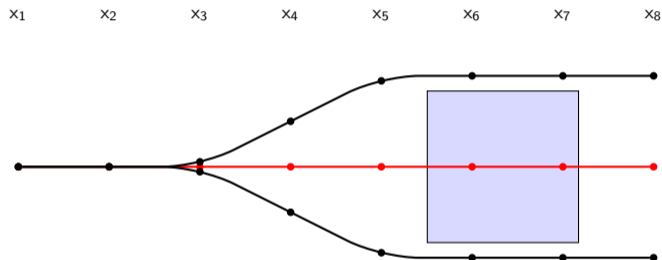


Machine Learning on Sequences

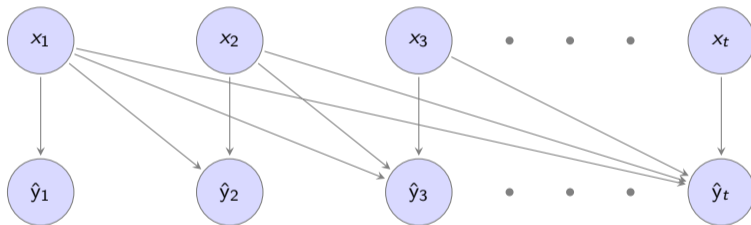
- ▶ GNNs are architectures specialized in learning data defined over graph supports
- ▶ Several processes have a **sequential** nature. To learn from them, we need dedicated architectures

- ▶ Often, we want to learn **properties of a sequence** \Rightarrow Is the particle entering the forbidden area?



- ▶ This problem is not just a simple **sequence of classifications** $\Rightarrow y_t = \phi(x_t)$
- ▶ It is a (sequence of) **classifications of a sequence** $\Rightarrow y_t = \phi(x_{1:t}) = \phi(x_t, x_{t-1}, \dots, x_1)$

- Predictions on a sequence depends on observation histories $\Rightarrow \hat{y}_t = \Phi(x_t, x_{t-1}, \dots, x_1)$



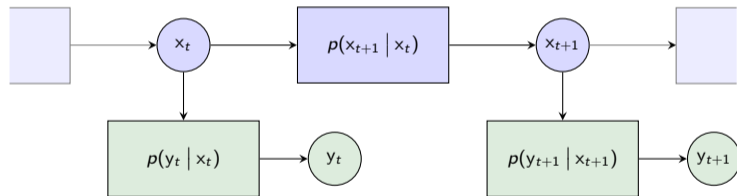
- Recurrent neural networks (RNNs) estimate a **hidden state** to avoid this **unbounded memory growth**

- ▶ A stochastic process (random sequence) is said to be **Markov** or **memoryless** if

$$p(x_{t+1} | \mathbf{x}_{1:t}) = p(x_{t+1} | x_t)$$

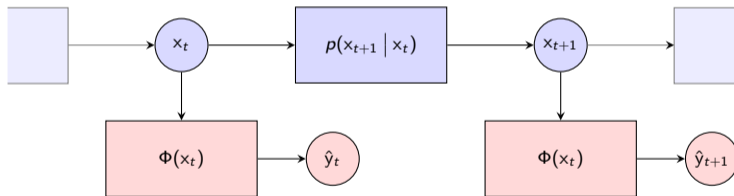
- ▶ It is the same to condition on the current value x_t or conditioning on the whole trajectory $x_{0:t}$
 - ⇒ The **future**, given the present, is **independent of the past**
 - ⇒ For predicting the future, knowledge of the past is irrelevant
- ▶ Outputs (e.g, trajectory categories) are **conditionally independent** ⇒ $p(y_t | x_t) = p(y_t | \mathbf{x}_{1:t})$

- ▶ In a **memoryless** Markov Process, learning **is equivalent** reduce to a **sequence of learning** problems
- ▶ State evolution is a chain of memoryless transitions. And outputs depend on the current state only



- ▶ An AI to predict y_t **mimics the conditional distribution of the observations**. The past is irrelevant

- ▶ In a **memoryless** Markov Process, learning **is equivalent** reduce to a **sequence of learning** problems
- ▶ State evolution is a chain of memoryless transitions. And outputs depend on the current state only

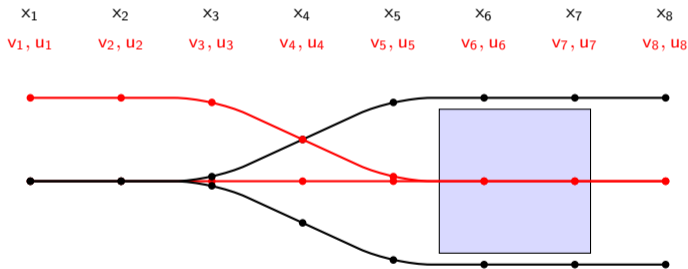


- ▶ An AI to predict y_t **mimics the conditional distribution of the observations**. The past is irrelevant

Recurrent Neural Networks

- ▶ Machine Learning in stochastic processes that are **not Markov** \Rightarrow The past is relevant in learning

- ▶ The evolution of the trajectory is not a Markov Process if we observe positions only
- ▶ But it is Markov if we have access to velocities and accelerations \Rightarrow Hidden (unobserved) states



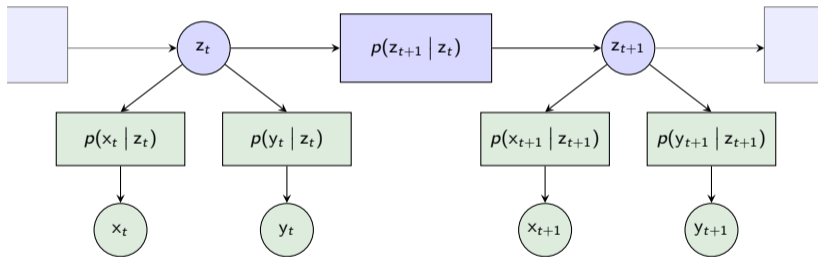
- ▶ All systems are Markov \Rightarrow We often lack enough information to observe their Markov structure

- ▶ Stochastic process x_t follows a **hidden Markov model** if there exists a **process z_t** such that

$$p(z_{t+1} | \mathbf{z}_{1:t}) = p(z_{t+1} | z_t) \quad \text{and} \quad p(x_t | z_t) = p(x_t | \mathbf{z}_{0:t})$$

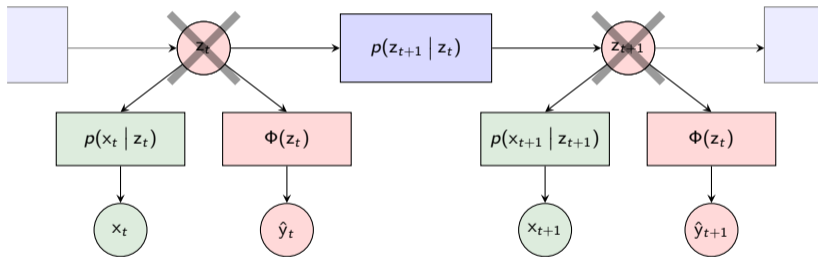
- ▶ The **hidden state z_t** is a **memoryless Markov** stochastic process
- ▶ The **observed state x_t** is **conditionally independent**. Depends only on the current hidden state z_t
- ▶ Outputs are also **conditionally independent** $\Rightarrow p(y_t | z_t) = p(y_t | \mathbf{z}_{1:t})$

- ▶ In a **hidden Markov model** learning is **not equivalent** to a **sequence of learning problems**
- ▶ The AI can try to mimic the conditional distribution $p(y_t | z_t)$. But we don't have access to z_t



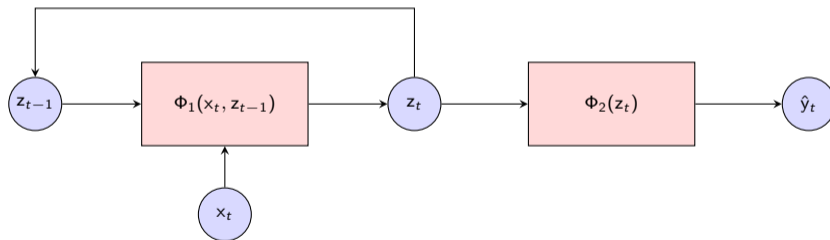
- ▶ Recurrent Neural Network (RNN) \Rightarrow Use observed state x_t to estimate hidden state z_t

- ▶ In a **hidden Markov model** learning is **not equivalent** to a **sequence of learning problems**
- ▶ The AI can try to mimic the conditional distribution $p(y_t | z_t)$. But we don't have access to z_t



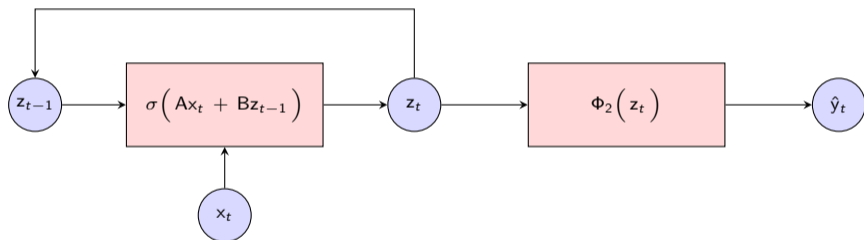
- ▶ Recurrent Neural Network (RNN) \Rightarrow Use observed state x_t to estimate hidden state z_t

- ▶ A recurrent neural network is made up of two separate learning parametrizations
 - ⇒ $\Phi_1(x_t, z_{t-1})$ ⇒ From observed state x_t ⇒ and hidden state z_{t-1} ⇒ to hidden state update z_t
 - ⇒ $\Phi_2(z_t)$ ⇒ From updated hidden state z_t ⇒ to output estimate \hat{y}_t



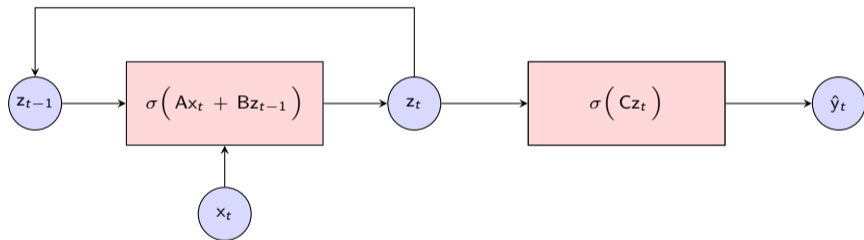
- ▶ It is a **recurrent** neural network because hidden states are **fed-back** as inputs for the next time step

- ▶ Use a **perceptron** for the AI that updates the hidden state $\Rightarrow \Phi_1(x_t, z_{t-1}) = \sigma(Ax_t + Bz_{t-1})$



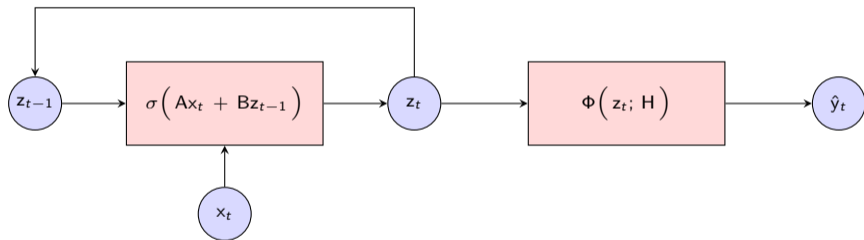
- ▶ Number of learnable parameters \equiv Entries of A and B \Rightarrow Does not depend on the time index t

- ▶ Use another **perceptron** for the AI that predicts the output $\Rightarrow \Phi_1(z_t) = \sigma(Cz_t)$



- ▶ We can also use a multi-layer neural network for the output prediction AI $\Rightarrow \Phi_1(z_t) = \Phi(z_t; H)$

- ▶ Use another **perceptron** for the AI that predicts the output $\Rightarrow \Phi_1(z_t) = \sigma(Cz_t)$

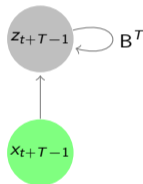


- ▶ We can also use a multi-layer neural network for the output prediction AI $\Rightarrow \Phi_1(z_t) = \Phi(z_t; H)$

Time Gating

- ▶ We discuss the problem of **vanishing/exploding gradients** in recurrent neural networks
- ▶ We introduce **gating mechanisms** in the form of long short-term memories and gated recurrent units

- ▶ In some tasks, the RNN may have to learn how to model **long term dependencies** of length T
- ▶ This poses a challenge \Rightarrow the Jacobian $\partial z_T / \partial \mathbf{B}$ will depend on a chain of multiplications by \mathbf{B}



- ▶ If eigenvalues of $\mathbf{B} \ll 1$, the gradients tend to **vanish**, leading to **exponentially smaller** weights \mathbf{B}
- ▶ If eigenvalues of $\mathbf{B} \gg 1$, the gradients tend to **explode**, leading to **exponentially larger** weights \mathbf{B}

- ▶ Consider a simplification of the RNN where we omit the nonlinear function $\sigma(\cdot)$ and the inputs x_t

$$z_t = Bz_{t-1}$$

- ▶ At time $t = T$, the state variable z_T depends on the **T th** power of the matrix B

$$z_T = B^T z_{t-T}$$

- ▶ If B admits an **eigendecomposition** $B = Q\Lambda Q^T$, the recurrence can be rewritten as

$$z_T = Q\Lambda^T Q^T z_{t-T}$$

⇒ Eigenvalues less than one will **vanish** and eigenvalues greater than one will **explode**

⇒ Any component of z_{t-T} not aligned with the largest eigenvalues will be **discarded**

- ▶ To address the issue of vanishing gradients, we add a **gating mechanism** to RNNs
- ▶ Gates are scalars in $[0, 1]$ acting on the **current input** and on the **previous state**
 - ⇒ Control how much of the input and past time information should be taken into account
- ▶ The value of each gate is **updated at every step of the sequence**
 - ⇒ Allows creating paths through time with derivatives that **neither vanish nor explode**
 - ⇒ Creates dependency paths that allow encoding both short and long term dependencies

- ▶ The most popular gated RNN architecture is the Long Short-Term Memory (LSTM) cell
- ▶ Three gates: a **forget gate** $f_t \in [0, 1]$, an **input gate** $g_t \in [0, 1]$, and a cell **output gate** $q_t \in [0, 1]$
- ▶ Let x_t be the input, z_t the state, and define the **internal memory** s_t of the LSTM cell
- ▶ Memory s_t updated by applying the **forget gate** to s_{t-1} and the **input gate** to the state update

$$s_t = f_t s_{t-1} + g_t \sigma(Ax_t + Bz_{t-1})$$

- ▶ State z_t updated by applying the cell **output gate** to the internal cell memory s_t

$$z_t = q_t \sigma(s_t)$$

- ▶ The Gated Recurrent Unit (GRU) is a second popular gated version of the RNN
- ▶ Slight variation of LSTM \Rightarrow **single gate** $u_t \in [0, 1]$ plays the role of input and forget gates

$$z_t = u_t z_{t-1} + (1 - u_t) \sigma(Ax_t + r_t Bz_{t-1})$$

- \Rightarrow **Reset gate** $r_t \in [0, 1]$ controls contribution of previous state z_{t-1} to updated state
- ▶ Besides the LSTM and the GRU, many more variants of gating mechanisms for RNNs exist

- ▶ In the LSTM and the GRU, the gates themselves are **calculated as the outputs of RNNs**
- ▶ For example, the forget gate f_t of the LSTM has its **own state variable z'_t** (as do all the other gates)

$$z'_t = \sigma (A'x_t + B'z'_{t-1})$$

- ▶ The forget gate f_t is then calculated from the input x_t and the state z'_t as

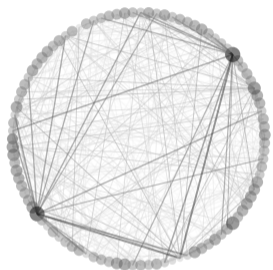
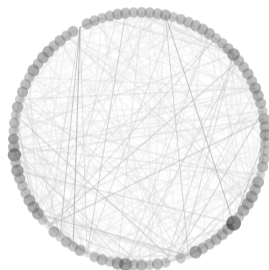
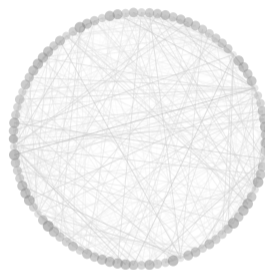
$$f_t = \text{sigmoid} (Ux_t + Wz'_t)$$

- ⇒ With U and W linear layers mapping the input and state features to a **single scalar**
- ⇒ And the sigmoid activation function ensuring **gate values in the [0, 1] interval**

Graph Recurrent Neural Networks

- ▶ We define Graph Recurrent Neural Networks (GRNNs) as particular cases of RNNs

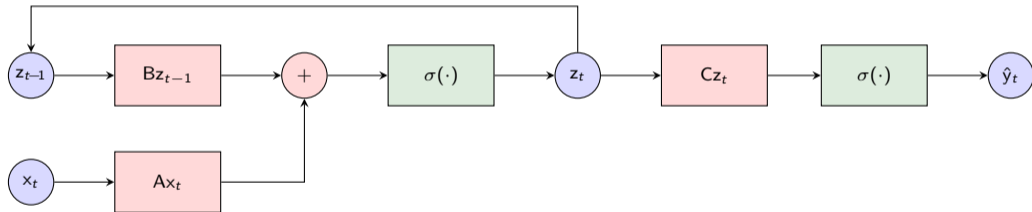
- ▶ Consider a **time varying process** x_t in which each of the signals is **supported on shift operator S**

 x_{t-2}  x_{t-1}  x_t

- ▶ A **graph recurrent** neural network (GRNN) combines

⇒ A **GNN** because x_t is supported on a graph. ⇒ An **RNN** because x_t is a sequence

- ▶ An RNN has a hidden state z_t updated with the perceptron $\Rightarrow z_t = \sigma(Ax_t + Bz_{t-1})$
- ▶ An it has an output prediction \hat{y}_t given by the perceptron $\Rightarrow \hat{y}_t = \sigma(Cz_t)$

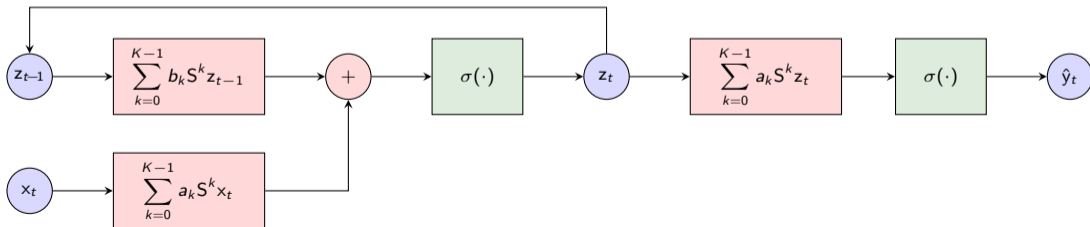


- ▶ The observed state x_t and the output y_t are graph signals supported on the graph shift operator S
 - \Rightarrow The **hidden state z_t** is constructed to be a graph signal **supported on the graph shift operator S**

- ▶ Hidden and observed state are **propagated through graph filters** to update the hidden state

$$A = A(S) = \sum_{k=0}^{K-1} a_k S^k \quad B = B(S)x = \sum_{k=0}^{K-1} b_k S^k$$

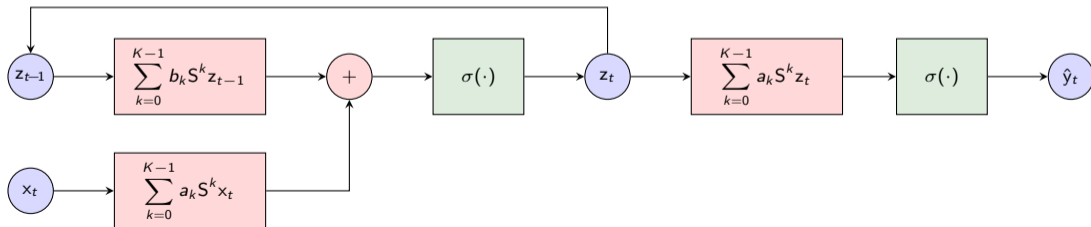
- ▶ The state update is $\Rightarrow z_t = \sigma \left[A(S)x_t + B(S)z_{t-1} \right] = \sigma \left[\sum_{k=0}^{K-1} a_k S^k x_t + \sum_{k=0}^{K-1} b_k S^k z_{t-1} \right]$



- ▶ The hidden state z_t is **propagated through a graph filter** to make a prediction \hat{y}_t of the output y_t

$$C = C(S) = \sum_{k=0}^{K-1} c_k S^k$$

- ▶ The prediction of the output y_t is given by $\Rightarrow \hat{y}_t = \sigma \left[C(S)z_t \right] = \sigma \left[\sum_{k=0}^{K-1} c_k S^k z_t \right]$



- ▶ A GRNN is made up a hidden state update perceptron and an output prediction perceptron

$$\mathbf{z}_t = \sigma \left[\sum_{k=0}^{K-1} \mathbf{a}_k \mathbf{s}^k \mathbf{x}_t + \sum_{k=0}^{K-1} \mathbf{b}_k \mathbf{s}^k \mathbf{z}_{t-1} \right] \quad \hat{\mathbf{y}}_t = \sigma \left[\sum_{k=0}^{K-1} \mathbf{c}_k \mathbf{s}^k \mathbf{z}_t \right]$$

- ▶ Each of these filters can be replaced by a **MIMO filter** to yield a GRNN with multiple features

$$\mathbf{Z}_t = \sigma \left[\sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{X}_t \mathbf{A}_k + \sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{Z}_{t-1} \mathbf{B}_k \right] \quad \hat{\mathbf{Y}}_t = \sigma \left[\sum_{k=0}^{K-1} \mathbf{S}^k \mathbf{Z}_t \mathbf{C}_k \right]$$

- ▶ Multiple-feature hidden state \mathbf{Z}_t permits larger dimensionality relative to observed states

Spatial Gating

- ▶ We extend **time gating** to GRNNs to handle the problem of **vanishing/exploding gradients**
- ▶ We discuss **long range graph dependencies** and introduce **node and edge gating**

- ▶ Like RNNs, GRNNs may also experience the problem of **vanishing/exploding gradients**

⇒ Happens when eigenvalues of $B(S)$ are **much smaller/larger than 1**

- ▶ Similarly to what we did for RNNs, we address it by adding gating operators to GRNNs

$$Z_t = \sigma \left(\hat{Q} \{A_S(X_t)\} + \check{Q} \{B_S(Z_{t-1})\} \right)$$

- ▶ **Input gate** operator $\hat{Q} : \mathbb{R}^{N \times H} \rightarrow \mathbb{R}^{N \times H} \Rightarrow$ controls the importance of the input X_t at time t

- ▶ **Forget gate** operator $\check{Q} : \mathbb{R}^{N \times H} \rightarrow \mathbb{R}^{N \times H} \Rightarrow$ controls the importance of the state Z_t at time t

- ▶ First type of gating for GRNNs is **time gating** \Rightarrow simple extension of input and forget gates of RNNs
- ▶ In the **Time-Gated GRNN**, the input and forget gate operators are expressed as

$$\hat{Q}\{A_S(X_t)\} = \hat{q}_t A_S(X_t), \quad \check{Q}\{B_S(Z_t)\} = \check{q}_t B_S(Z_t)$$

- ▶ Time gating multiplies the input and the state by scalar gates $\hat{q}_t \in [0, 1]$ and $\check{q}_t \in [0, 1]$
- ▶ A single scalar gate is applied to the whole graph signal \Rightarrow same gate value for all nodes

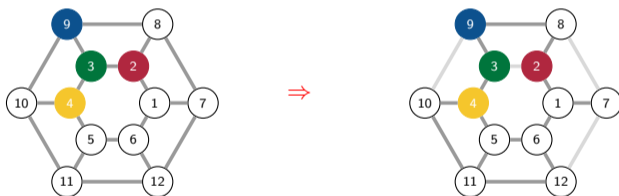
- ▶ Even if eigenvalues of $B(S) \sim 1$ **spatial imbalances** can cause gradients to vanish in **space**
 - ⇒ Some nodes/paths might get assigned more importance than others in long range exchanges
- ▶ **Example:** graphs with community structure, where some nodes are highly connected within clusters
 - ⇒ Gradients of Z_T depend on successive products of $B(S)$ ⇒ successive products of S
 - ⇒ For large T , the matrix entries in S^T with **highly connected nodes** will get densely populated
 - ⇒ **Overshadows community structure** ⇒ can't **encode long processes** that are **local** on the graph

- ▶ Node and edge structure of the graph allows for other forms of gating \Rightarrow **spatial gating**
- ▶ **Node gating** \Rightarrow one input and one forget gate for **each node of the graph**



- ▶ Spatial gating strategies help encode long range spatial dependencies in graph processes

- ▶ Node and edge structure of the graph allows for other forms of gating \Rightarrow **spatial gating**
- ▶ **Edge gating** \Rightarrow one input and one forget gate for **each edge of the graph**



- ▶ Spatial gating strategies help encode long range spatial dependencies in graph processes

- ▶ In the Node-Gated GRNN, the **input gate** and **forget gate** operators are expressed as

$$\hat{Q}\{\mathcal{A}_S(X_t)\} = \text{diag}(\hat{q}_t)\mathcal{A}_S(X_t), \quad \check{Q}\{\mathcal{B}_S(Z_t)\} = \text{diag}(\check{q}_t)\mathcal{B}_S(Z_t)$$

- ▶ Gating operators correspond to **multiplication of the input and state by diagonal matrices**

⇒ The diagonals are the input and forget vector gates $\hat{q}_t \in [0, 1]^N$ and $\check{q}_t \in [0, 1]^N$

- ▶ A scalar gate applied to each nodal component of the signal ⇒ different gate values for each node

- ▶ In the Edge-Gated GRNN, the **input gate** and **forget gate** operators are expressed as

$$\hat{Q} \{ \mathcal{A}_S(X_t) \} = \mathcal{A}_{S \odot \hat{Q}_t}(X_t), \quad \check{Q} \{ \mathcal{B}_S(Z_t) \} = \mathcal{B}_{S \odot \check{Q}_t}(Z_t)$$

- ▶ Gating operators correspond to **elementwise multiplication of the shift operator by gate matrices**

⇒ The matrices multiplying the GSOs are the input and forget matrix gates \hat{Q}_t and $\check{Q}_t \in [0, 1]^{N \times N}$

- ▶ Separate gate for **each edge** ⇒ control the amount of information transmitted across edges

- ▶ Parameters of input and forget gate operators are the outputs of GRNNs themselves
- ▶ **Input** and **forget** gate states are expressed as

$$\hat{\mathbf{Z}}_t = \hat{\sigma} \left(\hat{\mathcal{A}}_S(\mathbf{X}_t) + \hat{\mathcal{B}}_S(\hat{\mathbf{Z}}_{t-1}) \right) \quad \check{\mathbf{Z}}_t = \check{\sigma} \left(\check{\mathcal{A}}_S(\mathbf{X}_t) + \check{\mathcal{B}}_S(\check{\mathbf{Z}}_{t-1}) \right)$$

⇒ Gate computation takes **different forms** depending on the **type of gating**

- ▶ In the case of **time gating**, the gates are calculated as

$$\hat{\mathbf{q}}_t = \text{sigmoid}(\hat{\mathbf{c}}^T \text{vec}(\hat{\mathbf{Z}}_t)) \quad \check{\mathbf{q}}_t = \text{sigmoid}(\check{\mathbf{c}}^T \text{vec}(\check{\mathbf{Z}}_t))$$

⇒ Where $\hat{\mathbf{c}} \in \mathbb{R}^{\hat{H}N}$ and $\check{\mathbf{c}} \in \mathbb{R}^{\check{H}N}$ are fully connected layers and the sigmoid ensures gates in $[0, 1]$

- ▶ Parameters of input and forget gate operators are the outputs of GRNNs themselves
- ▶ **Input** and **forget** gate states are expressed as

$$\hat{\mathbf{Z}}_t = \hat{\sigma} \left(\hat{\mathcal{A}}_S(\mathbf{X}_t) + \hat{\mathcal{B}}_S(\hat{\mathbf{Z}}_{t-1}) \right) \quad \check{\mathbf{Z}}_t = \check{\sigma} \left(\check{\mathcal{A}}_S(\mathbf{X}_t) + \check{\mathcal{B}}_S(\check{\mathbf{Z}}_{t-1}) \right)$$

⇒ Gate computation takes **different forms** depending on the **type of gating**

- ▶ In the case of **node gating**, the gates are calculated as

$$\hat{\mathbf{q}}_t = \text{sigmoid} \left(\hat{\mathcal{C}}_S(\hat{\mathbf{Z}}_t) \right) \quad \check{\mathbf{q}}_t = \text{sigmoid} \left(\check{\mathcal{C}}_S(\check{\mathbf{Z}}_t) \right)$$

⇒ Where $\hat{\mathcal{C}}_S$ and $\check{\mathcal{C}}_S$ are graph convolutions and the sigmoid ensures gates in $[0, 1]^N$

- ▶ Parameters of input and forget gate operators are the outputs of GRNNs themselves
- ▶ **Input** and **forget** gate states are expressed as

$$\hat{Z}_t = \hat{\sigma} \left(\hat{A}_S(X_t) + \hat{B}_S(\hat{Z}_{t-1}) \right) \quad \check{Z}_t = \check{\sigma} \left(\check{A}_S(X_t) + \check{B}_S(\check{Z}_{t-1}) \right)$$

⇒ Gate computation takes **different forms** depending on the **type of gating**

- ▶ In the case of **edge gating**, the gates are calculated as

$$[\hat{Q}_t]_{ij} = \text{sigmoid} \left(\hat{c}^T [\delta_i^T \hat{Z}_t \hat{C} \parallel \delta_j^T \hat{Z}_t \hat{C}]^T \right) \quad [\check{Q}_t]_{ij} = \text{sigmoid} \left(\check{c}^T [\delta_i^T \check{Z}_t \check{C} \parallel \delta_j^T \check{Z}_t \check{C}]^T \right)$$

⇒ Where δ_i and δ_j are N -dimensional Dirac deltas; $\hat{C} \in \mathbb{R}^{\hat{H} \times \hat{H}'}$, $\check{C} \in \mathbb{R}^{\check{H} \times \check{H}'}$ are linear layers

⇒ And $\hat{c} \in \mathbb{R}^{2\hat{H}' \times 1}$ and $\check{c} \in \mathbb{R}^{2\check{H}' \times 1}$ are f.c. layers applied to concatenation \parallel of features of i and j

Stability of GRNNs

- ▶ GRNNs can be seen as a time extension of GNNs, therefore they inherit their stability properties

Definition (Relative perturbation matrices)

Given GSOs S and \tilde{S} , we define the set of relative perturbation matrices modulo permutation as

$$\mathcal{E}(S, \tilde{S}) = \left\{ E \in \mathbb{R}^{N \times N} : P^T \tilde{S} P = S + ES + SE^T, P \in \mathcal{P} \right\} \quad (1)$$

where $\mathcal{P} = \{P \in \{0, 1\}^{N \times N} : P\mathbf{1} = \mathbf{1}, P^T\mathbf{1} = \mathbf{1}\}$.

- ▶ We consider that the distance between two graphs S and \tilde{S} is given by $d(S, \tilde{S}) = \min_{E \in \mathcal{E}(S, \tilde{S})} \|E\|$
- ▶ Notice that if \tilde{S} is a permutation of the shift matrix S , then we have $d(S, \tilde{S}) = 0$

Definition (Integral Lipschitz filters)

A filter $A(S) = \sum_{k=0}^{K-1} a_k S^k$ is integral Lipschitz if there exists $C > 0$ such that $a(\lambda) = \sum_{k=0}^{K-1} a_k \lambda^k$ satisfies

$$|a(\lambda_2) - a(\lambda_1)| \leq C \frac{|\lambda_2 - \lambda_1|}{|\lambda_1 + \lambda_2|/2} \quad (2)$$

for all $\lambda_1, \lambda_2 \in \mathbb{R}$.

- ▶ Integral Lipschitz filters also satisfy $|\lambda a'(\lambda)| \leq C$, where $a'(\lambda)$ is the derivative of $a(\lambda)$
- ▶ Recall that the frequency response of integral Lipschitz filters becomes **flat for large λ**

- ▶ We consider a GRNN with $F_X = 1$ input feature, $F_Z = 1$ state feature, and $F_Y = 1$ output feature

$$z_t = \sigma(A(S)x_t + B(S)z_{t-1}) \quad \hat{y}_t = \rho(C(S)z_t)$$

- (A1) A, B and C are integral Lipschitz with constants C_A , C_B and C_C and $\|A\| = \|B\| = \|C\| = 1$
- (A2) Nonlinearities σ and ρ satisfy: $|\sigma(b) - \sigma(a)| \leq |b - a|$ for all $a, b \in \mathbb{R}$, $\sigma(0) = \rho(0) = 0$
- (A3) Initial hidden state is identically zero, i.e., $z_0 = 0$, and the x_t satisfy $\|x_t\| \leq \|x\| = 1$ for all t

Theorem (Stability of GRNNs)

Let $S = V\Lambda V^H$ and \tilde{S} be the GSOs of the original and perturbed graph, and let $E = U\mu U^H \in \mathcal{E}(S, \tilde{S})$ such that $d(S, \tilde{S}) \leq \|E\| \leq \varepsilon$. Let y_t and \tilde{y}_t be the outputs of the GRNNs running on S and \tilde{S} respectively, and satisfying assumptions (A1)-(A3). Then,

$$\min_{P \in \mathcal{P}} \|y_t - P^T \tilde{y}_t\| \leq C(1 + \sqrt{N}\delta)(t^2 + 3t)\varepsilon + \mathcal{O}(\varepsilon^2) \quad (3)$$

where $C = \max\{C_A, C_B, C_C\}$ and $\delta = (\|U - V\| + 1)^2 - 1$.

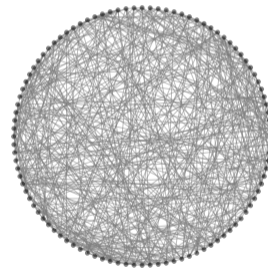
- ▶ GRNNs are stable to relative perturbations with constant $C(1 + \sqrt{N}\delta)(T^2 + 3T)$, T process length
- ▶ C could be set at a **fixed value** or **learned** from data through A , B and $C \Rightarrow$ design parameter
- ▶ The term $(1 + \delta\sqrt{N})$ is a property of the graph perturbation \Rightarrow **cannot be controlled by design**
- ▶ Eigenvector misalignment $\delta = (\|U - V\| + 1)^2 - 1$ **measures commutativity** of matrices S and E
- ▶ Polynomial dependence on $T \Rightarrow$ due to recurrence relationship in the computation of x_t

Epidemic Modeling with GRNNs

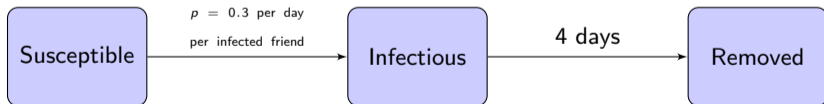
- ▶ We use a GRNN, a GNN, and a RNN to track an epidemic on a high school friendship network

- ▶ Model the spread of an infectious disease over a friendship network as a **graph process**
- ▶ Graph is a symmetric **friendship network** corresponding to a high school in France
- ▶ Model the spread of the disease on the graph using **Susceptible-Infectious-Removed (SIR)** model
- ▶ Compare the performance of a **GRNN, a RNN, and a GNN** in predicting infections after 8 days

- ▶ Real-world friendship network corresponding to 134 students from a high school in Marseille
- ▶ Each **node of the graph** represents a **student**
- ▶ **Friendships** are modeled as **symmetric unweighted edges**
- ▶ Isolated nodes are removed to make the graph fully connected
- ▶ **Assumption:** friends are likely to **be in contact** with each other



- ▶ Process starts with random seed infections on day 0 \Rightarrow probability $p_{seed} = 0.05$
- ▶ Each person is in one of the three SIR states \Rightarrow updated each day with the following rules
- ▶ **Susceptible**: can get the disease from an infected friend with probability $p_{inf} = 0.3$
- ▶ **Infectious**: can spread the disease for 4 days after being infected, after which they recover
- ▶ **Removed**: have overcome the disease and can no longer spread it or contract it



- ▶ **Problem:** given the node states, goal is to **predict whether each node will be infected in 8 days**
- ▶ **Input:** graph process x_t where, at each time t , $[x_t]_i$ is given by

$$[x_t]_i = \begin{cases} 0, & \text{if student } i \text{ is susceptible} \\ 1, & \text{if student } i \text{ is infected} \\ 2, & \text{if student } i \text{ is removed} \end{cases}$$

- ▶ **Output:** binary graph process $y_t \Rightarrow$ our goal is only to track infections

$$[y_t]_i = \begin{cases} 0, & \text{if student } i \text{ is susceptible or removed} \\ 1, & \text{if student } i \text{ is infected} \end{cases}$$

- ▶ Given $x_t, x_{t+1}, \dots, x_{t+7}$, we want to predict $y_{t+8}, y_{t+9}, \dots, y_{t+15} \Rightarrow$ **binary node classification**

- ▶ Accuracy is not a good performance metric \Rightarrow does not distinguish true positives and true negatives
- ▶ In epidemic tracking, true positives are more important than true negatives \Rightarrow maximize F1 score

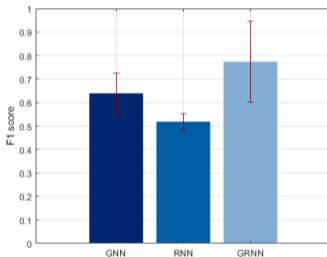
$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

- ▶ **Precision** = True Positive/Predicted Positive
 \Rightarrow Proportion of correct positive predictions
- ▶ **Recall** = True Positive/All Actual Positive
 \Rightarrow Proportion of correctly predicted positives

		Actual	
		Positive	Negative
Predicted	Positive	True Positive	False Positive
	Negative	False Negative	True Negative

- ▶ Loss function we minimize is $1 - F1 \Rightarrow$ trade-off between minimizing FPs and FNs

- ▶ We compare a GRNN with a GNN and a RNN, all with roughly the same number of parameters
 - ⇒ In the GNN, the **time instants** become input features ⇒ parameters depend on T
 - ⇒ In the RNN, the **nodal components** become input features ⇒ parameters depend on N



- ▶ GRNN improves upon RNN and GNN ⇒ exploits **both spatial and temporal structure** of the data