

Operator Convergence of Diffusion Maps and the Bistochastic Normalisation

Caroline L. Wormell

Joint work with Sebastian Reich

May 25, 2021

Introduction

Diffusion maps: on a random point sample, create matrix approximation of semigroup of weighted Laplacian.

- ▶ Eigendata of Laplacian (e.g. for dimensionality reduction, visualisation. . .)
- ▶ Non-parametric forecasting
- ▶ Approximation of more complex operators (e.g. Berry '18)

Diffusion maps

- ▶ Sample of M points $x^i \sim \rho$ abs. cts on domain $\mathbb{D} = (\mathbb{R}/L\mathbb{Z})^d$.
- ▶ Construct $M \times M$ kernel matrix K

$$K_{ij} = \frac{1}{M} g_\epsilon(x^i - x^j)$$

where g_ϵ is Gaussian kernel of *variance* ϵ .

- ▶ With appropriate weight vectors u and $v := 1/(Ku)$, construct Markov matrix

$$P = \text{diag } v \ K \ \text{diag } u$$

- ▶ As $M \rightarrow \infty$ and $\epsilon \rightarrow 0$ appropriately, P is approximation of $e^{\epsilon \mathcal{L}}$ where

$$\mathcal{L} = \frac{1}{2} \Delta + \log p \cdot \nabla \phi$$

Diffusion maps: convergence rates

Expect in general:

$$\left\| f(P^{t/\epsilon}) - f(e^{t\mathcal{L}}) \right\| = \mathcal{O} \left(\underbrace{M^{-\frac{1}{2}} \epsilon^{-\frac{d}{4} - \frac{1}{2}} \log(\dots)}_{\text{"variance error"}} + \underbrace{\epsilon^\theta}_{\text{"bias error"}} \right)$$

Know rigorously this works for

- ▶ f = pointwise evaluation of functions (von Luxburg *et al.* '08)
- ▶ f = eigendata of graph Laplacian (Calder and Trillos '20)

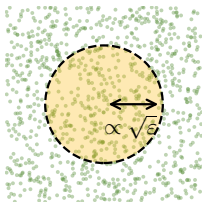


Figure: Effective support of g_ϵ contains $\mathcal{O}(M\epsilon^{d/2})$ data points.

Questions

Some mysteries we will investigate:

1. How does an operator defined on a random point cloud converge *as an operator* to a continuous kernel? (At the rate seen in practice?)
2. What is the (best possible) exponent in the bias error? How can we best choose weight vectors?

Kernel operator interpolation

The following operator on $C^0(\mathbb{D})$ matches kernel matrix K at sample points:

$$\mathcal{K}_\epsilon^M \phi = \sum_{i=1}^M \frac{1}{M} g_\epsilon(\cdot - x^i) \phi(x^i) = g_\epsilon * [\rho^M \phi]$$

As $M \rightarrow \infty$, expect \mathcal{K}_ϵ^M to converge to continuous kernel operator

$$\mathcal{K}_\epsilon \phi := g_\epsilon * [\rho \phi],$$

ideally in some Banach space $\mathcal{B}_\epsilon \subseteq C^0$.

Kernel operator interpolation

Because $g_\epsilon = g_{\epsilon/2} * g_{\epsilon/2}$ we can try for

$$\begin{aligned} \mathcal{K}_\epsilon^M - \mathcal{K}_\epsilon &= \underbrace{g_{\epsilon/2}}_{\text{bd. } C^0 \rightarrow \mathcal{B}_\epsilon} * \underbrace{(\mathcal{K}_{\epsilon/2}^M - \mathcal{K}_{\epsilon/2})}_{\text{small } \mathcal{B}_\epsilon \rightarrow C^0} \\ &= \text{small } \mathcal{B}_\epsilon \rightarrow \mathcal{B}_\epsilon \end{aligned}$$

Choice of \mathcal{B}_ϵ

As $\epsilon \rightarrow 0$, convolution by $g_\epsilon * \phi \rightarrow \phi$, so we expect $\mathcal{B}_0 = C^0$.
Let the complex domain

$$\mathbb{D}_\epsilon = \mathbb{D} + B_{\mathbb{C}}(\sqrt{\epsilon/2}).$$

One of the “smallest” candidates is

$$\mathcal{B}_\epsilon(\mathbb{D}) := \{\text{ct's analytic functions on } \mathbb{D}_\epsilon\}$$

endowed with $C^0(\mathbb{D}_\epsilon)$ norm.

This is good because

$$\|g_{\epsilon/2} * \phi\|_{\mathcal{B}_\epsilon} = \|g_{\epsilon/2}\|_{L^1(\partial\mathbb{D}_\epsilon)} \|\phi\|_{C^0} = e^{1/2} \|\phi\|_{C^0}$$

which gives us the bounded norm $C^0 \rightarrow \mathcal{B}_{\epsilon/2}$.

Kernel operator interpolation

Want to show that, up to log terms,

$$\delta := \|\mathcal{K}_{\epsilon/2}\phi - \mathcal{K}_{\epsilon/2}^M\phi\|_{\mathcal{B}_\epsilon \rightarrow C^0} \approx \text{pointwise bound} = \mathcal{O}(M^{-1/2}\epsilon^{-d/4})$$

We know that* for fixed ϕ and x ,

$$\left| (\mathcal{K}_\epsilon\phi - \mathcal{K}_\epsilon^M\phi)(x) \right| \leq \frac{C\epsilon^{-d/4}}{M^{1/2}} |\mathcal{N}(0, 1)|,$$

i.e error is $\mathcal{O}(M^{-1/2}\epsilon^{-d/4})$

How to extend efficiently to uniform bounds for all $\phi \in \mathcal{B}_\epsilon$, $x \in \mathbb{D}$?

* except for large deviations

Kernel operator interpolation

Want to show that, up to log terms,

$$\delta := \|\mathcal{K}_{\epsilon/2}\phi - \mathcal{K}_{\epsilon/2}^M\phi\|_{\mathcal{B}_\epsilon \rightarrow C^0} \approx \text{pointwise bound} = \mathcal{O}(M^{-1/2}\epsilon^{-d/4})$$

We know that* for fixed ϕ and x ,

$$\left| (\mathcal{K}_\epsilon\phi - \mathcal{K}_\epsilon^M\phi)(x) \right| \leq \frac{C\epsilon^{-d/4}}{M^{1/2}} |\mathcal{N}(0, 1)|,$$

i.e error is $\mathcal{O}(M^{-1/2}\epsilon^{-d/4})$

How to extend efficiently to uniform bounds for all $\phi \in \mathcal{B}_\epsilon$? Say,

$$\sup_{\|\phi\|_{\mathcal{B}_\epsilon}=1} \left| (\mathcal{K}_\epsilon\phi - \mathcal{K}_\epsilon^M\phi)(x) \right| \sim \frac{C\epsilon^{-d/4}}{M^{1/2}} \times \text{log terms}$$

* except for large deviations

Naive idea (Glivenko-Cantelli)

We have (bad) a priori estimate

$$\|\mathcal{K}_\epsilon - \mathcal{K}_\epsilon^M\|_{C^0} \leq 2 \sup g_\epsilon = C\epsilon^{-d/2}.$$

The unit ball in \mathcal{B}_ϵ is compact in C^0 , so we can cover the unit ball with a finite number of C^0 balls, i.e. there is a collection of $\#(\mathcal{B}_\epsilon, \xi)$ functions ϕ_n so that every ϕ with $\|\phi\|_{\mathcal{B}_\epsilon} \leq 1$ has $\|\phi_n - \phi\| \leq \xi$ for some n .

Naive idea (Glivenko-Cantelli)

Maximising over the ϕ_n ,

$$\sup_n \left| (\mathcal{K}_\epsilon \phi_n - \mathcal{K}_\epsilon^M \phi_n)(x) \right| \leq \frac{C\epsilon^{-d/4}}{M^{1/2}} \mathcal{N}_{\#(\mathcal{B}_\epsilon, \xi)},$$

where the maximum absolute value of T (non-ind.) standard normal distributions is $\mathcal{N}_T = \mathcal{O}(\sqrt{\log T})$. Thus,

$$\sup_{\|\phi\|_{\mathcal{B}_\epsilon}=1} \left| (\mathcal{K}_\epsilon \phi - \mathcal{K}_\epsilon^M \phi)(x) \right| \leq \frac{C\epsilon^{-d/4}}{M^{1/2}} \mathcal{N}_{\#(\mathcal{B}_\epsilon, \xi)} + C\epsilon^{-d/2}\xi.$$

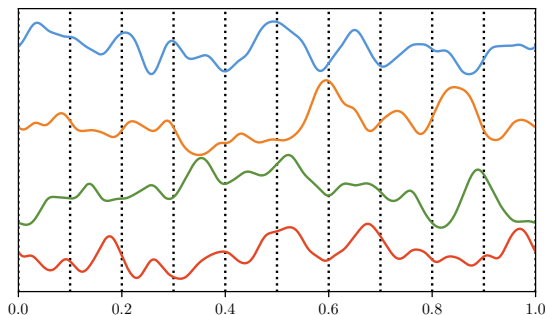
Want $\sqrt{\log \#(\mathcal{B}_\epsilon, \xi)}$ to grow sub-polynomially with $\epsilon, \xi \rightarrow 0$.

Naive idea (Glivenko-Cantelli)

In practice, if $X \subset \mathbb{R}^d$ is a hypercube of length L then

$$\log \#(C^0(X), B_\epsilon(X), \xi) = \mathcal{O}\left((L\epsilon^{-1/2} \log \xi^{-1})^d\right)$$

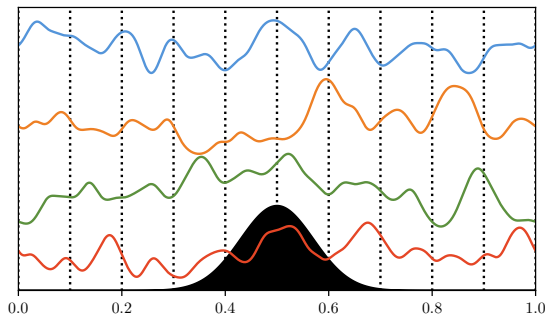
This gives us problems when $\epsilon^{1/2} \ll \text{diam } \mathbb{D}$.



Local Glivenko-Cantelli

However, we only see ϕ on a small part of the domain!

$$(\mathbf{g}_{\epsilon/2} * \psi)(x) = \mathbf{g}_{\epsilon/2} * (\mathbb{1}_{B(x, l\sqrt{\epsilon})}\psi) + \mathcal{O}(e^{-Cl^2})\|\psi\|_{L^1}.$$



Local Glivenko-Cantelli

We really just want a set of radius $l\sqrt{\epsilon}$, where l grows logarithmically.

$$\mathcal{B}_\epsilon^{x,l} := \{\text{bd. analytic functions on } B_{\mathbb{R}}(x, l\sqrt{\epsilon}) + B_{\mathbb{C}}(0, \sqrt{\epsilon}/2)\} \supset \mathcal{B}_\epsilon.$$

with

$$\log \#(\mathcal{B}_\epsilon^{x,l}, \xi) = \mathcal{O}\left((l \log \xi^{-1})^d\right)$$

and we are in business:

$$\begin{aligned} \sup_{\|\phi\|_{\mathcal{B}_\epsilon}=1} \left| (\mathcal{K}_\epsilon \phi - \mathcal{K}_\epsilon^M \phi)(x) \right| &\leq \frac{C\epsilon^{-d/4}}{M^{1/2}} \mathcal{N}_{\#(\mathcal{B}_\epsilon^{x,l}, \xi)} + C\epsilon^{-d/2}\xi + Ce^{-Cl^2} \\ &= \mathcal{O}\left(\epsilon^{-d/4} M^{-1/2} (\log M\epsilon^{-1})^{d-1/2}\right) \end{aligned}$$

Local Glivenko-Cantelli

We can use an easier compactness argument to extend to a supremum over all x , giving

$$\begin{aligned}\delta &:= \left\| (\mathcal{K}_{\epsilon/2}\phi - \mathcal{K}_{\epsilon/2}^M)\phi \right\|_{\mathcal{B}_\epsilon \rightarrow C^0} = \mathcal{O}\left(\epsilon^{-d/4} M^{-1/2} (\log M \epsilon^{-1})^{d-1/2}\right) \\ &= \text{appropriately small}\end{aligned}$$

All discretisation errors are then controlled by δ ! In particular,

$$\|\mathcal{P}_\epsilon^M - \mathcal{P}_\epsilon\|_{\mathcal{B}_\epsilon} = \mathcal{O}(\delta).$$

Comments

Result: convergence of spectral data, complex operator problems, etc. at near-pointwise rates.

- ▶ Requires very smooth kernel with exponentially decaying tails.
- ▶ Will generalise nicely to curved manifolds.
- ▶ Argument not based on Markov normalisation.
- ▶ Specialisation to Markov kernels would improve by $\mathcal{O}(\epsilon^{1/2})$ factor (Singer '06, Calder and Trillos '20).

Bias error analysis

Our weight vectors u, v are interpolated by functions $U_\epsilon^M, V_\epsilon^M$ which converge to U_ϵ, V_ϵ as $M \rightarrow \infty$.

Have infinite limit

$$\mathcal{P}_\epsilon \phi = V_\epsilon \mathcal{K}_\epsilon [U_\epsilon \phi].$$

Want to show that as $\epsilon \rightarrow 0$

$$\mathcal{P}_\epsilon \rightarrow e^{\epsilon \mathcal{L}}.$$

Bias error analysis

Our weight vectors u, v are interpolated by functions $U_\epsilon^M, V_\epsilon^M$ which converge to U_ϵ, V_ϵ as $M \rightarrow \infty$.

Have infinite limit

$$\mathcal{P}_\epsilon \phi = V_\epsilon \mathcal{K}_\epsilon [U_\epsilon \phi].$$

Want to show that as $\epsilon \rightarrow 0$

$$\mathcal{P}_\epsilon^{t/\epsilon} \rightarrow e^{t\mathcal{L}}.$$

Bias error analysis

Know \mathcal{L} is generator of SDE for invariant density p

$$dX = -\frac{1}{2}\nabla p dt + dW_t$$

We can study $\mathcal{P}_\epsilon^{t/\epsilon}$ as the evolution operator of a (time-varying) SDE.

Bias error: SDE formulation

Let

$$e^{s_t} = g_t * (\rho U_\epsilon) = e^{t\Delta/2}(\rho U_\epsilon).$$

Then $\rho U_\epsilon = e^{s_0}$ and $V_\epsilon = e^{-s_\epsilon}$.

$$\mathcal{P}_\epsilon \phi := V_\epsilon g_\epsilon \star (\rho U_\epsilon \phi) = e^{-s_\epsilon} e^{\epsilon\Delta/2} e^{s_0} \phi$$

is time- ϵ operator of forward equation of SDE

$$dX_t = -\nabla s_t dt + dW_t$$

So $\mathcal{P}_\epsilon^{t/\epsilon}$ is the time- t operator of

$$dX_t = \underbrace{-\nabla s_{\epsilon\{t/\epsilon\}}}_{\text{fast, periodic}} dt + dW_t$$

Bias error: SDE formulation

Time-average with $\mathcal{O}(t\epsilon^2)$ error:

$$dX_t \approx -\nabla \bar{s} dt + dW_t$$

$$\begin{aligned}\bar{s} &= \frac{1}{\epsilon} \int_0^\epsilon s_t dt \\ &= \frac{1}{2}(s_0 + s_\epsilon) + \mathcal{O}(\epsilon^2) \\ &= \frac{1}{2} \log(\rho U_\epsilon / V_\epsilon) + \mathcal{O}(\epsilon^2)\end{aligned}$$

Bias error: SDE formulation

Time-average with $\mathcal{O}(t\epsilon^2)$ error:

$$dX_t \approx -\nabla \bar{s} dt + dW_t$$

$$\begin{aligned}\bar{s} &= \frac{1}{\epsilon} \int_0^\epsilon s_t dt \\ &= \frac{1}{2}(s_0 + s_\epsilon) + \mathcal{O}(\epsilon^2) \\ &= \underbrace{\frac{1}{2} \log(\rho U_\epsilon / V_\epsilon)}_{\text{want} = \frac{1}{2} \log \rho} + \mathcal{O}(\epsilon^2)\end{aligned}$$

Bias error: SDE formulation

Time-average with $\mathcal{O}(t\epsilon^2)$ error:

$$dX_t \approx -\nabla \bar{s} dt + dW_t$$

$$\begin{aligned}\bar{s} &= \frac{1}{\epsilon} \int_0^\epsilon s_t dt \\ &= \frac{1}{2}(s_0 + s_\epsilon) + \mathcal{O}(\epsilon^2) \\ &= \underbrace{\frac{1}{2} \log(\rho U_\epsilon / V_\epsilon)}_{\text{want} = \frac{1}{2} \log \rho} + \mathcal{O}(\epsilon^2)\end{aligned}$$

- ▶ Typically we fit $e^{s_0}/\rho = U_\epsilon \approx p^{1/2}/\rho$. Since $s_\epsilon = s_0 + \mathcal{O}(\epsilon)$, get $\mathcal{O}(\epsilon)$ error (for $\rho \in C^{3/2+\alpha}$).

Bias error: SDE formulation

Time-average with $\mathcal{O}(t\epsilon^2)$ error:

$$dX_t \approx -\nabla \bar{s} dt + dW_t$$

$$\begin{aligned}\bar{s} &= \frac{1}{\epsilon} \int_0^\epsilon s_t dt \\ &= \frac{1}{2}(s_0 + s_\epsilon) + \mathcal{O}(\epsilon^2) \\ &= \underbrace{\frac{1}{2} \log(\rho U_\epsilon / V_\epsilon)}_{\text{want} = \frac{1}{2} \log \rho} + \mathcal{O}(\epsilon^2)\end{aligned}$$

- ▶ Typically we fit $e^{s_0}/\rho = U_\epsilon \approx p^{1/2}/\rho$. Since $s_\epsilon = s_0 + \mathcal{O}(\epsilon)$, get $\mathcal{O}(\epsilon)$ error (for $\rho \in C^{3/2+\alpha}$).
- ▶ Optimally accurate approximation is $\mathcal{O}(\epsilon^2)$, obtained via fitting weight ratio: $U_\epsilon/V_\epsilon = p/\rho$.

Sinkhorn problem

Since by Markov constraint $V = 1/(\mathcal{K}U)$, this means solving symmetric Sinkhorn problem for U :

$$U \times (\mathcal{K}U) = p/\rho.$$

- ▶ Only need $\rho, p \in C^{2+\alpha}$ for $\mathcal{O}(\epsilon^2)$ eigendata convergence.
- ▶ Fast iterative algorithm to compute U .

Sinkhorn problem

Since by Markov constraint $V = 1/(\mathcal{K}U)$, this means solving symmetric Sinkhorn problem for U :

$$U \times (\mathcal{K}U) = \rho/\rho.$$

- ▶ Only need $\rho, p \in C^{2+\alpha}$ for $\mathcal{O}(\epsilon^2)$ eigendata convergence.
- ▶ Fast iterative algorithm to compute U .

In paper: $p = \rho$, i.e. \mathcal{L} generates Langevin diffusion on ρ .

- ▶ P symmetric ($U = V$)
- ▶ P bistochastic (i.e. gives reversible Markov chain)

Comments

- ▶ In practice variance error $\mathcal{O}(M^{-1/2}\epsilon^{-d/4-1/2})$ will dominate bias error $\mathcal{O}(\epsilon^2)$!
- ▶ Expect convergence speed-up to work for symmetric kernels with correct 4th moments
- ▶ Only expect $\mathcal{O}(\epsilon)$ convergence on curved domains

Paper

We give, albeit in fairly specific setting, operator convergence with:

- ▶ Near-pointwise convergence rates for variance error
- ▶ Optimal convergence rates/choice of weights for bias error

In paper: proof of spectral convergence rates for standard and bistochastic normalisations.

Wormell, Caroline L., and Sebastian Reich. “Spectral convergence of diffusion maps: improved error bounds and an alternative normalisation.” *arxiv:2006.02037, to appear in SINUM* (2021).