal weights $\hat{\mathbf{w}}_n$ summing to 1: Bayesian integration

$$\hat{\mathbf{w}}_n = \frac{\tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^\top \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$$

$$\begin{cases} \hat{\mathbf{w}}_n &= \text{weights} \\ \sigma^2 \mathcal{E}_K(\hat{\xi}_n - \mu) &= \text{variance} \end{cases} \quad \text{for Bayesian integration w.r.t. } \mu$$

in the model

$$f(\mathbf{x}) = \beta + Z_x \quad \text{with } E\{Z_x\} = 0 \text{ and } E\{Z_x Z_{x'}\} = \sigma^2 K(\mathbf{x}, \mathbf{x}')$$

- Posterior mean: $\hat{I}_n = \hat{\mathbf{w}}_n^\top \mathbf{y}_n$, with $\mathbf{y}_n = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^\top$
- Posterior variance $\sigma^2 s_n^2 = \sigma^2 \mathcal{E}_K(\hat{\xi}_n - \mu) = \sigma^2 \mathcal{E}_{\tilde{\mathbf{K}}_\mu}(\hat{\xi}_n) = \frac{\sigma^2}{\mathbf{1}_n^\top \tilde{\mathbf{K}}_\mu^{-1} \mathbf{1}_n}$

► all weights = $1/n$

... ?!? initially we had

$$\begin{aligned}\hat{\mathbf{w}}_n &= \left(\mathbf{K}_n^{-1} - \frac{\mathbf{K}_n^{-1} \mathbf{1}_n \mathbf{1}_n^T \mathbf{K}_n^{-1}}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n} \right) \mathbf{p}_n(\mu) + \frac{\mathbf{K}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n} \\ s_n^2 &= \mathcal{E}_{\mathbf{K}}(\mu) - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{p}_n(\mu) + \frac{(1 - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}\end{aligned}$$

...?!? initially we had

$$\hat{\mathbf{w}}_n = \left(\mathbf{K}_n^{-1} - \frac{\mathbf{1}_n \mathbf{1}_n^T \mathbf{K}_n^{-1}}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n} \right) \mathbf{p}_n(\mu) + \frac{\mathbf{K}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

$$s_n^2 = \mathcal{E}_K(\mu) - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{p}_n(\mu) + \frac{(1 - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

$\mathfrak{p} \triangleq$ orthogonal projection of $L^2(\mathcal{X}, \mu)$ onto the linear space spanned by 1

$$f(\mathbf{x}) = \beta + \underbrace{\mathfrak{p} Z_x + (\text{Id}_{L^2} - \mathfrak{p}) Z_x}_{\rightarrow \sigma^2 K} = \underbrace{\beta' + (\beta + \mathfrak{p} Z_x)}_{\beta + \mathfrak{p} Z_x} + \underbrace{(\text{Id}_{L^2} - \mathfrak{p}) Z_x}_{\rightarrow \sigma^2 \tilde{K}_\mu}$$

New trend β' , but the predictions are not modified when $K \rightarrow \tilde{K}_\mu$

... ?!? initially we had

$$\hat{\mathbf{w}}_n = \left(\mathbf{K}_n^{-1} - \frac{\mathbf{K}_n^{-1} \mathbf{1}_n \mathbf{1}_n^T \mathbf{K}_n^{-1}}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n} \right) \mathbf{p}_n(\mu) + \frac{\mathbf{K}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

$$s_n^2 = \mathcal{E}_K(\mu) - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{p}_n(\mu) + \frac{(1 - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

$\mathbf{p} \triangleq$ orthogonal projection of $L^2(\mathcal{X}, \mu)$ onto the linear space spanned by $\mathbf{1}$

$$f(\mathbf{x}) = \beta + \underbrace{Z_x}_{\mathbf{p} Z_x + (\text{Id}_{L^2} - \mathbf{p}) Z_x} \rightarrow \sigma^2 K \quad \underbrace{\beta'}_{\beta + \mathbf{p} Z_x} + \underbrace{\tilde{Z}_x}_{(\text{Id}_{L^2} - \mathbf{p}) Z_x} \rightarrow \sigma^2 \tilde{K}_\mu$$

New trend β' , but the predictions are not modified when $K \rightarrow \tilde{K}_\mu$... and
 $\mathcal{E}_{\tilde{K}_\mu}(\mu) = 0$, $\tilde{\mathbf{p}}_n(\mu) = \mathbf{0}$

$$\hat{\mathbf{w}}_n = \left(\tilde{\mathbf{K}}_n^{-1} - \frac{\tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n \mathbf{1}_n^T \tilde{\mathbf{K}}_n^{-1}}{\mathbf{1}_n^T \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n} \right) \tilde{\mathbf{p}}_n(\mu) + \frac{\tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^T \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$$

$$s_n^2 = \mathcal{E}_{\tilde{K}_\mu}(\mu) - \tilde{\mathbf{p}}_n^T(\mu) \tilde{\mathbf{K}}_n^{-1} \tilde{\mathbf{p}}_n(\mu) + \frac{(1 - \tilde{\mathbf{p}}_n^T(\mu) \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$$

... ?!? initially we had

$$\hat{\mathbf{w}}_n = \left(\mathbf{K}_n^{-1} - \frac{\mathbf{K}_n^{-1} \mathbf{1}_n \mathbf{1}_n^T \mathbf{K}_n^{-1}}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n} \right) \mathbf{p}_n(\mu) + \frac{\mathbf{K}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

$$s_n^2 = \mathcal{E}_K(\mu) - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{p}_n(\mu) + \frac{(1 - \mathbf{p}_n^T(\mu) \mathbf{K}_n^{-1} \mathbf{1}_n)^2}{\mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n}$$

$\mathfrak{p} \triangleq$ orthogonal projection of $L^2(\mathcal{X}, \mu)$ onto the linear space spanned by 1

$$f(\mathbf{x}) = \beta + \underbrace{\mathfrak{p} Z_x + (\text{Id}_{L^2} - \mathfrak{p}) Z_x}_{\rightarrow \sigma^2 K} = \underbrace{\beta'}_{\beta + \mathfrak{p} Z_x} + \underbrace{\widetilde{Z}_x}_{(\text{Id}_{L^2} - \mathfrak{p}) Z_x \rightarrow \sigma^2 \widetilde{K}_\mu}$$

New trend β' , but the predictions are not modified when $K \rightarrow \widetilde{K}_\mu$... and
 $\mathcal{E}_{\widetilde{K}_\mu}(\mu) = 0$, $\widetilde{\mathbf{p}}_n(\mu) = \mathbf{0}$

$$\hat{\mathbf{w}}_n = \frac{\widetilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}{\mathbf{1}_n^T \widetilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$$

$$s_n^2 = \frac{1}{\mathbf{1}_n^T \widetilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$$

Equivalence with estimation in a location model

$$\left\{ \begin{array}{l} \hat{\mathbf{w}}_n = \text{weights} \\ \sigma^2 \mathcal{E}_K(\hat{\xi}_n - \mu) = \text{variance} \end{array} \right. \text{ for (discrete) BLUE of } \beta' \text{ in the model}$$

$f(\mathbf{x}) = \beta' + \tilde{Z}_{\mathbf{x}}$ with $E\{\tilde{Z}_{\mathbf{x}}\} = 0$, $E\{\tilde{Z}_{\mathbf{x}} \tilde{Z}_{\mathbf{x}'}\} = \sigma^2 \tilde{K}_{\mu}(\mathbf{x}, \mathbf{x}')$

- $\hat{I}_n = \hat{\beta}'^n =$ (discrete) BLUE of β'
- $\sigma^2 s_n^2 = \text{var}[\hat{\beta}'^n]$

(design for) Bayesian integration ⇔ (design for) parameter estimation
 in a model with errors having
 a modified covariance

Equivalence with estimation in a location model

$\left\{ \begin{array}{l} \hat{\mathbf{w}}_n \\ \sigma^2 \mathcal{E}_K(\hat{\xi}_n - \mu) \end{array} \right. \begin{array}{l} = \text{weights} \\ = \text{variance} \end{array}$ for **(discrete) BLUE** of β' in the model

$$f(\mathbf{x}) = \beta' + \tilde{Z}_{\mathbf{x}} \quad \text{with } E\{\tilde{Z}_{\mathbf{x}}\} = 0, E\{\tilde{Z}_{\mathbf{x}} \tilde{Z}_{\mathbf{x}'}\} = \sigma^2 \tilde{K}_{\mu}(\mathbf{x}, \mathbf{x}')$$

→ $\hat{I}_n = \hat{\beta}'^n =$ (discrete) BLUE of β'

→ $\sigma^2 s_n^2 = \text{var}[\hat{\beta}'^n]$

(design for) Bayesian integration ⇔ (design for) parameter estimation
 in a model with errors having
 a modified covariance

$[\mu = \delta_{\mathbf{x}_0}$ yields a similar equivalence

between function approximation (at \mathbf{x}_0) and parameter estimation]

1/ Fix the weights to $1/n$, optimise w.r.t. \mathbf{X}_n

→ choose \mathbf{X}_n that minimises $\mathcal{E}_K(\xi_n - \mu) = \mathbf{1}_n^\top \tilde{\mathbf{K}}_n \mathbf{1}_n / n^2$

→ *Projection kernel* of Mak and Joseph (2017);
Support points of Mak and Joseph (2018) based on energy distance;
 L_2 discrepancy of QMC (symmetric, centred, wrap-around...) based
on tensorised kernels related to variants of Brownian-motion covariance (Hickernell, 1998)

2/ Nested designs (extensible point sequences)

Choose n points sequentially among a finite set $\mathcal{X}_Q = \{\mathbf{s}_1, \dots, \mathbf{s}_Q\} \subset \mathcal{X}$
A measure ξ on $\mathcal{X}_Q \Leftrightarrow$ a vector $\omega \in \mathbb{R}^Q$

2/ Nested designs (extensible point sequences)

Choose n points sequentially among a finite set $\mathcal{X}_Q = \{\mathbf{s}_1, \dots, \mathbf{s}_Q\} \subset \mathcal{X}$
 A measure ξ on $\mathcal{X}_Q \Leftrightarrow$ a vector $\omega \in \mathbb{R}^Q$

Vertex Direction (=“kernel herding” in machine learning):

→ minimise $\mathcal{E}_K(\xi - \mu)$, quadratic in ω ,

using Frank-Wolfe conditional gradient (Wynn’s algorithm):

$$\xi^{(n+1)} = \frac{n}{n+1} \xi^{(n)} + \frac{1}{n+1} \delta_{\mathbf{x}_{n+1}}$$

with $\mathbf{x}_{n+1} \in \operatorname{Arg} \min_{\mathbf{x} \in \mathcal{X}_Q} \left[\sum_{i=1}^n w_i^{(n)} K(\mathbf{x}, \mathbf{x}_i) - P_\mu(\mathbf{x}) \right]$
 (well defined also for singular kernels, e.g., Riesz kernels)

2/ Nested designs (extensible point sequences)

Choose n points sequentially among a finite set $\mathcal{X}_Q = \{\mathbf{s}_1, \dots, \mathbf{s}_Q\} \subset \mathcal{X}$
 A measure ξ on $\mathcal{X}_Q \Leftrightarrow$ a vector $\omega \in \mathbb{R}^Q$

Vertex Direction (=“kernel herding” in machine learning):

→ minimise $\mathcal{E}_K(\xi - \mu)$, quadratic in ω ,

using Frank-Wolfe conditional gradient (Wynn’s algorithm):

$$\xi^{(n+1)} = \frac{n}{n+1} \xi^{(n)} + \frac{1}{n+1} \delta_{\mathbf{x}_{n+1}}$$

with $\mathbf{x}_{n+1} \in \operatorname{Arg} \min_{\mathbf{x} \in \mathcal{X}_Q} \left[\sum_{i=1}^n w_i^{(n)} K(\mathbf{x}, \mathbf{x}_i) - P_\mu(\mathbf{x}) \right]$

(well defined also for singular kernels, e.g., Riesz kernels)

- $\xi^{(1)} = \delta_{\mathbf{x}_1} \rightarrow$ all weights $w_i^{(n)}$ in $\xi^{(n)}$ equal $1/n$
- Rate of decrease of $\mathcal{E}_K(\xi^{(n)} - \mu) = \mathcal{O}(1/n)$

Cost $\mathcal{O}(Q)$ at each iteration ($\mathcal{O}(nQ)$ for n iterations)

→ \mathbf{X}_n^{VD}

Minimum-Norm-point algorithm of (Bach et al., 2012):

replace $\xi^{(n)}$ (uniform on its support) by $\hat{\xi}^{(n)}$ having
the same support but optimal weights, positive with sum = 1

► We use \hat{w}_n from Bayesian quadrature

(extra comput. cost $\rightarrow \mathcal{O}(n^2 Q)$ for n iterations (P., 2021)) $\rightarrow \mathbf{X}_n^{MN}$

Minimum-Norm-point algorithm of (Bach et al., 2012):

replace $\xi^{(n)}$ (uniform on its support) by $\hat{\xi}^{(n)}$ having the same support but optimal weights, positive with sum = 1

► We use $\hat{\mathbf{w}}_n$ from Bayesian quadrature

(extra comput. cost $\rightarrow \mathcal{O}(n^2 Q)$ for n iterations (P., 2021)) $\rightarrow \mathbf{X}_n^{MN}$

Comparison with:

- (scrambled) Sobol' sequence $\rightarrow \mathbf{X}_n^{sS}$
- Extensible Lattice sequence \mathbf{X}_n^{EL} : $\mathbf{x}_k = \{k\mathbf{g}\}$ (\mathbf{g} has irrational components independent over the rationals, $\{\cdot\}$ = fractional part)
We use $\mathbf{g} = (1/\varphi_d, 1/\varphi_d^2, \dots, 1/\varphi_d^d)^\top$, with φ_d the unique positive root of $x^{d+1} = x + 1$, see <http://extremelearning.com.au/> $\rightarrow \mathbf{X}_n^{EL}$
(but not competitive for large d)

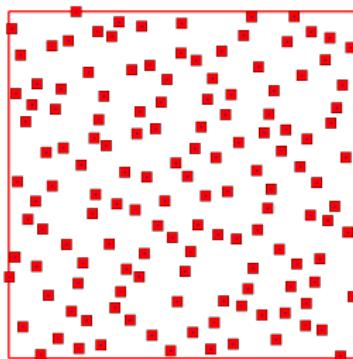
Example: $d = 2$, $\mathcal{X} = [0, 1]^2$, \mathcal{X}_Q = regular grid 64×64

$n_{\max} = 140$, $\mathbf{X}_1 = (0.5, 0.5)$, $\xi_{(1)} = \delta_{\mathbf{X}_1}$

K = tensor product of Matérn 3/2:

$$K_{3/2, \theta}(x, x') = (1 + \sqrt{3}\theta|x - x'|) \exp(-\sqrt{3}\theta|x - x'|), \theta = 10 \quad (n_{\max}^{1/d} \simeq 11.8)$$

scrambled Sobol' LDS \mathbf{X}_{140}^{ss}



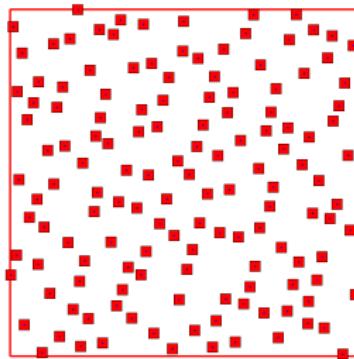
Example: $d = 2$, $\mathcal{X} = [0, 1]^2$, \mathcal{X}_Q = regular grid 64×64

$n_{\max} = 140$, $\mathbf{X}_1 = (0.5, 0.5)$, $\xi_{(1)} = \delta_{\mathbf{X}_1}$

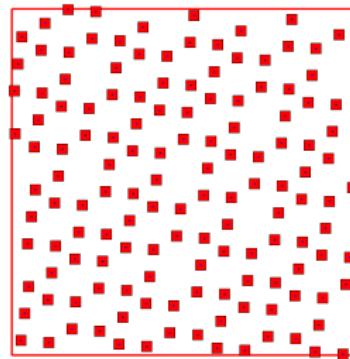
K = tensor product of Matérn 3/2:

$$K_{3/2, \theta}(x, x') = (1 + \sqrt{3}\theta|x - x'|) \exp(-\sqrt{3}\theta|x - x'|), \quad \theta = 10 \quad (n_{\max}^{1/d} \simeq 11.8)$$

scrambled Sobol' LDS \mathbf{X}_{140}^{sS}



Extensible Lattice sequence \mathbf{X}_{140}^{EL}



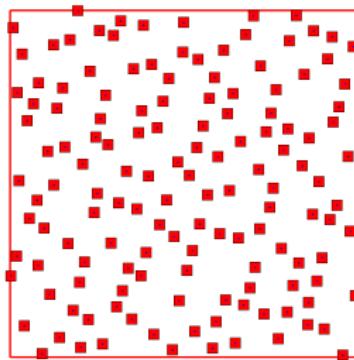
Example: $d = 2$, $\mathcal{X} = [0, 1]^2$, \mathcal{X}_Q = regular grid 64×64

$n_{\max} = 140$, $\mathbf{X}_1 = (0.5, 0.5)$, $\xi_{(1)} = \delta_{\mathbf{X}_1}$

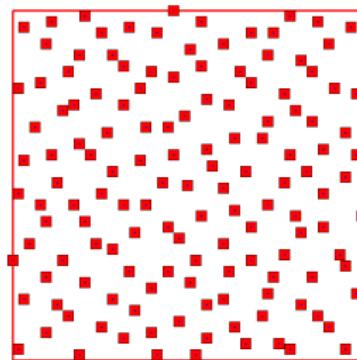
K = tensor product of Matérn 3/2:

$$K_{3/2, \theta}(x, x') = (1 + \sqrt{3}\theta|x - x'|) \exp(-\sqrt{3}\theta|x - x'|), \quad \theta = 10 \quad (n_{\max}^{1/d} \simeq 11.8)$$

scrambled Sobol' LDS \mathbf{X}_{140}^{sS}



Vertex Direction \mathbf{X}_{140}^{VD}



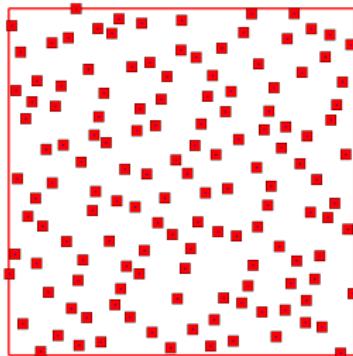
Example: $d = 2$, $\mathcal{X} = [0, 1]^2$, \mathcal{X}_Q = regular grid 64×64

$n_{\max} = 140$, $\mathbf{X}_1 = (0.5, 0.5)$, $\xi_{(1)} = \delta_{\mathbf{X}_1}$

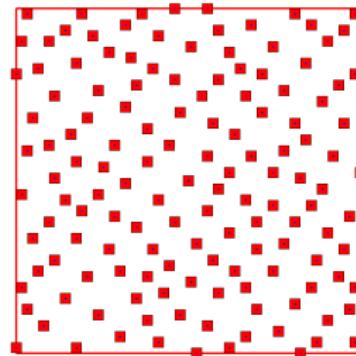
K = tensor product of Matérn 3/2:

$$K_{3/2, \theta}(x, x') = (1 + \sqrt{3}\theta|x - x'|) \exp(-\sqrt{3}\theta|x - x'|), \quad \theta = 10 \quad (n_{\max}^{1/d} \simeq 11.8)$$

scrambled Sobol' LDS \mathbf{X}_{140}^{ss}



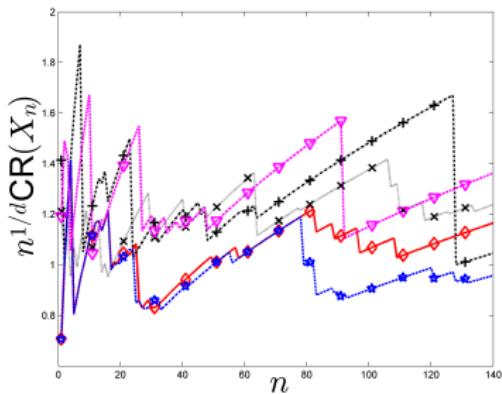
Minimum-Norm \mathbf{X}_{140}^{MN}



Space-filling performance

$$n^{1/d} \text{CR}(\mathbf{X}_n)$$

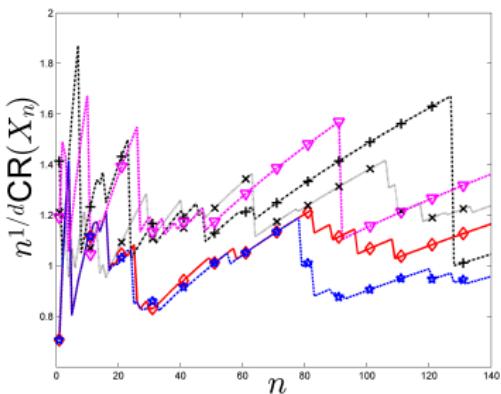
(smaller is better)



\mathbf{X}_n^S , \mathbf{X}_n^{sS} , \mathbf{X}_n^{EL} , \mathbf{X}_n^{VD} , \mathbf{X}_n^{MN}

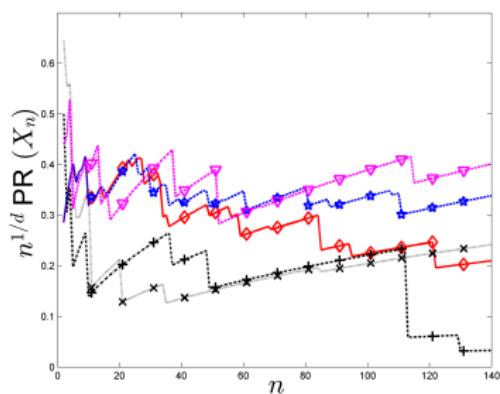
Space-filling performance

$n^{1/d} \text{CR}(\mathbf{X}_n)$
(smaller is better)



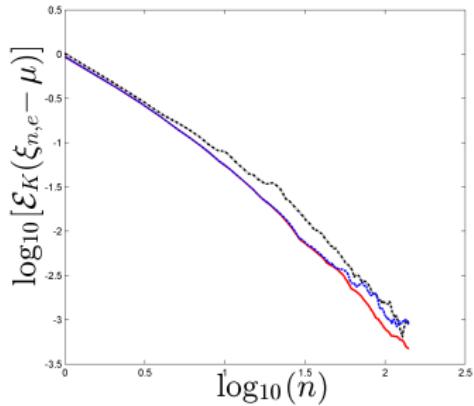
\mathbf{X}_n^S , \mathbf{X}_n^{sS} , \mathbf{X}_n^{EL} , \mathbf{X}_n^{VD} , \mathbf{X}_n^{MN}

$n^{1/d} \text{PR}(\mathbf{X}_n)$
(larger is better)

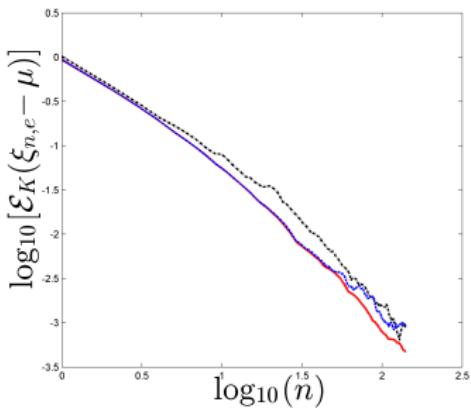


\mathbf{X}_n^S , \mathbf{X}_n^{sS} , \mathbf{X}_n^{EL} , \mathbf{X}_n^{VD} , \mathbf{X}_n^{MN}

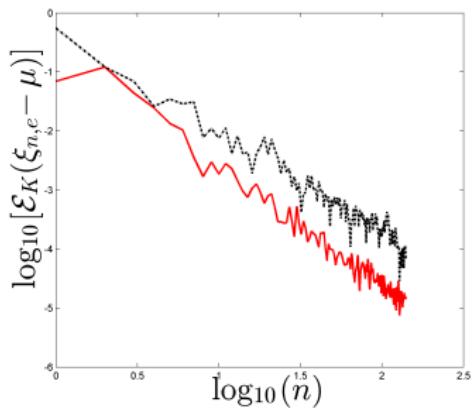
Decrease of $\mathcal{E}_{K_{3/2,10}}(\xi_n - \mu)$
 $(\mathbf{X}_n^S, \mathbf{X}_n^{VD}, \mathbf{X}_n^{MN})$



Decrease of $\mathcal{E}_{K_{3/2,10}}(\xi_n - \mu)$
 $(\mathbf{X}_n^S, \mathbf{X}_n^{VD}, \mathbf{X}_n^{MN})$



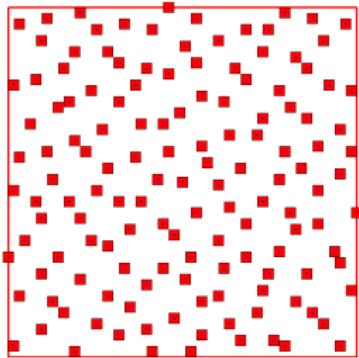
Decrease of $\mathcal{E}_{K_{3/2,1}}(\xi_n - \mu)$
 $(\mathbf{X}_n^S, \mathbf{X}_n^{VD})$



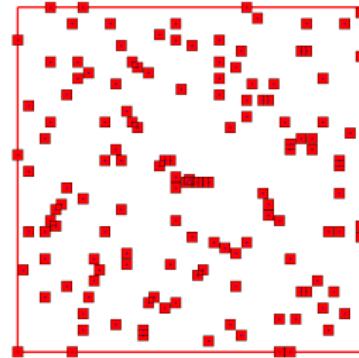
For a large enough correlation length the rate of decrease can reach $\mathcal{O}(1/n^2)$ (no proof available)

How meaningful is the rate of convergence?

“kernel herding” with $K_{3/2, 10}$
 \mathbf{X}_{140}^{VD} (rate $\lesssim 1/k$)



“kernel herding” with $K_{3/2, 1}$
 \mathbf{X}_{140}^{VD} (rate $\lesssim 1/k^2$)



A faster rate does not mean that points are more evenly distributed!

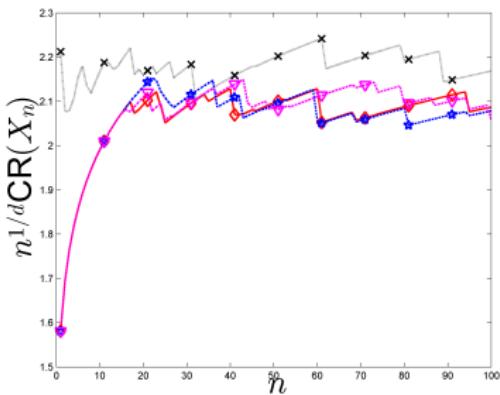
Example: $d = 10$, $\mathcal{X}_Q = 2^{12}$ points of scrambled Sobol' in $\mathcal{X} = [0, 1]^2$

$n_{\max} = 140$, $\mathbf{X}_1 = \mathbf{1}_d/2$, $\xi_{(1)} = \delta_{\mathbf{x}_1}$

K = tensor product of Matérn 3/2:

$$K_{3/2, \theta}(x, x') = (1 + \sqrt{3}\theta|x - x'|) \exp(-\sqrt{3}\theta|x - x'|), \quad \theta = n_{\max}^{1/d}$$

\mathbf{X}_n^S , \mathbf{X}_n^{VD} , \mathbf{X}_n^{MN} , $\mathbf{X}_n^{MN-\log}$



$$[K_{\log}(x_i, x'_i) = -\log |x_i - x'_i|]$$

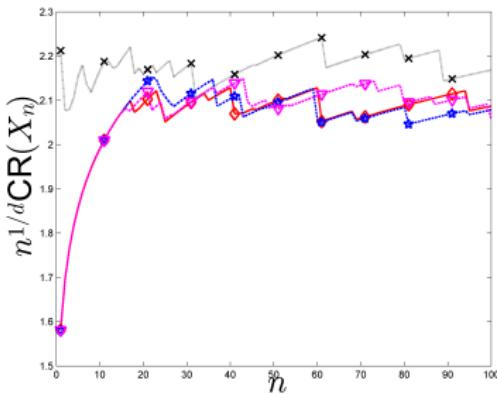
Example: $d = 10$, $\mathcal{X}_Q = 2^{12}$ points of scrambled Sobol' in $\mathcal{X} = [0, 1]^d$

$n_{\max} = 140$, $\mathbf{X}_1 = \mathbf{1}_d/2$, $\xi_{(1)} = \delta_{\mathbf{X}_1}$

K = tensor product of Matérn 3/2:

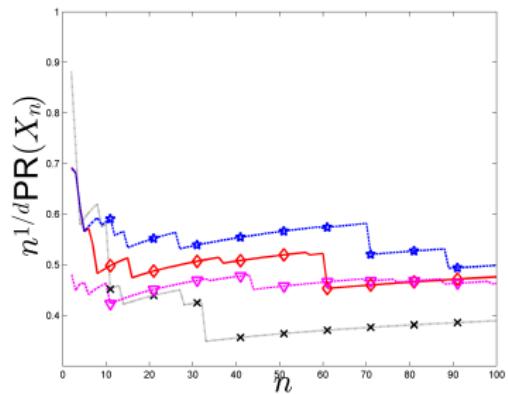
$$K_{3/2, \theta}(x, x') = (1 + \sqrt{3}\theta|x - x'|) \exp(-\sqrt{3}\theta|x - x'|), \quad \theta = n_{\max}^{1/d}$$

\mathbf{X}_n^S , \mathbf{X}_n^{VD} , \mathbf{X}_n^{MN} , $\mathbf{X}_n^{MN-\log}$



$$[K_{\log}(x_i, x'_i) = -\log |x_i - x'_i|]$$

\mathbf{X}_n^S , \mathbf{X}_n^{VD} , \mathbf{X}_n^{MN} , $\mathbf{X}_n^{MN-\log}$



3/ n fixed: minimise $s_n^2 = \mathcal{E}_K(\hat{\xi}_n - \mu) = \frac{1}{\mathbf{1}_n^\top \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$ w.r.t. \mathbf{X}_n

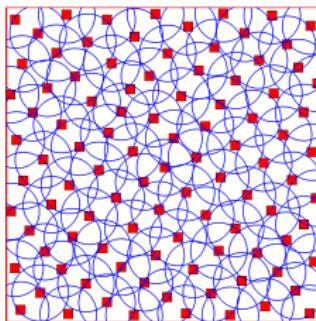
a) First (initialisation), use a greedy algorithm, stop at n

$$\rightarrow \mathbf{X}_n^{VD}$$

b) Then, maximise $\mathbf{1}_n^\top \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n$ viewed as a function of \mathbf{X}_n

$$\rightarrow \mathbf{X}_n^{VDopt}$$

a) Initialisation at \mathbf{X}_{100}^{VD} generated by Vertex-Direction method
with $CR(\mathbf{X}_{100}^{VD}) \simeq 0.0925$, $2PR(\mathbf{X}_{100}^{VD}) \simeq 0.0523$



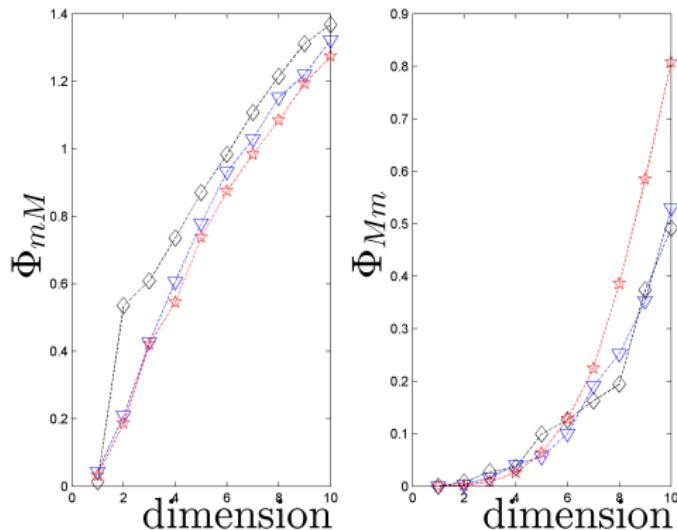
\rightarrow b) $CR(\mathbf{X}_{100}^{VDopt}) \simeq 0.0897$, $2PR(\mathbf{X}_{100}^{VDopt}) \simeq 0.0738$

Example: $d = 10$, $\mathcal{X} = [0, 1]^{10}$, $n_{\max} = 100$

$$K = \bigotimes K_{3/2, \theta}(x, x'), \theta = \lfloor n_{\max}^{1/d} \rfloor = 1$$

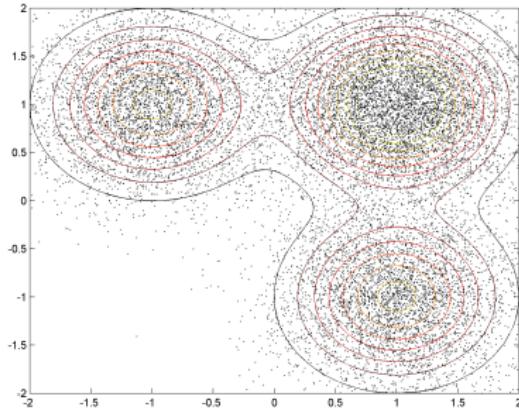
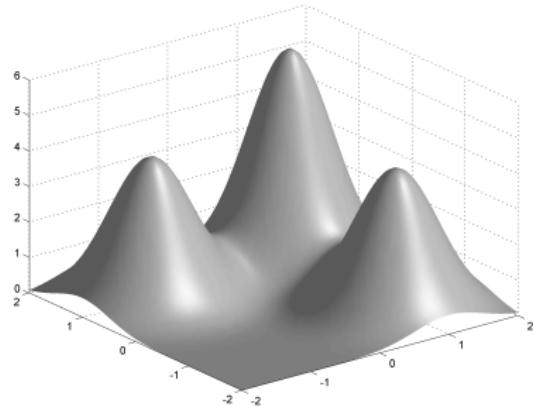
a) VD \mathbf{X}_n^{VD} , $n = 1, \dots, 100$, then b) \mathbf{X}_{100}^{VDopt}

$\max \text{CR} = \max \Phi_{mM}$ and $2 \min \text{PR} = \min \Phi_{Mm}$ on projections
 $\diamond \mathbf{X}_{100}^S$, \triangledown : \mathbf{X}_{100}^{VD} , \star : \mathbf{X}_{100}^{VDopt}



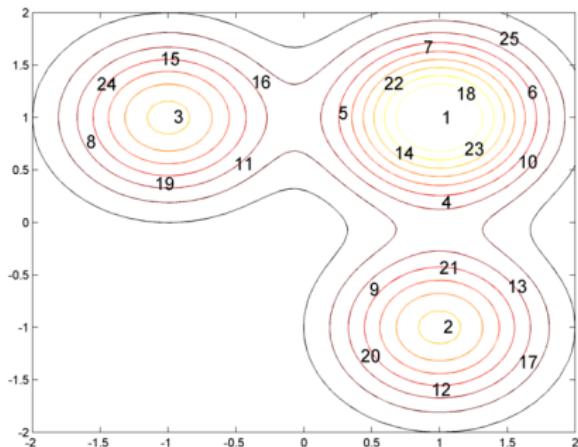
Kernel herding is not restricted to μ being uniform:

Example: Gaussian mixture $\mu = \sum_{j=1}^3 \beta_j \mu_N(\mathbf{a}_j, \sigma_j)$, $Q = 2^{14} = 16\,384$
 (for $K_\theta(\mathbf{x}, \mathbf{x}') = \exp -(\theta \|\mathbf{x} - \mathbf{x}'\|^2)$, we know $P_\mu(\cdot)$ and $\mathcal{E}_K(\mu)$)

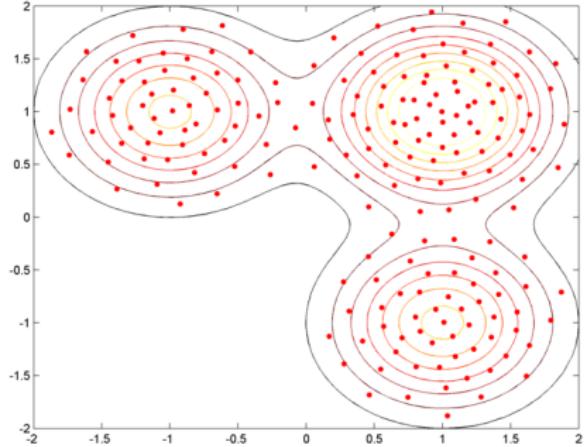


Kernel herding is not restricted to μ being uniform:

Example: Gaussian mixture $\mu = \sum_{j=1}^3 \beta_j \mu_N(\mathbf{a}_j, \sigma_j)$, $Q = 2^{14} = 16\,384$
 (for $K_\theta(\mathbf{x}, \mathbf{x}') = \exp -(\theta \|\mathbf{x} - \mathbf{x}'\|^2)$, we know $P_\mu(\cdot)$ and $\mathcal{E}_K(\mu)$)



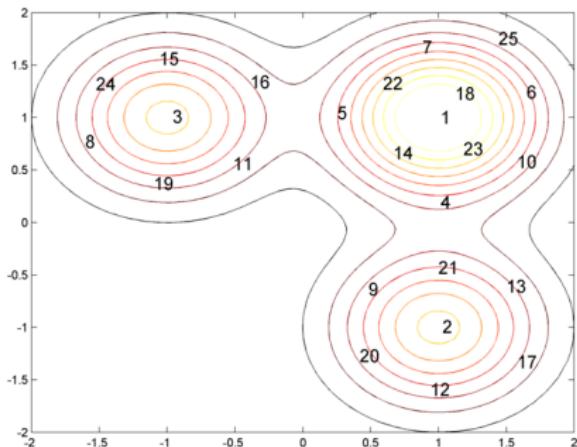
$n = 25$



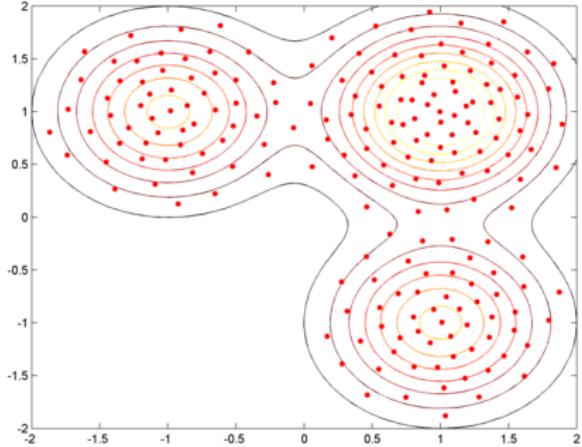
$n = 200$

Kernel herding is not restricted to μ being uniform:

Example: Gaussian mixture $\mu = \sum_{j=1}^3 \beta_j \mu_N(\mathbf{a}_j, \sigma_j)$, $Q = 2^{14} = 16\,384$
 (for $K_\theta(\mathbf{x}, \mathbf{x}') = \exp -(\theta \|\mathbf{x} - \mathbf{x}'\|^2)$, we know $P_\mu(\cdot)$ and $\mathcal{E}_K(\mu)$)



$n = 25$



$n = 200$

- Comparison between kernel herding, greedy MMD minimisation and Sequential Bayesian Quadrature (P., 2021)
- Extension to Stein discrepancy (Teymur et al., 2020) (K'_μ such that $P_\mu(\cdot) \equiv 0$ and $\mathcal{E}_{K'_\mu}(\mu) = 0$ without knowing the normalising constant in μ)
- Singular kernels (via completely monotone functions) (P. & Zhigljavsky, 2021)

4 Conclusions

- Maximum Mean Discrepancy connects Bayesian quadrature, potential theory, design for parameter estimation, space filling design
 - For many kernels K : metric for the weak cv. of prob. measures, easier to compute than optimal transport (= Wasserstein)
 - (design for) Bayesian integration \Leftrightarrow (design for) parameter estimation in a model with modified covariance

4 Conclusions

- Maximum Mean Discrepancy connects Bayesian quadrature, potential theory, design for parameter estimation, space filling design
 - For many kernels K : metric for the weak cv. of prob. measures, easier to compute than optimal transport (= Wasserstein)
 - (design for) Bayesian integration \Leftrightarrow (design for) parameter estimation in a model with modified covariance
 - → minimising $s_n^2 = \frac{1}{\mathbf{1}_n^\top \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$ as an alternative to minimising
$$\text{IMSPE}(\mathbf{X}_n) = \int_{\mathcal{X}} \rho_n^2(\mathbf{x}) d\mu(\mathbf{x})$$

4 Conclusions

- Maximum Mean Discrepancy connects Bayesian quadrature, potential theory, design for parameter estimation, space filling design
 - For many kernels K : metric for the weak cv. of prob. measures, easier to compute than optimal transport (= Wasserstein)
 - (design for) Bayesian integration \Leftrightarrow (design for) parameter estimation in a model with modified covariance
 - → minimising $s_n^2 = \frac{1}{\mathbf{1}_n^\top \tilde{\mathbf{K}}_n^{-1} \mathbf{1}_n}$ as an alternative to minimising

$$\text{IMSPE}(\mathbf{X}_n) = \int_{\mathcal{X}} \rho_n^2(\mathbf{x}) d\mu(\mathbf{x})$$
- Optimisation algorithm for n fixed:
 Any nonlinear programming method (unconstrained) can be used
 (a Maj.-Min. approach is used in (Mak and Joseph, 2017, 2018))

A few more general questions and open issues:

- Space-filling performance? (discrepancy $\not\Rightarrow$ dispersion)
Efficiency bounds for CR and PR?
- Which K ? (wide choice, but does not seem to be crucial)

A few more general questions and open issues:

- Space-filling performance? (discrepancy $\not\Rightarrow$ dispersion)
Efficiency bounds for CR and PR?
- Which K ? (wide choice, but does not seem to be crucial)
- Steepest descent or Greedy optimisation?
→ To construct a sequence of embedded designs with small $CR(\mathbf{X}_n)$,
is kernel herding (= Frank-Wolfe applied to MMD minimisation)
preferable to methods exploiting submodularity, such as (**Nogales
Gómez et al., 2021**)?

A few more general questions and open issues:

- Space-filling performance? (discrepancy $\not\Rightarrow$ dispersion)
Efficiency bounds for CR and PR?
- Which K ? (wide choice, but does not seem to be crucial)
- Steepest descent or Greedy optimisation?
→ To construct a sequence of embedded designs with small $CR(\mathbf{X}_n)$,
is kernel herding (= Frank-Wolfe applied to MMD minimisation)
preferable to methods exploiting submodularity, such as (**Nogales
Gómez et al., 2021**)?

Thank you for your attention !

References I

- Bach, F., Lacoste-Julien, S., Obozinski, G., 2012. On the equivalence between herding and conditional gradient algorithms. In: Proc. 29th Annual International Conference on Machine Learning. pp. 1355–1362.
- Damelin, S., Hickernell, F., Ragozin, D., Zeng, X., 2010. On energy, discrepancy and group invariant measures on measurable subsets of Euclidean space. *J. Fourier Anal. Appl.* 16, 813–839.
- Grenander, U., 1950. Stochastic processes and statistical inference. *Arkiv för Matematik* 1 (3), 195–277.
- Hájek, J., 1956. Linear estimation of the mean value of a stationary random process with convex correlation function. *Czechoslovak Mathematical Journal* 6 (81), 94–117.
- Hickernell, F., 1998. A generalized discrepancy and quadrature error bound. *Mathematics of Computation* 67 (221), 299–322.
- Mak, S., Joseph, V., 2017. Projected support points, with application to optimal MCMC reduction. arXiv preprint arXiv:1708.06897.
- Mak, S., Joseph, V., 2018. Support points. *Annals of Statistics* 46 (6A), 2562–2592.
- Näther, W., 1985. Effective Observation of Random Fields. Vol. 72. Teubner-Texte zur Mathematik, Leipzig.

References II

- Nogales Gómez, A., Pronzato, L., Rendas, M.-J., 2021. Incremental space-filling design based on coverings and spacings: improving upon low discrepancy sequences. *Journal of Statistical Theory and Practice* (to appear).
- O'Hagan, A., 1991. Bayes-Hermite quadrature. *Journal of Statistical Planning and Inference* 29 (3), 245–260.
- Pronzato, L., 2021. Performance analysis of greedy algorithms for minimising a maximum mean discrepancy. hal-03114891, arXiv:2101.07564.
- Pronzato, L., Zhigljavsky, A., 2019. Measures minimizing regularized dispersion. *Journal of Scientific Computing* 78 (3), 1550–1570.
- Pronzato, L., Zhigljavsky, A., 2020. Bayesian quadrature, energy minimization and space-filling design. *SIAM/ASA J. Uncertainty Quantification* 8 (3), 959–1011.
- Pronzato, L., Zhigljavsky, A., 2021. Minimum-energy measures for singular kernels. *Journal of Computational and Applied Mathematics* 382, (16 pages, to appear) hal-02495643.
- Sejdinovic, S., Sriperumbudur, B., Gretton, A., Fukumizu, K., 2013. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *The Annals of Statistics* 41 (5), 2263–2291.
- Sriperumbudur, B., Gretton, A., Fukumizu, K., Schölkopf, B., Lanckriet, G., 2010. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research* 11, 1517–1561.

References III

- Szabó, Z., Sriperumbudur, B., 2018. Characteristic and universal tensor product kernels. *Journal of Machine Learning Research* 18, 1–29.
- Székely, G., Rizzo, M., 2013. Energy statistics: A class of statistics based on distances. *Journal of Statistical Planning and Inference* 143 (8), 1249–1272.
- Teymur, O., Gorham, J., Riabiz, M., Oates, C., 2020. Optimal quantisation of probability measures using maximum mean discrepancy. arXiv preprint arXiv:2010.07064v1.