

2/26/2024

Geometry
@ uw
mmp

kernel Thinning

\subseteq finding core-sets

\neq data thinning !!! \rightarrow sample splitting

Problem: approximate distribution μ on X by a sample
 \xrightarrow{n} a large sample μ_N by smaller sample

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

empirical distribution

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

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

Plan

- Kernel Thinning problem
- Ideas for solving it
 - Herting
 - K. Thinning
 - weighted Quad.
- (Basic) RKHS facts ← geometric
Important/useful
- The methods : H, KT, WQ (if time)

Problem: approximate distribution μ on \mathcal{X} by a sample
— a large sample μ_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** $\mu f = \mu_N f \approx \mu_Q f$ for all $f \in \mathcal{H}_k$

RKHS

uniformly

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

More precisely

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

Maximum Mean Discrepancy

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

$$\text{MMD}_k(\mu_N, \mu_Q) < \varepsilon_n$$

$$\text{MMD}_k(\mu, \mu_Q) < \varepsilon_n$$

Why RKHS?

- compact specification of $\{f's\}$
- no need to enumerate!
- nice, elegant, well understood math

Wahba

$$\text{polynomials } k = (1 + x^T x)^d$$

kernel reg., classification

$k = \text{Gaussian, Matérn, ...}$

string, tree kernels

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

Why core-sets / thinning?

want to do ML with large data

- non-parametric
- f not known yet

⇒ save computation time
tractability

so Kernel $k(\cdot)$ must match
program !!

Why expectations $E_\mu[f]$?

many functionals (criteria of quality) are expectations

$P[\text{err}]$
 $E[\text{err}^2]$
 quantile

Problem (rephrased)

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

• kernel $k \equiv RKBHS$ \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 ε_n ?

- $\mu_N \sim iid \mu$

$\mu_N \approx \mu \Rightarrow$ want same rate for $\mu_Q \approx \mu_N$

- $E_\mu[x] = \mu_n^x$

\mathcal{H}_k

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i \rightarrow \text{Var } \hat{\mu} \propto \frac{1}{N}$$

$$\text{std } \hat{\mu} \propto \frac{1}{\sqrt{N}} = N^{-\frac{1}{2}}$$

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



benchmark rate

Problem (rephrased)

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

• kernel $k \equiv R\text{KHS}$ the

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

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

Solutions

- grid on $X \leftarrow$ optimal N^{-2r}
- MC = iid sampling $N^{-\frac{1}{2}}$
- MCMC $\sim N^{-\alpha} \quad \alpha < \frac{1}{2}$
- Quasi MC $\sim N^{-\alpha'} \quad \alpha' > \alpha$
- Herding
- DPP (Determinantal Point Process)
- Thinning
- Weighted Quadrature

dependent (not iid)
neg. corr.

Problem (rephrased)

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

• kernel $k \equiv \text{RKHS}$

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]
RKHS approximation
+ cubature algorithm (Caratheodory)

Solutions

• grid on $X \leftarrow$ 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 Point Process)
[Belhadji ...]

Thinning [Dwivedi, Mackey '21]

Weighted Quadrature

What makes thinning pb. hard? core-set $\varepsilon_n \sim \frac{C}{n^{\alpha}}$

For ex:

μ compact support

μ subgaussian

μ heavy tailed

rate of decay of
kernel, μ

- similar for k : compact support, ...

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

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

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

Kernel Thinning

kernel decay radius

- $R_{k,n} = \inf \underline{r}$ such that $\sup_{\|x-y\|_2 \geq \underline{r}} |k(x,y)| \leq \frac{1}{n} \|k\|_\infty$

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

Ex: $k(x,y) = N(x-y; 0, \sigma^2 I) \Rightarrow \bar{\phi}_k(r) = \int e^{-\frac{\|x-(x-y)\|^2}{2\sigma^2}} dy = \text{tail prob for } N(0, \frac{\sigma^2}{2} I)$ at $r \cdot \sqrt{2}$

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

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

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

It helps if these are known analytically \leftarrow otherwise

or exact, tractable

- kernel tail probabilities $Z_k(r)$ computable they must be approximated, majorized, ..
decay radii $R_{k,n}$

(+ decay rates of μ) (Thinning)

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

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

- (Bochner) Fourier repn of $k(\cdot, \cdot)$ $\xrightarrow{\text{Thinning}} k_{1/2}$

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

\hookrightarrow weighted Quadrature

Bochner

Caratheodory

Calculating MMD

Weight optimization

Computing MMD_k

Bochner

Let $k(x, y) \equiv k(x-y)$ continuous kernel on \mathbb{R}^d
shift invariant

Then

$$k(z) = \int_{\mathbb{R}^d} p(\omega) e^{-i\omega^T z} d\omega$$
 Fourier transform of
non-negative measure p

For ex: $k(z)$ $p(\omega)$ up to normalization

$$e^{-\|z\|^2/2}$$

$$e^{-\|\omega\|^2/2}$$

Gaussian

$$e^{-\|z\|_1}$$

$$\prod_{j=1}^d \frac{1}{1+\omega_j^2}$$

Laplace

$$\prod_{j=1}^d \frac{1}{1+z_j^2}$$

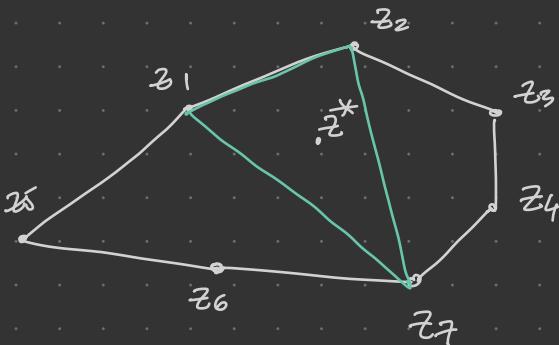
$$e^{-\|\omega\|_1}$$

Product kernel

Caratheodory

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

$$z'_{1:d} \subset \{z_{1:N}\}$$



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

$d=2 \Rightarrow 3$ points sufficient

For us:

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

points $y_{1:N}$

$$\bar{z} = \left[\frac{1}{N} \sum_{i=1}^N \varphi_j(y_i) \right] \bar{z}_j$$

$\Rightarrow \exists x_{1:n} \in \{y_{1:N}\}$ so that $\bar{z}_j = \sum_{i=1} w_{ij} \varphi_j(x_{ij})$ for $j=1:n-1$

$$w_{1:n} \geq 0, \sum w_{ij} = 1$$

Calculating MMD

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

[Koksma-Havka inequality]

straightforward from

$$\|\mu\|_{\mathcal{H}} = \sup_{f \in \mathcal{H}} \frac{|\mu f|}{\|f\|_{\mathcal{H}}}$$

From e.g. [Gretton]

$$\text{MMD}_k(\mu, \nu) = \sup_{\|f\|_{\mathcal{H}} \leq 1} |\mu f - \nu f| = \|\mu - \nu\|_{\mathcal{H}} = \|\mu \varphi - \nu \varphi\|_{\mathcal{H}} \xrightarrow{\varphi(x) \in \mathcal{H}_k} \mathbb{E}_{x \sim \nu} [\varphi(x)] - \mathbb{E}_{x \sim \mu} [\varphi(x)]$$

$$\|\mu \varphi\|_{\mathcal{H}}^2 = \langle \mathbb{E}_{x \sim \mu} [\varphi(x)], \mathbb{E}_{y \sim \mu} [\varphi(y)] \rangle = \mathbb{E}_{\substack{x \sim \mu \\ y \sim \mu}} \langle \varphi(x), \varphi(y) \rangle = \mathbb{E}_{\mu \otimes \mu} k(x, y)$$

$$\begin{aligned} \mu &\leftarrow \mu_N \quad \Rightarrow \quad \|\mu_N - \mu_Q\|^2 = \frac{1}{n^2} \mathbf{1}^T k(Q, Q) \mathbf{1} + \frac{1}{N^2} \mathbf{1}^T k(Y_{1:N}, Y_{1:N}) \mathbf{1} - \\ &\quad - \frac{2}{nN} \mathbf{1}^T k(Q, Y_{1:N}) \end{aligned}$$

Gram matrix Gram matrix

Computing MMD_k

from Weighted Quadrature (14)

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

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

$$= w^\top k(X, X) w - 2E_y [w^\top k(X, y)] + E_{y, y'} [k(y, y')]$$

↑
weights $w_{1:n}$
or $1/n$

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

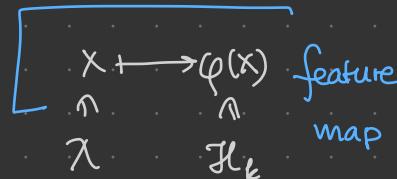
Weight optimization

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

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

Kernel Herding [Chen, Welling, Smola^[12]]

- Assume $k(x, x) = 1$ e.g. Gaussian
 $\| \varphi(x) \|_{\mathcal{H}_k}^2$



① define loss $G: \mathcal{H}_k \rightarrow \mathbb{R}$

$$G(z) = \langle z, \bar{\varphi} \rangle - \max_x \langle z, x \rangle$$

\uparrow
 concave \uparrow
 linear in z

$\bar{\varphi} = \underset{x \in \mu}{\mathbb{E}} [\varphi(x)] = \mu \varphi$
 ↪ piecewise linear, convex

② gradient

$$\nabla_z G = \bar{\varphi} - \varphi(x)$$

\uparrow
 $x = \underset{x}{\operatorname{argmax}} \langle z, \varphi(x) \rangle$

③ "gradient ascent" (step size = 1) :

given z_t

$$\left\{ \begin{array}{l} x_{t+1} = \underset{x}{\operatorname{argmax}} \langle z_t, \varphi(x) \rangle \\ z_{t+1} = z_t + \bar{\varphi} - \varphi(x_{t+1}) \end{array} \right.$$

\uparrow
 need to calculate!

need to solve!

\Rightarrow sequence x_1, x_2, \dots, x_T

= correct Q

But is it any good?

Kernel Herding
- 2 -

④ let's first look at $z_{1,2,\dots}$

converges to
 $\bar{\varphi}$?

$$z_T = z_{T-1} + \bar{\varphi} - \varphi(x_T) = z_0 + T\bar{\varphi} - \sum_{t=1}^T \varphi(x_t) \Rightarrow \boxed{\frac{1}{T} \sum_t \varphi(x_t) = \bar{\varphi} + \frac{z_0 - z_T}{T}}$$

\uparrow
...
 $z_{T-2} + \bar{\varphi} - \varphi(x_{T-1})$

remember $\|z_0\| = \|z_T\|_{\mathcal{H}} = 1 \Rightarrow \frac{\|z_0 - z_T\|}{T} \leq 2T^{-1}$ nice!

⑤ now let's look at $x_{1,2,\dots}$

$$\begin{aligned} x_{T+1} &= \underset{x}{\operatorname{argmax}} \left\langle \bar{\varphi} + T\bar{\varphi} - \sum_{t=1}^T \varphi(x_t), \varphi(x) \right\rangle = \\ &= \underset{x \in \mathcal{X}}{\operatorname{argmax}} \left\{ (T+1) \underset{x' \sim p}{\mathbb{E}} k(x, x') - \sum_{t=1}^T k(x, x_t) \right\} \\ &= \underset{x \in \mathcal{X}}{\operatorname{argmax}} \left\{ \underset{x' \sim p}{\mathbb{E}} k(x, x') - \frac{1}{T+1} \sum_{t=1}^T k(x, x_t) \right\} \end{aligned}$$

⑥ and now look at squared error

$$\mathcal{E}_T^2 = \left\| \bar{\varphi} - \frac{1}{T} \sum_{t=1}^T \varphi(x_t) \right\|_{\mathcal{H}}^2 = \left\langle \bar{\varphi} - \frac{1}{T} \sum \varphi(x_t), \dots \right\rangle$$

find \underline{x}_{T+1} that minimizes $\sum_{t=1}^T \|\varphi(x_t)\|^2$ for $x_{1:T}$ fixed

$$\left\langle \bar{\varphi} - \frac{1}{T+1} \sum_{t=1}^T \varphi(x_t), \bar{\varphi} - \frac{1}{T+1} \sum_{t=1}^T \varphi(x_t) \right\rangle =$$

$$= -2 \left\langle \bar{\varphi}, \frac{1}{T+1} \varphi(\underline{x}_{T+1}) \right\rangle + \frac{1}{(T+1)^2} \left\langle \varphi(\underline{x}_{T+1}), \varphi(\underline{x}_{T+1}) + \sum_{t=1}^T \varphi(x_t) \right\rangle + \dots$$

$$= -\frac{2}{T+1} \left\{ E_{\mu} k(x, \underline{x}_{T+1}) - \frac{1}{T+1} \underbrace{\left(k(\underline{x}_{T+1}, \underline{x}_{T+1}) + \sum_{t=1}^T k(\underline{x}_{T+1}, x_t) \right)}_1 \right\}$$

Same as (5)

$\Rightarrow x_{T+1}$ greedily minimizes $\sum_{t=1}^T \|\varphi(x_t)\|^2$

Remarks:

- sequence x_1, x_2, \dots NOT Markov
- elegant! but restrictive $k(x, x) = 1$, needs (estimate of)
(not all kernels and μ 's)

$$\langle \bar{\varphi}, z \rangle, \arg \max_z \langle \varphi(x), z \rangle$$

- how about MMD?

Calculating MMD

Weighted Kernel Quadrature [Hayakawa, Oberhauser, Lyons '21]

Given $y_{1:N} \in \mathcal{X}$, find $x_{1:n} \subset y_{1:N}$, weights $w_{1:n}$ with $w_i \geq 0$, $\sum_{i=1}^n w_i = 1$
so that $\mu_Q = \sum_{i=1}^n w_i \delta_{x_i}$ and $MMD_k(\mu_X, \mu_Q) \rightarrow$ with n fast

Idea

- ① approximate \mathcal{H}_k by some \mathcal{H}_{k_0} \leadsto finite dimensional decay
- ② then match \mathcal{H}_{k_0} exactly \leadsto Caratheodory

① Representer Theorem: $k(x, y) = \sum_{j=1}^{\infty} \lambda_j e_j(x) e_j(y)$
(Mercer)

Then, let $\underline{k}_0(x, y) = \sum_{j=1}^{n-1} \lambda_j e_j(x) e_j(y)$ $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$

Nice: $K \succeq k_0 \succeq 0$ \leadsto but must know spectral decomp of K

1.2 Nystrom to the rescue

Weighted Quadrature

- 2 -

1. sample ℓ points $\tilde{y}_{1:\ell}$ (different from $y_{1:n}$)

2. construct $W = \left[k(\tilde{y}_i, \tilde{y}_j) \right]_{i,j=1:\ell}^{\text{assume rank } \ell}$

3. SVD: $W = U \Lambda V^T$

4. keep components 1:s $k_0 \equiv \tilde{k}_s(x, y) = \sum_{j=1}^s \frac{1}{\lambda_j} \underbrace{U_j^T k(\tilde{y}, x)}_{\text{TR}} \underbrace{U_j k(\tilde{y}, y)}_{\mathcal{R}}$
 $s \ll \ell$

Note: $k \approx \tilde{k}(x, y) = k(x, \tilde{y}) W^T k(\tilde{y}, y) \Rightarrow K \gtrsim k_0 \gtrsim 0$

1.3 Theorem If $\bullet k - k_0 = k_1$ is Mercer kernel then $\exists \mu_Q$ (as above) with

$\bullet \dim \mathcal{H}_{k_0} \leq n-1$

$$\boxed{\mathbb{E}_{y_{1:N}} [\text{MMSE}_k^2(\mu, \mu_Q)] \leq 8 \left[\underbrace{\int k_1(x, x) d\mu(x)}_{\text{residual kernel}} + \frac{1}{N} C_{k_1, \mu} \right]}$$

↑ how diagonal is k
 ↓ within x $\int k(x, x) d\mu - \int k(x, y) d\mu_x d\mu_y$
 must be also $\frac{1}{N} \sim \frac{1}{n^2}$

(2) given $k, y_{1:N}, k_0$, functions $\varphi_{1:n-1}$ weighted Quadrature

-3-

with $\mathcal{H}_{k_0} = \text{span} \varphi_{1:n-1} \}$

find $x_{1:n}, w_{1:n}$

$$\text{Alg 1. } \tilde{\varphi}_{1:n-1} = \varphi_{1:n-1} \begin{bmatrix} y_1 \\ y_N \end{bmatrix}, \quad \tilde{\varphi}_0 = \begin{bmatrix} k_1(y_1, y_0) \end{bmatrix}$$

now supported only on
 $y_{1:N}$, measure is μ_N

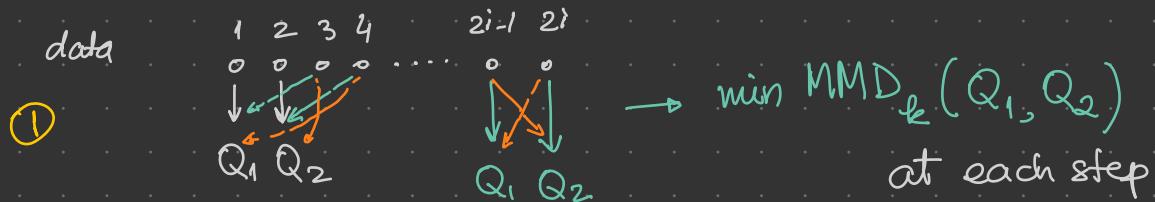
2. "Recombination algo" \Rightarrow obtain $\tilde{x}_{1:n+1}, v_{1:n+1}$ so that $\tilde{\varphi}_{0:n-1}^T(\tilde{x}) v = \tilde{\varphi}^T$

3. ...

$$\begin{array}{c} \text{---} \\ n+1 \\ \vdots \\ n+1 \end{array} \begin{array}{c} \text{---} \\ N \\ \vdots \\ N \end{array} = \frac{1}{N} \begin{array}{c} \text{---} \\ n \\ \vdots \\ 1 \end{array} \begin{array}{c} \text{---} \\ N \\ \vdots \\ 1 \end{array} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}^N$$

Kernel Thinning [Arivazh & Mackey '21]

Idea: Kernel halving



$$1:N = Q_1 \cup Q_2$$

$$Q_1 \leftarrow Q_1 \cup x \quad \text{with } x, x'y = \{x_{2i-1}, x_{2i}\}$$

$$Q_2 \leftarrow Q_2 \cup x'$$

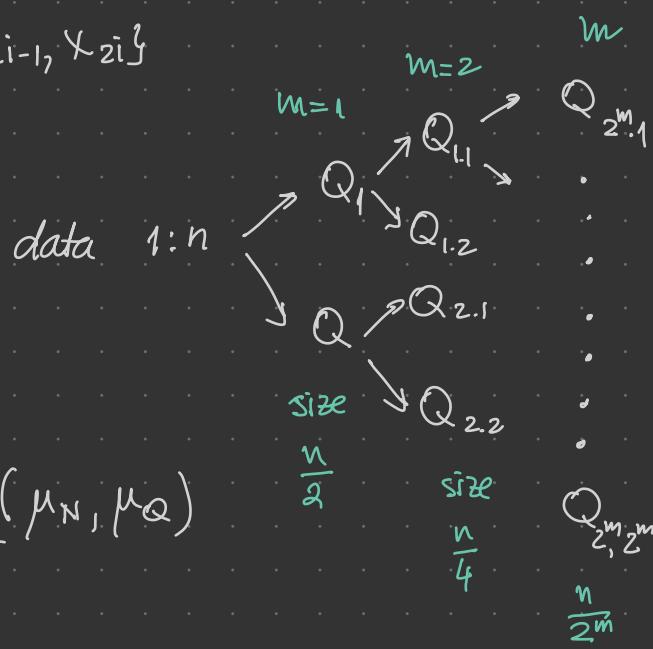
② Recursive Halving = kT-split

- Alg that runs kT-Halving in parallel at all m levels

③ kT-SWAP

$$Q^{\text{prelim}} \leftarrow \underset{l=1:2^m}{\operatorname{argmin}} \text{MMD}_k(\mu_N, \mu_Q)$$

+ some more refinement $\Rightarrow Q$



- What should m be?

$$\text{want } |Q| = \sqrt{n} = \frac{n}{2^m} \Rightarrow n = 2^{2m} \Rightarrow m = \frac{1}{2} \log_2 n$$

Kernel Thinning
- 2 -

- Run time $\sim n^2$

- Proof Ideas 1. Vector balancing

too beam

bounds $\|\mu_{\mathbb{N}^k} - \mu_Q k\|_\infty$

- 2. Square root kernel trick
to get bound on MMD $_k$

- $k_{1/2}$ is square root for k iff

$$k(x, y) = \int_x k_{1/2}(x, z) k_{1/2}(y, z) dz$$

- Theorem $MMD_k(\mu, \nu) \leq C \|\mu k_{1/2} - \nu k_{1/2}\|_\infty + \text{tail terms}$

- Algorithm: KT-Split uses $k_{1/2}$

Data $x_1, x_2, \dots, x_n \dots$ sequential

at time t , choose $s_t \in \pm 1$ sign so that

$$\left\| \sum_{t=1}^n s_t x_t \right\|_\infty \text{ is small}$$

How small? $\sim \sqrt{\log d/\delta \cdot \log n/\delta}$ w.p. $1-\delta$

worst case $\sim n$ (when $\|x_i\| \leq 1$)

Square root kernel

Kernel Thinning
- 2.1 -

If $k(x, y) = k(x-y)$ and even

and k_{k_2} — " — " —

then

$$\hat{k}_{k_2} = (\hat{k})^{1/2}$$

Fourier
transforms

Proof

$$k(x, y) = \int k_{k_2}(z-x) k_{k_2}(y-z) dz = \int_x^x k_{k_2}(t) k_{k_2}(y-x-t) dt = (\hat{k}_{k_2} * \hat{k}_{k_2})(y-x)$$

" x t $y-x-t$

$\hat{k}(y-x)$

convolution

$$\text{Fourier } \ast \mapsto \cdot \Rightarrow \hat{k} = \hat{k}_{k_2} \hat{k}_{k_2}$$