# Kernel Thinning

$\subseteq$ finding core-sets

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

<u>Problem</u> : approximate distribution $\mu$ on $\mathcal{X}$ by a sample

—"— a large sample $\mu_N$ by smaller sample

$y_{1:N} \sim iid \; \mu$

empirical distribution

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

<u>so that</u> $\mu_Q \approx \mu_N \approx \mu$

# Plan

- Kernel Thinning problem
- Ideas for solving it $\left\{\begin{array}{l}\text{Herding}\\ \text{K. Thinning}\\ \text{weighted Quad.}\end{array}\right.$
- (Basic) RKHS facts $\longleftarrow$ 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

$y_{1:N} \sim$ iid $\mu$

empirical distribution

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

**More precisely**

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

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

RKHS

uniformly

More precisely

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

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
  - normed
  - "Euclidean" (intuitive)
  - tractable computations
  - non-parametric
  - general
  - . . . .

*Wahba*

- for ex:

    polynomials $k = (1 + x^T x')^d$

    kernel regr, classification

       $k = $ Gaussian, Matérn, ...

  string, tree kernels

# Why core-sets / thinning?

   want to do ML with large data
  - non-parametric
  - $f$ not known yet

$\Rightarrow$ save computation time
    tractability

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

# Why expectations $E_\mu[f]$

   many functionals (criteria of quality) are expectations

$P[err]$
$E[err^2]$
quantile

## Problem (rephrased)

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

wanted
- $x_{1:n} \equiv \mu_Q$  <u>small sample</u>

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

<u>What is a reasonable $\varepsilon_n$?</u>

- $\mu_N \sim iid \; \mu$
  $\mu_N \approx \mu \Rightarrow$ want same rate for $\mu_Q \approx \mu_N$

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

- $E_\mu [x] = \mu \underset{\mathcal{H}_k}{\overset{\wedge}{x}}$

  $\hat{\mu} = \frac{1}{N} \sum\limits_{i=1}^{N} y_i \longrightarrow$ Var $\hat{\mu} \propto \frac{1}{N}$

  std $\hat{\mu} \propto \frac{1}{\sqrt{N}} = N^{-1/2}$  $\longleftarrow$ benchmark rate

# 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 $\mathrm{MMD}_k(\mu_N, \mu_Q) \leq \varepsilon_n$

# Solutions

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

• MC = iid sampling $N^{-\frac{1}{2}}$

• 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 $k \equiv$ 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

recursive vector
balancing
(paired comparisons)
+ square root kernel

# Solutions

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

• MC = iid sampling $N^{-1/2}$

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

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

sequential greedy →

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

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

Thinning [Dwivedi, Mackey '21]

Weighted Quadrature
weighted sample [Hayakawa, Oberhauser, Lyons '21]
RKHS approximation ⋆
+ cubature algorithm (Caratheodory) ⋆

What makes thinning pb. hard? $\varepsilon_n \sim \underline{c\, n^{-1}}$

core set

For ex:

$\mu$ compact support

$\wedge$

$\mu$ subgaussian

$\wedge$

$\mu$ heavy tailed

rate of decay of
kernel, $\mu$

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

weighted
Quadrature : $\left( \sigma_{1,2,\ldots}, \, e_{1,2,\ldots}(x) \right) = $ e-values, e-functions of operator

$$\left( k\, f \right)(x) = \int_\mathcal{X} f(x')\, k(x', x)\, d\mu(x')$$

typically
convolution
op.

$\Rightarrow$ err depends on $\sigma_n, \; r_{n+1} = \sum\limits_{j=n+1}^{\infty} \sigma_j$ (residual spectrum)

# Kernel thinning

kernel decay radius

- $R_{k,n} = \inf\ r$

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

- $\mathcal{T}_k(r) = \sup_{x} \int_{\|y\|_2 \geq r} k^2(x, x-y)\,dy \longleftarrow$ tail weight of $k^2()$ at $r$

  Ex: $k(x,y) = N\left(x-y;\ 0, \sigma^2 I\right) \Rightarrow \mathcal{T}_k(r) = \int_{\|y\| \geq r} e^{-2\frac{\|x-(x-y)\|^2}{2\sigma^2}}\,dy = $ tail pdas
  
  Gaussian kernel
  
  for $N\left(0, \frac{\sigma^2}{2} I\right)$ at $r \cdot \sqrt{2}$

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

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

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

It helps if these are known <u>analytically</u> ← otherwise

or <u>exact, tractably</u> they must be approximated,

- kernel tail probabilities $\overline{\tau}_k(r)$ <u>computable</u> majorized, ...

decay radii $R_{k,n}$

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

- $\overline{\varphi} = \mathbb{E}_\mu[\varphi(x)]$  or  $\langle z, \overline{\varphi} \rangle_k \in \mathbb{R}$

  $\in \mathcal{H}_k$  $\mathcal{H}_k$

  $\max_x \langle z, \varphi(x) \rangle$  (Herding)

  $\mathcal{H}_k$

  Thinning $k_{1/2}$

- (Bochner) Fourier repn of $k(,)$

  spectrum of operator $k$  $(v_{1,2,..m}, e_{1,2,..m}(x))$

  $\hookrightarrow$ weighted Quadrature

Bochner

Caratheodory

Universal kernel

Weight optimization

Computing $MMD_k$

Bochner

$\boxed{\text{Caratheodory}}$

$\underline{\text{Thm}}$  $z_{1:N} \in \mathbb{R}^{d-1}$, $\quad z^* \in \text{convex hull}(z_{1:N}) \implies$

$\boxed{\text{Caratheodory}}$

$\qquad z^* \in \text{conv. hull}(z_1', \cdots z_d')$  where

$\qquad\qquad z_{1:d}' \subset \{z_{1:N}\}$



$\longrightarrow z^* \in \text{conv}(z_1, z_2, z_7)$

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

$\underline{\text{For us:}}$

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

points $y_{1:N}$

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

$\Rightarrow \begin{bmatrix} \exists\ x_{1:n} \in \{y_{1:N}\} \text{ so that } \underline{\bar{z}_j = \sum_{i=1}^{n} w_i\ \varphi_j(x_i)}\ \text{ for } j = 1:n-1 \\ w_{1:n} \geq 0\ ,\ \sum w_i = 1 \end{bmatrix}$

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

[ Koksma Hlavka inequality ]

## Computing MMD$_k$

from Weighted Quadrature (14)

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

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

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

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

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

## Weight optimization

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

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