

Today: → popular GNN archits. that can be expressed in conv. form

- MPNN
- GCN
- ChebNet
- GraphSAGE

→ GATs

Recap: GNNs

A GNN is a sequence of L layers of the form:

$$U_\ell = \sum_{k=0}^{K-1} S^k X_{\ell-1} H_{\ell,k}$$

$$X_\ell = \sigma(U_\ell)$$

where: S is the GSO (e.g. $S = A, L, \dots$)
 $X_{l-1} \in \mathbb{R}^{n \times d_{l-1}}$
 $H_{l,k} \in \mathbb{R}^{d_{l-1} \times d_l}$
 $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ nonlinear (e.g. ReLU)
 and: $X_0 = X; \quad Y = \bigoplus_{\mathcal{H}} (X, S) = X_L$

$$\mathcal{H} = \{ H_{l,k} \}_{l,k}$$

This is a convolutional GNN: it is based on convolutional filters (from GSP)

↳ GNNs inherit amazing theoretical properties, which explain their practical success, from them!

But it is general enough to encompass a majority of the most popular architectures in use these days.

► Message-passing neural networks (MPNNs) (Gilmer et al.)

Each layer consists of 2 operations:

$$1) \text{ message: } [m_e]_i = \sum_{j \in \mathcal{N}(i)} M_e([x_e]_i, [x_e]_j, A_{ij})$$

$(i \in V)$

$$2) \text{ update: } [x_{e+1}]_i = U_e([x_e]_i, [m_e]_i)$$

\Rightarrow As long as M_e is linear, "message" can be expressed as a graph convolution

e.g.:

$$M_e([x_e]_i, [x_e]_j, A_{ij}) = \alpha [x_e]_i + \beta \sum_{j \in \mathcal{N}(i)} A_{ij} [x_e]_j$$



$$m_e = \alpha x_e + \beta A \cdot x_e \rightarrow \text{graph convolution with } K=2, S=A, h_0=\alpha, h_1=\beta$$

\Rightarrow as long as U_e is the composition of a pointwise nonlinearity σ with a linear fcn on $[m_e]_i, [x_e]_i$, can be expressed as GNN layer

e.g.: $U_e([x_e]_i, [m_e]_i) = \sigma(\alpha' [x_e]_i + \beta' [m_e]_i)$

$$x_{e+1} = \sigma \left(\underbrace{(\alpha' + \beta' \alpha)}_{\substack{\text{pointwise} \\ \text{NL}}} x_e + \underbrace{\beta' \beta A}_{\substack{\text{graph conv. w/ } S=A, K=2}} x_e \right)$$

$h_0 = \alpha' + \beta' \alpha$
 $h_1 = \beta' \beta$

Ex.: from the perspective of learning the weights (h_0, h_1 vs. $\alpha, \alpha', \beta, \beta'$), what is different?

► Graph convolutional networks (Kipf & Welling, 2017)

A GCN layer is given by:

$$[x_e]_i = \sigma \left(\sum_{j \in \mathcal{N}(i)} \frac{[x_{e-1}]_j}{|\mathcal{N}(i)|} H_e \right)$$

$$x_e = \underbrace{\sigma}_{\substack{\text{pointwise} \\ \text{nonlinearity}}} \left(\underbrace{S x_{e-1} H_e}_{\substack{\text{graph conv. with} \\ K=2, H_{e0}=0}} \right)$$

$S = D_{\text{bin}}^{-\frac{1}{2}} A_{\text{bin}} D_{\text{bin}}^{-\frac{1}{2}}, H_{e1} = H_e$

(A_{bin} is the binary
 $A_{\text{bin}} = (A > 0)$ adj.
↳ in Python

↳ advantages:

- only one local diffusion step
- simple / interpretable
- normalization by degree avoids exploding / vanishing gradients

↳ disadvantages:

- only one diffusion step
- no edge weights
- only supports $S=A$
- no "self loops" (unless present in A)
- even if self-loops, no ability to weight $[X_{t-1}]_i$ & $[X_{t-1}]_j, j \in N(i)$ differently



Chebnet (Defferrard et al., 2017)

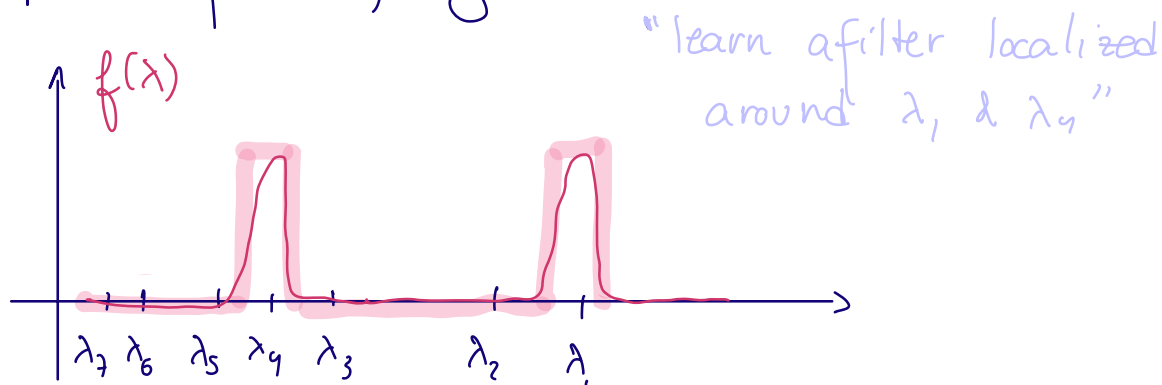
"Fast" implementation of:

$$h(\lambda) = \sum_{k=0}^{K-1} a_k \lambda^k(\underline{L})$$

↳ Laplacian

which we already know is the spectral domain representation / frequency response of the graph convolution \rightarrow the weights h_k are the same

Learning in the spectral domain provides an "inductive bias" for learning localized spectral filters, e.g.

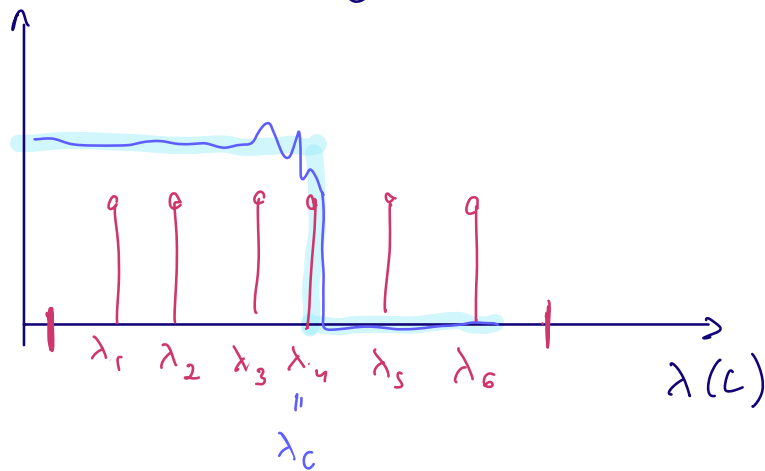


However, translating back to the spectral domain is expensive:

- we need to compute $\text{GFT}(x) = \hat{x} = V^T x$ & the inverse GFT of $h(\lambda)\hat{x}$, $y = V h(\lambda)\hat{x}$
- ↳ cost of 2 additional mv multiplications + cost of eigendecomposition of L

Solution: Chebyshev polynomials. In conventional SP, filters based on these polynomials

have better cutoff behavior in the spectral domain
(compared with polynomial filters)
I.e., say we want to approximate
the following lowpass filter:



Using $h(\lambda) = \sum_{k=0}^{K-1} h_k \lambda^k$ requires more coeffs.

for the same quality of approx. than the Chebyshev polynomial approx.:

$$h(\lambda) = \sum_{k=0}^{K'-1} h_k T_k(\lambda)$$

i.e., for the same approx. quality, $K' < K$

The Chebyshev polynomials can be efficiently calculated using the following recurrence:

$$\begin{cases} T_0(\lambda) = 1 \\ T_1(\lambda) = \lambda \\ T_{k+1}(\lambda) = 2\lambda T_k(\lambda) - T_{k-1}(\lambda) \end{cases}$$

satisfy $T_n(\cos\theta) = \cos(n\theta)$

Why is Chebyshev expansion better than Taylor expansion in this case?

→ Chebyshev polynomial expansion minimizes L_∞ norm ($\|x\|_\infty = \max_i |x_i|$) over the approximation interval
 ↳ while Taylor series is a local approximation around some $\tilde{\lambda}$

↳ check the Chebyshev equioscillation theorem for the formal statement & proof.

Further: the polynomials T_0, T_1, T_2, \dots are orthogonal, so the orthonormal basis comes for free!

↳ with respect to the inner product:

$$\langle f, g \rangle = \int_{-1}^1 f(x) g(x) \frac{dx}{\sqrt{1-x^2}}$$

Pf. $\int_{-1}^1 T_n(x) \cdot T_m(x) \frac{dx}{\sqrt{1-x^2}} \Rightarrow \begin{matrix} x = \cos \theta \\ \frac{dx}{d\theta} = -\sin \theta \end{matrix}$

$$\Rightarrow \int_{\pi}^{2\pi} T_n(\cos \theta) T_m(\cos \theta) \frac{dx}{\sqrt{1-x^2}} \quad \text{where } \frac{dx}{d\theta} = -\sin \theta$$

$$= \int_{\pi}^{2\pi} \cos(n\theta) \cos(m\theta) d\theta \quad (T_n = \cos(n\theta))$$

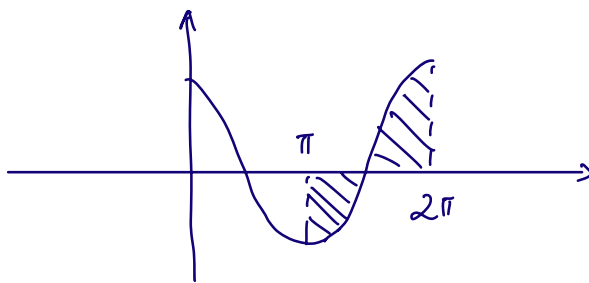
$$= \int_{\pi}^{2\pi} \left(\frac{e^{in\theta} + e^{-in\theta}}{2} \right) \left(\frac{e^{im\theta} + e^{-im\theta}}{2} \right) d\theta$$

$$= \frac{1}{4} \int_{\pi}^{2\pi} e^{-i(n+m)\theta} + e^{-i(n-m)\theta} + e^{-i(m-n)\theta} + e^{i(n+m)\theta} d\theta$$

•) if $n=m=0$: $\frac{1}{4} \int_{\pi}^{2\pi} 4 d\theta = \left. \theta \right|_{\pi}^{2\pi} = \pi$

•) if $n=m \neq 0$: $\frac{1}{4} \int_{\pi}^{2\pi} 2 d\theta + \frac{1}{4} \int_{\pi}^{2\pi} 2 \cos((n+m)\theta) d\theta$

$$= \frac{\pi}{2} + 0 = \frac{\pi}{2}$$



•) if $n \neq m$, then we get

$$\frac{1}{4\pi} \int_0^{2\pi} 2 \cos((m+n)\theta) + 2 \cos((n-m)\theta) d\theta = \boxed{0}$$

Hence: we get orthogonality for free!

no need to implement in the spectral domain

Hence: we get orthogonality for free!

no need to implement in the spectral domain

$$\left\{ \begin{array}{l} \text{Chebnet layer: } X_L = \sigma \left(\sum_{k=0}^{K-1} T_k(\bar{L}) X_{L-1} H_{L,k} \right) \\ \bar{L} = \frac{2L}{\lambda_{\max}} - I \quad (\text{to ensure spectra btw } [-1, 1]) \end{array} \right.$$

↳ advantages

- fast, cheap
- localized spectral filters

↳ disadvantages

- unstable
- S restricted to L
- difficult to interpret in node domain

► Graph SAGE (Hamilton-Ying-Leskovec, 2017)

Each SAGE layer implements the following operations:

$$[U_e]_i = \text{AGGREGATE}_e(\{[X_{e-1}]_j, j \in \mathcal{N}(i)\})$$

$$[X_e]_i = \sigma(\text{CONCAT}([X_{e-1}]_i, [U_e]_i) H_e)$$

$$[X_e]_i = \frac{[X_e]_i}{\|[X_e]_i\|_2}$$

The standard AGGREGATE operation is an average over $\mathcal{N}(i)$. Letting $H = \begin{bmatrix} H_0 \\ H_1 \end{bmatrix}$, we get:

$$[U_e]_i = \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} [X_{e-1}]_j \Rightarrow \underline{U_e = S X_{e-1}}$$

$S = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$
binary

$$X_e = \sigma([X_{e-1} \parallel U_e] \begin{bmatrix} H_0 \\ H_1 \end{bmatrix}) = \sigma(X_{e-1} H_0 + S X_{e-1} H_1)$$

↳ convolutional GNN w/
 $K=2$

↳ equivalent to GCN! The only difference is the
node wise normalization $\frac{[X_e]_i}{\|[X_e]_i\|_2}$, which helps
in some cases (empirically)

→ SAGE implementations allow for a variety of
AGGREGATE functions, including max, LSTM, \rightarrow no longer
look at $N(i)$ as a sequence \leftarrow convolutional
↳ lose perm. equiv.

→ same advantages & disadvantages as GCN

Obs.: the authors of SAGE popularized the AGGREGATE-
UPDATE representation of GNNs ($K=2$):

$$X_e = \sigma \left(S X_{e-1} \parallel H_e \right)$$

↳ aggregate \rightarrow update