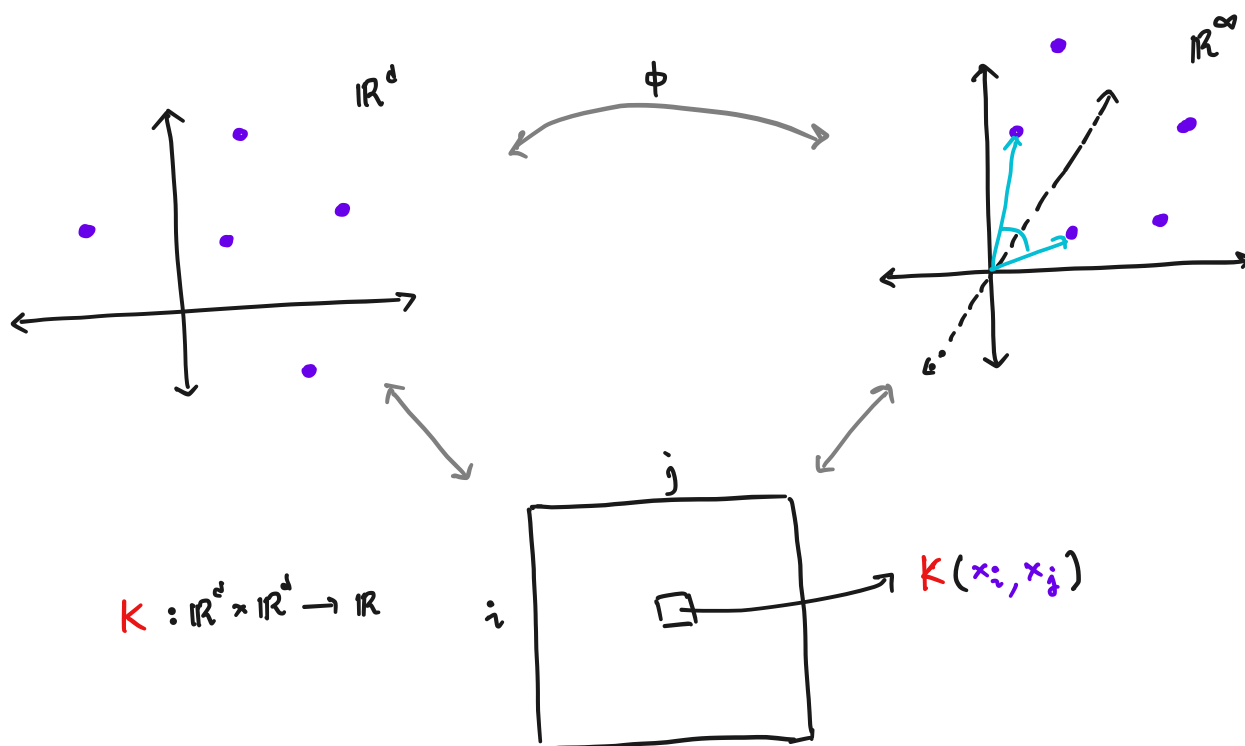


A Quasi-Monte Carlo Data Structure For Smooth Kernel Evaluations

Joint with Moses Charikar (STANFORD)
Michael Kapralov (EPFL)

RESEARCH DIRECTION: How to compute in
(HIGH DIMENSIONAL) GEOMETRIC spaces.

"Kernel Trick":



MOTIVATING QUESTIONS:

1. Algorithmic tradeoffs? Approximation, dimensionality....etc.
2. Are different kernels easier/harder?
3. ... (near)-linear time?

Problem Setup:

□ Kernel function $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$
(PSD, "radially decaying")

□ Dataset of vectors $p_1, \dots, p_n \in \mathbb{R}^d$.

□ Goal: Produce data structure s.t on
input $q \in \mathbb{R}^d$, can output

$$\frac{1}{n} \sum_{i=1}^n K(p_i, q)$$

(approximately + w.h.p.).

Ex: Gaussian kernel: $K(p, q) = \exp(-\|p-q\|_2^2/2)$

t-Student kernel: $K(p, q) = \frac{1}{1 + \|p-q\|_2^2}$.

Approx. Guarantees:

□ accuracy $\epsilon > 0$, failure probability $\delta > 0$

□ Fix query $q \in \mathbb{R}^d$:

$$\mu = \text{true answer} \triangleq \frac{1}{n} \sum_{i=1}^n K(p_i, q)$$

$$\hat{\mu} = \text{r.v. output}$$

$$\Pr \left[(1-\epsilon)\mu \leq \hat{\mu} \leq (1+\epsilon)\mu \right] \geq 1-\delta.$$

Parameters

- Dataset size n in dimension d
- Accuracy ϵ , Failure probability δ
- True answer μ
- Query time + Space complexity.

Curse of dimensionality

[Greenland-Rokhlin '83, Alman-Chu-Schild-Song '20]

FMM: $\log^{\Theta(d)}(1/\epsilon\mu)$ query time.

$n \log^{\Theta(d)}(1/\epsilon\mu)$ space complexity

SETH: require $n^{1-\Omega(1)}$ query time when $d \geq \log n$

for any $\text{poly}(n)$ -space complexity for
"high" accuracy $\log(1/\epsilon)$ -dependence.

Theorem: [This work]

For radially decaying smooth kernels $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0,1]$

□ Query time: $\frac{\text{poly}(d \log(1/\mu))}{\epsilon}$

□ Space complexity: $n \cdot \frac{\text{poly}(d \log(\frac{1}{\mu}))}{\epsilon}$

↳ [Charikar - Siminelakis '17, Backurs - Charikar - Indyk - Siminelakis '18,
... Phillips - Tai '20, ...]

Smoothness: $K(x, y)$ decays polynomially
with $\|x - y\|_2$

ex: 2-Student kernel $\leftrightarrow K(x, y) = \frac{1}{1 + \|x - y\|_2^2}$

non-ex: Gaussian kernel $\leftrightarrow K(x, y) = \exp\left(-\frac{\|x - y\|_2^2}{2}\right)$.

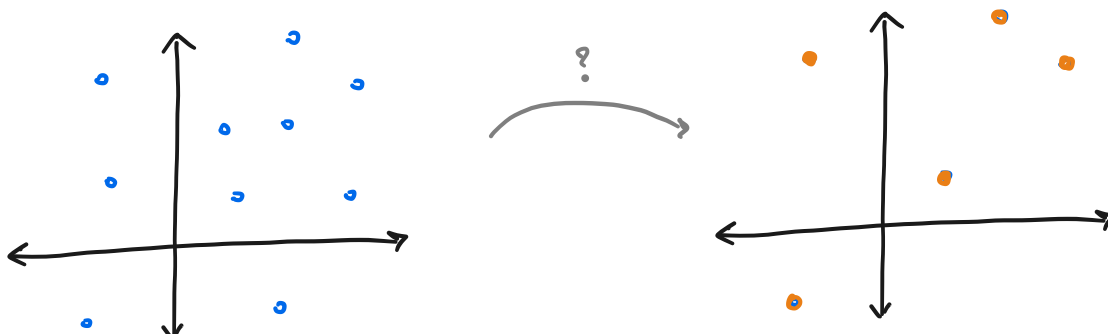
[Backurs - Charikar - Indyk - Siminelakis '18]

Smoothness is algorithmically useful primitive.

↳ $\frac{\text{poly}(\log(\frac{1}{\mu}))}{\epsilon^2}$ vs. $\frac{\mu^{-0.173}}{\epsilon^2}$

Sparsification, Coresets, Sampling ...

Which are my important data points?



$$\mu = \frac{1}{n} \sum_{i=1}^n k(p_i, q)$$

$\approx \pm \epsilon \mu$

$$\sum_{j=1}^s w_j \cdot k(p_i, q)$$

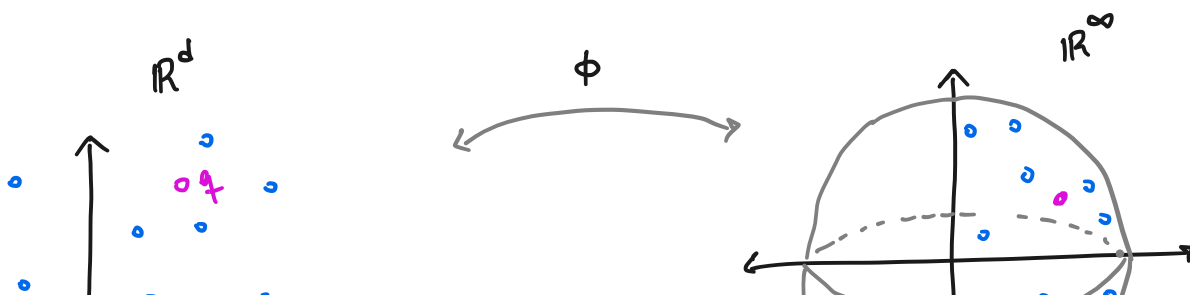
Outline

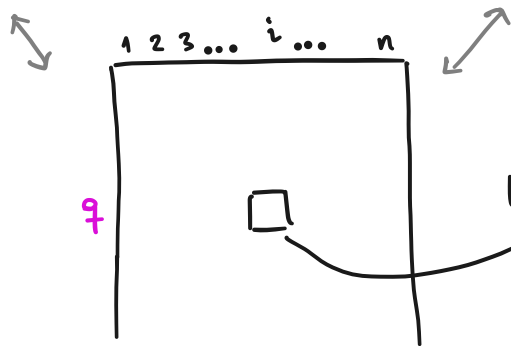
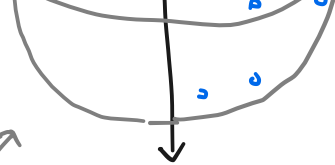
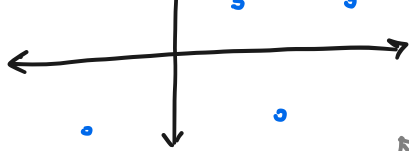
Benchmark 0: Random Sampling $\frac{d}{\epsilon^2 \mu}$

Benchmark 1: Discrepancy [Phillips-Tai '20] $\frac{d}{\epsilon \mu}$

Our Approach: LSH + Discrepancy $\frac{\text{poly}(d \log(1/\mu))}{\epsilon}$

Discrepancy for Cosets [Phillips-Tai '20]





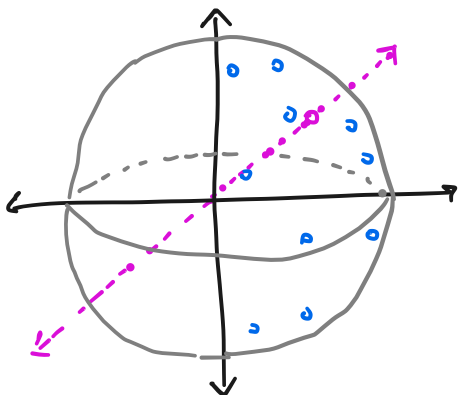
$$K(p_i, q) = \langle \phi(p_i), \phi(q) \rangle$$

⋮

Discrepancy: $x = \begin{bmatrix} \square & \square & \square & \dots & \square \end{bmatrix} \in \{-1, 1\}^n$
s.t.

$$\left| \sum_{i=1}^n x_i \cdot K(p_i, q) \right| = \left| \sum_{x_i=1} K(p_i, q) - \sum_{x_i=-1} K(p_i, q) \right| \leq c. \quad \text{w.h.p}$$

[Banerjee '98]



Vector Balancing Question

$$v = \sum_{i=1}^n x_i \phi(p_i)$$

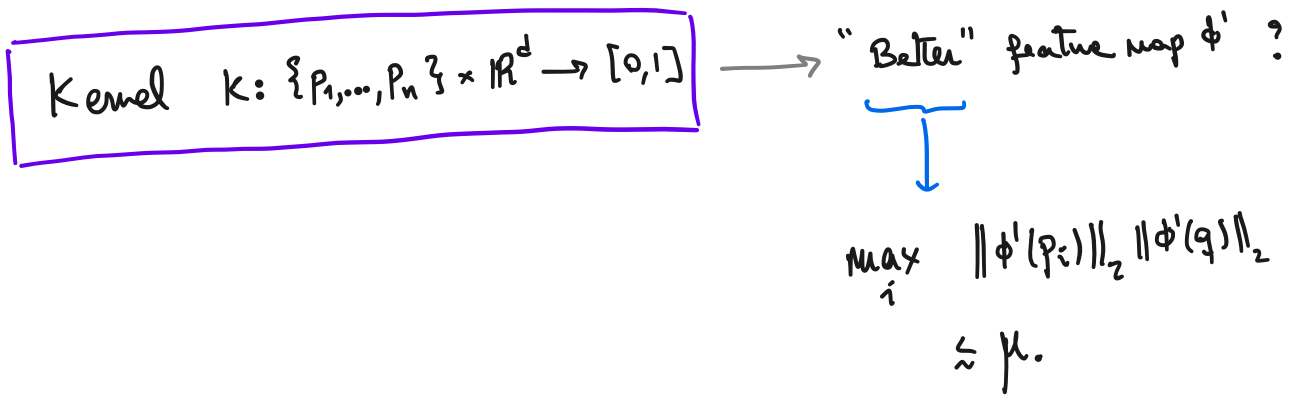
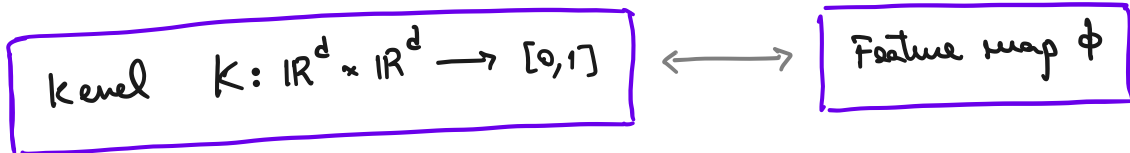
s.t. $|\langle v, \phi(q) \rangle|$ is small.

If q is fixed, becomes 1D problem

"Self-Balancing Walk" [Alweiss-Liu-Sawhney '20]

$$\text{discrepancy} \lesssim \max_{i \in [n]} \|\phi(p_i)\|_2 \|\phi(q)\|_2.$$

Conceptual Question:



What we show:

[Main Technical Claim - Informal] For dataset P , smooth kernel K

\exists random partition of $\mathbb{R}^d = R_1 \cup \dots \cup R_\ell$

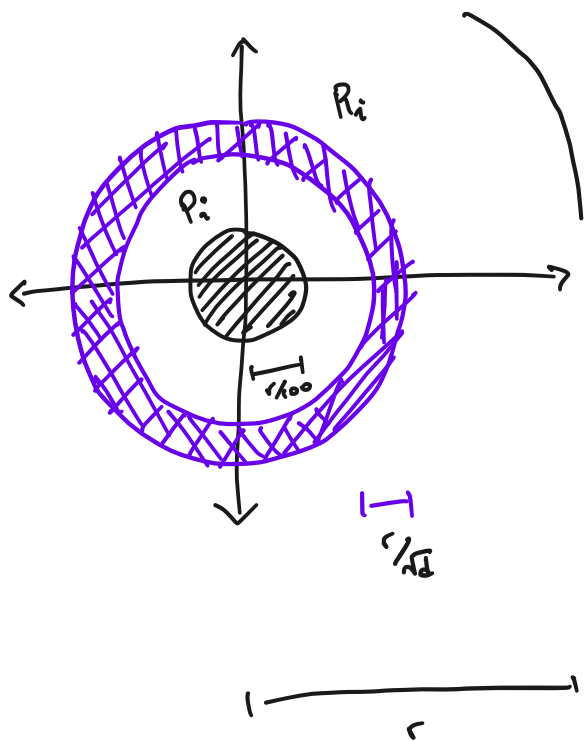
$P = P_1 \cup \dots \cup P_\ell$

and decomposition

s.t. efficiently identifiable + stored (via decision tree in time $\text{poly}(d \log(1/\mu))$)

$K: P_i \times R_i \rightarrow [0,1]$ has feature map ϕ' s.t

$$\|\phi'(p)\|_2 \|\phi'(q)\|_2 \leq \text{poly}(d \log(1/\mu)) * K(p,q).$$



Sweet Spot

1. Algorithmically decompose \mathbb{R}^d into few regions.
2. Admit better feature maps.

Tradeoff: decay + separation

Open:

- Gaussian i.e. "non-smooth" kernels.
- More complex kernel computations?
- Overcoming curse of dimensionality?

Thanks!