

A Source Separation Approach to Temporal Graph Modelling for Computer Networks

Corentin Larroche

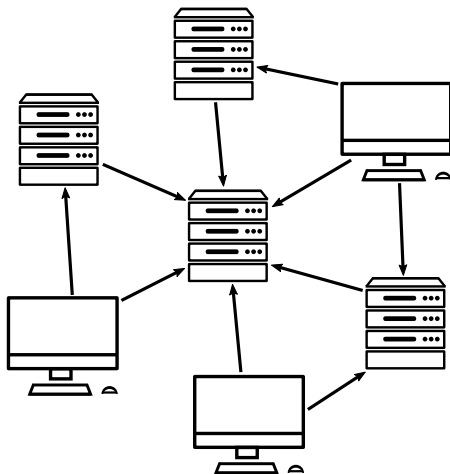
`corentin.larroche@ssi.gouv.fr`



French National Cybersecurity Agency (ANSSI), Paris, France

MLCS '23, Torino, Italy

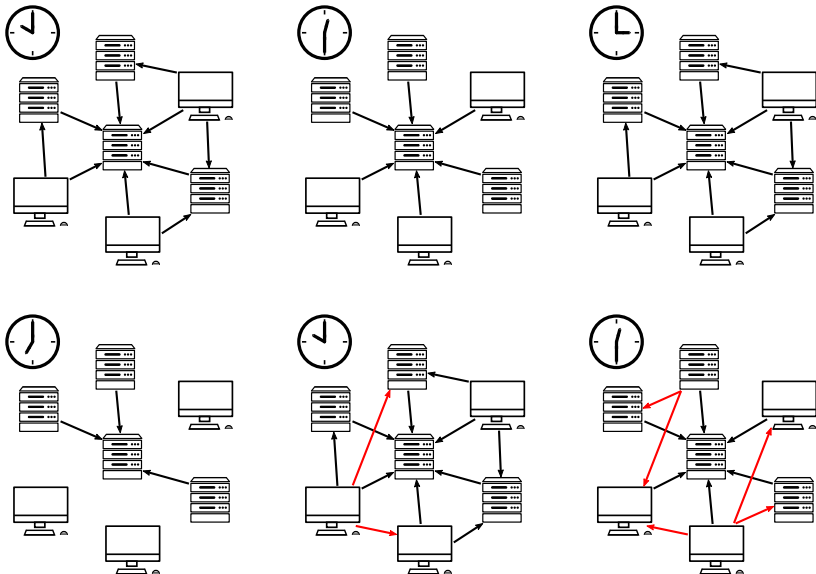
September 22nd, 2023

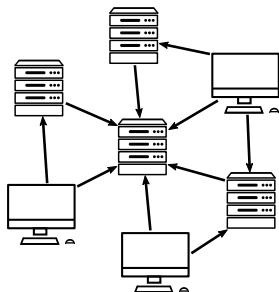


Network monitoring

- ▶ Traffic within an **enterprise network** analyzed to detect **malicious activity**
- ▶ Network traffic represented as a **directed graph**
- ▶ **Goal:** detect anomalous edges

Computer network monitoring – Temporal perspective





Definitions

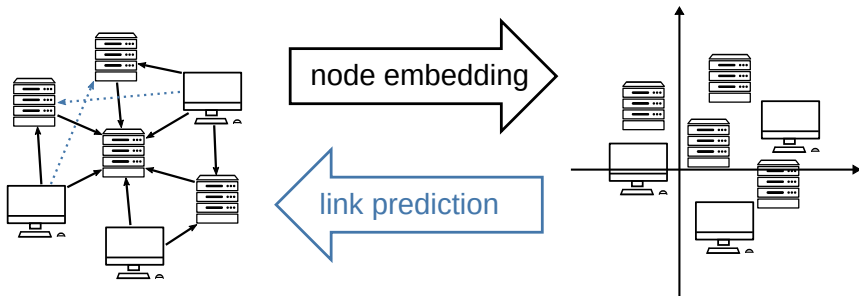
- ▶ $(\mathcal{G}_t)_{t \geq 1}$ is a sequence of directed graphs with **shared node set** $\mathcal{V} = \{1, \dots, N\}$
- ▶ \mathbf{A}_t is the (binary, non-symmetric) adjacency matrix of \mathcal{G}_t

- ▶ **Problem:** given normal graphs $\mathcal{G}_1, \dots, \mathcal{G}_T$, detect anomalous edges in new graphs \mathcal{G}_t for $t > T$
- ▶ Equivalent to **temporal link prediction**: if we can predict normal edges, we can detect anomalous ones

Latent space models (aka MF, graph embedding, GNNs)

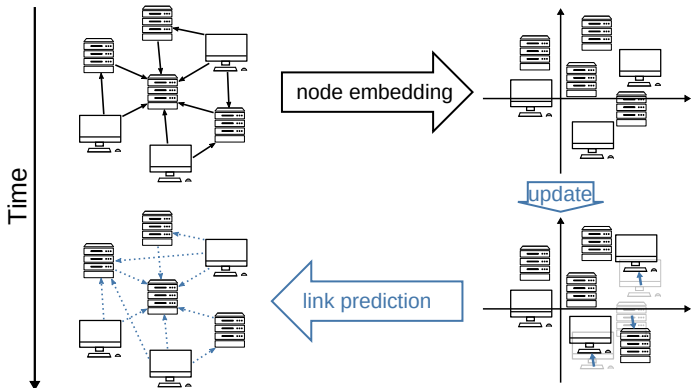
Workflow of most link prediction methods:

- ▶ Learn **node embeddings** $\{\mathbf{u}_i \in \mathbb{R}^K\}_{i \in \mathcal{V}}$ and **link predictor** $g : \mathbb{R}^K \times \mathbb{R}^K \rightarrow [0, 1]$ that maximize the probability of observed edges, $p(i, j) = g(\mathbf{u}_i, \mathbf{u}_j)$
- ▶ Predict new edges (i', j') using $g(\mathbf{u}_{i'}, \mathbf{u}_{j'})$



Generalizing latent space models to temporal graphs

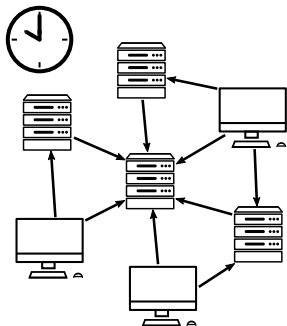
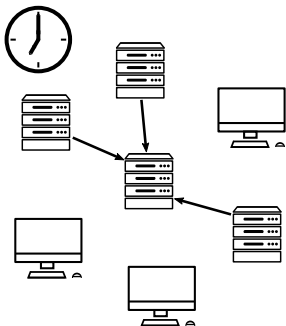
- ▶ Just replace \mathbf{u}_i with a sequence $(\mathbf{u}_{i,t})_{t \geq 1}$
- ▶ Predict future embeddings using recursive Bayesian estimation [Lee et al., 2022] or RNNs [King and Huang, 2022]



The problem with dynamic latent space models

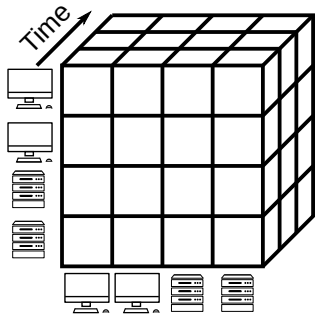
Dynamic LSM vs. network traffic dynamics

- ▶ In dynamic LSMs, temporal dynamics are **node-driven**: local, smooth evolution of the graph
- ▶ In enterprise networks, observed traffic undergoes **sharp, global variations**



Temporal graphs as tensors [Dunlavy et al., 2011]

The sequence $(\mathbf{A}_t)_{t=1,\dots,T}$ can be seen as a **3-mode tensor**, with modes standing for time step, origin and destination nodes.



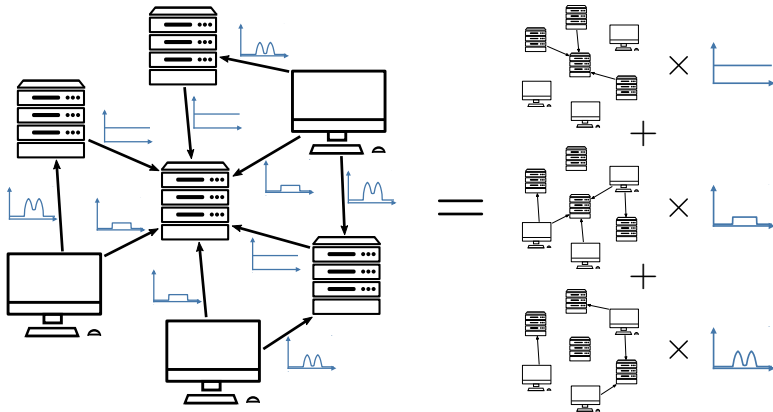
CANDECOMP/PARAFAC 3-mode tensor factorization:

$$\hat{\mathbf{A}}_t = \sum_{k=1}^K w_{tk} \mathbf{U}_k \mathbf{V}_k^{\top},$$

with $\mathbf{U}_k, \mathbf{V}_k \in \mathbb{R}^{N \times 1}$ origin and destination embedding matrices.

Tensor factorization – Source separation interpretation

$$\hat{\mathbf{A}}_t = \sum_{k=1}^K w_{tk} \mathbf{U}_k \mathbf{V}_k^\top$$



An improved source separation approach

Tensor factorization models each source as a **rank-one matrix**.

- ▶ **Idea:** use richer models for the activity sources

Our approach: Superposed Nonnegative Matrix Factorization (SNMF)

- ▶ For each of L activity sources, define origin and destination embedding matrices $\mathbf{U}_\ell, \mathbf{V}_\ell \in \mathbb{R}_+^{N \times K}$
- ▶ For each time step t , define mixing coefficients $w_{t\ell}$ such that $\mathbf{A}_t = \sum_{\ell=1}^L w_{t\ell} \mathbf{U}_\ell \mathbf{V}_\ell^\top$

Model inference: Given adjacency matrices $(\mathbf{A}_1, \dots, \mathbf{A}_T)$, find $\mathbb{U} = (\mathbf{U}_\ell)_{\ell=1}^L, \mathbb{V} = (\mathbf{V}_\ell)_{\ell=1}^L, \mathbf{W} = (w_{t\ell}) \in \mathbb{R}_+^{T \times L}$ minimizing

$$J(\mathbb{U}, \mathbb{V}, \mathbf{W}) = \frac{1}{2} \sum_{t=1}^T \left\| (\mathbf{1}_N - \mathbf{I}_N) \odot \left(\mathbf{A}_t - \sum_{\ell=1}^L w_{t\ell} \mathbf{U}_\ell \mathbf{V}_\ell^\top \right) \right\|_F^2 \\ + \lambda_1 \|\mathbf{W}\|_1 + \frac{\lambda_2}{2} \sum_{\ell=1}^L \left(\|\mathbf{U}_\ell\|_F^2 + \|\mathbf{V}_\ell\|_F^2 \right)$$

After inference ($t > T$), given a new adjacency matrix \mathbf{A}_t :

- ▶ Predict mixing coefficients $\hat{\mathbf{w}}_t$ using a **seasonal model**:

$$\hat{\mathbf{w}}_t = \text{MEAN}(\mathbf{w}_{t'}; t' \in [t-1] : t \equiv t' \pmod{\tau}),$$

with period $\tau > 0$ (one week here)

- ▶ Predict adjacency matrix $\hat{\mathbf{A}}_t$ using $\hat{\mathbf{w}}_t$ and compute **anomaly score matrix** $\mathbf{A}_t - \hat{\mathbf{A}}_t$
- ▶ Compute **true mixing coefficients** \mathbf{w}_t to update the seasonal model

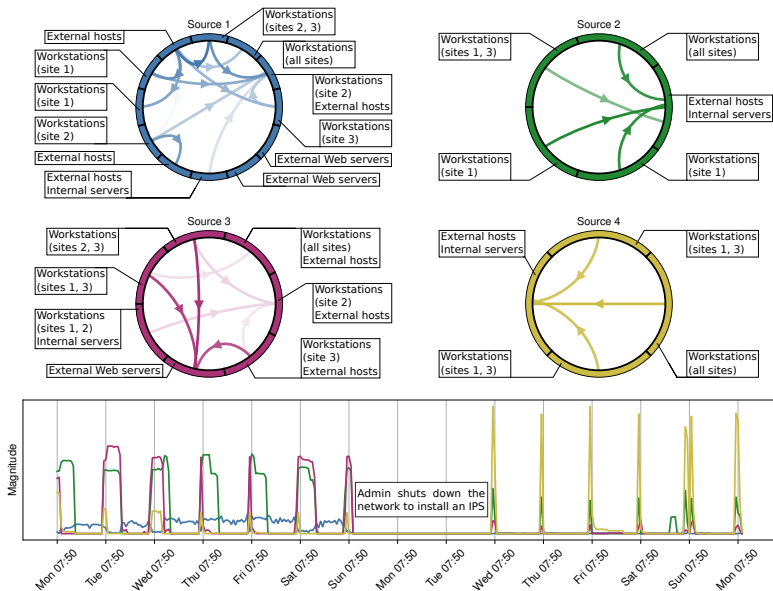
Dataset: VAST 2013 Mini-Challenge 3

- ▶ Simulated network traffic (enterprise network + external hosts)
- ▶ 1.4K hosts, 14 days
- ▶ Ground truth description of host roles and observed traffic

Goal: analyze the activity sources inferred by SNMF.

- ▶ For each source:
 - ▶ Cluster nodes based on their embeddings
 - ▶ Display predicted edges between clusters
- ▶ Plot mixing coefficients over time

Qualitative evaluation – Results



Dataset: LANL CMSCSE

- ▶ Real-world enterprise network
- ▶ 12.7K hosts, 30 days
- ▶ Labelled malicious edges

Evaluation tasks

- ▶ Anomaly detection
- ▶ Temporal link prediction with sampled negative edges

Baseline methods

- ▶ Dynamic link predictors:
 - ▶ Tensor factorization (**PTF** [Eren et al., 2020])
 - ▶ Latent space model (**BME** [Lee et al., 2022])
- ▶ Static link predictors:
 - ▶ Matrix factorization (**HPF** [Sanna Passino et al., 2022])
 - ▶ Graph embedding (**GL-GV** [Bowman et al., 2020])
- ▶ Generic/naive baselines: **EDGEBANK** [Poursafaei et al., 2022], **SEDANSPOT** [Eswaran and Faloutsos, 2018]

SNMF outperforms baselines in terms of AUC on anomaly detection and historical link prediction tasks.

Method	Anomaly	Random	Historical	Inductive
SNMF	99.1\pm0.1	98.4 \pm 0.0	76.9\pm0.2	98.2 \pm 0.1
PTF	97.6 \pm 0.9	98.6 \pm 0.0	68.5 \pm 0.2	96.7 \pm 0.3
BME	90.3 \pm 0.2	98.5 \pm 0.0	73.3 \pm 0.0	96.2 \pm 0.0
HPF	97.7 \pm 0.3	99.1\pm0.0	69.6 \pm 0.1	97.5 \pm 0.0
GL-GV	87.0 \pm 2.4	95.8 \pm 1.0	61.2 \pm 0.5	74.6 \pm 1.9
SEDANSPOT	63.6 \pm 7.3	51.2 \pm 2.2	54.7 \pm 1.4	53.2 \pm 2.7
EDGEBANK $_{\infty}$	96.2 \pm 0.0	97.2 \pm 0.0	56.0 \pm 0.0	98.3\pm0.0
EDGEBANK $_w$	96.0 \pm 0.0	97.0 \pm 0.0	58.0 \pm 0.0	98.1 \pm 0.0

Conclusion: The nature of temporal dynamics matters

Key takeaways

- ▶ Activity within enterprise networks has **specific temporal dynamics**
- ▶ Temporal graph models introduced in other domains are **not well-suited** to this context
- ▶ However, a **simple source separation approach** fits these dynamics rather well

Future research directions:

- ▶ Extension to **more complex data representations** than graphs
- ▶ **Better prediction** of the mixing coefficients
- ▶ Including **long-term dynamics** (concept drift, new nodes)

- [Bowman et al., 2020] Bowman, B., Laprade, C., Ji, Y., and Huang, H. H. (2020). Detecting lateral movement in enterprise computer networks with unsupervised graph AI. In *RAID*.
- [Dunlavy et al., 2011] Dunlavy, D. M., Kolda, T. G., and Acar, E. (2011). Temporal link prediction using matrix and tensor factorizations. *ACM Trans. Knowl. Discov. Data*, 5(2):1–27.
- [Eren et al., 2020] Eren, M. E., Moore, J. S., and Alexandro, B. S. (2020). Multi-dimensional anomalous entity detection via Poisson tensor factorization. In *ISI*.
- [Eswaran and Faloutsos, 2018] Eswaran, D. and Faloutsos, C. (2018). Sedanspot: Detecting anomalies in edge streams. In *ICDM*.
- [King and Huang, 2022] King, I. J. and Huang, H. H. (2022). EULER: Detecting network lateral movement via scalable temporal link prediction. In *NDSS*.

- [Lee et al., 2022] Lee, W., McCormick, T. H., Neil, J., Sodja, C., and Cui, Y. (2022). Anomaly detection in large-scale networks with latent space models. *Technometrics*, 64(2):241–252.
- [Poursafaei et al., 2022] Poursafaei, F., Huang, S., Pelrine, K., and Rabbany, R. (2022). Towards better evaluation for dynamic link prediction. In *NeurIPS*.
- [Sanna Passino et al., 2022] Sanna Passino, F., Turcotte, M. J., and Heard, N. A. (2022). Graph link prediction in computer networks using Poisson matrix factorisation. *Ann. Appl. Stat.*, 16(3):1313–1332.