

- Today:**
- a note on transferability rates
  - learning by transference - leveraging transferability in ERM

(THM) Given  $(g_n, x_n) \sim (\omega, x)$  and convolutions

$$h(g_n) = \sum_{k=0}^{\infty} h_k \left( \frac{g_n}{n} \right)^k \text{ and } T_h = \sum_{k=0}^{\infty} h_k T_w^{(k)}, \text{ under}$$

assumptions 3 & 4, we have:

$$\|y_n - y\| \leq \left( L + \frac{\pi n_c}{\delta_c} \right) \|T_w - T_{w_n}\| + (L_c + 2) \|x - x_n\| + 2\ell_c$$

↳ Actual transferability rates ( $\in \mathbb{N}$ ) depend on  $\|T_w - T_{w_n}\|$  and  $\|x - x_n\|$

- For  $\|T_w - T_{w_n}\|$ , we know that the relation b/w the operator norm and the cut norm,

combined with the graphon sampling lemma,  
give us the bound:

$$\|T_w - T_{w_n}\| \leq \sqrt{8} \|w - w_n\|_{\square} \leq \sqrt{\frac{\epsilon \cdot \alpha_2}{\sqrt{\log n}}} = \frac{4 \sqrt{11}}{\sqrt[4]{\log n}}$$

- For  $\|\chi - \chi_n\|_2$ , we have  $\|\chi - \chi_n\|_2 \leq \|\chi - \chi_n\|_1$

Further, Levie et al. (2023) defined a cut distance for graphon signals (akin to the graphon cut distance), which is equivalent to the  $L_1$  norm in the following sense:

$$\|\chi\|_{\square} \leq \|\chi\|_1 \leq \alpha \|\chi\|_{\square}$$

Leveraging these together with the graphon signal sampling lemma (Levie et al., 2023), we get:

$$\|\chi - \chi_n\|_2 \leq \frac{2.15}{\sqrt{\log n}} = \frac{30}{\sqrt{\log n}}$$

$\Rightarrow$  transferability rate of order  $\mathcal{O}((\log n)^{-\frac{1}{2}})$

→ Though applicable to general  $(w, \chi)$ , these rates are extremely slow.

↳ We can obtain better rates by making stricter assumptions on the graphon & graphon signal

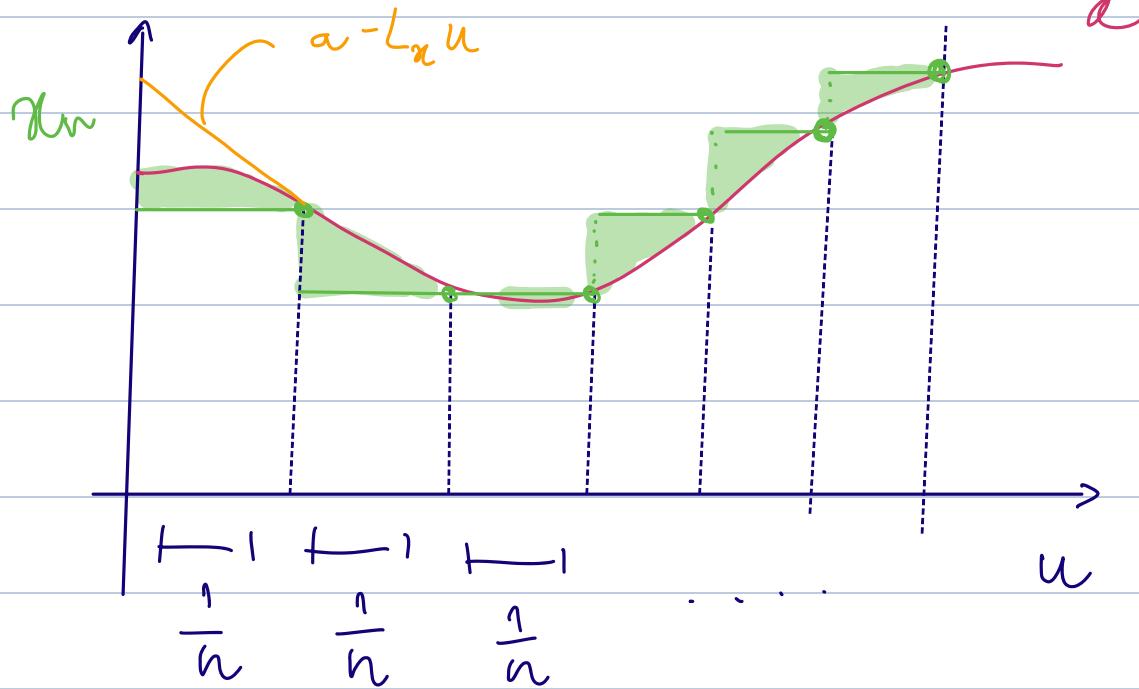
- For  $L_w$ -Lipschitz  $w$  &  $L_\chi$ -Lipschitz  $\chi$ , we can get

$$\|T_{w_n} - T_w\| \leq \|w - w_n\|_{HS} \leq \frac{L_w}{n}$$

$$\|\chi - \chi_n\| \leq \frac{L_\chi}{n}$$

(For deterministic  $G_n, x_n$ )

E.g. i



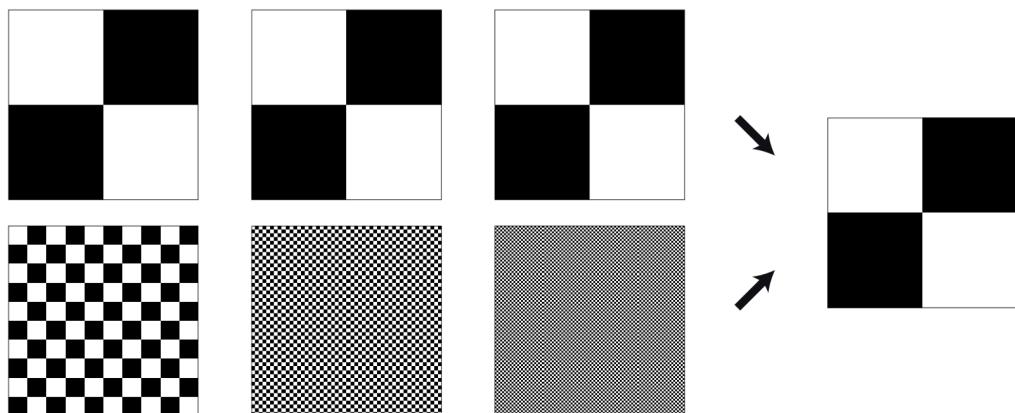
$$\begin{aligned}\|x - x_n\|^2 &= \int_a^b (x(u) - x_n(u))^2 du = \\ &= \sum_{I_i} \int_{I_i} (x(u) - x_n(u))^2 du \\ &\leq \|x\| \cdot \frac{1}{n} \left(\frac{L}{n}\right)^2 \\ \Rightarrow \|x - x_n\| &\leq \frac{L}{n}\end{aligned}$$

(similar example for graphian)

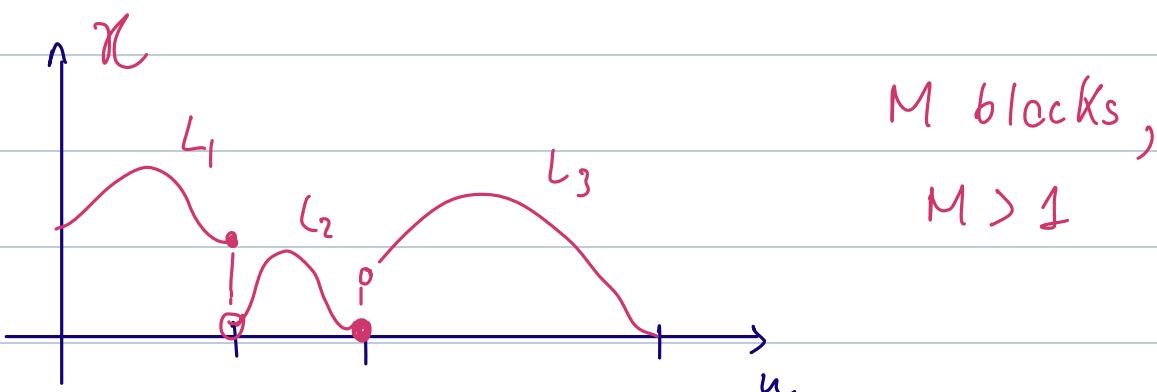
$\Rightarrow$  much better transferability rate, of order  $O\left(\frac{1}{n}\right)$

↳ check Ruiz, Chamorro, Ribeiro (2023) for stochastic.

However, Lipschitz assumptions on  $W$  and  $\chi$  can be problematic because we lose exchangeability. Conversely, the Lipschitz property is not preserved under measure preserving bijections (not "invariant to permutations")



↳ One compromise is to consider piecewise continuous  $W, \chi$ .



⇒ In this case, we get rates of order  $\Theta\left(\sqrt{\frac{M}{n}}\right)$   
(for node-deterministic  $G_m, X_n$ )

Check Arella-Medina et al. (2018) for stochastic.

→ Learning by transference

Instead of :

Sampling  $\rightarrow$  training  $\Phi$   $\rightarrow$  running  $\Phi$   
on  $G_n$  on  $G_n$  on  $G_m$ ,  $m > n$

We do :

sample  $\rightarrow$  train  $\Phi$   
on  $G_n$  for  $\Delta$  steps

increase  
 $n = n + \delta n$



run  $\Phi$   
on  $G_m$   
 $m > \max(n)$

- We do this to leverage graph convergence in the training process of the GNN.

- Helpful because, ideally, we would like to solve the ERM problem on the graph:

$$\min_{\mathcal{W}} \sum_{i=1}^{|G|} \ell(y^{(i)}, \Phi_{\mathcal{W}}(x^{(i)}, w))$$

with gradient updates:

$$\mathcal{R}_{\text{Rfi}} = \mathcal{H}_k - \gamma_k \frac{1}{|G|} \sum_{i=1}^{|G|} \nabla_{\mathcal{W}} \ell(y^{(i)}, \Phi_{\mathcal{W}}(x^{(i)}, w))$$

step size

- But in practice can only compute gradients on graphs:

$$\mathcal{R}_{\text{Rfi}} = \mathcal{H}_k - \gamma_k \frac{1}{|G|} \sum_{i=1}^{|G|} \nabla_{\mathcal{W}} \ell(y_n^{(i)}, \Phi_{\mathcal{W}}(x_n^{(i)}, g_n))$$

Still, if the graph gradients

$$\nabla_{f_G} \ell(y_n^{(i)}, \Phi_{f_G}(x_n^{(i)}, g_n))$$

converge to the graphon gradients

$$\nabla_{f_G} \ell(y^{(i)}, \Phi_{f_G}(x^{(i)}, w))$$

we can hope to learn  $f_G$ 's that are stationary points of the graphon ERM.

(THM) Let  $(g_n, x_n) \rightarrow (w, x)$  in the usual sense. Then, under the same assumptions needed for GNN transferability, and assuming the loss  $\ell$  and its gradients are Lipschitz, and  $\ell \geq 0$ , we have:

$$E\left(\|\nabla_{f_G} \ell(y_n^{(i)}, \Phi_{f_G}(x_n^{(i)}, g_n))\right)$$

$$- \nabla_{f_G} \ell(y^{(i)}, \Phi_{f_G}(x^{(i)}, w))\|$$

$$= C_1 \cdot \epsilon + C_2 \cdot \|w - w_n\|$$

where  $C_1$  &  $C_2$  are constants depending on  $D, F$ , the Lipschitz constants, and the graphon spectrum.

I.e., the graph gradients converge to the graphon gradient.

The proof follows easily (modulo some algebra from Lipschitz continuity of  $\ell$  and  $\nabla \ell$  and transferability of  $\Phi$ ).

- gradient convergence implies that algorithm  $(*)$  converges to the solution of the graphon ERM:

(THM) Under the previous assumptions, and further assuming that the ANN  $\Phi$  and  $\nabla_{\mathcal{X}} \Phi$  are stable to perturbations of  $\mathcal{X}$ , we have that, if at each step  $k$  Alg.  $(*)$  satisfies:

$$\mathbb{E}(\|\nabla_{\mathcal{X}_k} \Phi(w_k) - \nabla_{\mathcal{X}_k} \Phi(w)\|) \leq \|\nabla_{\mathcal{X}_k} \Phi(w)\| + \varepsilon,$$

then under the appropriate step sizes

$\eta_k$ , Alg. (\*) converges to an  $\epsilon$ -neighborhood of the solution of the graphon ERM in  $k^* = \mathcal{O}\left(\frac{1}{\epsilon^2}\right)$  iterations, at most.

Pf.: The proof relies on the following lemma proved in (Cervinó, 2023):

(Lemma) If at step  $k$   $n$  is such that:

$$\begin{aligned} \mathbb{E}\left(\|\nabla_{\delta \ell_k} \ell(y, \phi_{\delta \ell_k}(w)) - \nabla_{\delta \ell_k} \ell(y_n, \phi_{\delta \ell_k}(w_n))\|\right) \\ - \|\nabla_{\delta \ell_k} \ell(y, \phi_{\delta \ell_k}(w))\| \leq \epsilon \end{aligned}$$

then under the appropriate step size  $\eta_k$ , the next iterate  $\delta \ell_{k+1}$  satisfies:

$$\mathbb{E}[\ell(y, \phi_{\delta \ell_{k+1}}(w))] \leq \ell(y, \phi_{\delta \ell_k}(w)) - \frac{\eta_k (\delta \nabla \ell)}{2} \stackrel{< 0}{\curvearrowleft}$$

where  $\delta \nabla \ell = \|\nabla_{\delta \ell_k} \ell(y, \phi_{\delta \ell_k}(w))\|^2$

$$- \mathbb{E}\left(\|\nabla_{\delta \ell_k} \ell(y, \phi_{\delta \ell_k}(w)) - \nabla_{\delta \ell_k} \ell(y_n, \phi_{\delta \ell_k}(w_n))\|\right)^2 \geq \epsilon^2$$

Proof of theorem:

For every  $\epsilon > 0$ , we define the stopping time:

$$k^* = \min_{k \geq 0} \left\{ \mathbb{E} (\|\nabla_{\theta_k} \ell(y, \phi_{\theta_k}(w))\|) \leq C_1 + \epsilon \right\}$$

We have:

$$\mathbb{E} [\ell(y, \Phi_{\theta_k}(w)) - \ell(y, \Phi_{\theta^{*}}(w)) | k^*]$$

$$\mathbb{E} \left[ \sum_{k=1}^{k^*} \ell(y, \Phi_{\theta_{k-1}}(w)) - \ell(y, \Phi_{\theta_k}(w)) | k^* \right]$$

$$\Rightarrow \mathbb{E} [\ell(y, \Phi_{\theta_k}(w)) - \ell(y, \Phi_{\theta^{*}}(w))]$$

$$= \sum_{t=0}^{\infty} \mathbb{E} \left[ \sum_{k=1}^t \ell(y, \Phi_{\theta_{k-1}}(w)) - \ell(y, \Phi_{\theta_k}(w)) \right] P(k^* = t)$$

From Lemma 1, for any  $k \leq k^*$ , we have:

$$\mathbb{E} \left[ \sum_{k=1}^t \ell(y, \Phi_{\theta_{k-1}}(w)) - \ell(y, \Phi_{\theta_k}(w)) \right] \geq \frac{\eta}{2} \epsilon^2$$

And so:

$$\mathbb{E}[\ell(y, \Phi_{\mathcal{H}_0}(w)) - \ell(y, \Phi_{\mathcal{H}^*}(w))] \geq \frac{\eta}{2} \epsilon^2 \sum_{t=0}^{\infty} t P(k^*=t)$$
$$= \frac{\eta}{2} \epsilon^2 E(k^*)$$

Since the loss function is non-negative,

$$\frac{\mathbb{E}[\ell(y, \Phi_{\mathcal{H}_0}(w))]}{\frac{\eta}{2} \epsilon^2} \geq E(k^*)$$

$$\Rightarrow k^* = O\left(\frac{1}{\epsilon^2}\right) \blacksquare$$