

(THM) Non-asymptotic convergence of graph convolutions.

Given $(G_n, x_n) \sim (w, z)$ and convolutions

$$H(G_n) = \sum_{k=0}^{\infty} h_k \left(\frac{A_n}{n} \right)^k \text{ and } T_H = \sum_{k=0}^{\infty} h_k T_w^{(k)}, \text{ under}$$

assumptions 3 & 4, we have:

$$\|y_n - y\| \leq \left(L + \frac{\pi n c}{\delta_c} \right) \|T_w - T_{w_n}\| + (Lc + 2) \|z - x_n\| + 2lc$$

\Rightarrow Convergence $(G_n, x_n) \rightarrow (w, z)$ (w/ appropriate node labeling) means approximation improves w/ n as expected

\Rightarrow Convergence - discriminability tradeoff is explicit;
larger L & smaller c (= more discriminative filters)
lead to higher error bound.

\Rightarrow In the finite-sample regime, unless $l=0$, there is always leftover "nontransferable energy" z_{lc} corresponding to spectral components w/ $|\lambda_i| < c$, which do not converge.

Pf. We give a simplified version here, but the ideas are the same.

The proof is constructive; we leverage the fact that any filter can be written as the sum of a lowpass and a highpass filters.

Specifically,

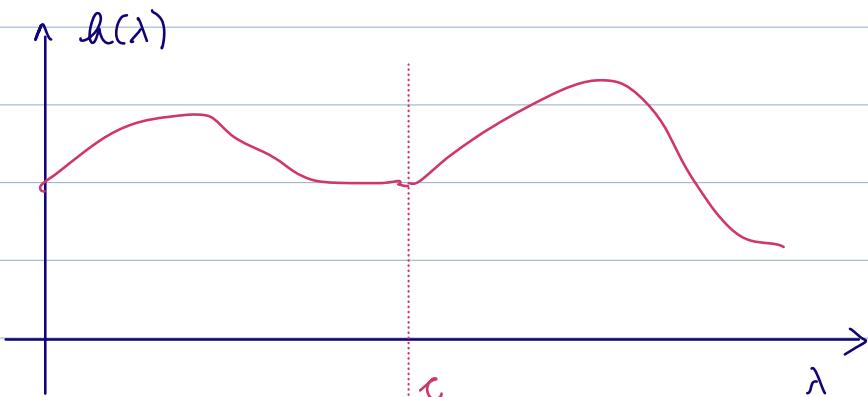
$$h(\lambda) = h_{LP}(\lambda) + h_{HP}(\lambda)$$

where

$$h_{LP}(\lambda) = \begin{cases} h(\lambda), & |\lambda| \geq c \\ 0 & \text{o.w.} \end{cases}$$

$$h_{HP}(\lambda) = \begin{cases} 0, & |\lambda| \geq c \\ h(\lambda) & \text{o.w.} \end{cases}$$

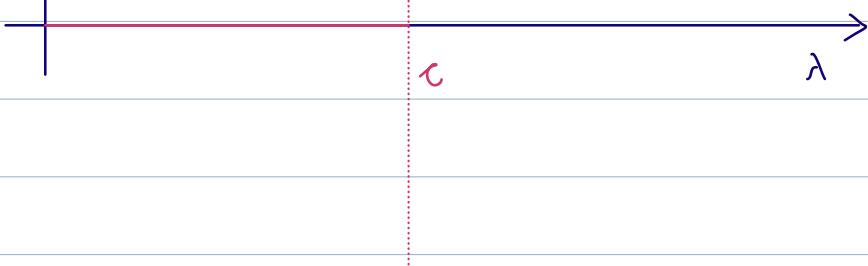
E.g.:



$h_{HP}(\lambda)$



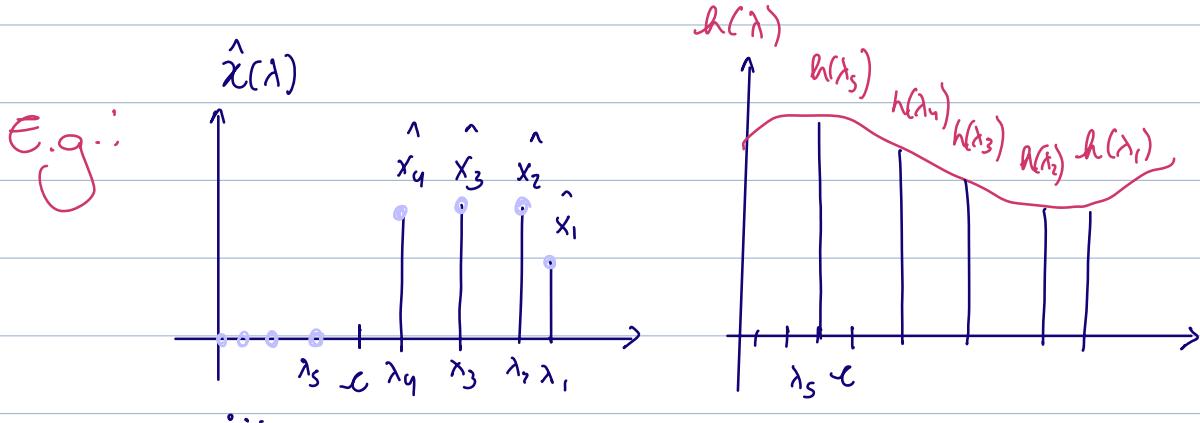
$h_{LP}(\lambda)$



→ The approximation error of h can then be upper bounded by bounding the approximation errors of h_{LP} & h_{HP}

- We start by upper bounding the approximation error of the low pass component.

To do that, we note that convergence of convolutions for LP filters is equivalent to convergence of general filters for band limited input signals.



$$\hat{y}_1 = \hat{x}_1 \cdot h(\lambda_1)$$

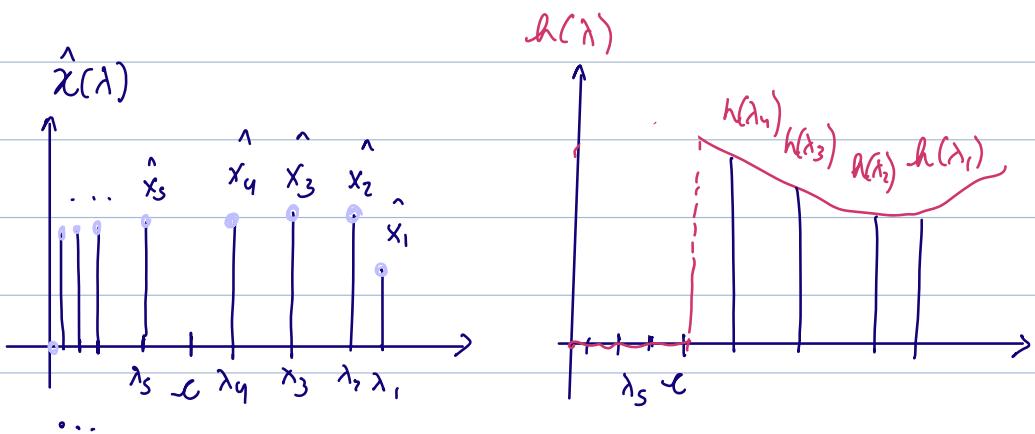
$$y_2 = \hat{x}_2 \cdot h(\lambda_2)$$

.

$$\begin{aligned} \hat{y}_s &= \hat{x}_s \cdot h(\lambda_s) = 0 \\ \hat{y}_6 &= \hat{x}_6 \cdot h(\lambda_6) = 0 \end{aligned} \quad \text{regardless of } h(\lambda)$$

O

This is equivalent to:



Therefore, we can leverage the following step from the proof of asymptotic convergence:

$$\|y - y_n\| \leq \left\| \sum_{j \in E} a(\lambda_j) \hat{x}_j \varphi_j \right. \\ \left. - \sum_{j \in E} a(\lambda_j^n) (\hat{x}_j \varphi_j - [\hat{x}_n]_j \varphi_j^n) \right\| \quad (4)$$

Adding and subtracting $\sum_{j \in E} a(\lambda_j^n) \hat{x}_j \varphi_j$, we get:

$$\|y - y_n\| \leq \left\| \sum_{j \in E} (a(\lambda_j) - a(\lambda_j^n)) \hat{x}_j \varphi_j \right\| \\ + \left\| \sum_{j \in E} a(\lambda_j^n) (\hat{x}_j \varphi_j - [\hat{x}_n]_j \varphi_j^n) \right\|$$

$\leq L \cdot \|T_w - T_{w_n}\| \|x\|^1$

$$+ \begin{pmatrix} |\lambda_i - \lambda'_i| \\ \leq \|T - T'\| \end{pmatrix}$$

$$+ \sum_{j \in E} \|a(\lambda_j)\| \left(\|\hat{x}_j\| \|\varphi_j - \varphi_j^n\| \right) \leq 1$$

$$+ \left\| \sum_{j \in E} a(\lambda_j) \varphi_j^n (\hat{x}_j - [\hat{x}_n]_j) \right\|$$



Davis - Kahan

$$\|\varphi_j - \varphi_j^n\| \leq \frac{\pi \|T_w - T_{w_n}\|}{2\delta_c} \quad \forall j \in \mathcal{C}$$



$$\left\| \sum_{j \in \mathcal{C}} \alpha(\lambda_j) \varphi_j^n (\hat{x}_j - (\hat{x}_n)_j) \right\| = \left\| \sum_{j \in \mathcal{C}} \alpha(\lambda_j) \varphi_j^n \langle \varphi_j - \varphi_j^n, x \rangle \right\|$$

$$+ \sum_{j \in \mathcal{C}} \alpha(\lambda_j) \varphi_j^n \langle (x - x_n), \varphi_j^n \rangle \|$$

$$\leq n_c \|\varphi_j - \varphi_j^n\| + \|x - x_n\|$$

Leading to:

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

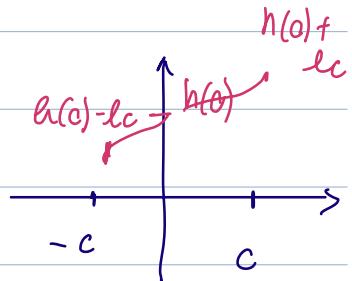
(a)

- Next, we consider the high-pass component. We reuse the argument from the proof of asymptotic convergence.

$$\left\| \sum_{j \in \mathbb{Z} \setminus \{0\}} a(x_j) \hat{x}_j \varphi_j - \sum_{j \in \mathbb{Z} \setminus \{0\}} a(x_j^n) (\hat{x}_j^n)_j \varphi_j^n \right\|$$

$$\leq \left\| (a(c) + l_c) x - (a(c) - l_c) x_n \right\|$$

$$\leq |a(c)| \|x - x_n\| + l_c \|x + x_n\|$$



$$\leq \|x - x_n\| + l_c \|2x + x_n - x\| \quad (b)$$

$$\leq (1 + l_c) \|x - x_n\| + \alpha \|x\| l_c$$

Combining (a) & (b) concludes the proof.

⇒ From convergence to transferability

In practice, we are not trying to approximate a graphon convolution, but to trans-

fer it; e.g., design on a smaller / offline graph & transfer it to a larger / online (dynamic) graph.

Transferability error bound can be derived from approx. error bound via Δ inequality.

↳ given $(G_n, x_n); (G_m, x_m) \sim (w, x)$

and a graph convolution w/ coeffs.

h_0, h_1, h_2, \dots , we have:

$$\|y_n - y_m\| \leq \|y_n - y\| + \|y_m - y\|$$

induced graphon
signals, in

sum of the approx. error
bounds on G_n & G_m

order to compare

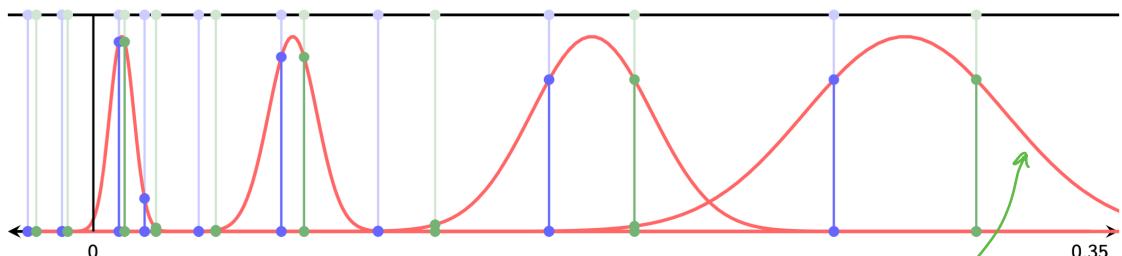
graph signals of \neq sizes

The theorem is essentially the same; since it will often be dominated by the approximation

error of the smallest graph, twice the approximation error for $\min\{n, m\}$ is a safe upper bound for the transferability error.

↳ Takeaways are the same

The most important: tradeoff between transf.
& discriminability



sharper filters are more discriminative but less
transferable

(↑L)

smoother filters are more transferable but less
discriminative

To improve the tradeoff, we need GNNs.

→ Transferability of GNNs

GNNs inherit the transferability property from graph filters.

(AS5) The nonlinearities σ are normalized Lipschitz ($L=1$). \rightarrow satisfied by ReLU, sigmoid, tanh ...

(THM) GNN transferability

Given $(G_n, x_n) \sim (w, x)$ and $(G_m, x_m) \sim (w, x)$
+ GNNs $\tilde{\phi}_{\theta}$ (G_n, x_n) and $\tilde{\phi}_{\theta} (G_m, x_m)$ with
the same set of weights θ_j ; under assumptions
AS3, AS4, AS5:

$$\|y_n - y_m\| \leq 2DF^{D-1} \left(L + \frac{\pi n_c}{\min_{n,m}(\delta_c)} \right) \max_{j=n,m} \|w - w_j\|$$

$$+ 2(L_c + 2) \max_{j=n,m} \|x - x_j\| + 4DF^{D-1} l_c$$

where D is depth and F the width of the GNN.

↳ The proof follows easily from the proof of transferability of graph convolutions.

- Note that if we have:

$$y = \delta(T_H x) \quad \text{and} \quad y_n = \delta(T_H^n x_n),$$

$$\|y - y_n\| \leq L_\delta \cdot \|T_H x - T_H^n x_n\|$$

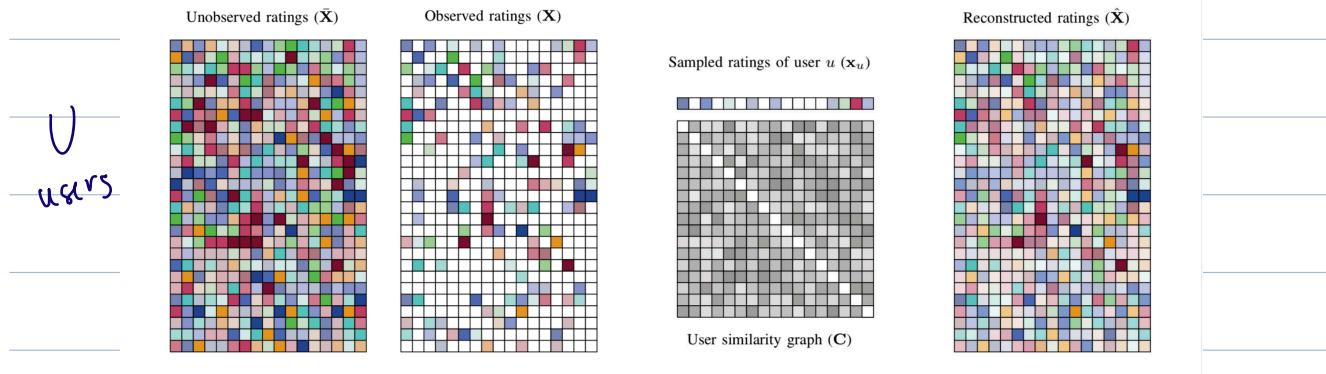
↑ by
AS5

↑ transf. of convolution

- D and F show up by recursively bounding $\|x_L - x_{L,n}\|$ in terms of $\|x_{L-1} - x_{L-1,n}\|$ for $1 \leq L \leq D$.

E.g.: Transferability in recommender systems

M_{movies}



→ Based on the observed rating matrix X , we can build either a user or movie similarity graph from the users'/movies' correlations

→ In the case where the graph connects similar movies (i.e., with similar ratings):

$$\delta_{m_1, m_2} = \frac{1}{|\mathcal{U}_{m_1, m_2}|} \sum_{u \in \mathcal{U}_{m_1, m_2}} (x_{um_1} - \mu_{m_1, m_2})(x_{um_2} - \mu_{m_2, m_1})$$

where m_1, m_2 are movies, \mathcal{U}_{m_1, m_2} is the set users having rated m_1 & m_2 , and μ_{ij} is the rating rating of movie i by $u \in \mathcal{U}_{ij}$

-o The graph is defined as :

$$[A]_{m_1, m_2} = \frac{\sigma_{m_1, m_2}}{\sqrt{\sigma_{m_1, m_1} \sigma_{m_2, m_2}}}$$

The signals are defined as :

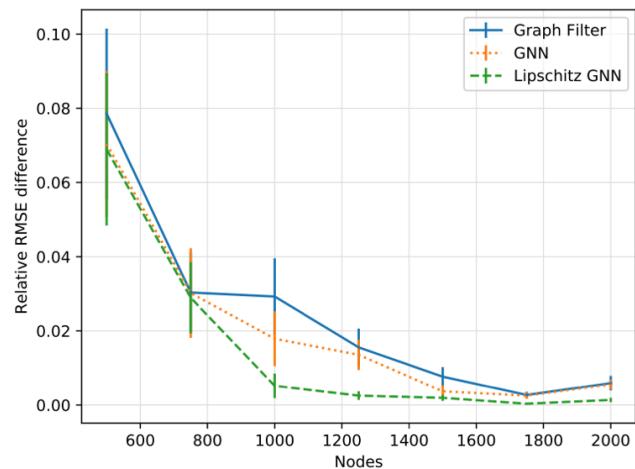
$$x_u \in \mathbb{R}^M \quad 1 \leq u \leq U$$

$$[x_u]_m = X_{um} \quad \text{if available (0 o.w.)}$$

-o We hold out some of the existing ratings as labels; i.e., for $u \in \mathcal{U}'$ & $m \in \mathcal{M}'$, we set $(x_u)_m = 0$ and $[y_u]_m = x_{um}$. We then compute the loss over $\mathcal{U}' \& \mathcal{M}'$

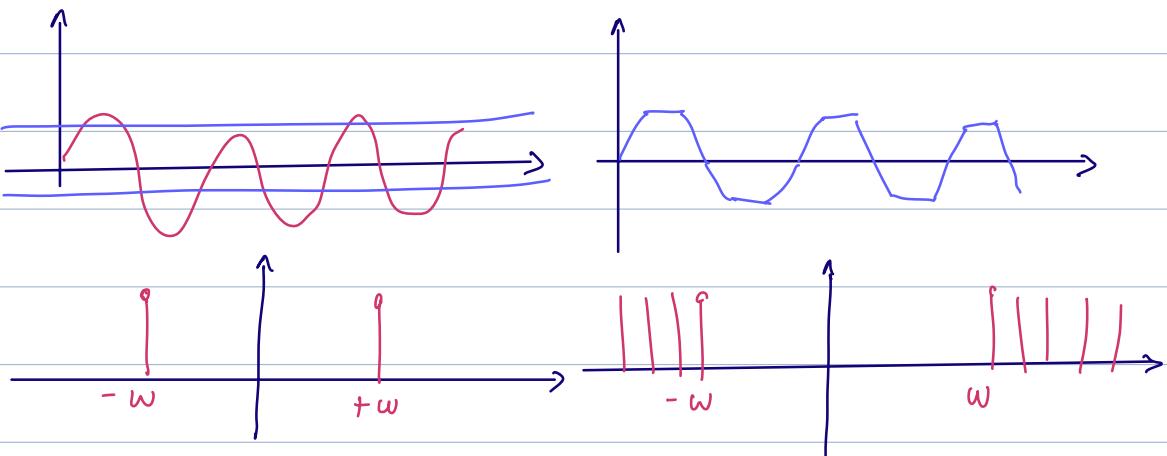
-o Then, we run a transferability experiment (like in HW3): train on G_n , transfer to G_m , $m > n$, while varying n and computing:

$$\underline{C} = \frac{\text{error}_{G_n} - \text{error}_{G_m}}{\text{error}_{G_n}}$$



If the transferability error of GNNs is essentially the same as that of graph convolutions (modulo scaling by $D_{\text{f}}^{\frac{1}{2}}$), why are they more transferable?

↳ For the same reason they are more stable, nonlinearities scatter energy $\sum_{i \in e} \|x_{e-1, i}\|^2$ to $\sum_{i \in e} \|x_{e, i}\|^2$, where components can be discriminated



↳ for the same level of transferability, better
discriminability, and vice-versa