## Learning with Graph Signals

▶ Almost ready to introduce GNNs. We begin with a short discussion of learning with graph signals

▶ In this course, machine learning (ML) on graphs $\equiv$ empirical risk minimization (ERM) on graphs.

▶ In ERM we are given:

⇒ A training set $\mathcal{T}$ containing observation pairs $(\mathbf{x}, \mathbf{y}) \in \mathcal{T}$. Assume equal length $\mathbf{x}, \mathbf{y}, \in \mathbb{R}^n$.

⇒ A loss function $\ell(\mathbf{y}, \hat{\mathbf{y}})$ to evaluate the similarity between $\mathbf{y}$ and an estimate $\hat{\mathbf{y}}$

⇒ A function class $\mathcal{C}$

▶ Learning means finding function $\Phi^* \in \mathcal{C}$ that minimizes loss $\ell\Big(\mathbf{y}, \Phi(\mathbf{x})\Big)$ averaged over training set
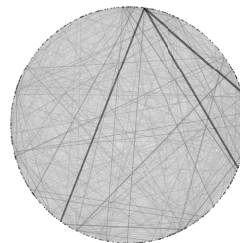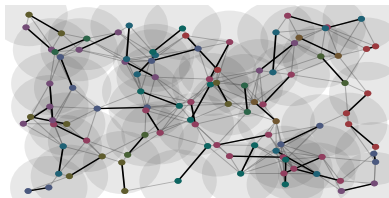
$$\Phi^* = \operatorname*{argmin}_{\Phi \in \mathcal{C}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \ell\Big(\mathbf{y}, \Phi(\mathbf{x}),\Big)$$

▶ We use $\Phi^*(\mathbf{x})$ to estimate outputs $\hat{\mathbf{y}} = \Phi^*(\mathbf{x})$ when inputs $\mathbf{x}$ are observed but outputs $\mathbf{y}$ are unknown

▶ In ERM, the function class $\mathcal{C}$ is the degree of freedom available to the system's designer

$$\Phi^* = \operatorname*{argmin}_{\Phi \in \mathcal{C}} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{T}} \ell\Big(\mathbf{y}, \Phi(\mathbf{x})\Big)$$

▶ Designing a Machine Learning $\equiv$ finding the right function class $\mathcal{C}$

▶ Since we are interested in graph signals, graph convolutional filters are a good starting point

- Input / output signals $\mathbf{x}$ / $\mathbf{y}$ are graph signals supported on a common graph with shift operator $\mathbf{S}$

- Function class $\Rightarrow$ graph filters of order $K$ supported on $\mathbf{S}$ $\Rightarrow$ $\Phi(\mathbf{x}) = \sum_{k=0}^{K-1} h_k \mathbf{S}^k \mathbf{x} = \Phi(\mathbf{x}; \mathbf{S}, \mathbf{h})$

$$\mathbf{x} \longrightarrow \boxed{\mathbf{z} = \sum_{k=0}^{K-1} h_k \mathbf{S}^k \mathbf{x}} \longrightarrow \mathbf{z} = \Phi(\mathbf{x}; \mathbf{S}, \mathbf{h})$$

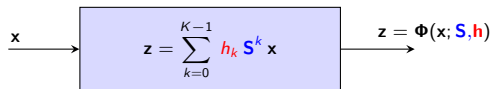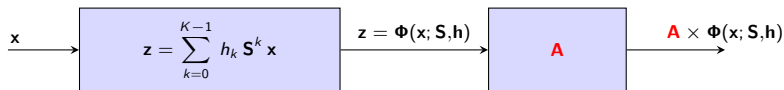- Learn ERM solution restricted to graph filter class $\Rightarrow$ $\mathbf{h}^* = \underset{\mathbf{h}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \ell\Big(\mathbf{y}, \Phi(\mathbf{x}; \mathbf{S}, \mathbf{h})\Big)$

  $\Rightarrow$ Optimization is over filter coefficients $\mathbf{h}$ with the graph shift operator $\mathbf{S}$ given

▶ Outputs $\mathbf{y} \in \mathbb{R}^m$ are not graph signals $\Rightarrow$ Add readout layer at filter's output to match dimensions

▶ Readout matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ yields parametrization $\Rightarrow$ $\mathbf{A} \times \Phi(\mathbf{x};\mathbf{S},\mathbf{h}) = \mathbf{A} \times \sum_{k=0}^{K-1} h_k \mathbf{S}^k \mathbf{x}$

$$\mathbf{x} \longrightarrow \boxed{\mathbf{z} = \sum_{k=0}^{K-1} h_k \mathbf{S}^k \mathbf{x}} \xrightarrow{\mathbf{z} = \Phi(\mathbf{x};\mathbf{S},\mathbf{h})} \boxed{\mathbf{A}} \xrightarrow{\mathbf{A} \times \Phi(\mathbf{x};\mathbf{S},\mathbf{h})}$$

▶ Making $\mathbf{A}$ trainable is inadvisable. Learn filter only. $\Rightarrow$ $\mathbf{h}^* = \underset{\mathbf{h}}{\operatorname{argmin}} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{T}} \ell\Big(\mathbf{y}, \mathbf{A} \times \Phi(\mathbf{x};\mathbf{S},\mathbf{h})\Big)$

▶ Readouts are simple. Read out node $i$ $\Rightarrow$ $\mathbf{A} = \mathbf{e}_i^T$. Read out signal average $\Rightarrow$ $\mathbf{A} = \mathbf{1}^T$.
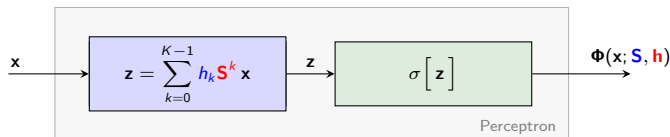
# Graph Neural Networks (GNNs)

▶ A pointwise nonlinearity is a nonlinear function applied componentwise. Without mixing entries

▶ The result of applying pointwise $\sigma$ to a vector $\mathbf{x}$ is $\Rightarrow \sigma \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \sigma(x_1) \\ \sigma(x_2) \\ \vdots \\ \sigma(x_n) \end{bmatrix}$

▶ A pointwise nonlinearity is the simplest nonlinear function we can apply to a vector

▶ ReLU: $\sigma(x) = \max(0, x)$. Hyperbolic tangent: $\sigma(x) = (e^{2x} - 1)/(e^{2x} + 1)$. Absolute value: $\sigma(x) = |x|$.

▶ Pointwise nonlinearities decrease variability. $\Rightarrow$ They function as demodulators.

▶ Graph filters have limited expressive power because they can only learn linear maps

▶ A first approach to nonlinear maps is the graph perceptron $\Rightarrow$ $\Phi(\mathbf{x}) = \sigma\left[\sum_{k=0}^{K-1} h_k \mathbf{S}^k \mathbf{x}\right] = \Phi(\mathbf{x}; \mathbf{S}, \mathbf{h})$



$$\sigma\left[\mathbf{x}\right] = \sigma\left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array}\right] = \left[\begin{array}{c} \sigma(x_1) \\ \sigma(x_2) \\ \vdots \\ \sigma(x_n) \end{array}\right]$$

▶ Optimal regressor restricted to perceptron class $\Rightarrow$ $\mathbf{h}^* = \underset{\mathbf{h}}{\mathrm{argmin}} \sum_{(\mathbf{x},\mathbf{y})\in\mathcal{T}} \ell\big(\mathbf{y}, \Phi(\mathbf{x}; \mathbf{S}, \mathbf{h})\big)$

$\Rightarrow$ Perceptron allows learning of nonlinear maps $\Rightarrow$ More expressive. Larger Representable Class

▶ To define a GNN we compose several graph perceptrons $\Rightarrow$ We layer graph perceptrons

▶ Layer 1 processes input signal $\mathbf{x}$ with the perceptron $\mathbf{h}_1 = [h_{10}, \ldots, h_{1,K-1}]$ to produce output $\mathbf{x}_1$

$$\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right] = \sigma\left[\sum_{k=0}^{K-1} h_{1k}\, \mathbf{S}^k\, \mathbf{x}\right]$$

▶ The Output of Layer 1 $\mathbf{x}_1$ becomes an input to Layer 2. Still $\mathbf{x}_1$ but with different interpretation

▶ Repeat analogous operations for $L$ times (the GNNs depth) $\Rightarrow$ Yields the GNN predicted output $\mathbf{x}_L$

▶ To define a GNN we compose several graph perceptrons $\Rightarrow$ We layer graph perceptrons

▶ Layer 2 processes its input signal $\mathbf{x}_1$ with the perceptron $\mathbf{h}_2 = [h_{20}, \ldots, h_{2,K-1}]$ to produce output $\mathbf{x}_2$

$$\mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right] = \sigma\left[\sum_{k=0}^{K-1} h_{2k} \mathbf{S}^k \mathbf{x}_1\right]$$

▶ The Output of Layer 2 $\mathbf{x}_2$ becomes an input to Layer 3. Still $\mathbf{x}_2$ but with different interpretation

▶ Repeat analogous operations for $L$ times (the GNNs depth) $\Rightarrow$ Yields the GNN predicted output $\mathbf{x}_L$

▶ A generic layer of the GNN, Layer $\ell$, takes as input the output $\mathbf{x}_{\ell-1}$ of the previous layer $(\ell-1)$

▶ Layer $\ell$ processes its input signal $\mathbf{x}_{\ell-1}$ with perceptron $\mathbf{h}_\ell = [h_{\ell 0}, \ldots, h_{\ell, K-1}]$ to produce output $\mathbf{x}_\ell$

$$\mathbf{x}_\ell = \sigma \left[ \mathbf{z}_\ell \right] = \sigma \left[ \sum_{k=0}^{K-1} h_{\ell k} \, \mathbf{S}^k \, \mathbf{x}_{\ell-1} \right]$$

▶ With the convention that the Layer 1 input is $\mathbf{x}_0 = \mathbf{x}$, this provides a recursive definition of a GNN

▶ If it has $L$ layers, the GNN output $\Rightarrow \mathbf{x}_L = \Phi\Big(\mathbf{x}; \mathbf{S}, \mathbf{h}_1, \ldots, \mathbf{h}_L\Big) = \Phi\Big(\mathbf{x}; \mathbf{S}, \mathcal{H}\Big)$

▶ The filter tensor $\mathcal{H} = [\mathbf{h}_1, \ldots, \mathbf{h}_L]$ is the trainable parameter. The graph shift is prior information
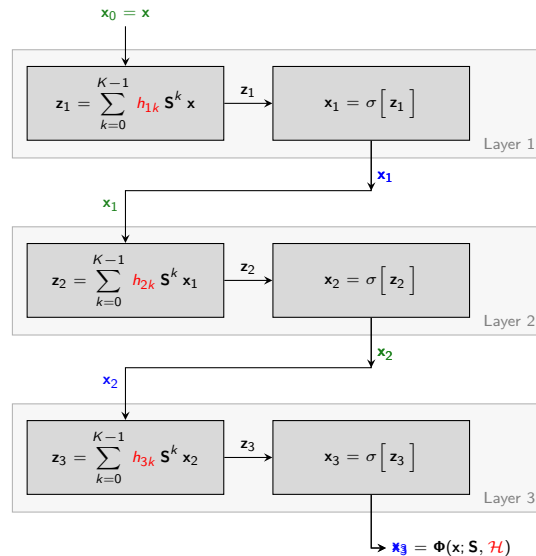
- Illustrate definition with a GNN with 3 layers

- Feed input signal $\mathbf{x} = \mathbf{x}_0$ into Layer 1

$$\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right] = \sigma\left[\sum_{k=0}^{K-1} h_{1k}\,\mathbf{S}^k\,\mathbf{x}_0\right]$$

- Last layer output is the GNN output $\Rightarrow \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

  $\Rightarrow$ Parametrized by filter tensor $\mathcal{H} = [\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3]$



$\mathbf{x}_0 = \mathbf{x}$

$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k}\,\mathbf{S}^k\,\mathbf{x} \quad \xrightarrow{\mathbf{z}_1} \quad \mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

Layer 1

$\mathbf{x}_1$

$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k}\,\mathbf{S}^k\,\mathbf{x}_1 \quad \xrightarrow{\mathbf{z}_2} \quad \mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

Layer 2

$\mathbf{x}_2$

$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k}\,\mathbf{S}^k\,\mathbf{x}_2 \quad \xrightarrow{\mathbf{z}_3} \quad \mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

Layer 3

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

- Illustrate definition with a GNN with 3 layers

- Feed Layer 1 output as an input to Layer 2

$$\mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right] = \sigma\left[\sum_{k=0}^{K-1} h_{2k}\,\mathbf{S}^k\,\mathbf{x}_1\right]$$

- Last layer output is the GNN output $\Rightarrow \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

  $\Rightarrow$ Parametrized by filter tensor $\mathcal{H} = [\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3]$
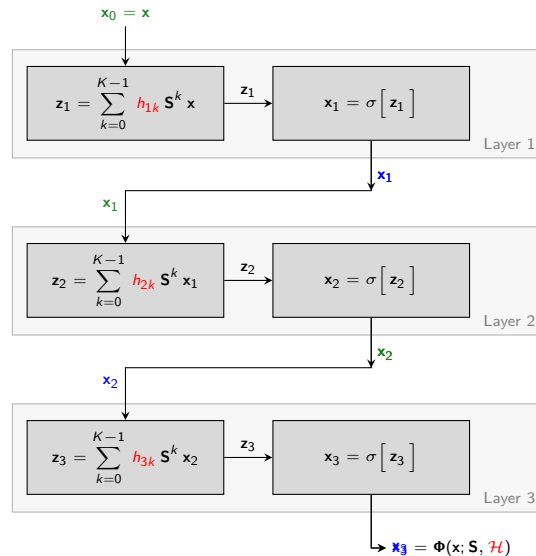
$\mathbf{x}_0 = \mathbf{x}$

$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k}\,\mathbf{S}^k\,\mathbf{x}$$  $\mathbf{z}_1$  $\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$

Layer 1

$\mathbf{x}_1$

$\mathbf{x}_1$

$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k}\,\mathbf{S}^k\,\mathbf{x}_1$$  $\mathbf{z}_2$  $\mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$

Layer 2

$\mathbf{x}_2$

$\mathbf{x}_2$

$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k}\,\mathbf{S}^k\,\mathbf{x}_2$$  $\mathbf{z}_3$  $\mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$

Layer 3

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

- Illustrate definition with a GNN with 3 layers

- Feed Layer 2 output as an input to Layer 3

$$\mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right] = \sigma\left[\sum_{k=0}^{K-1} h_{3k}\,\mathbf{S}^k\,\mathbf{x}_2\right]$$

- Last layer output is the GNN output $\Rightarrow \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

    $\Rightarrow$ Parametrized by filter tensor $\mathcal{H} = [\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3]$
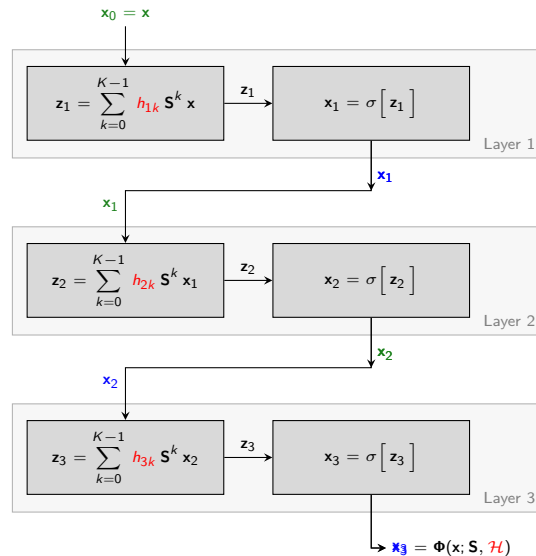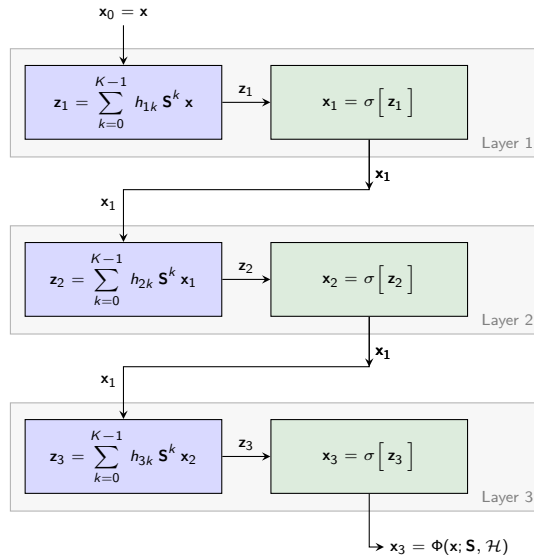


$\mathbf{x}_0 = \mathbf{x}$

Layer 1
$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k}\,\mathbf{S}^k\,\mathbf{x}$$
$\mathbf{z}_1$
$$\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

$\mathbf{x}_1$

Layer 2
$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k}\,\mathbf{S}^k\,\mathbf{x}_1$$
$\mathbf{z}_2$
$$\mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

$\mathbf{x}_2$

Layer 3
$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k}\,\mathbf{S}^k\,\mathbf{x}_2$$
$\mathbf{z}_3$
$$\mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

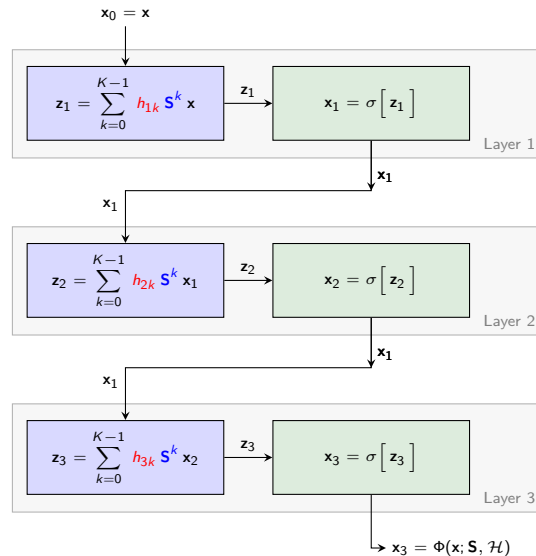Some Observations about Graph Neural Networks

- A GNN with $L$ layers follows $L$ recursions of the form

$$\mathbf{x}_\ell = \sigma\left[\mathbf{z}_\ell\right] = \sigma\left[\sum_{k=0}^{K-1} h_{\ell k}\,\mathbf{S}^k\,\mathbf{x}_{\ell-1}\right]$$

- A composition of $L$ layers. Each of which itself a...

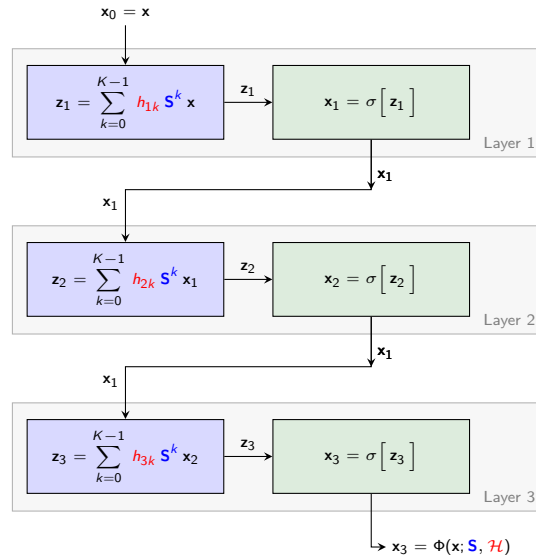  $\Rightarrow$ Compositions of Filters & Pointwise nonlinearities

- A GNN with $L$ layers follows $L$ recursions of the form

$$\mathbf{x}_\ell = \sigma\left[\mathbf{z}_\ell\right] = \sigma\left[\sum_{k=0}^{K-1} h_{\ell k}\,\mathbf{S}^k\,\mathbf{x}_{\ell-1}\right]$$

- Filters are parametrized by...

  $\Rightarrow$ Coefficients $h_{\ell k}$ and graph shift operators $\mathbf{S}$

$\mathbf{x}_0 = \mathbf{x}$

$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k}\,\mathbf{S}^k\,\mathbf{x} \qquad \mathbf{z}_1 \qquad \mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

Layer 1

$\mathbf{x}_1$

$\mathbf{x}_1$

$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k}\,\mathbf{S}^k\,\mathbf{x}_1 \qquad \mathbf{z}_2 \qquad \mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

Layer 2

$\mathbf{x}_1$

$\mathbf{x}_1$

$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k}\,\mathbf{S}^k\,\mathbf{x}_2 \qquad \mathbf{z}_3 \qquad \mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

Layer 3

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

- A GNN with $L$ layers follows $L$ recursions of the form

$$\mathbf{x}_\ell = \sigma\left[\mathbf{z}_\ell\right] = \sigma\left[\sum_{k=0}^{K-1} h_{\ell k}\, \mathbf{S}^k\, \mathbf{x}_{\ell-1}\right]$$

- Output $\mathbf{x}_L = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$ parametrized by...

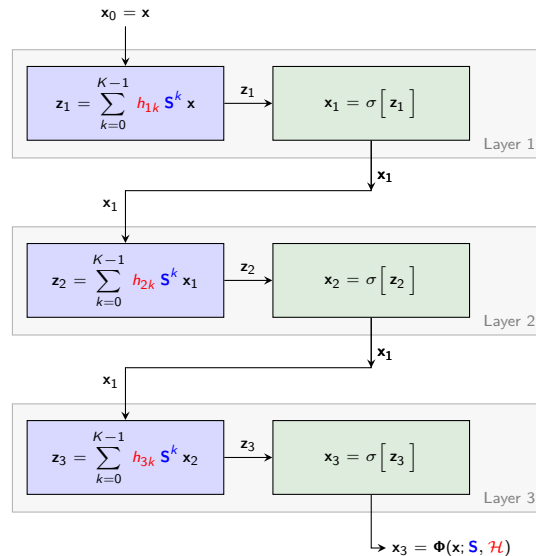  $\Rightarrow$ Learnable Filter tensor $\mathcal{H} = [\mathbf{h}_1, \ldots, \mathbf{h}_L]$



$$\mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k}\, \mathbf{S}^k\, \mathbf{x} \qquad \mathbf{z}_1 \qquad \mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

Layer 1

$$\mathbf{x}_1$$

$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k}\, \mathbf{S}^k\, \mathbf{x}_1 \qquad \mathbf{z}_2 \qquad \mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

Layer 2

$$\mathbf{x}_1$$

$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k}\, \mathbf{S}^k\, \mathbf{x}_2 \qquad \mathbf{z}_3 \qquad \mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

Layer 3

$$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$$

- Learn Optimal GNN tensor $\mathcal{H}^* = (\mathbf{h}_1^*, \mathbf{h}_2^*, \mathbf{h}_3^*)$ as
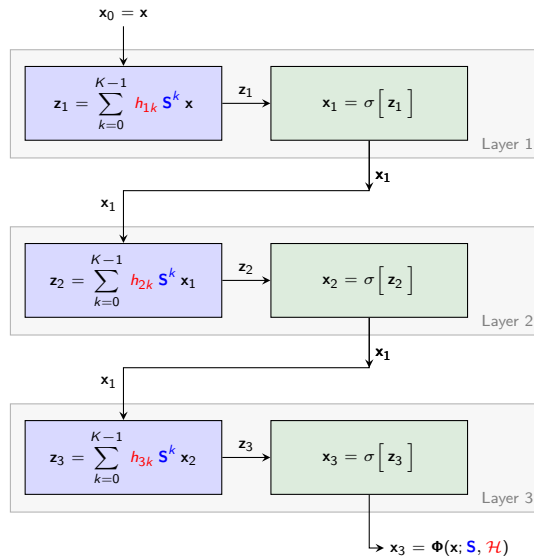
$$\mathcal{H}^* = \operatorname*{argmin}_{\mathcal{H}} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{T}} \ell\Big(\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}), \mathbf{y}\Big)$$

- Optimization is over tensor only. Graph $\mathbf{S}$ is given
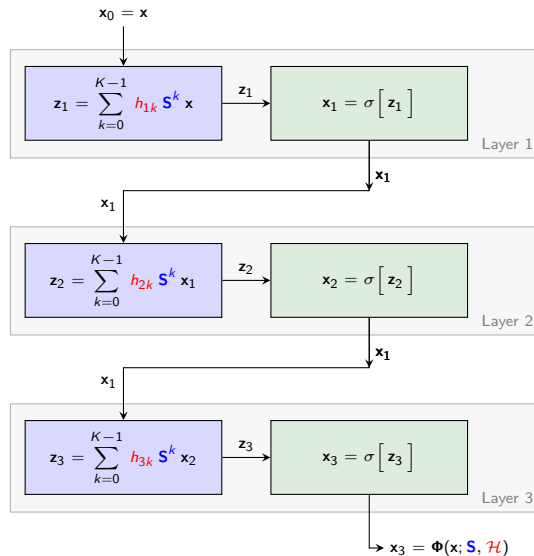
  $\Rightarrow$ Prior information given to the GNN

$\mathbf{x}_0 = \mathbf{x}$

$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k} \, \mathbf{S}^k \, \mathbf{x}$$

$\mathbf{z}_1$

$$\mathbf{x}_1 = \sigma\big[\mathbf{z}_1\big]$$

Layer 1

$\mathbf{x}_1$

$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k} \, \mathbf{S}^k \, \mathbf{x}_1$$

$\mathbf{z}_2$

$$\mathbf{x}_2 = \sigma\big[\mathbf{z}_2\big]$$

Layer 2

$\mathbf{x}_1$

$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k} \, \mathbf{S}^k \, \mathbf{x}_2$$

$\mathbf{z}_3$

$$\mathbf{x}_3 = \sigma\big[\mathbf{z}_3\big]$$

Layer 3

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

- ▶ GNNs are minor variations of graph filters

- ▶ Add pointwise nonlinearities and layer compositions

  ⇒ Nonlinearities process individual entries

  ⇒ Component mixing is done by graph filters only

- ▶ GNNs do work (much) better than graph filters

  ⇒ Which is unexpected and deserves explanation

  ⇒ Which we will attempt with stability analyses



$\mathbf{x}_0 = \mathbf{x}$

Layer 1:
$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k} \mathbf{S}^k \mathbf{x}$$
$\mathbf{z}_1$
$$\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

$\mathbf{x}_1$

Layer 2:
$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k} \mathbf{S}^k \mathbf{x}_1$$
$\mathbf{z}_2$
$$\mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

$\mathbf{x}_1$

Layer 3:
$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k} \mathbf{S}^k \mathbf{x}_2$$
$\mathbf{z}_3$
$$\mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$

- GNN Output depends on the graph $\mathbf{S}$.

- Interpret $\mathbf{S}$ as a parameter

  $\Rightarrow$ Encodes prior information. As we have done so far

$$\mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k}\, \mathbf{S}^k\, \mathbf{x} \qquad \mathbf{z}_1 \qquad \mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

Layer 1

$$\mathbf{x}_1$$

$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k}\, \mathbf{S}^k\, \mathbf{x}_1 \qquad \mathbf{z}_2 \qquad \mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

Layer 2

$$\mathbf{x}_1$$

$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k}\, \mathbf{S}^k\, \mathbf{x}_2 \qquad \mathbf{z}_3 \qquad \mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

Layer 3

$$\mathbf{x}_3 = \mathbf{\Phi}(\mathbf{x}; \mathbf{S}, \mathcal{H})$$

16

- But we can reinterpret $\mathbf{S}$ as an input of the GNN

    - $\Rightarrow$ Enabling transference across graphs

        $$\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}) \Rightarrow \Phi(\mathbf{x}; \tilde{\mathbf{S}}, \mathcal{H})$$

    - $\Rightarrow$ Same as we enable transference across signals

        $$\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}) \Rightarrow \Phi(\tilde{\mathbf{x}}; \mathbf{S}, \mathcal{H})$$

- A trained GNN is just a filter tensor $\mathcal{H}^*$



$$\mathbf{x}_0 = \mathbf{x}$$

Layer 1
$$\mathbf{z}_1 = \sum_{k=0}^{K-1} h_{1k} \mathbf{S}^k \mathbf{x} \quad \xrightarrow{\mathbf{z}_1} \quad \mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$$

$$\mathbf{x}_1$$

Layer 2
$$\mathbf{z}_2 = \sum_{k=0}^{K-1} h_{2k} \mathbf{S}^k \mathbf{x}_1 \quad \xrightarrow{\mathbf{z}_2} \quad \mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$$

$$\mathbf{x}_1$$

Layer 3
$$\mathbf{z}_3 = \sum_{k=0}^{K-1} h_{3k} \mathbf{S}^k \mathbf{x}_2 \quad \xrightarrow{\mathbf{z}_3} \quad \mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$$

$$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$$

- There is no difference between CNNs and GNNs

- To recover a CNN just particularize the shift operator the adjacency matrix of the directed line graph

$$S = \begin{bmatrix} & \vdots & \vdots & \vdots & \\ \cdots & 0 & 0 & 0 & \cdots \\ \cdots & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \end{bmatrix}$$



- GNNs are proper generalizations of CNNs

# Fully Connected Neural Networks

▶ We chose graph filters and graph neural networks (GNNs) because of our interest in graph signals

▶ We argued this is a good idea because they are generalizations of convolutional filters and CNNs

▶ We can explore this better if we go back to the road not taken ⇒ Fully connected neural networks

▶ Instead of graph filters, we choose arbitrary linear functions $\Rightarrow \Phi(\mathbf{x}) = \Phi(\mathbf{x}; \mathbf{H}) = \mathbf{H}\,\mathbf{x}$



$$\mathbf{x} \longrightarrow \boxed{\mathbf{z} = \mathbf{H}\,\mathbf{x}} \longrightarrow \mathbf{z} = \Phi(\mathbf{x}; \mathbf{H})$$

▶ Optimal regressor is ERM solution restricted to linear class $\Rightarrow \mathbf{H}^* = \underset{\mathbf{H}}{\mathrm{argmin}} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{T}} \ell\Big(\Phi(\mathbf{x}; \mathbf{H}), \mathbf{y}\Big)$

▶ We increase expressive power with the introduction of a perceptrons $\Rightarrow \boldsymbol{\Phi}(\mathbf{x}) = \boldsymbol{\Phi}(\mathbf{x}; \mathbf{H}) = \sigma\left[\mathbf{H}\mathbf{x}\right]$



▶ Optimal regressor restricted to perceptron class $\Rightarrow \mathbf{H}^* = \underset{\mathbf{H}}{\operatorname{argmin}} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{T}} \ell\Big(\boldsymbol{\Phi}(\mathbf{x}; \mathbf{H}), \mathbf{y}\Big)$

▶ A generic layer, Layer $\ell$ of a FCNN, takes as input the output $\mathbf{x}_{\ell-1}$ of the previous layer $(\ell-1)$

▶ Layer $\ell$ processes its input signal $\mathbf{x}_{\ell-1}$ with a linear perceptron $\mathbf{H}_\ell$ to produce output $\mathbf{x}_\ell$

$$\mathbf{x}_\ell = \sigma\left[\mathbf{z}_\ell\right] = \sigma\left[\mathbf{H}_\ell\,\mathbf{x}_{\ell-1}\right]$$

▶ With the convention that the Layer 1 input is $\mathbf{x}_0 = \mathbf{x}$, this provides a recursive definition of a GNN

▶ If it has $L$ layers, the FCNN output $\Rightarrow \mathbf{x}_L = \Phi\left(\mathbf{x};\mathbf{H}_1,\ldots,\mathbf{H}_L\right) = \Phi\left(\mathbf{x};\mathcal{H}\right)$

▶ The filter tensor $\mathcal{H} = [\mathbf{H}_1,\ldots,\mathbf{H}_L]$ is the trainable parameter.

# Fully Connected Neural Network Block Diagram

- Illustrate definition with an FCNN with 3 layers

- Feed input signal $\mathbf{x} = \mathbf{x}_0$ into Layer 1

$$\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right] = \sigma\left[\mathbf{H}_{1k}\,\mathbf{x}_0\right]$$

- Output $\Phi(\mathbf{x}; \mathcal{H})$ Parametrized by $\mathcal{H} = [\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3]$

$\mathbf{x}_0 = \mathbf{x}$

| | |
|---|---|
| $\mathbf{z}_1 = \mathbf{H}_1\,\mathbf{x}$ | $\mathbf{x}_1 = \sigma\left[\mathbf{z}_1\right]$ |

$\mathbf{z}_1$

Layer 1

$\mathbf{x}_1$

$\mathbf{x}_1$

| | |
|---|---|
| $\mathbf{z}_2 = \mathbf{H}_2\,\mathbf{x}_1$ | $\mathbf{x}_2 = \sigma\left[\mathbf{z}_2\right]$ |

$\mathbf{z}_2$

Layer 2

$\mathbf{x}_2$

$\mathbf{x}_2$

| | |
|---|---|
| $\mathbf{z}_3 = \mathbf{H}_3\,\mathbf{x}_2$ | $\mathbf{x}_3 = \sigma\left[\mathbf{z}_3\right]$ |

$\mathbf{z}_3$

Layer 3

$\mathbf{x}_3 = \Phi(\mathbf{x}; \mathcal{H})$

# Fully Connected Neural Network Block Diagram

- Illustrate definition with an FCNN with 3 layers

- Feed Layer 1 output as an input to Layer 2

$$\mathbf{x}_2 = \sigma \Big[ \mathbf{z}_2 \Big] = \sigma \Big[ \mathbf{H}_2 \, \mathbf{x}_1 \Big]$$

- Output $\Phi(\mathbf{x}; \mathcal{H})$ Parametrized by $\mathcal{H} = [\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3]$

# Fully Connected Neural Network Block Diagram

- Illustrate definition with an FCNN with 3 layers

- Feed Layer 2 output as an input to Layer 3

$$\mathbf{x}_3 = \sigma\Big[\mathbf{z}_3\Big] = \sigma\Big[\mathbf{H}_3\,\mathbf{x}_2\Big]$$

- Output $\Phi(\mathbf{x}; \mathcal{H})$ Parametrized by $\mathcal{H} = [\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3]$

Neural Networks vs Graph Neural Networks

▶ Since the GNN is a particular case of a fully connected NN, the latter attains a smaller cost

$$\min_{\mathcal{H}} \sum_{(\mathbf{x},\mathbf{y})\in\mathcal{T}} \ell\Big(\mathbf{\Phi}(\mathbf{x};\mathcal{H}),\mathbf{y}\Big) \leq \min_{\mathcal{H}} \sum_{(\mathbf{x},\mathbf{y})\in\mathcal{T}} \ell\Big(\mathbf{\Phi}(\mathbf{x};\mathbf{S},\mathcal{H}),\mathbf{y}\Big)$$

▶ The fully connected NN does better. But this holds for the training set

▶ In practice, the GNN does better because it generalizes better to unseen signals

⇒ Because it exploits internal symmetries of graph signals codified in the graph shift operator

▶ Suppose the graph represents a recommendation system where we want to fill empty ratings

▶ We observe ratings with the structure in the left. But we do not observe examples like the other two

▶ From examples like the one in the left, the NN learns how to fill the middle signal but not the right

▶ The GNN will succeed at predicting ratings for the signal on the right because it knows the graph

▶ The GNN still learns how to fill the middle signal. But it also learns how to fill the right signal

▶ The GNN exploits symmetries of the signal to effectively multiply available data

▶ This will be formalized later as the permutation equivariance of graph neural networks

# Graph Filter Banks

► Filters isolate features. When we are interested in multiple features, we use Banks of filters

▶ A graph filter bank is a collection of filters. Use $F$ to denote total number of filters in the bank

▶ Filter $f$ in the bank uses coefficients $\mathbf{h}^f = [h_1^f; \ldots; h_{K-1}^f] \Rightarrow$ Output $\mathbf{z}^f$ is a graph signal



▶ Filter bank output is a collection of $F$ graph signals $\Rightarrow$ Matrix graph signal $\mathbf{Z} = [\mathbf{z}^1, \ldots, \mathbf{z}^F]$

▶ The input of a filter bank is a single graph signal **x**. Rows of **x** are signals components $x_i$.

▶ Output matrix **Z** is a collection of signals $\mathbf{z}^f$. Rows of which are components $z_i^f$.

▶ Vector $\mathbf{z}_i$ supported at each node. Columns of **Z** are graph signals $\mathbf{z}^f$. Rows of **Z** are node features $\mathbf{z}_i$
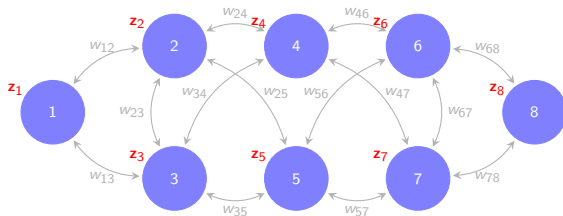
▶ The input of a filter bank is a single graph signal $\mathbf{x}$. Rows of $\mathbf{x}$ are signals components $x_i$.

▶ Output matrix $\mathbf{Z}$ is a collection of signals $\mathbf{z}^f$. Rows of which are components $z_i^f$.

▶ Vector $\mathbf{z}_i$ supported at each node. Columns of $\mathbf{Z}$ are graph signals $\mathbf{z}^f$. Rows of $\mathbf{Z}$ are node features $\mathbf{z}_i$



$$
\mathbf{Z} = \begin{bmatrix} z_1^1 & \cdots & z_1^f & \cdots & z_1^F \\ \vdots & & \vdots & & \vdots \\ z_i^1 & \cdots & z_i^f & \cdots & z_i^F \\ \vdots & & \vdots & & \vdots \\ z_n^1 & \cdots & z_n^f & \cdots & z_n^F \end{bmatrix} = \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_i \\ \vdots \\ \mathbf{z}_n \end{bmatrix}
$$

$$
= \begin{bmatrix} \mathbf{z}^1 & \cdots & \mathbf{z}^f & \cdots & \mathbf{z}^F \end{bmatrix}
$$

▶ The input of a filter bank is a single graph signal $\mathbf{x}$. Rows of $\mathbf{x}$ are signals components $x_i$.

▶ Output matrix $\mathbf{Z}$ is a collection of signals $\mathbf{z}^f$. Rows of which are components $z_i^f$.

▶ Vector $\mathbf{z}_i$ supported at each node. Columns of $\mathbf{Z}$ are graph signals $\mathbf{z}^f$. Rows of $\mathbf{Z}$ are node features $\mathbf{z}_i$



$$
\mathbf{Z} = \begin{bmatrix} z_1^1 & \cdots & z_1^f & \cdots & z_1^F \\ \vdots & & \vdots & & \vdots \\ z_i^1 & \cdots & z_i^f & \cdots & z_i^F \\ \vdots & & \vdots & & \vdots \\ z_n^1 & \cdots & z_n^f & \cdots & z_n^F \end{bmatrix} = \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_i \\ \vdots \\ \mathbf{z}_n \end{bmatrix}
$$

$$
= \begin{bmatrix} \mathbf{z}^1 & \cdots & \mathbf{z}^f & \cdots & \mathbf{z}^F \end{bmatrix}
$$

**Theorem (Output Energy of a Graph Filter)**

Consider graph filter $\mathbf{h}$ with coefficients $h_k$ and frequency response $\tilde{h}(\lambda) = \sum\limits_{k=0}^{\infty} h_k \lambda^k$ . The energy of the filter's output $\mathbf{z} = \sum\limits_{k=0}^{\infty} h_k \mathbf{S}^k \mathbf{x}$ is given by

$$\| \mathbf{z} \|^2 = \sum_{i=1}^{n} \left( \tilde{h}(\lambda_i) \, \tilde{x}_i \right)^2$$

where $\lambda_i$ are eigenvalues of symmetric $\mathbf{S}$ and $\tilde{x}_i$ are components of the GFT of $\mathbf{x}$, $\tilde{\mathbf{x}} = \mathbf{V}^H \mathbf{x}$ is

**Proof:** The GFT is a unitary transform that preserves energy. Indeed, with $\tilde{\mathbf{z}} = \mathbf{V}^H \mathbf{z}$ we have

$$\left\| \tilde{\mathbf{z}} \right\|^2 = \tilde{\mathbf{z}}^H \tilde{\mathbf{z}} = \left( \mathbf{V}^H \mathbf{z} \right)^H \left( \mathbf{V}^H \mathbf{z} \right) = \mathbf{z}^H \mathbf{V} \mathbf{V}^H \mathbf{z} = \mathbf{z}^H \mathbf{I} \mathbf{z} = \left\| \mathbf{z} \right\|^2$$
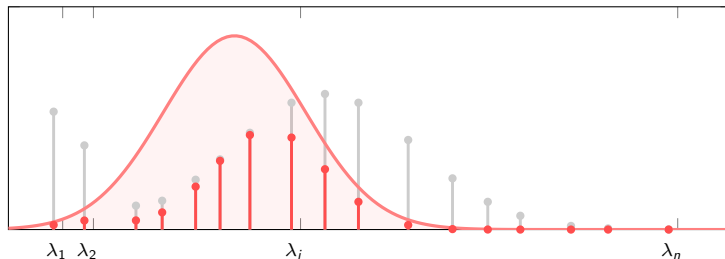
▶ We know that graph filters are pointwise in the frequency domain $\Rightarrow \tilde{z}_i = \tilde{h}(\lambda_i) \tilde{x}_i$

$$\left\| \tilde{\mathbf{z}} \right\|^2 = \tilde{\mathbf{z}}^H \tilde{\mathbf{z}} = \sum_{i=1}^{n} \tilde{z}_i^2 = \sum_{i=1}^{n} \left( \tilde{h}^f(\lambda_i) \tilde{x}_i \right)^2$$

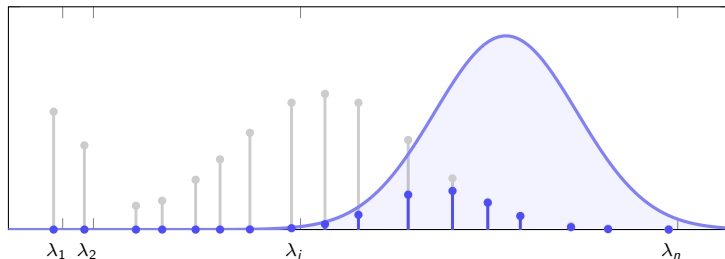▶ We have the energy expressed in the form we want. Except that it is in the frequency domain.

▶ But we have just seen the GFT preserves energy $\Rightarrow \left\| \mathbf{z} \right\|^2 = \left\| \tilde{\mathbf{z}} \right\|^2 = \sum_{i=1}^{n} \left( \tilde{h}(\lambda_i) \tilde{x}_i \right)^2$ ■

▶ The energy that graph filters let pass is a sort of "area under the frequency response curve."

▶ Graph Filter banks are helpful in identifying frequency signatures of different signals



▶ Filter banks scatter the energy of signal $\mathbf{x}$ into the signals $\mathbf{z}^f$ at the output of the filters.

⇒ Different signals concentrate energy on different outputs $\mathbf{z}^f$

▶ The energy that graph filters let pass is a sort of "area under the frequency response curve."

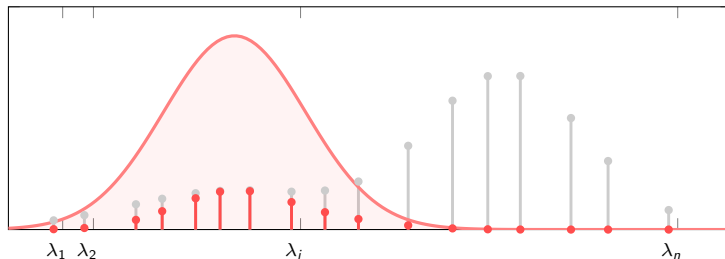▶ Graph Filter banks are helpful in identifying frequency signatures of different signals



▶ Filter banks scatter the energy of signal $\mathbf{x}$ into the signals $\mathbf{z}^f$ at the output of the filters.

⇒ Different signals concentrate energy on different outputs $\mathbf{z}^f$

▶ The energy that graph filters let pass is a sort of "area under the frequency response curve."

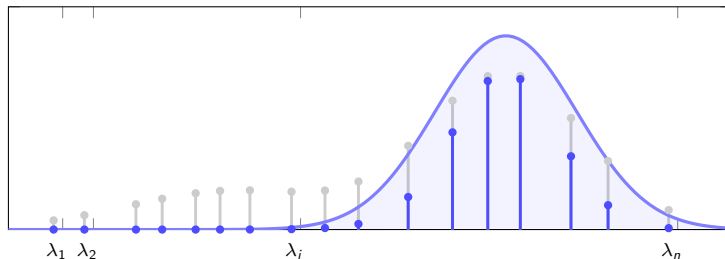▶ Graph Filter banks are helpful in identifying frequency signatures of different signals



▶ Filter banks scatter the energy of signal $\mathbf{x}$ into the signals $\mathbf{z}^f$ at the output of the filters.

⇒ Different signals concentrate energy on different outputs $\mathbf{z}^f$
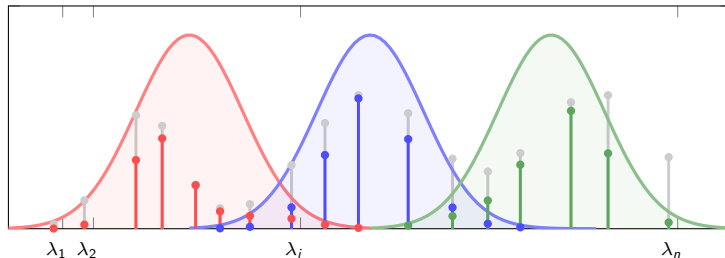
▶ The energy that graph filters let pass is a sort of "area under the frequency response curve."

▶ Graph Filter banks are helpful in identifying frequency signatures of different signals



▶ Filter banks scatter the energy of signal $\mathbf{x}$ into the signals $\mathbf{z}^f$ at the output of the filters.

⇒ Different signals concentrate energy on different outputs $\mathbf{z}^f$

► The filter bank isolates groups of frequency components

$\Rightarrow$ Energy of bank output $\mathbf{z}^f = \displaystyle\sum_{k=0}^{\infty} h_k^f \mathbf{S}^k \mathbf{x}$ is area under the curve $\Rightarrow \left\| \mathbf{z}^f \right\|^2 = \displaystyle\sum_{i=1}^{n} \left( \tilde{h}^f(\lambda_i) \, \tilde{x}_i \right)^2$



► We use the filter bank to identify signals with different spectral signatures.

▶ The GFT preserves energy $\Rightarrow$ It scatters information. But it doesn't loose information

▶ A filter bank is a frame if there exist constants $m \leq M \Rightarrow m\|\mathbf{x}\|^2 \leq \sum_{f=1}^{F} \|\mathbf{z}^f\|^2 \leq M\|\mathbf{x}\|^2$

▶ A filter banks is a tight frame if $m = M = 1 \Rightarrow \|\mathbf{x}\|^2 = \sum_{f=1}^{F} \|\mathbf{z}^f\|^2$

▶ No signal is vanquished by a frame. Energy is preserved by a tight frame

▶ Because filters are pointwise in the GFT domain, a frame must satisfy $\Rightarrow m \leq \sum_{f=1}^{F} \left[ \tilde{h}^f(\lambda) \right]^2 \leq M$

▶ All frequencies $\lambda$ must have at least one filter $\mathbf{h}^f$ with response $m \leq \left[ \tilde{h}^f(\lambda) \right]^2$



$\lambda_1 \ \lambda_2 \qquad\qquad \lambda_i \qquad\qquad\qquad \lambda_n$

▶ Likewise, a tight frame must be such that for all $\lambda \Rightarrow \sum_{f=1}^{F} \left[ \tilde{h}^f(\lambda) \right]^2 = 1$

▶ A Sufficient condition is that all frequencies accumulate unit energy when summing across all filters



▶ We will not design filter banks. We will learn them. But keeping them close to frames is good.

## Multiple Feature GNNs

► We leverage filter banks to create GNNs that process multiple features per layer

▶ Filter banks output a collection of multiple graph signals $\Rightarrow$ A matrix graph signal $\mathbf{Z} = [\mathbf{z}^1, \ldots, \mathbf{z}^F]$

▶ The $F$ graph signals $\mathbf{z}^f$ represent $F$ features per node. A vector $\mathbf{z}_i$ supported at each node



▶ We would now like to process multiple feature graph signals. Process each feature with a filterbank.

▶ Filter banks output a collection of multiple graph signals $\Rightarrow$ A matrix graph signal $\mathbf{Z} = [\mathbf{z}^1, \ldots, \mathbf{z}^F]$

▶ The $F$ graph signals $\mathbf{z}^f$ represent $F$ features per node. A vector $\mathbf{z}_i$ supported at each node



▶ We would now like to process multiple feature graph signals. Process each feature with a filterbank.

▶ Filter banks output a collection of multiple graph signals $\Rightarrow$ A matrix graph signal $\mathbf{Z} = [\mathbf{z}^1, \ldots, \mathbf{z}^F]$

▶ The $F$ graph signals $\mathbf{z}^f$ represent $F$ features per node. A vector $\mathbf{z}_i$ supported at each node



▶ We would now like to process multiple feature graph signals. Process each feature with a filterbank.

► Each of the $F$ features $\mathbf{x}^f$ is processed with $G$ filters with coefficients $h_k^{fg}$ $\Rightarrow$ $\mathbf{u}^{fg} = \sum_{k=0}^{K-1} h_k^{fg} \mathbf{S}^k \mathbf{x}^f$

▶ This Multiple-Input-Multiple-Output Graph Filter generates an output with $F \times G$ features



$$z^1 = u^{11} + u^{21} + \ldots + u^{F1} \qquad z^2 = u^{12} + u^{22} + \ldots + u^{F2} \qquad z^2 = u^{1G} + u^{2G} + \ldots + u^{FG}$$

▶ Reduce to $G$ outputs with sum over input features for given $g$ $\Rightarrow$ $\mathbf{z}^g = \sum_{f=1}^{F} \mathbf{u}^{fg} = \sum_{f=1}^{F} \sum_{k=0}^{K-1} h_k^{fg} \, \mathbf{S}^k \, \mathbf{x}^f$



$z^1 = \mathbf{u}^{11} + \mathbf{u}^{21} + \ldots + \mathbf{u}^{F1}$
$\qquad$ $z^2 = \mathbf{u}^{12} + \mathbf{u}^{22} + \ldots + \mathbf{u}^{F2}$
$\qquad$ $z^2 = \mathbf{u}^{1G} + \mathbf{u}^{2G} + \ldots + \mathbf{u}^{FG}$

▶ MIMO graph filters are cumbersome, not difficult. Just $F \times G$ filters. Or $F$ filter banks.

▶ Easier with matrices $\Rightarrow$ $G \times F$ coefficient matrix $\mathbf{H}_k$ with entries $\left( \mathbf{H}_k \right)_{fg} = h_k^{fg}$

$$\mathbf{Z} = \sum_{k=0}^{K-1} \mathbf{S}^k \times \mathbf{X} \times \mathbf{H}_k$$

▶ This is a more compact format of the MIMO filter. It is equivalent

$$\begin{bmatrix} \mathbf{z}^1 & \cdot\cdot & \mathbf{z}^g & \cdot\cdot & \mathbf{z}^G \end{bmatrix} = \sum_{k=0}^{K-1} \mathbf{S}^k \times \begin{bmatrix} \mathbf{x}^1 & \cdot\cdot & \mathbf{x}^f & \cdot\cdot & \mathbf{x}^F \end{bmatrix} \times \begin{bmatrix} h_k^{11} & \cdot\cdot & h_k^{1g} & \cdot\cdot & h_k^{1G} \\ & & \vdots & & \\ h_k^{f1} & \cdot\cdot & h_k^{fg} & \cdot\cdot & h_k^{fG} \\ & & \vdots & & \\ h_k^{F1} & \cdot\cdot & h_k^{Fg} & \cdot\cdot & h_k^{FG} \end{bmatrix}$$

► MIMO GNN stacks MIMO perceptrons $\Rightarrow$ Compose of MIMO filters with pointwise nonlinearities

► Layer $\ell$ processes input signal $\mathbf{X}_{\ell-1}$ with perceptron $\mathbf{H}_\ell = [\mathbf{H}_{\ell 0}, \ldots, \mathbf{H}_{\ell, K-1}]$ to produce output $\mathbf{X}_\ell$

$$\mathbf{X}_\ell = \sigma\Big[\mathbf{Z}_\ell\Big] = \sigma\left[\sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_{\ell-1} \mathbf{H}_{\ell k}\right]$$

► Denoting the Layer 1 input as $\mathbf{X}_0 = \mathbf{X}$, this provides a recursive definition of a MIMO GNN

► If it has $L$ layers, the GNN output $\Rightarrow \mathbf{X}_L = \Phi\Big(\mathbf{x}; \mathbf{S}, \mathbf{H}_1, \ldots, \mathbf{H}_L\Big) = \Phi\Big(\mathbf{x}; \mathbf{S}, \mathcal{H}\Big)$

► The filter tensor $\mathcal{H} = [\mathbf{H}_1, \ldots, \mathbf{H}_L]$ is the trainable parameter. The graph shift is prior information

▶ We illustrate with a MIMO GNN with 3 layers

▶ Feed input signal $\mathbf{X} = \mathbf{X}_0$ into Layer 1 ($F_0$ features)

$$\mathbf{X}_1 = \sigma\left[\mathbf{Z}_1\right] = \sigma\left[\sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_0 \mathbf{H}_{1k}\right]$$

▶ Last layer output is the GNN output $\Rightarrow \Phi(\mathbf{X}; \mathbf{S}, \mathcal{H})$

$\Rightarrow$ Parametrized by trainable tensor $\mathcal{H} = [\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3]$



$\mathbf{x}_0 = \mathbf{x}$

**Layer 1**
$\mathbf{Z}_1 = \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X} \mathbf{H}_{1k}$ $\quad \mathbf{Z}_1 \quad$ $\mathbf{X}_1 = \sigma\left[\mathbf{Z}_1\right]$

$\mathbf{x}_1$

**Layer 2**
$\mathbf{Z}_2 = \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_1 \mathbf{H}_{2k}$ $\quad \mathbf{Z}_2 \quad$ $\mathbf{X}_2 = \sigma\left[\mathbf{Z}_2\right]$

$\mathbf{x}_2$

**Layer 3**
$\mathbf{Z}_3 = \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_2 \mathbf{H}_{3k}$ $\quad \mathbf{Z}_3 \quad$ $\mathbf{X}_3 = \sigma\left[\mathbf{Z}_3\right]$

$\mathbf{x}_3 = \Phi(\mathbf{X}; \mathbf{S}, \mathcal{H})$

- We illustrate with a MIMO GNN with 3 layers

- Feed Layer 1 output as an input to Layer 2 ($F_1$ features)

$$\mathbf{X}_2 = \sigma\left[\mathbf{Z}_2\right] = \sigma\left[\sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_1 \mathbf{H}_{2k}\right]$$

- Last layer output is the GNN output $\Rightarrow \Phi(\mathbf{X}; \mathbf{S}, \mathcal{H})$

  $\Rightarrow$ Parametrized by trainable tensor $\mathcal{H} = [\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3]$

$\mathbf{x}_0 = \mathbf{x}$

Layer 1:
$$\mathbf{Z}_1 = \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X} \mathbf{H}_{1k}$$
$\mathbf{Z}_1$
$$\mathbf{X}_1 = \sigma\left[\mathbf{Z}_1\right]$$
$\mathbf{X}_1$

Layer 2:
$\mathbf{X}_1$
$$\mathbf{Z}_2 = \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_1 \mathbf{H}_{2k}$$
$\mathbf{Z}_2$
$$\mathbf{X}_2 = \sigma\left[\mathbf{Z}_2\right]$$
$\mathbf{X}_2$

Layer 3:
$\mathbf{X}_2$
$$\mathbf{Z}_3 = \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_2 \mathbf{H}_{3k}$$
$\mathbf{Z}_3$
$$\mathbf{X}_3 = \sigma\left[\mathbf{Z}_3\right]$$

$\mathbf{X}_3 = \Phi(\mathbf{X}; \mathbf{S}, \mathcal{H})$

- We illustrate with a MIMO GNN with 3 layers

- Feed Layer 2 output ($F_2$ features) as an input to Layer 3

$$\mathbf{X}_3 = \sigma \Big[ \mathbf{Z}_3 \Big] = \sigma \left[ \sum_{k=0}^{K-1} \mathbf{S}^k \, \mathbf{X}_2 \, \mathbf{H}_{3k} \right]$$
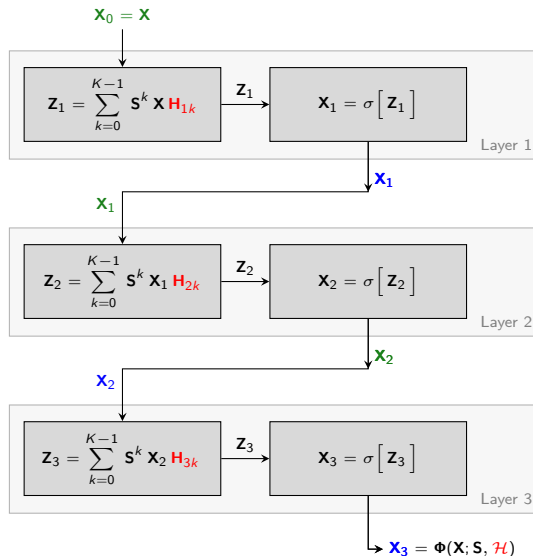
- Last layer output is the GNN output $\Rightarrow \Phi(\mathbf{X}; \mathbf{S}, \mathcal{H})$

  $\Rightarrow$ Parametrized by trainable tensor $\mathcal{H} = [\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3]$

Diagram:

$\mathbf{x}_0 = \mathbf{x}$

Layer 1: $\mathbf{z}_1 = \sum_{k=0}^{K-1} \mathbf{S}^k \, \mathbf{X} \, \mathbf{H}_{1k}$ $\xrightarrow{\mathbf{z}_1}$ $\mathbf{x}_1 = \sigma \Big[ \mathbf{z}_1 \Big]$

$\mathbf{x}_1$

Layer 2: $\mathbf{z}_2 = \sum_{k=0}^{K-1} \mathbf{S}^k \, \mathbf{x}_1 \, \mathbf{H}_{2k}$ $\xrightarrow{\mathbf{z}_2}$ $\mathbf{x}_2 = \sigma \Big[ \mathbf{z}_2 \Big]$

$\mathbf{x}_2$

Layer 3: $\mathbf{z}_3 = \sum_{k=0}^{K-1} \mathbf{S}^k \, \mathbf{x}_2 \, \mathbf{H}_{3k}$ $\xrightarrow{\mathbf{z}_3}$ $\mathbf{x}_3 = \sigma \Big[ \mathbf{z}_3 \Big]$

$\mathbf{x}_3 = \Phi(\mathbf{X}; \mathbf{S}, \mathcal{H})$