

# kernel Thinning

2/26/2024

Geometry  
@ UW  
mmp

$\subseteq$  finding core-sets

$\neq$  data thinning !!!  $\rightarrow$  sample splitting

Problem: approximate distribution  $\mu$  on  $\mathcal{X}$  by a sample  
—— " —— a large sample  $\mu_N$  by smaller sample

$y_{1:N} \sim \text{iid } \mu$  ————  
empirical distribution

$Q = \{x_{1:n} \in \mathcal{X}\} \rightarrow \mu_Q = \text{approximator}$

so that  $\mu_Q \approx \mu_N \approx \mu$

# Plan

- Kernel Thinning problem
- Ideas for solving it  $\leftarrow$ 

Herding

K. Thinning

weighted Quad.
- (Basic) RKHS facts  $\leftarrow$ 

geometric

Important/useful
- The methods : H, KT, wQ (if time)

Problem: approximate distribution  $\mu$  on  $\mathcal{X}$  by a sample  
 ————— a large sample  $\mu_N$  by smaller sample

$y_{1:N} \sim \text{iid } \mu$  —————  
 empirical distribution

$Q = \{x_{1:n} \in \mathcal{X}\} \rightarrow \mu_Q = \text{approximator}$

so that  $\mu_Q \approx \mu_N \approx \mu$  **More precisely**  $\mu f \approx \mu_N f \approx \mu_Q f$  for all  $f \in \mathcal{H}_k$

$$\mu f \equiv E_{\mu}[f(x)]$$

uniformly  
 RKHS

**More precisely**

•  $MMD_k(\nu, \mu) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} |\nu f - \mu f|$

Maximum Mean Discrepancy

wanted  $Q = \{x_{1:n}\}$   
 so that

$$MMD_k(\mu_N, \mu_Q) < \varepsilon_n$$

$$MMD_k(\mu, \mu_Q) < \varepsilon_n$$

## Why RKHS?

- compact specification of  $\{f's\}$  - no need to enumerate!

- nice, elegant, well understood math: • complete

- for ex:

polynomials  $k = (1 + x^T x)^d$

kernel reg, classification

$k = \text{Gaussian, Matérn, ...}$

string, tree kernels

• "Euclidean" (intuitive)

• tractable computations

• non-parametric

• general

.....

Wahba

## Why core-sets / thinning?

want to do ML-with large data

- non-parametric

-  $f$  not known yet

⇒ save computation time  
tractability

Kernel  $k(\cdot)$  must match  
problem!!

## Why expectations $E_\mu[f]$

many functionals (criteria of quality) are expectations



## Problem (rephrased)

- given •  $y_{1:n} \equiv \mu_N$  large sample  
• kernel  $k \equiv RKHS \mathcal{H}_k$

wanted •  $x_{1:n} \equiv \mu_Q$  small sample

so that  $MMD_k(\mu_N, \mu_Q) \leq \varepsilon_n$

What is a reasonable  $\varepsilon_n$ ?

- $\mu_N \sim iid \mu$   
 $\mu_N \approx \mu \Rightarrow$  want same rate for  $\mu_Q \approx \mu_N$
- $E_\mu[x] = \mu \int_{\mathcal{H}_k} x$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i \rightarrow \text{Var } \hat{\mu} \propto \frac{1}{N}$$
$$\text{std } \hat{\mu} \propto \frac{1}{\sqrt{N}} = N^{-1/2}$$

$MMD_k$  rate no faster than  $\hat{\mu}$  rate

benchmark rate

## Problem (rephrased)

given •  $y_{1:n} \equiv \mu_N$  large sample

• kernel  $k \equiv \text{RKHS } \mathcal{H}_k$

wanted •  $x_{1:n} \equiv \mu_Q$  small sample

so that  $\text{MMD}_k(\mu_N, \mu_Q) \leq \varepsilon_n$

## Solutions

• grid on  $\mathcal{X} \leftarrow \text{optimal } N^{-2r}$

• MC = iid sampling  $N^{-1/2}$

• MCMC  $\sim N^{-\alpha}$   $\alpha < \frac{1}{2}$

• Quasi MC  $\sim N^{-\alpha'}$   $\alpha' > \alpha$

• Herding

• DPP (Determinantal Point P)

• Thinning

• Weighted Quadrature

dependent (not iid)  
neg. correl.  
↓

## Problem (rephrased)

given •  $y_{1:N} \equiv \mu_N$  large sample

• kernel  $k \equiv \text{RKHS } \mathcal{H}_k$

wanted •  $x_{1:n} \equiv \mu_Q$  small sample

so that  $\text{MMD}_k(\mu_N, \mu_Q) \leq \varepsilon_n$

→ all select subsample  $Q$  of  $y_{1:N}$

→  $n = \sqrt{N}$

rate  $\sim n^{-1} = N^{-1/2}$  but with  $\sqrt{N}$  points

sequential  
greedy

recursive vector  
balancing  
(paired comparisons)

+ square root kernel

weighted sample [Hayakawa, Oberhauser, Lyons<sub>21</sub>]

RKHS approximation\*

+ cubature algorithm (Caratheodory)\*

## Solutions

• grid on  $\mathcal{X} \leftarrow \text{optimal } N^{-2r}$

• MC = iid sampling  $N^{-1/2}$

• MCMC  $\sim N^{-\alpha}$   $\alpha < \frac{1}{2}$

• Quasi MC  $\sim N^{-\alpha'}$   $\alpha' > \alpha$

• Herding [Yutian Chen, Welling, Smola'12]

• DPP (Determinantal [Belhadji ...] Point P)

Thinning [Dwivedi, Mackey'21]

Weighted Quadrature

What makes thinning pb. hard? con. set  $\varepsilon_n \sim \underline{\underline{C n^{-1}}}$

rate of decay of  
kernel,  $\mu$

For ex:

$\mu$  compact support

$\mu$  subgaussian

$\mu$  heavy tailed

- similar for  $k$ : compact support, ...

Weighted Quadrature:  $(\sigma_{1,2}, \dots, e_{1,2}, \dots(x)) = e\text{-values, } e\text{-functions of operator}$

$$(kf)(x) = \int_{\mathcal{X}} f(x') k(x', x) d\mu(x') \quad \leftarrow \text{typically convolution op.}$$

$\Rightarrow$  err depends on  $\sigma_n, k_{n+1} = \sum_{j=n+1}^{\infty} \sigma_j^2$  (residual spectrum)



## kernel thinning

kernel decay radius

$$R_{k,n} = \inf_{\substack{\dots \\ \dots}} r \sup_{\|x-y\|_2 \geq r} |k(x,y)| \leq \frac{1}{n} \|k\|_\infty$$

$$\phi_k(r) = \sup_x \int_{\|y\|_2 \geq r} k^2(x, x-y) dy \leftarrow \text{tail weight of } k^2(\cdot) \text{ at } r$$

Ex:  $k(x,y) = N(x-y; 0, \sigma^2 I) \Rightarrow \phi_k(r) = \int_{\|y\|_2 \geq r} e^{-2 \frac{\|x-(x-y)\|^2}{2\sigma^2}} dy = \text{tail probs for } N(0, \frac{\sigma^2}{2} I) \text{ at } r \cdot \sqrt{2}$

Gaussian kernel

for  $\sigma=1$ :  $\|k\|_\infty = \frac{1}{(2\pi)^{d/2}}$

$$\|x-y\| \geq r \Rightarrow k(x,y) \leq \frac{1}{(2\pi)^{d/2}} e^{-\frac{r^2}{2}}$$

$$\Rightarrow e^{-\frac{r^2}{2}} = \frac{1}{n} \Rightarrow R_{k,n} = \sqrt{2 \ln n}$$

It helps if these are known analytically  $\leftarrow$  otherwise

- kernel tail probabilities  $\tau_k(r)$  <sup>or exact, tractably</sup> computable they must be approximated, majorized, ...

decay radii  $\underline{R}_{k,n}$

(+ decay rates of  $\mu$ ) (Thinning)

- $\bar{\varphi} = \mathbb{E}_{\mu}[\varphi(x)]$  or  $\langle z, \bar{\varphi} \rangle_k \in \mathbb{R}$   
 $\cap \mathcal{H}_k$   $\cap \mathcal{H}_k$

$\max_x \langle z, \varphi(x) \rangle$  (Herdling)  
 $\cap \mathcal{H}_k$

- (Bochner) Fourier repr of  $k(\cdot, \cdot)$   $\nearrow$  Thinning  $k_{1/2}$

Spectrum of operator  $k$   $(\sigma_{1,2,\dots,m}, e_{1,2,\dots,m}(x))$

$\rightarrow$  weighted Quadrature

Bochner

Carathéodory

Universal kernel

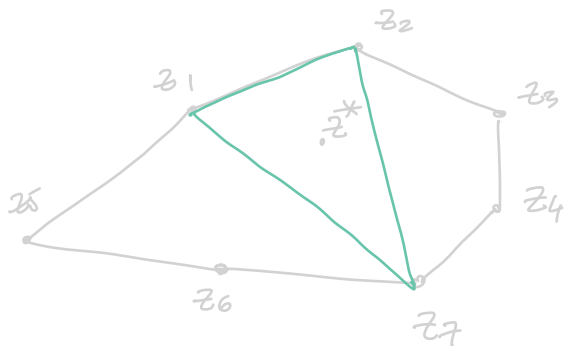
Weight optimization

Computing  $MMD_L$

Bochner

# Caratheodory

Thm  $z_{1:N} \in \mathbb{R}^{d-1}$ ,  $z^* \in \text{conv hull}(z_{1:N}) \implies$   
 $z^* \in \text{conv. hull}(z'_1, \dots, z'_d)$  where  
 $z'_{1:d} \subset \{z_{1:N}\}$



$$\implies z^* \in \text{conv}(z_1, z_2, z_4)$$

$d=2 \implies \underbrace{d+1}_{3} \text{ points sufficient}$

For us:

functions  $\varphi_{1:n-1}$   
 points  $y_{1:N}$  }  $\implies$  define  $z_i = \begin{bmatrix} \varphi_1(y_i) \\ \vdots \\ \varphi_{n-1}(y_i) \end{bmatrix} \in \mathbb{R}^{n-1}$

$$\bar{z} = \begin{bmatrix} \vdots \\ \frac{1}{N} \sum_{i=1}^N \varphi_j(y_i) \\ \vdots \end{bmatrix} \leftarrow \bar{z}_j$$

$\implies \left[ \begin{array}{l} \exists x_{1:n} \in \{y_{1:N}\} \text{ so that} \\ \underline{w_{1:n} \geq 0}, \underline{\sum w_i = 1} \end{array} \right] \underline{\bar{z}_j = \sum_{i=1} w_i \varphi_j(x_i) \text{ for } j=1:n-1}$

## Universal kernel

If  $k$  universal kernel  $\Rightarrow |\mu f - \nu f| \leq \|f\|_{\mathcal{H}} \|\mu - \nu\|_{\mathcal{H}}$

[Koksma-Hlawka inequality]

# Computing $\text{MMD}_k$

from Weighted Quadrature (14)

$$\text{MMD}_k(\mu_Q, \mu)^2 = \|\mu_Q \varphi - \mu \varphi\|_{\mathcal{H}}^2$$

$\varphi: \mathcal{X} \rightarrow \mathcal{H}$  feature map

$$= w^T k(X, X) w - 2 E_y [w^T k(X, y)] + E_{y, y'} [k(y, y')]$$

weights  $w_{1:n}$   
or  $1/n$

$\mu = \mu_N$   
or need to know analytically

## Weight optimization

Quadratic in  $w \in \mathbb{R}^n \Rightarrow \min_w \text{MMD}^2 \text{ s.t. } w \geq 0 \text{ with } x_{1:n}, y_{1:n} \text{ given}$   
 $\sum w_i = 1$

- optimizes any coe-set (e.g. from QMC, Thinning, ...)