# Graph Signal Processing

Alejandro Ribeiro

Dept. of Electrical and Systems Engineering

University of Pennsylvania

Email: aribeiro@seas.upenn.edu

Web: alelab.seas.upenn.edu

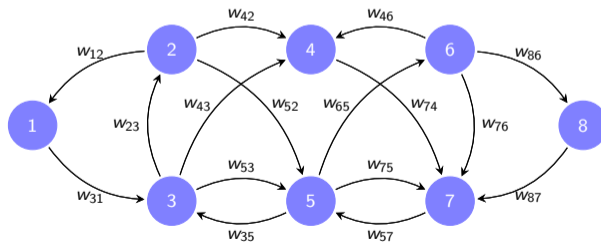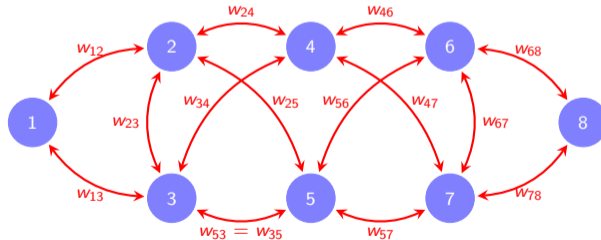May 17, 2021

# Graphs and Graph Signals

▶ A graph is a triplet $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$, which includes vertices $\mathcal{V}$, edges $\mathcal{E}$, and weights $\mathcal{W}$

⇒ Vertices or nodes are a set of n labels. Typical labels are $\mathcal{V} = \{1, \ldots, n\}$

⇒ Edges are ordered pairs of labels $(i, j)$. We interpret $(i, j) \in \mathcal{E}$ as "i can be influenced by j."

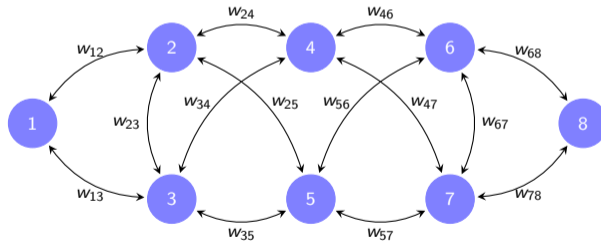⇒ Weights $w_{ij} \in \mathbb{R}$ are numbers associated to edges $(i, j)$. "Strength of the influence of j on i."

▶ A graph is symmetric or undirected if both, the edge set and the weight are symmetric

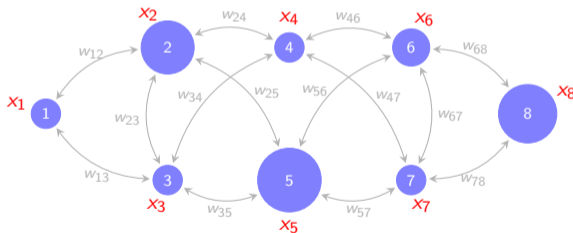⇒ Edges come in pairs ⇒ We have $(i, j) \in \mathcal{E}$ if and only if $(j, i) \in \mathcal{E}$

⇒ Weights are symmetric ⇒ We must have $w_{ij} = w_{ji}$ for all $(i, j) \in \mathcal{E}$

► Graphs can be directed or symmetric. Separately, they can be weighted or unweighted.

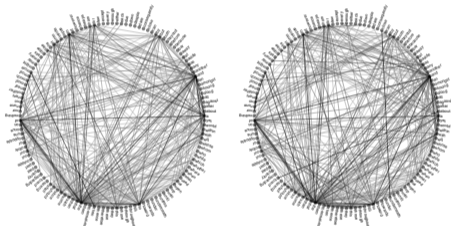► Most of the graphs we encounter in practical situations are symmetric and weighted

▶ Consider a given graph $\mathcal{G}$ with $n$ nodes, edge set $\mathcal{E}$ and weights $\mathcal{W}$

▶ A graph signal is a vector $x \in \mathbb{R}^n$ in which component $x_i$ is associated with node $i$

▶ To emphasize that the graph is intrinsic to the signal we may write the signal as a pair $\Rightarrow (\mathcal{G}, x)$



▶ The graph is an expectation of proximity or similarity between components of the signal $x$

▶ Graphs are generic models of signal structure that can help to learn in several practical problems
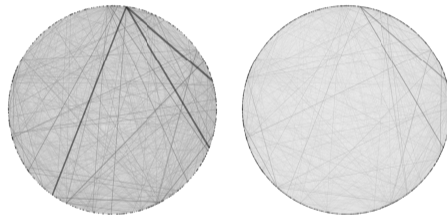
Authorship Attribution



Identify the author of a text of unknown provenance

Segarra et al '16, arxiv.org/abs/1805.00165

Recommendation Systems
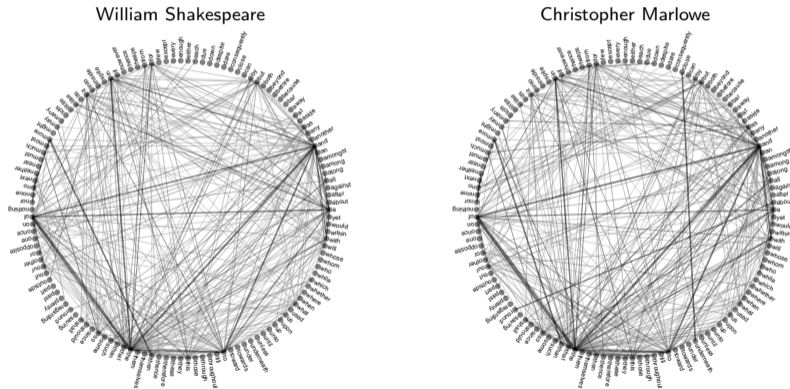


Predict the rating a customer would give to a product

Ruiz et al '18, arxiv.org/abs/1903.12575

▶ In both cases there exists a graph that contains meaningful information about the problem to solve

▶ **Nodes** represent different **function words** and **edges how often words appear close** to each other

⇒ A proxy for the different ways in which different authors use the English language grammar
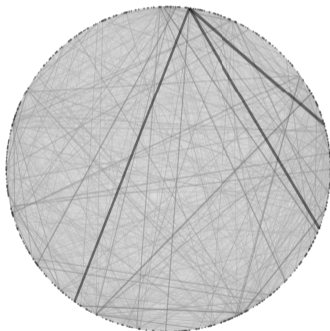
William Shakespeare

Christopher Marlowe



▶ WAN differences **differentiate the writing styles of Marlowe and Shakespeare** in, e.g., Henry VI

▶ **Nodes** represent different **customers** and **edges** their average **similarity in product ratings**

$\Rightarrow$ The graph informs the completion of ratings when some are unknown and are to be predicted

Variation Diagram for Original (sampled) ratings

Variation Diagram for Reconstructed (predicted) ratings



▶ Variation energy of reconstructed signal is (much) smaller than variation energy of sampled signal

Ruiz-Gama-Marques-Ribeiro, *Invariance-Preserving Localized Activation Functions for Graph Neural Networks*, arxiv.org/abs/1903.12575

▶ Graphs are more than data structures ⇒ They are models of physical systems with multiple agents

Decentralized Control of Autonomous Systems

Wireless Communications Networks



Coordinate a team of agents without central coordination

Tolstaya et al '19, arxiv.org/abs/1903.10527

Manage interference when allocating bandwidth and power

Eisen-Ribeiro '19, arxiv.org/abs/1909.01865

▶ The graph is the source of the problem ⇒ Challenge is that goals are global but information is local

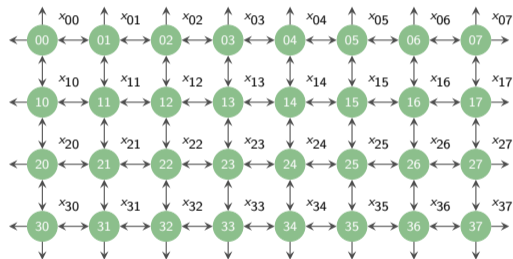▶ We can describe discrete time and space using graphs that support time or space signals

Description of time with a directed line graph

Description of images (space) with a grid graph



▶ Line graph represents adjacency of points in time. Grid graph represents adjacency of points in space

▶ A covariance matrix $\Sigma$ with entries $((\Sigma))_{ij} = \sigma_{ij}$ is also representable with a graph

⇒ One that has self loops to represent the variances $\sigma_{ii}$

▶ A realization $x$ of a random signal X is a signal supported on the covariance matrix graph

► Time and Space are pervasive and important, but still a (very) limited class of signals

► Use graphs as generic descriptors of signal structure with signal values associated to nodes and edges expressing expected similarity between signal components

A signal supported on a graph

Another signal supported on another graph



► Nodes are customers. Signal values are product ratings. Edges are cosine similarities of past scores

▶ Time and Space are pervasive and important, but still a (very) limited class of signals

▶ Use graphs as generic descriptors of signal structure with signal values associated to nodes and edges expressing expected similarity between signal components

A signal supported on a graph

Another signal supported on another graph



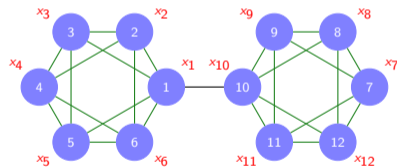▶ Nodes are drones. Signal values are velocities. Edges are sensing and communication ranges

▶ Time and Space are pervasive and important, but still a (very) limited class of signals

▶ Use graphs as generic descriptors of signal structure with signal values associated to nodes and edges expressing expected similarity between signal components

A signal supported on a graph

Another signal supported on another graph



▶ Nodes are transceivers. Signal values are QoS requirements. Edges are wireless channels strength

▶ Time and Space are pervasive and important, but still a (very) limited class of signals

▶ Use graphs as generic descriptors of signal structure with signal values associated to nodes and edges expressing expected similarity between signal components
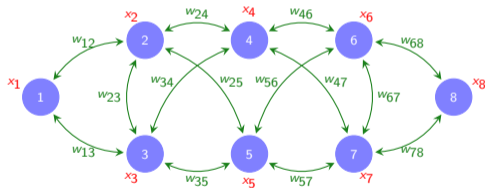
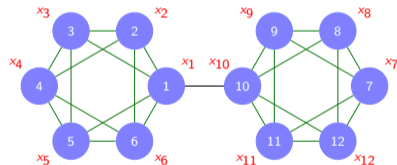A signal supported on a graph

Another signal supported on another graph



▶ Nodes are points in time. Signal values. Edges denote time causality

► Time and Space are pervasive and important, but still a (very) limited class of signals

► Use graphs as generic descriptors of signal structure with signal values associated to nodes and edges expressing expected similarity between signal components

A signal supported on a graph

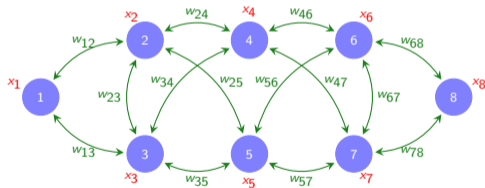Another signal supported on another graph



► Nodes are pixels. Signal values are luminances. Edges denote spatial proximities

- Time and Space are pervasive and important, but still a (very) limited class of signals

- Use graphs as generic descriptors of signal structure with signal values associated to nodes and edges expressing expected similarity between signal components

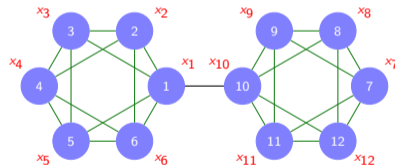A signal supported on a graph

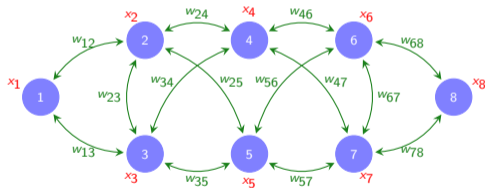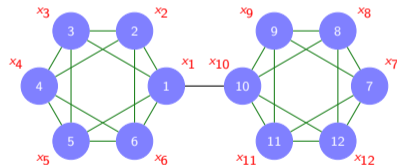Another signal supported on another graph



- Nodes are entries. Signal values. Edges denote crosscovariances

► Techniques to process signals on graphs that...

   ⇒ Generalize techniques developed for time, space, and random signals

   ⇒ Recover techniques developed for time, space, and random signals as particular cases

► Graph Fourier transform ⇒ Recovers DFT, 2D-DFT and PCA as particular cases

► Graph Convolutional Filters ⇒ Recovers time and spatial convolutions as particular cases

# Graph Shift Operators

▶ Graphs have matrix representations. Which in this course, we call graph shift operators (GSOs)

▶ The adjacency matrix of graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$ is the sparse matrix $A$ with nonzero entries

$$A_{ij} = w_{ij}, \text{ for all } (i,j) \in \mathcal{E}$$

▶ If the graph is symmetric, the adjacency matrix is symmetric $\Rightarrow A = A^T$. As in the example



$$A = \begin{bmatrix} 0 & w_{12} & w_{13} & 0 & 0 \\ w_{21} & 0 & w_{23} & w_{24} & 0 \\ w_{31} & w_{32} & 0 & 0 & w_{35} \\ 0 & w_{42} & 0 & 0 & w_{45} \\ 0 & 0 & w_{53} & w_{54} & 0 \end{bmatrix}.$$

▶ For the particular case in which the graph is <span style="color:red">unweighted</span>. Weights interpreted as units

$$A_{ij} = 1, \qquad \text{for all} \qquad (i,j) \in \mathcal{E}$$



$$A = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}.$$

▶ The neighborhood of node $i$ is the set of nodes that influence $i$ $\Rightarrow$ $n(i) := \{j : (i,j) \in \mathcal{E}\}$

▶ Degree $d_i$ of node $i$ is the sum of the weights of its incident edges $\Rightarrow$ $d_i = \sum_{j \in n(i)} w_{ij} = \sum_{j:(i,j) \in \mathcal{E}} w_{ij}$



▶ Node 1 neighborhood $\Rightarrow$ $n(1) = \{2, 3\}$

▶ Node 1 degree $\Rightarrow$ $n(1) = w_{12} + w_{13}$

▶ The degree matrix is a diagonal matrix D with degrees as diagonal entries $\Rightarrow D_{ii} = d_i$

▶ Write in terms of adjacency matrix as $D = \text{diag}(A1)$. Because $(A1)_i = \sum_j w_{ij} = d_i$



$$D = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix}$$

▶ The Laplacian matrix of a graph with adjacency matrix A is ⇒ $L = D - A = \text{diag}(A1) - A$

▶ Can also be written explicitly in terms of graph weights $A_{ij} = w_{ij}$

⇒ Off diagonal entries ⇒ $L_{ij} = -A_{ij} = -w_{ij}$

⇒ Diagonal entries ⇒ $L_{ii} = d_i = \sum_{j \in n(i)} w_{ij}$

$$L = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 \\ -1 & 3 & -1 & -1 & 0 \\ -1 & -1 & 3 & 0 & -1 \\ 0 & -1 & 0 & 2 & -1 \\ 0 & 0 & -1 & -1 & 2 \end{bmatrix}$$

▶ Normalized adjacency and Laplacian matrices express weights relative to the nodes' degrees

▶ Normalized adjacency matrix $\Rightarrow \bar{A} := D^{-1/2} A D^{-1/2} \Rightarrow$ Results in entries $(\bar{A})_{ij} = \dfrac{w_{ij}}{\sqrt{d_i d_j}}$

▶ The normalized adjacency is symmetric if the graph is symmetric $\Rightarrow \bar{A}^T = \bar{A}$.

▶ **Normalized Laplacian** matrix $\Rightarrow \bar{L} := D^{-1/2}LD^{-1/2}$. Same normalization of adjacency matrix

▶ Given definitions normalized representations $\Rightarrow \bar{L} = D^{-1/2}\Big(D - A\Big)D^{-1/2} = I - \bar{A}$

$\Rightarrow$ The normalized Laplacian and adjacency are essentially the same linear transformation.

▶ Normalized operators are more homogeneous. The entries in the vector A1 tend to be similar.

▶ The Graph Shift Operator S is a stand in for any of the matrix representations of the graph

Adjacency Matrix $\qquad$ Laplacian Matrix $\qquad$ Normalized Adjacency $\qquad$ Normalized Laplacian

$$S = A \qquad\qquad S = L \qquad\qquad S = \bar{A} \qquad\qquad S = \bar{L}$$

▶ If the graph is symmetric, the shift operator S is symmetric $\Rightarrow S = S^T$

▶ The specific choice matters in practice but most of results and analysis hold for any choice of S

# Laplacians and Graph Signal Variability

▶ The variability of a graph signal has to be measured with respect to the structure of the graph

▶ The quadratic form of the graph's Laplacian provides this measure

▶ We are given a graph signal x and a <span style="color:red">symmetric</span> graph with edge set $\mathcal{E}$ and edge weights $w_{ij}$

---

**Definition (Total Variation Energy)**

The total variation energy of the signal x with respect to the graph $\mathcal{G}$ is defined as

$$\text{TV}(\mathsf{x}) := \frac{1}{2} \sum_{(i,j) \in \mathcal{E}} w_{ij}(x_i - x_j)^2$$

---

▶ $(x_i - x_j)^2 \Rightarrow$ Energy of difference between the signal values $x_i$ and $x_j$ observed at <span style="color:blue">node $i$</span> and <span style="color:green">node $j$</span>

▶ Weighted by the edge weight $w_{ij}$ and summed across all edges

▶ In the total variation energy $\text{TV}(x) := \dfrac{1}{2} \displaystyle\sum_{(i,j)\in\mathcal{E}} w_{ij}(x_i - x_j)^2$ there is a term associated to each edge



▶ The factor 2 appear because the graph is symmetric. Each arrow counts for two edges

▶ We are given a graph signal x and a symmetric graph with Laplacian L

▶ The Laplacian quadratic form is the function $\Rightarrow x^T L x$ (row $\times$ matrix $\times$ column $=$ scalar)

**Theorem (Laplacian Quadratic Form)**

The Laplacian quadratic form of graph signal x is equal to its total variation energy

$$x^T L x \;=\; TV(x) \;=\; \frac{1}{2} \sum_{(i,j) \in \mathcal{E}} w_{ij}(x_i - x_j)^2$$

▶ The Laplacian quadratic form measures the variability of different graph signals

**Proof:**

▶ This is an annoying algebraic calculation

▶ Isolate an edge $e = (ij) \in \mathcal{E}$ and define a symmetric graph with edge $e$. It's Laplacian satisfies

$$((\mathsf{L}_e))_{ij} = ((\mathsf{L}_e))_{ji} = -w_{ij} \qquad ((\mathsf{L}_e))_{ii} = ((\mathsf{L}_e))_{jj} = w_{ij}$$

▶ Since the matrix $\mathsf{L}_e$ has only four nonzero entries, the quadratic form $\mathsf{x}^T \mathsf{L}_e \mathsf{x}$ satisfies

$$\mathsf{x}^T \mathsf{L}_e \mathsf{x} = x_i w_{ij} x_i + x_j w_{ij} x_j - x_i w_{ij} x_j - x_j w_{ij} x_i = w_{ij}\big(x_i - x_j\big)^2$$

▶ To conclude notice that we have $\mathsf{L} = \dfrac{1}{2} \displaystyle\sum_{(i,j)\in\mathcal{E}} \mathsf{L}_e$ and therefore $\Rightarrow \mathsf{x}^T \mathsf{L} \mathsf{x} = \dfrac{1}{2} \displaystyle\sum_{(i,j)\in\mathcal{E}} \mathsf{x}^T \mathsf{L}_e \mathsf{x}$ ∎

▶ We say $v_k$ is an eigenvector of L with associated eigenvalue $\lambda_k$ if we have $Lv_k = \lambda_k v_k$

**Corollary (Variability of Laplacian Eigenvectors)**

The total variation energy of eigenvector $v_k$ is its associated eigenvalue $\Rightarrow TV(v_k) = \lambda_k$

**Proof:** As per the Laplacian quadratic form theorem $\Rightarrow TV(v_k) = v_k^T L v_k = v_k^T \lambda_k v_k = \lambda_k$ ■

▶ Eigenvectors of the Laplacian represent different rates of variability $\Rightarrow$ A (graph) Fourier transform

# Graph Fourier Transform

▶ The Graph Fourier Transform (GFT) is a tool for analyzing graph information processing systems

▶ We work with symmetric graph shift operators $\Rightarrow S = S^H$

▶ Introduce eigenvectors $v_i$ and eigenvalues $\lambda_i$ of graph shift operator S $\Rightarrow Sv_i = \lambda_i v_i$

$\Rightarrow$ For symmetric S eigenvalues are real. We have ordered them $\Rightarrow \lambda_0 \leq \lambda_1 \leq \ldots \leq \lambda_n$

▶ Define eigenvector matrix $V = [v_1, \ldots, v_n]$ and eigenvalue matrix $\Lambda = \text{diag}([\lambda_1; \ldots; \lambda_n])$

**Theorem (Eigenvectors Orthogonality of Symmetric Matrices)**

Consider a symmetric shift operator (matrix) S, with eigenvalues v and u associated with different

eigenvalues $\lambda$ and $\mu$. The eigenvectors are orthogonal

$$v^H u = 0.$$

▶ The eigenvectors of a symmetric shift operator can be used to define a unitary transform

**Proof:**

▶ Since eigenvectors v and u are respectively associated with eigenvalues $\lambda$ and $\mu$, we have that

$$\mathsf{Sv} = \lambda \mathsf{v}, \qquad \mathsf{Su} = \mu \mathsf{u}$$

▶ Since the matrix S is symmetric and real we have that $\mathsf{S}^H = \mathsf{S}$. For here, it follows that

$$\left( \mathsf{u}^H \mathsf{Sv} \right)^H \;=\; \mathsf{v}^H \mathsf{S}^H \mathsf{u} \;=\; \mathsf{v}^H \mathsf{Su}$$

▶ Substitute $\mathsf{Sv} = \lambda \mathsf{v}$ on the leftmost side. Substitute $\mathsf{Su} = \mu \mathsf{u}$ on the rightmost side.

$$\left( \lambda \mathsf{u}^H \mathsf{v} \right)^H \;=\; \left( \mathsf{u}^H \mathsf{Sv} \right)^H \;=\; \mathsf{v}^H \mathsf{S}^H \mathsf{u} \;=\; \mathsf{v}^H \mathsf{Su} \;=\; \mu \mathsf{v}^H \mathsf{u}$$

▶ For this to be true with $\lambda \neq \mu$ we must have that $\mathsf{v}^H \mathsf{u} = 0$

▶ The $k$th column of the eigenvector matrix $V = [v_1, \ldots, v_n]$ is the $k$th eigenvector $v_k$ of the shift S

▶ Since the eigenvectors $v_k$ are orthonormal, the eigenvector matrix V is unitary $\Rightarrow T^H T = I$

$$V^H V = \begin{bmatrix} v_1^H \\ \vdots \\ v_k^H \\ \vdots \\ v_n^H \end{bmatrix} \begin{bmatrix} v_1^H v_1 & \cdots & v_1^H v_k & \cdots & v_1^H v_n \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ v_k^H v_1 & \cdots & v_k^H v_k & \cdots & v_k^H v_n \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ v_n^H v_n & \cdots & v_n^H v_k & \cdots & v_n^H v_n \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 1 \end{bmatrix}$$

where the top row is $\begin{bmatrix} v_1 & \cdots & v_k & \cdots & v_n \end{bmatrix}$

▶ The eigenvalue matrix $\Lambda$ is a diagonal matrix with diagonal entries equal to eigenvalues of S

$$\Lambda = \begin{bmatrix} \lambda_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_k & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & \lambda_n \end{bmatrix} = \text{diag}(\lambda_1, \ldots, \lambda_n)$$

▶ Eigenvalue decomposition $\Rightarrow$ We can write the shift operator as $S = V\Lambda V^H$. Indeed,

$$SV = S[v_1, \ldots, v_n] = [Sv_1, \ldots, Sv_n] = [\lambda_1 v_1, \ldots, S\lambda_n] = V\Lambda$$

▶ Multiply from the right by $V^H$ and use the fact that V is unitary to eliminate $VV^H = I$

**Graph Fourier Transform**

Given a graph shift operator $S = V \Lambda V^H$, the graph Fourier transform (GFT) of graph signal $x$ is

$$\tilde{x} \;=\; V^H x$$

▶ GFT $\equiv$ projection on the eigenspace of the graph shift operator $\;\Rightarrow\; \tilde{x}_k \;=\; v_k^H x \;=\; \langle x, v_k \rangle$

▶ We say $\tilde{x}$ is a graph frequency representation of $x$. A representation in the graph frequency domain

**Inverse Graph Fourier Transform**

Given a graph shift operator $S = V \Lambda V^H$, the inverse graph Fourier transform (iGFT) of GFT $\tilde{x}$ is

$$\tilde{\tilde{x}} = V \tilde{x}$$

▶ Given that $V^H V = I$, the iGFT of the GFT of signal $x$ recovers the signal $x$

$$\tilde{\tilde{x}} = V \tilde{x} = V \left( V^H x \right) = I x = x$$

**Theorem (The GFT Preserves Energy)**

The energy $\|x\|^2$ of a signal and the energy of its GFT $\|\tilde{x}\|^2$ are the same $\Rightarrow \|x\|^2 = \|\tilde{x}\|^2$

▶ Given that $V^H V = I$, we have the chain of equalities

$$\|\tilde{x}\|^2 = \tilde{x}^H \tilde{x} = x^H V V^H x = x^H I x = x^H x = \|x\|^2$$

▶ Because of inverse theorem, we can write graph signals as $\Rightarrow \mathbf{x} = \mathbf{V}\tilde{\mathbf{x}} = \sum_{k=1}^{n} \tilde{x}_k \mathbf{v}_k$

▶ Because of Parseval, the energy $\left|\tilde{x}_k\right|^2$ of the $k$th coefficient is the energy $\mathbf{v}_k$ contributes to $\mathbf{x}$

▶ Use the Laplacian as shift operator $\Rightarrow \mathbf{S} = \mathbf{L}$

$\Rightarrow$ Total variation energy of Laplacian eigenvectors $\Rightarrow \mathrm{TV}(\mathbf{v}_k) = \lambda_k = \dfrac{1}{2} \sum_{(i,j) \in \mathcal{E}} w_{ij}(x_i - x_j)^2$

$\Rightarrow$ Eigenvectors are sorted according to their variability $\Rightarrow \mathrm{TV}(\mathbf{v}_1) \leq \mathrm{TV}(\mathbf{v}_2) \leq \ldots \leq \mathrm{TV}(\mathbf{v}_n)$

▶ The Laplacian GFT decomposes signals $\mathbf{x}$ into components of progressively higher variability

▶ This variability interpretation is true for Laplacian shift operators only

▶ Adjacency matrix $\Rightarrow$ If $S = A$ this is sort of true if the node degrees are similar

▶ Normalized Laplacian $\Rightarrow$ If $S = \bar{L}$, analogous interpretation holds for normalized variation energy

$$\bar{T}V(v_k) = \frac{1}{2} \sum_{(i,j) \in \mathcal{E}} w_{ij} \left( \frac{x_i}{\sqrt{d_i}} - \frac{x_j}{\sqrt{d_j}} \right)^2$$

▶ Normalized Adjacency $\Rightarrow$ If $S = \bar{A}$ the same holds because eigenvectors coincide $\Rightarrow \bar{L} = I - \bar{L}$

# The GFT of Discrete Time Signals

▶ We can describe discrete time signals as signals supported on a directed line graph

Description of time with a directed line graph



The adjacency "matrix" of a directed line graph

$$S = A = \begin{bmatrix} & \cdot & \cdot & \cdot & \\ \cdot & 0 & 0 & 0 & \cdot \\ \cdot & 1 & 0 & 0 & \cdot \\ \cdot & 0 & 1 & 0 & \cdot \\ \cdot & 0 & 0 & 1 & \cdot \\ & \cdot & \cdot & \cdot & \end{bmatrix}$$

▶ This adjacency "matrix" has a GFT associated with it. Is it related to the DTFT?

▶ A time shifting of a time signal means moving the signal up on the time line $\Rightarrow$ Follow the arrows



▶ Time shift is reinterpreted as multiplication by the adjacency matrix S of the line graph

$$
S x \;=\; \begin{bmatrix} & \cdot & \cdot & \cdot & \\ \cdot & 0 & 0 & 0 & \cdot \\ \cdot & 1 & 0 & 0 & \cdot \\ \cdot & 0 & 1 & 0 & \cdot \\ \cdot & 0 & 0 & 1 & \cdot \\ & \cdot & \cdot & \cdot & \end{bmatrix} \begin{bmatrix} \cdot \\ x_0 \\ x_1 \\ x_2 \\ x_3 \\ \cdot \end{bmatrix} \;=\; \begin{bmatrix} \cdot \\ x_{-1} \\ x_0 \\ x_1 \\ x_2 \\ \cdot \end{bmatrix}
$$

▶ Product Sx is such that $\big( Sx \big)_n = x_{n-1}$ $\Rightarrow$ Signal components move up on the time line

▶ Particularize to the case in which the graph signal x is a complex exponential $\Rightarrow x_n = e^{j2\pi fnTs}$

▶ Moving components up in the time line for this particular signal yields

$$(Sx)_n = e^{j2\pi f(n-1)Ts} = e^{j2\pi f(-1)Ts} e^{j2\pi fnTs} = e^{-j2\pi fTs} x_n$$

▶ Complex exponential x is an eigenvector of the shift "matrix" S with associated eigenvalue $e^{-j2\pi fTs}$

▶ Let $e_{fT_s}$ be a discrete time complex exponential with components $x_n = e^{j2\pi f n T_s}$

**Theorem (GFT of a Directed Line Graph)**

The components of the GFT of a discrete time signal x are $\Rightarrow \tilde{x}_k = \langle x, e_{fT_s} \rangle = \sum\limits_{-\infty}^{+\infty} x_n e^{-j2\pi f n T_s}$

▶ Which is the exact same definition of the DTFT of the signal x

▶ DFT ≡ GFT of directed cycle graph (connect node $n - 1$ to node 1)

▶ 2D-DFT ∼ GFT of grid graph ⇒ In fact, it's complicated. But true enough

▶ PCA *equiv* GFT of covariance marix graph ⇒ Self evident. Same definition

# Graph Convolutional Filters

▶ Graph convolutional filters are the tool of choice for the linear processing of graph signals

▶ Convolutional filters process signals in time by leveraging the time shift operator



▶ The time convolution is a linear combination of time shifted inputs $\Rightarrow y_n = \sum_{k=0}^{K-1} h_k x_{n-k}$

▶ Time signals are representable as graph signals supported on a line graph S ⇒ The pair (S, x)



▶ Time shift is reinterpreted as multiplication by the adjacency matrix S of the line graph

$$S^3 x = S \left[ S^2 x \right] = S \left[ S \left( S x \right) \right] = \begin{bmatrix} \cdots & \vdots & \vdots & \vdots & \cdots \\ \cdots & 0 & 0 & 0 & \cdots \\ \cdots & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \end{bmatrix} \begin{bmatrix} \vdots \\ x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ x_{-3} \\ x_{-2} \\ x_{-1} \\ x_0 \\ \vdots \end{bmatrix}$$

▶ Components of the shift sequence are powers of the adjacency matrix applied to the original signal

⇒ We can rewrite convolutional filters as polynomials on S, the adjacency of the line graph

▶ The convolution operation is a linear combination of shifted versions of the input signal

▶ But we now know that time shifts are multiplications with the adjacency matrix S of line graph



▶ Time convolution is a polynomial on adjacency matrix of line graph $\Rightarrow y = h \star x = \sum_{k=0}^{K-1} h_k S^k x$
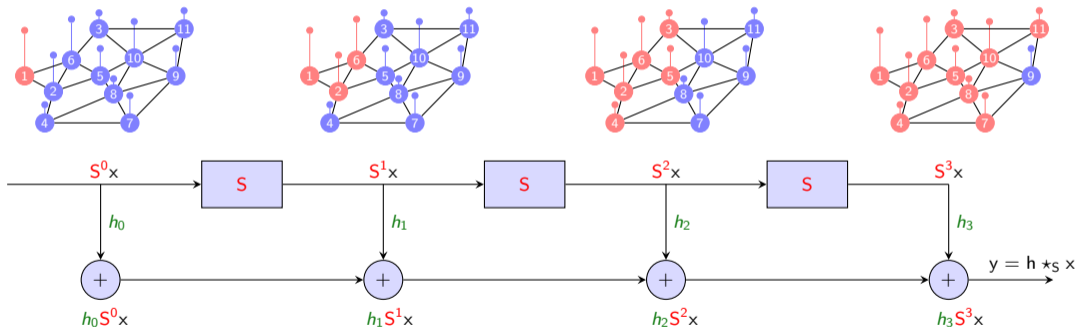
▶ The convolution operation is a linear combination of shifted versions of the input signal

▶ But we now know that time shifts are multiplications with the adjacency matrix S of line graph



▶ Time convolution is a polynomial on adjacency matrix of line graph $\Rightarrow y = h \star x = \sum_{k=0}^{K-1} h_k S^k x$

► If we let S be the shift operator of an arbitrary graph we recover the graph convolution

▶ Given graph shift operator S and coefficients $h_k$, a graph filter is a polynomial (series) on S

$$H(S) = \sum_{k=0}^{\infty} h_k S^k$$

▶ The result of applying the filter $H(S)$ to the signal x is the signal

$$y = H(S)x = \sum_{k=0}^{\infty} h_k S^k x$$

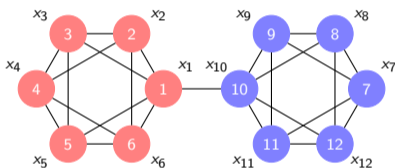▶ We say that $y = h \star_S x$ is the graph convolution of the filter $h = \{h_k\}_{k=0}^{\infty}$ with the signal x

▶ Graph convolutions aggregate information growing from local to global neighborhoods

▶ Consider a signal x supported on a graph with shift operator S. Along with filter $h = \{h_k\}_{k=0}^{K-1}$



▶ Graph convolution output $\Rightarrow y = h \star_S x = h_0 S^0 x + h_1 S^1 x + h_2 S^2 x + h_3 S^3 x + \ldots = \sum_{k=0}^{K-1} h_k S^k x$

▶ The same filter h = $\{h_k\}_{k=0}^{\infty}$ can be executed in multiple graphs ⇒ We can transfer the filter
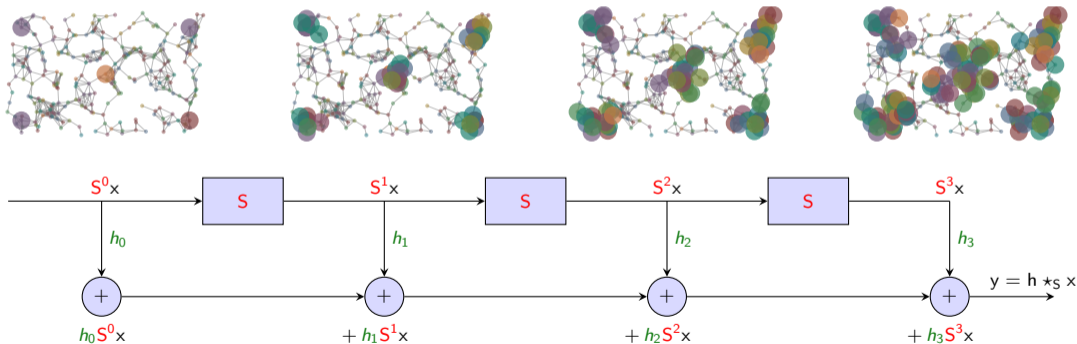
Graph Filter on a Graph

Same Graph Filter on Another Graph



▶ Graph convolution output ⇒ $y = h \star_S x = h_0 S^0 x + h_1 S^1 x + h_2 S^2 x + h_3 S^3 x + \ldots = \sum_{k=0}^{\infty} h_k S^k x$

▶ Output depends on the filter coefficients h, the graph shift operator S and the signal x

▶ A graph convolution is a weighted linear combination of the elements of the diffusion sequence

▶ Can represent graph convolutions with a shift register ⇒ Convolution ≡ Shift. Scale. Sum

## Graph Frequency Response of Graph Filters

▶ Graph filters admit a pointwise representation when projected into the shift operator's eigenspace

**Theorem (Graph frequency representation of graph filters)**

Consider graph filter h with coefficients $h_k$, graph signal x and the filtered signal $y = \sum_{k=0}^{\infty} h_k S^k x$.

The GFTs $\tilde{x} = V^H x$ and $\tilde{y} = V^H y$ are related by

$$\tilde{y} = \sum_{k=0}^{\infty} h_k \Lambda^k \tilde{x}$$

▶ The same polynomial but on different variables. One on S. The other on eigenvalue matrix Λ

**Proof:** Since $S = V \Lambda V^H$, can write shift operator powers as $S^k = V \Lambda^k V^H$. Therefore filter output is

$$y = \sum_{k=0}^{\infty} h_k S^k x = \sum_{k=0}^{\infty} h_k V \Lambda^k V^H x$$

▶ Multiply both sides by $V^H$ on the left $\Rightarrow V^H y = V^H \sum_{k=0}^{\infty} h_k V \Lambda^k V^H x$

▶ Copy and identify terms. Output GFT $V^H y = \tilde{y}$. Input GFT $V^H x = \tilde{x}$. Cancel out $V^H V$

$$V^H y = V^H \sum_{k=0}^{\infty} h_k V \Lambda^k V^H x \qquad \Rightarrow \qquad \tilde{y} = \sum_{k=0}^{\infty} h_k \Lambda^k \tilde{x} \qquad \blacksquare$$

▶ In the graph frequency domain graph filters are a **diagonal** matrices $\Rightarrow \tilde{y} = \displaystyle\sum_{k=0}^{\infty} h_k \Lambda^k \tilde{x}$

▶ Thus, graph convolutions are **pointwise in the GFT domain** $\Rightarrow \tilde{y}_i = \displaystyle\sum_{k=0}^{\infty} h_k \lambda_i^k \tilde{x}_i = \tilde{h}(\lambda_i)\tilde{x}_i$
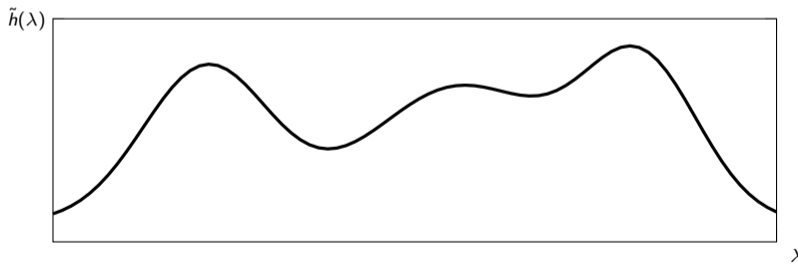
**Definition (Frequency Response of a Graph Filter)**

Given a graph filter with **coefficients** $h = \{h_k\}_{k=1}^{\infty}$, the graph frequency response is the polynomial

$$\tilde{h}(\lambda) = \sum_{k=0}^{\infty} h_k \lambda^k$$

**Definition (Frequency Response of a Graph Filter)**

Given a graph filter with coefficients h $= \{h_k\}_{k=1}^{\infty}$, the graph frequency response is the polynomial

$$\tilde{h}(\lambda) = \sum_{k=0}^{\infty} h_k \lambda^k$$

▶ Frequency response is the same polynomial that defines the graph filter $\Rightarrow$ but on scalar variable $\lambda$

▶ Frequency response is independent of the graph $\Rightarrow$ Depends only on filter coefficients

▶ The role of the graph is to determine the eigenvalues on which the response is instantiated

▶ Graph filter frequency response is a <span style="color:red">polynomial on a scalar variable</span> $\lambda \Rightarrow \tilde{h}(\lambda) = \sum_{k=0}^{\infty} h_k \lambda^k$

▶ Completely <span style="color:red">determined by the filter coefficients</span> $\mathsf{h} = \{h_k\}_{k=1}^{\infty}$ . The Graph has nothing to do with it

▶ A given (another) graph instantiates the response on its given (different) specific eigenvalues $\lambda_i$

▶ Eigenvectors do not appear in the frequency response. They determine the meaning of frequencies.

# Learning with Graph Signals

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

▶ Machine learning (ML) on graphs (or not) ≡ empirical risk minimization (ERM) on graphs (or not)

▶ In ERM we are given:

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

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

⇒ A function class $\mathcal{C}$

▶ Learning means finding function $\Phi^* \in \mathcal{C}$ that minimizes loss $\ell\left(y, \Phi(x)\right)$ averaged over training set

$$\Phi^* = \underset{\Phi \in \mathcal{C}}{\text{argmin}} \sum_{(x,y) \in \mathcal{T}} \ell\left(y, \Phi(x),\right)$$

▶ We use $\Phi^*(x)$ to estimate outputs $\hat{y} = \Phi^*(x)$ when inputs $x$ are observed but outputs $y$ are unknown
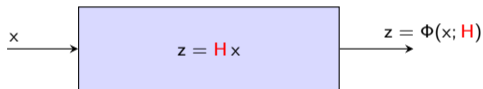
▶ In ERM, the function class $\mathcal{C}$ is the degree of freedom available to the system's designer

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

▶ Designing a Machine Learning ≡ 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 x / y are graph signals supported on a common graph with shift operator S

▶ Function class ⇒ Generic Linea function mapping inputs to oouputs ⇒ $\Phi(x) = Hx = \Phi(x; H)$



▶ Learn ERM solution restricted to graph filter class ⇒ $h^* = \underset{h}{\mathrm{argmin}} \sum_{(x,y)\in\mathcal{T}} \ell\Big(y, \Phi(x; H)\Big)$

⇒ Optimization is over matrices H. It does not take advantage of the graph

▶ Input / output signals x / y are graph signals supported on a common graph with shift operator S

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

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

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

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

# Learning Ratings in Recommendation Systems

▶ Formulate recommendation systems as ERM problems that predict ratings that users give to items

▶ In a recommendation system, we want to predict the rating a user would give to an item

▶ Collect ratings that some users give to some items ⇒ These are rating histories

▶ Exploit product similarities to predict ratings of unseen user-item pairs

▶ Example 1 ⇒ In an online store items are products and users are customers

▶ Example 2 ⇒ In a movie repository items are movies and users are watchers

▶ For all items $i$ and users $u$ there exist ratings $\Rightarrow y_{ui}$

  $\Rightarrow$ User rating vector $\mathbf{y}_u$ has entries $y_{ui}$

▶ We only observe a subset of ratings $\Rightarrow x_{ui}$

  $\Rightarrow$ Observed user rating vector $\mathbf{x}_u$ has entries $x_{ui}$

  $\Rightarrow$ We assume $x_{ui} = 0$ if item $i$ is unrated by user $u$

▶ For all items $i$ and users $u$ there exist ratings $\Rightarrow y_{ui}$

$\Rightarrow$ User rating vector $\mathbf{y}_u$ has entries $y_{ui}$

▶ We only observe a subset of ratings $\Rightarrow x_{ui}$

$\Rightarrow$ Observed user rating vector $\mathbf{x}_u$ has entries $x_{ui}$

$\Rightarrow$ We assume $x_{ui} = 0$ if item $i$ is unrated by user $u$

▶ Construct product similarity graph with weights $w_{ij}$ represent likelihood of similar scores

▶ Interpret vector of ratings $y_u$ of user $u$ as a graph signal supported on the product similarity graph

▶ The observed ratings $x_u$ of user $u$ are a subsampling of this graph signal.

▶ Our goal is to learn to reconstruct the rating graph signal $y_u$ from the observed ratings $x_u$

▶ Build similarity graph using available ratings. Use of expert knowledge is common as well

▶ Consider pair of products $i$ and $j$. Restrict attention to set of users that rated both products $\Rightarrow \mathcal{U}_{ij}$

▶ Mean ratings restricted to users that rated products $i$ and $j$

$$\mu_{ij} = \frac{1}{\#(\mathcal{U}_{ij})} \sum_{u \in \mathcal{U}_{ij}} x_{ui} \qquad \mu_{ji} = \frac{1}{\#(\mathcal{U}_{ij})} \sum_{u \in \mathcal{U}_{ji}} x_{uj}$$

▶ Similarity score = correlation restricted to users in $\mathcal{U}_{ij}$

$$\sigma_{ij} = \frac{1}{\#(\mathcal{U}_{ij})} \sum_{u \in \mathcal{U}_{ij}} \left( x_{ui} - \mu_{ij} \right) \left( x_{uj} - \mu_{ji} \right)$$

▶ Weights = normalized correlations $\Rightarrow w_{ij} = \sigma_{ij} \bigg/ \sqrt{\sigma_{ii}\sigma_{jj}}$

▶ Given observed ratings $x_u$ the AI produces estimates $\Phi(x_u)$. We want $\Phi(x_u)$ to approximate $y_u$

$$\ell\Big(y_u, \Phi(x_u)\Big) = \frac{1}{2}\Big\| y_u - \Phi(x_u) \Big\|^2$$

▶ In reality, we want to predict the rating of specific item $i$

$$\ell\Big(y_u, \Phi(x_u)\Big) = \frac{1}{2}\Big( e_i^T y_u - e_i^T \Phi(x_u) \Big)^2$$

▶ Where $e_i$ is a vector in the canonical basis $\Rightarrow (e_i)_i = 1$, $(e_i)_j = 0$ for $j \neq i$

▶ For each item $i$ let $\mathcal{U}_i$ be the set of users that have rated $i$. Construct training pairs $(x, y)$ with

$$y = \left(e_i^T x_u\right) e_i \qquad x = x_u - y \qquad \text{for all } u \in \mathcal{U}_i, \text{ for all } i$$

▶ Extract the rating $x_{ui}$ of item $i$. Record into graph signal $y$

▶ Remove rating $x_{ui}$ from $x_u$. Record to graph signal $x$

▶ Repeat for all users in the set $\mathcal{U}_i$ of users that rated $i$

▶ Repeat for all items $\Rightarrow$ Training set $\mathcal{T}$

▶ Parametrized AI $\Phi(x_u) = \Phi(x_u; \mathcal{H})$. We want to find solution of the ERM problem

$$\mathcal{H}^* = \operatorname*{argmin}_{\mathcal{H}} \sum_{(x,y) \in \mathcal{T}} \left( e_i^T y - e_i^T \Phi(x; \mathcal{H}) \right)^2$$
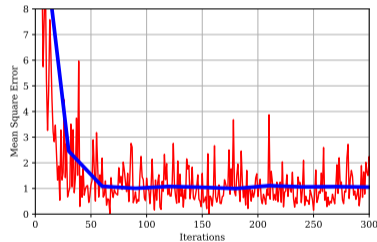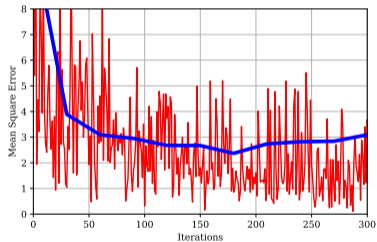
▶ A bad idea $\Rightarrow$ Linear regression with a generic linear function.

▶ A good idea $\Rightarrow$ Graph filters.

# Learning Ratings with Graph Filters

▶ We use graph filters to learn ratings in recommendation systems

▶ We contrast with the use of linear regression with a generic linear function

▶ Use MovieLens-100k as benchmark $\Rightarrow 10^6$ ratings given by $U = 943$ users to $M = 1,682$ movies

▶ The ratings for each movie are between 1 and 5. From one star to five starts

▶ Train and test several machine learning parametrizations.

▶ We predict ratings using AI that results from solving the ERM problem

$$\mathcal{H}^* \;=\; \underset{\mathcal{H}}{\operatorname{argmin}} \quad \sum_{(x,y)\in\mathcal{T}} \left( e_i^T y - e_i^T \Phi(x; \mathcal{H}) \right)^2$$

▶ Parameterizations that ignore data structure= ⇒ Linear regression. Fully connected NNs

▶ Parameterizations that leverage data structure= ⇒ Graph filters. Graph NNs

▶ Linear regression reduces training MSE to about 2. Quite bad for ratings that vary from 0 to 5

▶ Graph filter reduces training MSE to about 1. Not too good. Humans are not that predictable



▶ Graph filter outperforms linear regression ⇒ Leverages underlying permutation symmetries

▶ Linear regression works even worse in the test set

▶ The test MSE of the graph filter is about the same as the training MSE. It generalizes



▶ Graph filter outperforms linear regression ⇒ Leverages underlying permutation symmetries

# Permutation Equivariance of Graph Filters
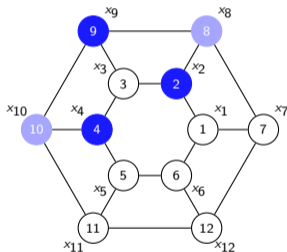
▶ We will show that graph convolutional filters are equivariant to permutations

**Definition (Permutation matrix)**

A square matrix $P$ is a permutation matrix if it has binary entries so that $P \in \{0,1\}^{n \times n}$ and it further satisfies $P1 = 1$ and $P^T 1 = 1$.

▶ The product $P^T x$ reorders the entries of the vector $x$.

▶ The product $P^T S P$ is a consistent reordering of the rows and columns of $S$

**Definition (Permutation matrix)**

A square matrix P is a permutation matrix if it has binary entries so that $P \in \{0, 1\}^{n \times n}$ and it further satisfies $P1 = 1$ and $P^T 1 = 1$.

▶ Since $P1 = P^T 1 = 1$ with binary entries $\Rightarrow$ Exactly one nonzero entry per row and column of P

▶ Permutation matrices are unitary $\Rightarrow P^T P = I$. Matrix $P^T$ undoes the reordering of matrix P

▶ If $(S, x)$ is a graph signal, $(P^T S P, P^T x)$ is a relabeling of $(S, x)$. Same signal. Different names

Graph signal $x$ Supported on $S$



Graph signal $\hat{x} = P^T x$ supported on $\hat{S} = P^T S P$



▶ Processing should be label-independent $\Rightarrow$ Permutation equivariance of graph filters and GNNs

▶ Graph filter H(S) is a polynomial on shift operator S with coefficients $h_k$. Outputs given by

$$H(S)x = \sum_{k=0}^{K-1} h_k S^k x$$

▶ We consider running the same filter on $(S, x)$ and permuted (relabeled) $(\hat{S}, \hat{x}) = (P^T S P, P^T x)$

$$H(S)x = \sum_{k=0}^{K-1} h_k S^k x \qquad\qquad H(\hat{S})\hat{x} = \sum_{k=0}^{K-1} h_k \hat{S}^k \hat{x}$$

▶ Filter $H(S)x \Rightarrow$ Coefficients $h_k$. Input signal $x$. Instantiated on shift $S$

▶ Filter $H(\hat{S})\hat{x} \Rightarrow$ Same Coefficients $h_k$. Permuted Input signal $\hat{x}$. Instantiated on permuted shift $\hat{S}$

**Theorem (Permutation equivariance of graph filters)**

Consider consistent permutations of the shift operator $\hat{S} = P^T S P$ and input signal $\hat{x} = P^T x$. Then

$$H(\hat{S})\,\hat{x} \;=\; H(P^T S P)\,(P^T x) \;=\; P^T H(S)\, x$$

▶ Graph filters are equivariant to permutations ⇒ Permute input and shift ≡ Permute output

**Proof:** Write filter output in polynomial form. Use permutation definitions $\hat{S} = P^T S P$ and $\hat{x} = P^T x$

$$H(\hat{S})\hat{x} \; = \; \sum_{k=0}^{K-1} h_k \hat{S}^k \hat{x} \; = \; \sum_{k=0}^{K-1} h_k \left( P^T S P \right)^k P^T x$$

▶ In the powers $\left( P^T S P \right)^k$, $P$ and $P^T$ undo each other $(P^T P = I) \;\Rightarrow\; \left( P^T S P \right)^k = P^T \left( S \right)^k P$

▶ Substitute this into filter's output expression. Cancel remaining $PP^T = I$ product. Factor $P^T$

$$H(\hat{S})\hat{x} \; = \; \sum_{k=0}^{K-1} h_k P^T S^k PP^T x \; = \; \sum_{k=0}^{K-1} h_k P^T S^k Ix \; = \; P^T \sum_{k=0}^{K-1} h_k S^k x \; = \; P^T H(S)x \qquad \blacksquare$$

▶ We request signal processing independent of labeling ⇒ Graph filters fulfill this request

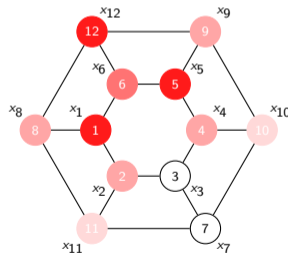⇒ Permute input and shift ≡ Relabel input ⇒ Permute output ≡ Relabel output

Graph signal x Supported on S

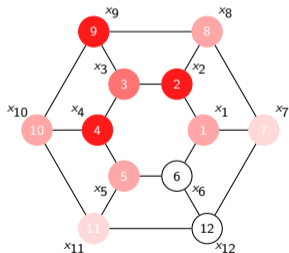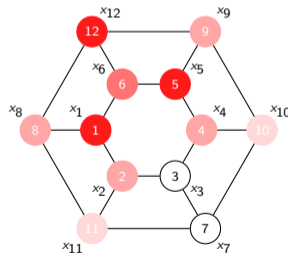Graph signal $\hat{x} = P^T x$ supported on $\hat{S} = P^T S P$

▶ We request signal processing independent of labeling $\Rightarrow$ Graph filters fulfill this request

$\Rightarrow$ Permute input and shift $\equiv$ Relabel input $\Rightarrow$ Permute output $\equiv$ Relabel output

Filter's output H(S)x Supported on S

Filter's Output H(Ŝ)x̂ supported on Ŝ

▶ We request signal processing independent of labeling $\Rightarrow$ Graph filters fulfill this request

$\Rightarrow$ Permute input and shift $\equiv$ Relabel input $\Rightarrow$ Permute output $\equiv$ Relabel output
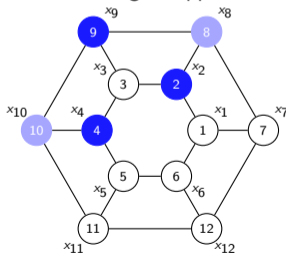


Filter's output H(S)x Supported on S



Equivariance theorem $\Rightarrow H(\hat{S})\hat{x} = P^T H(S)x$
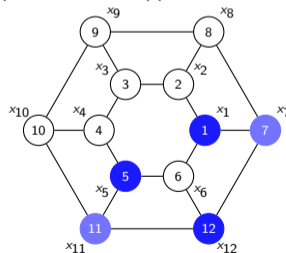
▶ Equivariance to permutations allows GNNs to exploit symmetries of graphs and graph signals

▶ By symmetry we mean that the graph can be permuted onto itself $\Rightarrow S = P^T S P$

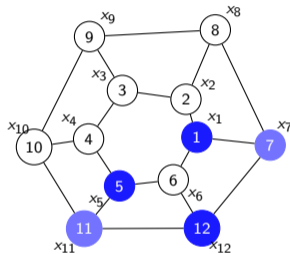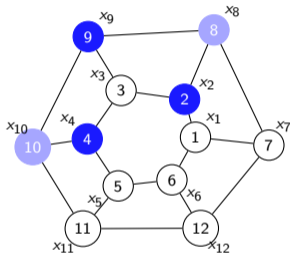▶ Equivariance theorem implies $\Rightarrow H(S)(P^T x) = H(P^T S P)(P^T x) = P^T H(S)(x)$

From observing x supported on S



Learn to process $P^T x$ supported on $S = P^T S P$

► Graph not symmetric but close to symmetric ⇒ perturbed version of a permutation of itself



► It can be shown that graph filters can lack stability to deformations ⇒ Graph Neural Networks

⇒ But this is a story for another day ⇒ Register for ESE 514. Or visit gnn.seas.upenn.edu