

Scaling limit of stochastic optimization over large networks

Zaid Harchaoui^{1,3}, Sewoong Oh^{1,4}, Soumik Pal², Raghav Somani¹ and Raghav Tripathi²

¹UW CSE, ²UW Math, ³UW Statistics & ⁴Google

Brown University, May 11, 2023



- Motivation: Why optimize over graphs?
- Detour: Interacting particle system
- Optimization on graphons: Setup and Results
- Proof Sketches
- Future directions

Motivation (Extremal graph theory)

Mantel-Turán Problem

Among all graphs on n vertices **containing no triangles**, maximize the number of edges.

Since we are interested in large n , we normalize. Let's define

$$t(K_3, G) = \frac{\text{No. of triangles in } G}{n^3}, \quad t(K_2, G) = \frac{\text{No. of edges in } G}{n^2}.$$

Problem

Maximize $t(K_2, G)$ subject to the constraint $t(K_3, G) = 0$.

Mantel-Turán Theorem

$$t(K_2, G) > \frac{1}{2} \implies t(K_3, G) > 0.$$

See Aigner and Ziegler '14.

Motivation (Statistical physics)

- Let G be a weighted graph on n vertices with weighted adjacency matrix A .
- Let F be a finite simple graph on m vertices.
- We define *homomorphism density* of F into G

$$t(F, G) = \frac{1}{n^m} \sum_{i_1, i_2, \dots, i_m} \prod_{\{u, v\} \in E(F)} A(i_u, i_v).$$

Ising model on graphs (See Lovasz' book *Large Networks and Graph Limits*)

- $F :=$ A graph on m vertices. Every vertex may have a state $1, 2, \dots, q$.
- Between two neighboring vertices with states i, j , there is an interaction energy J_{ij} .
- A configuration is a map $\sigma : V(F) \rightarrow [q]$.
- The partition function is given by

$$Z = \sum_{\sigma: V(F) \rightarrow [q]} \exp \left(- \sum_{uv \in E(F)} J_{\sigma(u), \sigma(v)} \right) = \sum_{\sigma: V(G) \rightarrow [q]} \prod_{uv \in E(F)} \beta_{\sigma(u), \sigma(v)},$$

where $\beta_{ij} = \exp(-J_{ij})$.

- Minimizing Z is equivalent to minimizing $t(F, K_q^\beta)$, where K_q^β is complete graph with edge weights β_{ij} .

Motivation (Exponential random graph models)

- ERGM is an exponential family of models on simple graphs. E.g.,

$$P(G_n = G) \propto \exp \left(n^2 \sum_{i=1}^k \beta_i t(F_i, G) \right),$$

where F_1, \dots, F_k are simple graphs and β_1, \dots, β_k are real parameters.

- Can be used to fit a random graph to some empirical homomorphism densities.
- Given complete or partial data on the graph, say the edge density or the degree distribution, the MLE is an optimization on graphs.
- A similar optimization problem appears in the large deviation limit. See Chatterjee and Diaconis '13, Chatterjee '17.

Summary

- There are interesting optimization problems on graphs.
- Some of these optimization problem may not admit solutions in the space of finite graphs.

Plan

- Fill in the holes in the space of graphs, that is, take a completion of the space of all finite graphs.
- Try solving optimization problem on the complete space.
- These optimization problems have rich symmetries, invariance under relabeling of vertices. Can we exploit that?

Detour: Interacting Diffusion (McKean, Kac, Snitzman, Otto ...)

Consider the following example of interacting diffusions

$$dX_t^{i,N} = \frac{1}{N} \sum_{j=1}^N \nabla b(X_t^{i,N} - X_t^{j,N}) dt + dW_t^i, \quad i = 1, \dots, N$$

$$X_0^{i,N} = x_0^i.$$

Let $\mu_t^N := N^{-1} \sum_{i=1}^N \delta_{X_t^{i,N}}$. Then, $\mu_t^N \rightarrow \mu_t$ weakly where μ_t is a gradient flow with respect to 2-Wasserstein metric, given as

$$\partial_t \mu_t(x) = -\operatorname{div}_x [\mu_t(x) \cdot (\nabla b * \mu_t)(x)] + \frac{1}{2} \Delta \mu_t(x). \quad (1.1)$$

Equation (1.1) is a Wasserstein gradient flow of

$$\rho \mapsto \iint b(x-y) \rho(dx) \rho(dy) + \operatorname{Ent}(\rho).$$

Detour continued...

Interacting particles system converges to McKean-Vlasov

Suppose $X_0^{i,N}$ are i.i.d. with distribution μ_0 . As $N \rightarrow \infty$, each $X^{i,N}$ has a natural limit \bar{X}^i . Each \bar{X}^i is an independent copy of following McKean-Vlasov process

$$dX_t = (\nabla b * \mu_t)(X_t) dt + dB_t, X_{t=0} = X_0 \sim \mu_0.$$

- Think of each particle $X^{i,N}$ as doing a (noisy) gradient flow.
- Drift of the particle $X^{i,N}$ depends on ‘itself’ $X^{i,N}$ and ‘on the ensemble’ $N^{-1} \sum_{i=1}^N \delta_{X^{i,N}}$ in a symmetric way.
- Then, ‘the ensemble limit’ also performs a gradient flow in suitable sense.
- And, the evolution of a typical particle can be described by a McKean-Vlasov equation.

Introduction

Objective

Study large scale optimization problems over dense weighted graphs.

Let $G = (V, E)$ be a graph and let A be an adjacency matrix of G .

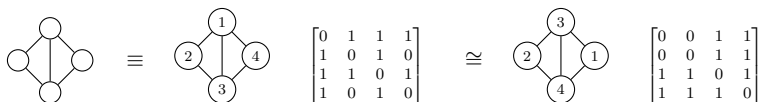


Figure: Symmetry in unlabeled graphs.

Invariant functions

A function $F: \mathcal{M}_n \rightarrow \mathbb{R}$ is said to be *invariant function/graph function* if $F(A) = F(A^\sigma)$ for all permutations $\sigma \in S_n$ and $A \in \mathcal{M}_n$, where $A^\sigma(i, j) = A(\sigma(i), \sigma(j))$.

Examples of functions

- Edge density: $h_{\square}(G) = (\# \text{ of edges in } G)/n^2$.
- Triangle density: $h_{\triangle}(G) = (\# \text{ of } \triangle \text{ in } G)/n^3$.

Plan and analogies with interacting diffusion

Objective

Let F be graph function. Our goal is to minimize F over large graphs.

Can perform gradient descent on finite graphs/symmetric matrices.

Exploiting the symmetry

- Think of the problem as an optimization problem on the space of ‘graphons’.
- Hope-Pray-Prove! The gradient descent process on finite graphs/symmetric matrices converge to a limit as $n \rightarrow \infty$.
- Can we show that the limit of GD is a **gradient flow on graphons**?
- What natural Markov processes on graphs converge to the gradient flow or related processes?

Graphons vs Wasserstein space

- Given a graph on n vertices is akin to particle ensemble
- Think of every edge as a *particle* and edge-weights are evolving

Setup and Results

Graphons

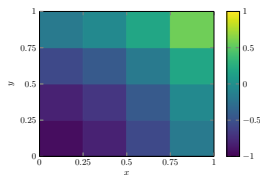
Kernels \mathcal{W}

A kernel is a measurable function $W : [0, 1]^2 \rightarrow [-1, 1]$ such that $W(x, y) = W(y, x)$.

- Adjacency matrix \equiv *kernel*.

$$\frac{1}{16} \begin{bmatrix} -16 & -15 & -12 & -7 \\ -15 & -14 & -11 & 1 \\ -12 & -11 & -6 & 4 \\ -7 & 1 & 4 & 9 \end{bmatrix}$$

Symmetric matrix A



Kernel representation of A

- Identify adjacency matrix/kernel up to ‘permutations’.
- Identify $W_1 \cong W_2$ if one can be obtained by ‘relabeling’ the vertices of the other, i.e.,

$$W_1(\varphi(x), \varphi(y)) = W_2(x, y), \quad \text{where } \varphi : [0, 1] \rightarrow [0, 1] \text{ is a measure preserving map.}$$

Graphons

Graphons $\widehat{\mathcal{W}}$ (Lovász & Szegedy, 2006): $\widehat{\mathcal{W}} := \mathcal{W}/\cong$

Cut metric :: Weak convergence

- Cut metric, δ_{\square} , metrizes graph convergence.
- $(\widehat{\mathcal{W}}, \delta_{\square})$ is **compact**.

Invariant L^2 metric δ_2 :: 2-Wasserstein metric \mathbb{W}_2

- Stronger than the cut metric (i.e., δ_{\square} convergence $\not\Rightarrow \delta_2$ convergence).
- **Gromov-Wasserstein distance** between $([0, 1], \text{Leb}, W_1)$ and $([0, 1], \text{Leb}, W_2)$.

We show (Oh, P., Somani, Tripathi, '21)

- The metric δ_2 is **geodesic** (just like \mathbb{W}_2). Geodesic convexity on $(\widehat{\mathcal{W}}, \delta_2)$.
- Notion of ‘gradient’ on $(\widehat{\mathcal{W}}, \delta_2)$ called ‘Frechét-like derivative’!
- Construction of ‘gradient flows’ on $(\widehat{\mathcal{W}}, \delta_2)^1$.

¹Gradient Flows: In Metric Spaces and in the Space of Probability Measures - Ambrosio, Gigli, Savaré, 2008

Existence of gradient flow on Graphons

Theorem [OPST '21]

If $R: \widehat{\mathcal{W}} \rightarrow \mathbb{R}$

- has a Fréchet-like derivative,
- is geodesically semiconvex in δ_2 ,

then starting from any $W_0 \in \widehat{\mathcal{W}}$, $\exists!$ gradient flow curve $(W_t)_{t \in \mathbb{R}_+}$ for R satisfying

$$W_t := W_0 - \int_0^t DR(W_s) ds + \text{boundary terms}, \quad t \in \mathbb{R}_+,$$

inside $\widehat{\mathcal{W}}$. At the boundary $\{-1, 1\}$ of $\widehat{\mathcal{W}}$, add constraints to contain it.

Scaling limits of GD [OPST '21 + HOPST '22]

Euclidean GD/SGD of R_n over $n \times n$ symmetric matrices, converges to the 'gradient flow' of R on the metric space of graphons.

Example

For $p \in [0, 1]$, define the entropy function $I(p) = p \log(p) + (1 - p) \log(p)$. In the following we assume $W : [0, 1]^2 \rightarrow [0, 1]$.

For a kernel W , define

$$I(W) := \iint I(W(x, y)) \, dx \, dy.$$

Gradient flow of $F(W) = t(K_3, W) + \beta I(W)$

$$W_t(x, y) = W_0(x, y) - 3 \int_0^t \int W_s(x, z) W_s(z, y) \, dz \, ds - \beta \int_0^t \log \left(\frac{W_s(x, y)}{1 - W_s(x, y)} \right) \, ds .$$

Finite dimensional gradient descent

$$W_t^{(n)}(i, j) = W_0^{(n)} - 3n^2 \int_0^t \frac{1}{n^3} \left(W_s^{(n)} \right)^2(i, j) \, ds - \beta \int_0^t \log \left(\frac{W_s^{(n)}(i, j)}{1 - W_s^{(n)}(i, j)} \right) \, ds .$$

Markov Chain converging to gradient flow

Suppose we want to construct a Markov process on graphs that converges to the gradient flow of triangle density $t(K_3, \cdot)$.

- Start with $G_{n,0}$.
- At each time step τ_n , all the edges in $G_{n,k}$ flip (or don't flip according to following rule).
 - If $\{i, j\}$ is not an edge in $G_{n,k}$ then $\{i, j\}$ remains a non-edge in $G_{n,k+1}$.
 - If $\{i, j\}$ is an edge in $G_{n,k}$ then drop it with probability

$$p_{ij} = \tau_n \frac{\Delta_{ij}}{n},$$

where Δ_{ij} = Number of triangles with containing $\{i, j\}$.

- Take the step-size $\tau_n = \frac{1}{n^2}$.
- One can also analyze variants of Metropolis-Hastings MC whose invariant measures are Gibbs measures.
- As $n \rightarrow \infty$, paths of such processes also converge. Current work with **Athreya-P.-Somani-Tripathi**.

Scaling limit of Noisy SGD

For $n \in \mathbb{N}$, let $\nabla R_n(A) = \mathbb{E}_\xi[\nabla \ell_n(A; \xi)]$ for $A \in \mathcal{M}_n$.

Stochastic Gradient Descent (SGD)

Given the k -th iterate $W_k^{(n)} \in \mathcal{M}_n$, sample ξ ,

$$W_{k+1}^{(n)} = W_k^{(n)} - \tau_n \cdot n^2 \underbrace{\nabla \ell_n(W_k^{(n)}; \xi)}_{\text{stochastic Euclidean gradient}}$$

Scaling limit of Noisy SGD

For $n \in \mathbb{N}$, let $\nabla R_n(A) = \mathbb{E}_\xi[\nabla \ell_n(A; \xi)]$ for $A \in \mathcal{M}_n$.

Noisy SGD

Given the k -th iterate $W_k^{(n)} \in \mathcal{M}_n$, sample ξ ,

$$W_{k+1}^{(n)} = W_k^{(n)} - \tau_n \cdot n^2 \underbrace{\nabla \ell_n(W_k^{(n)}; \xi)}_{\text{stochastic Euclidean gradient}} + \tau_n^{1/2} \cdot \underbrace{\xi_k \sim N(0, I)}_{\text{independent GOE noise}}$$

Scaling limit of Noisy SGD

For $n \in \mathbb{N}$, let $\nabla R_n(A) = \mathbb{E}_\xi[\nabla \ell_n(A; \xi)]$ for $A \in \mathcal{M}_n$.

Noisy SGD

Given the k -th iterate $W_k^{(n)} \in \mathcal{M}_n$, sample ξ ,

$$W_{k+1}^{(n)} = Proj \left(W_k^{(n)} - \tau_n \cdot n^2 \underbrace{\nabla \ell_n(W_k^{(n)}; \xi)}_{\text{stochastic Euclidean gradient}} + \tau_n^{1/2} \cdot \underbrace{\xi_k \sim N(0, I)}_{\text{independent GOE noise}} \right)$$

Scaling limit of Noisy SGD

For $n \in \mathbb{N}$, let $\nabla R_n(A) = \mathbb{E}_\xi[\nabla \ell_n(A; \xi)]$ for $A \in \mathcal{M}_n$.

Noisy SGD

Given the k -th iterate $W_k^{(n)} \in \mathcal{M}_n$, sample ξ ,

$$W_{k+1}^{(n)} = Proj \left(W_k^{(n)} - \tau_n \cdot n^2 \underbrace{\nabla \ell_n(W_k^{(n)}; \xi)}_{\text{stochastic Euclidean gradient}} + \tau_n^{1/2} \cdot \underbrace{\xi_k \sim N(0, I)}_{\text{independent GOE noise}} \right)$$

If $W_0^{(n)} \xrightarrow{\delta_2} W_0$, and $\tau_n \rightarrow 0$, as $n \rightarrow \infty$, then a.s.

$$W^{(n)} \xrightarrow{\delta_\square} \Gamma, \quad \text{as } n \rightarrow \infty,$$

where $\Gamma: t \mapsto \Gamma(t)$ is the curve described by the McKean-Vlasov equation.

McKean-Vlasov equation

- Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space with a Brownian Motion $B(t)$, and $(U, V) \stackrel{\text{i.i.d.}}{\sim} \text{Uni}[0, 1]$.
- R is a function on graphons and DR its \mathbf{L}^2 “gradient”.
- Consider the process $(\mathbf{X}(t), \mathbf{\Gamma}(t))$ such that on $\{U = u, V = v\}$,

$$d\mathbf{X}(t) = -(DR)(\mathbf{\Gamma}(t))(u, v) dt + dB(t) \underbrace{+ dL^-(t) - dL^+(t)}_{\text{constrain in } [-1, 1]}, \quad (\text{McKean-Vlasov})$$

$$\mathbf{\Gamma}(t)(x, y) = \mathbb{E}[\mathbf{X}(t) \mid (U, V) = (x, y)], \quad \forall (x, y) \in [0, 1]^2.$$

- $(\mathbf{\Gamma}(t), t \geq 0)$ is a process of kernels, given the initial random labeling.
- Expected to arise as limit of large number of graph dynamics.

A novel notion of “mean-field”

- “Mean-field interaction”: For any **edge-weight**, the effect of **all others edge-weights** on **its** evolution is invariant under vertex relabeling.
- “Propagation of chaos”: Every edge-weight between a set of m randomly chosen vertices evolves independently in the limit.

Proof Sketch: Scaling limits of gradient flow

- We show that the cut topology is *consistent* with the invariant L^2 metric δ_2^2 .
- At every $n \in \mathbb{N}$, consider *implicit Euler update* rule with positive a step size τ_n .
- The limit is obtained by showing Γ -convergence.

²Gradient Flows: In Metric Spaces and in the Space of Probability Measures - Ambrosio, Gigli, Savaré '08

Proof Sketch: Scaling limits of noisy SGD

- The existence of the deterministic limit Γ is obtained as a limit of a Picard iterations.
- Independently sample a sequence of vertices.
 - From SGD iterations $W^{(n)}(t)$, sample a random $m \times m$ submatrix process $W^{(n)}(t)[m]$.
 - Couple and get matrix processes $X(t)[m]$ & $\Gamma(t)[m]$ from McKean-Vlasov type SDEs.
- Use concentration estimates to show that as curves,

$$W^{(n)}[m] \xrightarrow{\delta \square} \Gamma, \quad \text{as } n \rightarrow \infty, \text{ and } m \rightarrow \infty, \quad \text{a.s.}$$

We recover the scaling limit of SGD (without added noise) as a corollary.

Upcoming work

Cut convergence gives limited information

- What can we infer if $W_n \rightarrow W$ in cut topology?
- We can infer the convergence of $t(F, W_n) \rightarrow t(F, W)$ for any finite graphs.
- Unfortunately, we can't say $\iint W_n(x, y)^2 dx dy \rightarrow \iint W(x, y)^2 dx dy$.

Cut topology is not good for weighted graphs

- Let $G(n, p)$ be the Erdős-Rényi graph.
- $G(n, p) \rightarrow W_p \equiv p$.
- Let $K(n, p)$ be the complete weighted graph with edge weights p .
- $K(n, p) \rightarrow W_p$.
- We would want to say $G(n, p)$ converges to an infinite exchangeable array $G(\infty, p)$ with i.i.d. Bernoulli random variables.
- And, $K(n, p)$ converges to an infinite (deterministic) array $K(\infty, p)$.
- Stronger but natural topology? Measure-valued graphons. **Upcoming paper- Athreya-P.-Somani-Tripathi**. Also analyze Metropolis-Hastings using gradient flows.

Simulations

- Turán's theorem: The n -vertex triangle-free graph with the maximum number of edges is a complete bipartite graph.

Q. Can we recover this theorem through an optimization problem on graphons?

$$F(W) = t(K_3, W) - \frac{1}{10}t(K_2, W) .$$

(a) GD ($n = 7$)

(b) GD ($n = 32$)

(c) GD ($n = 256$)

Future directions

- Extension to **Deep NNs**. Use a graphon for each layer (bipartite graph), respecting all joint layerwise permutation symmetries - **In progress**.

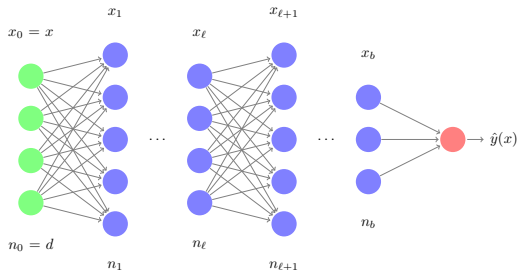


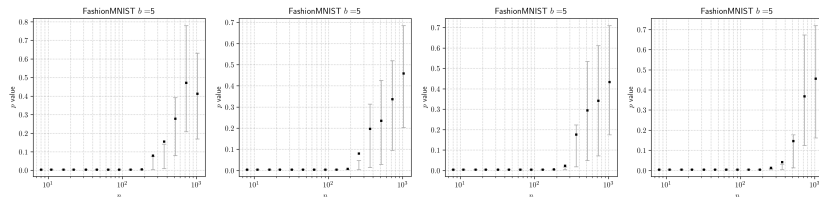
Figure: A b -layer NN.

- How does data distribution propagate across depth? Control theory, optimal transport - **Open**.

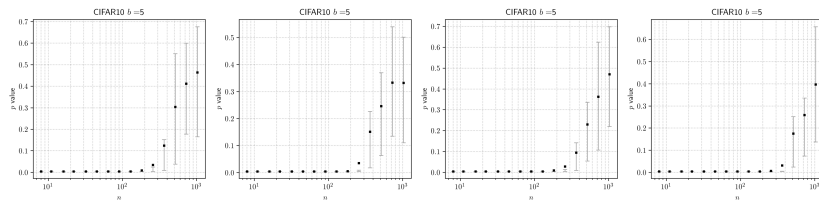
Propagation of Chaos experiments

- SGD training of a 5 layer deep feedforward ReLU networks.
- Test joint independence of elements in random 2×2 submatrices.
- Null hypothesis: All the 4 random variables are jointly independent.

$$\sigma: x \mapsto \max\{0, x\}.$$



(a) FashionMNIST

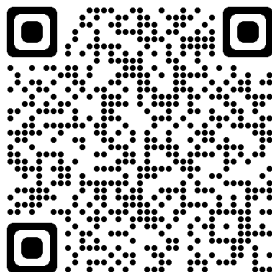


(b) Dataset: CIFAR10. x -axis: n , y -axis: p -value with interquartile range.

Thank you!

Thank you!

ArXiv version³: <https://arxiv.org/abs/2210.00422>



³Stochastic optimization on matrices and a graphon McKean-Vlasov limit - Harchaoui, Oh, Pal, Somani, Tripathi, 2022