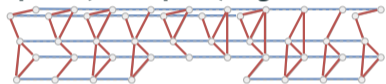


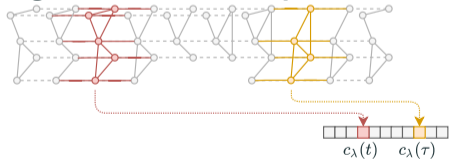
Latent graph learning

Discovering insightful correlation patterns

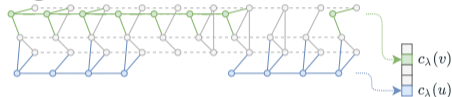
Spatial (or temporal) edges



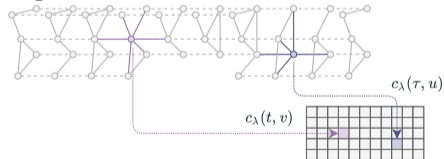
Edges related to a time step



Edges related to a node

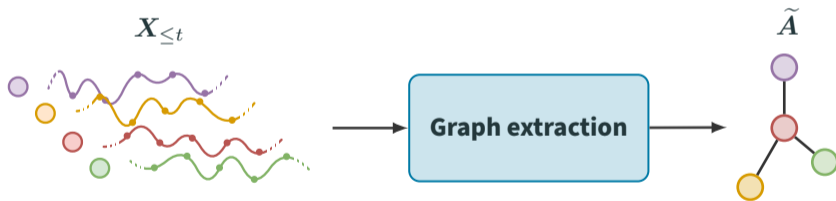


Edges related to a node



Learning and adjacency matrix

- ☹ Relational information is **not** always **available**
- ☹ or might be **ineffective** in capturing spatial dynamics.
- 😊 Relational architectural **biases** can nonetheless be exploited
→ **extract a graph** from the time series or node attributes

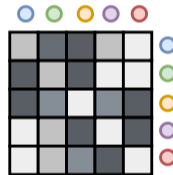


- It can be interpreted as **regularizing a spatial attention** operator.

Time-series similarities

Probably, the simplest approach to extract a graph from the time series is by computing **time series similarity scores**.

- Pearson correlation
- Correntropy
- Granger causality
- Kernels for time series
- ...



→ Thresholding might be necessary to obtain binary and sparse graphs.

Latent graph learning

An integrated approach: **learn** the **relations end-to-end** with the downstream task

- as a function of the **input** data,
- as trainable **parameters** of the model,
- or **both**.

This problem is known as **latent graph learning** (or latent graph inference).

Two different approaches:

1. learning directly an **adjacency matrix** $\tilde{\mathbf{A}} \in \mathbb{R}^{N \times N}$;
2. learning a **probability distribution over graphs** p_{Φ} generating $\tilde{\mathbf{A}}$.

! One key challenge is keeping both $\tilde{\mathbf{A}}$ and the subsequent computations **sparse**.
→ challenging with gradient-based optimization.

Direct approach

A direct approach consists in learning $\tilde{\mathbf{A}}$ as function $\xi(\cdot)$ of edge scores $\Phi \in \mathbb{R}^{N \times N}$ as

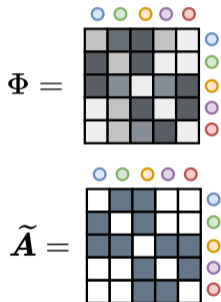
$$\tilde{\mathbf{A}} = \xi(\Phi)$$

Edge scores Φ

- can be a table of **learnable** model **parameters**,
- obtained as a **function** of the **inputs** and/or other parameters.

Function $\xi(\cdot)$ is a nonlinear activation

- it can be exploited to **make** $\tilde{\mathbf{A}}$ **sparse**.



Direct approach: factorization methods

Many of the methods directly learning $\tilde{\mathbf{A}}$, learn a **factorization** of the former to amortize the cost of the inference:

$$\tilde{\mathbf{A}} = \xi(\Phi)$$

$$\Phi = \mathbf{Z}_s \mathbf{Z}_t^\top$$

with

- $\mathbf{Z}_s \in \mathbb{R}^{N \times d}$ **source** node embeddings
- $\mathbf{Z}_t \in \mathbb{R}^{N \times d}$ **target** node embeddings

\mathbf{Z}_s and \mathbf{Z}_t can be learned as tables of (local) parameters or **as a function of the input window**.

Pro & Cons of the direct approach

- 😊 Easy to implement.
- 😊 Many possible parametrizations.
- 😊 Edge scores are usually easy to learn end-to-end.
- 😞 It often results in dense computations with $\mathcal{O}(N^2)$ complexity.
- 😞 Sparsifying \tilde{A} results in sparse gradients.
- 😞 Encoding prior structural information requires smart parametrizations.

Probabilistic methods


In this context, probabilistic methods aim at learning a parametric distribution p_{Φ} for $\tilde{\mathbf{A}}$ by minimizing

$$\mathcal{L}(\Phi) = \mathbb{E}_{\hat{\mathbf{A}} \sim p_{\Phi}} \left[\ell \left(\widehat{\mathbf{X}}_{t:t+H}, \mathbf{X}_{t:t+H} \right) \right]. \quad (15)$$

- Again, we can factorize Φ and make p_{Φ} input dependent, e.g.,

$$\Phi = \xi \left(\mathbf{Z}_s \mathbf{Z}_t^{\top} \right) \quad \tilde{\mathbf{A}} \sim p_{\Phi} \left(\mathbf{A} | \mathbf{X}_{<t}, \mathbf{U}_{<t}, \mathbf{V} \right)$$

- Different parametrizations of p_{Φ} allow for embedding **sparsity priors** on the sampled graphs [22].

 Gradient-based optimization requires $\nabla_{\Phi} \mathcal{L}(\Phi)$
 → it can be **challenging** and **computationally expensive**.

[22] A. Cini *et al.*, “Sparse graph learning from spatiotemporal time series”, JMLR 2023.

Monte Carlo gradient estimators

💡 One approach is to **reparametrize** $\tilde{\mathbf{A}} \sim p_{\Phi}(\mathbf{A})$ as: $\tilde{\mathbf{A}} = g(\Phi, \varepsilon)$, $\varepsilon \sim p(\varepsilon)$
 decoupling parameters Φ from the random component ε : $\nabla_{\Phi} \mathcal{L}(\Phi) = \mathbb{E}_{\varepsilon} \left[\nabla_{\Phi} \ell(\widehat{\mathbf{X}}, \mathbf{X}) \right]$.

😊 Practical and **easy** to implement,

😞 rely on **continuous relaxations** and make subsequent computations scale with $\mathcal{O}(N^2)$.

💡 Conversely, **score-function** (SF) gradient estimators rely on the relation

$$\nabla_{\Phi} \mathbb{E}_{p_{\Phi}} \left[\ell(\widehat{\mathbf{X}}, \mathbf{X}) \right] = \mathbb{E}_{p_{\Phi}} \left[\ell(\widehat{\mathbf{X}}, \mathbf{X}) \nabla_{\Phi} \log p_{\Phi} \right]$$

😞 suffer from **high variance** (use variance reduction techniques),

😊 allow to **keep computations sparse**.

→ we can use **Monte Carlo** estimator.

[24] T. Kipf *et al.*, “Neural relational inference for interacting systems”, ICML 2018.

[22] A. Cini *et al.*, “Sparse graph learning from spatiotemporal time series”, JMLR 2023.