



DATA-DRIVEN NETWORK INFERENCE FOR INFRASTRUCTURE SYSTEMS

Estéban Nocet-Binois, supervised by Jürgen Hackl

Department of Civil & Environmental Engineering | Princeton University, NJ, USA

MOTIVATION

The objective is to learn a basis that captures the dependencies and global modes of variation (called graph harmonics, or *frequencies*) present in the data. Analyzing how system states (like frequency deviations) decompose into these modes reveals which patterns are most likely to persist or amplify after a disturbance. Thus, given a *meaningful* basis, we can analyze how signals propagate, diffuse, or concentrate, and how oscillations or disturbances (such as power outages) manifest across the entire grid.



NOTATIONS

Graph Laplacian

- The Laplacian $L \in \mathbb{R}^{n \times n}$ encodes how values at each node deviate from its neighborhood. We denote the eigendecomposition:

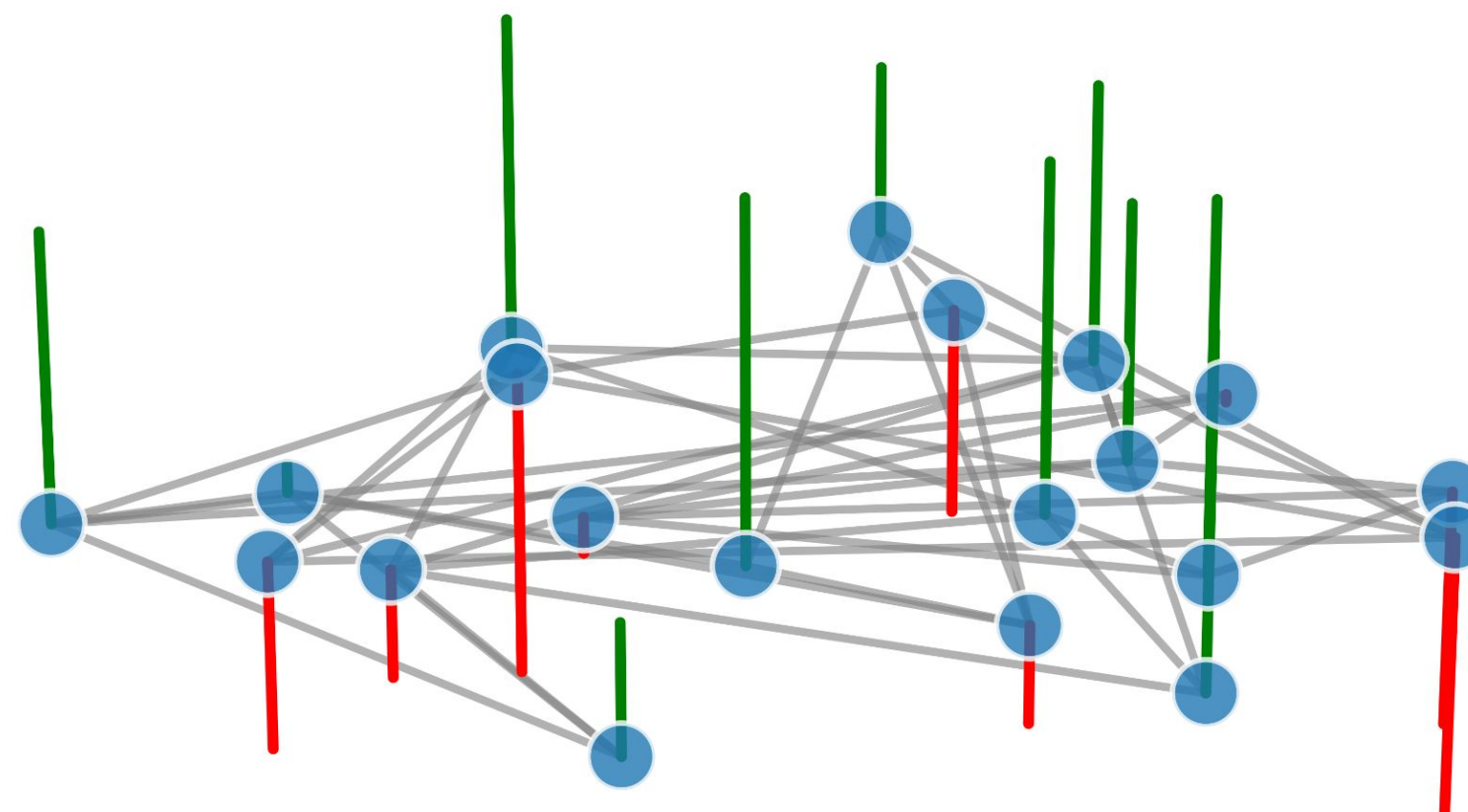
$$L = U \Lambda U^\top, \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

- We note the pseudoinverse:

$$L^\dagger = \sum_{k=1}^{n-1} \frac{1}{\lambda_k} \mathbf{u}_k \mathbf{u}_k^\top$$

Signal Matrix

- A graph signal is a n -valued tensor, collecting states across all nodes.



- Collection $X \in \mathbb{R}^{m \times n}$ of m signals (rows) over n nodes (columns)
- Sample covariance:

$$\Sigma = \frac{1}{m} X^\top X$$

The objective is, given a series of graphs signals, and optional adjacencies, to infer the dependency structure as edge weights.

LAPLACIAN LEARNING

L^\dagger encodes graph topology via $\mathbb{E}[\Sigma]$.

Edge Interpretation

- $\Omega = L$: Precision matrix
- $\Omega_{i,j} = -W_{i,j}$ (partial correlations)
- $\Omega_{i,i} = D_i$ (degree)

Conditional Independence

$$x_i \perp x_j \mid X_{\setminus\{i,j\}} \iff L_{i,j} = 0$$

Maximum Likelihood

$$\max_{L \in \mathcal{L}} \log \det(L) - \text{tr}(L \Sigma)$$

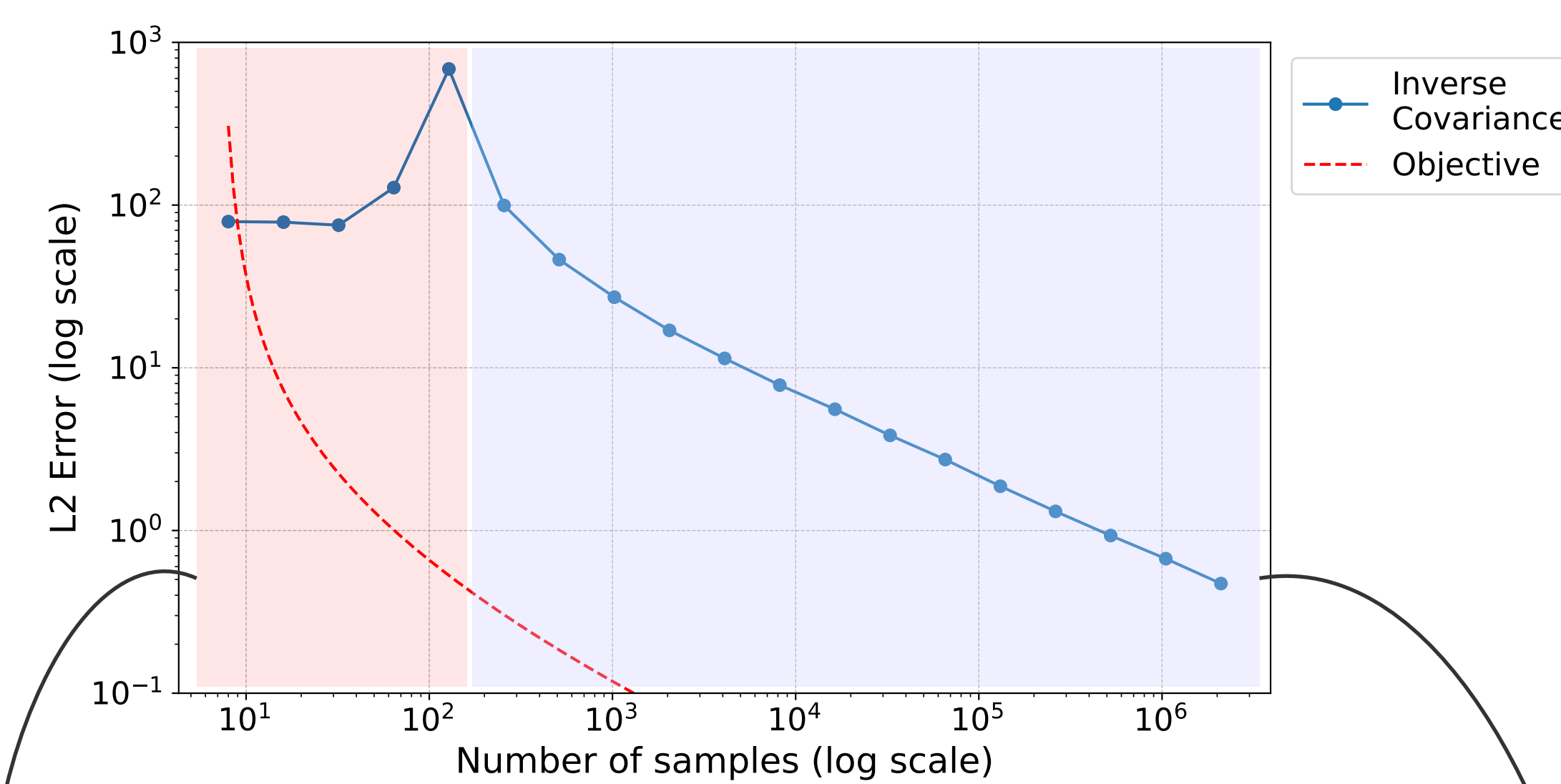
Promotes large eigenvalues (smooth surrogate for rank)
Penalizes deviation from data covariance Σ

$$\mathcal{L} = \{L : L \succeq 0, L \mathbf{1} = 0, L_{i \neq j} \leq 0\}$$

L is positive semidefinite
Rows sum to zero (centering)
Off-diagonal entries are non-positive

Note. We add the ridge term δI to ensure $L + \delta I \succ 0$.

Sampling Complexity



Spectral Relaxation

$$\hat{L} = \arg \min_{L \in \mathcal{L}} \left\{ \text{tr}(L \Sigma) + \alpha \|L\|_F^2 \right\}, \quad \alpha > 0$$

- Strongly convex due to Frobenius norm, while implicitly enforcing spectral properties.
- Stabilizes estimation when Σ is rank-deficient ($m < n$).

Sparsity Penalty

$$\min_{L \in \mathcal{L}} \left\{ \text{tr}(L \Sigma) + \alpha \psi(L) + \beta \sum_{i \neq j} \rho(L_{ij}) \right\}, \quad \alpha, \beta > 0$$

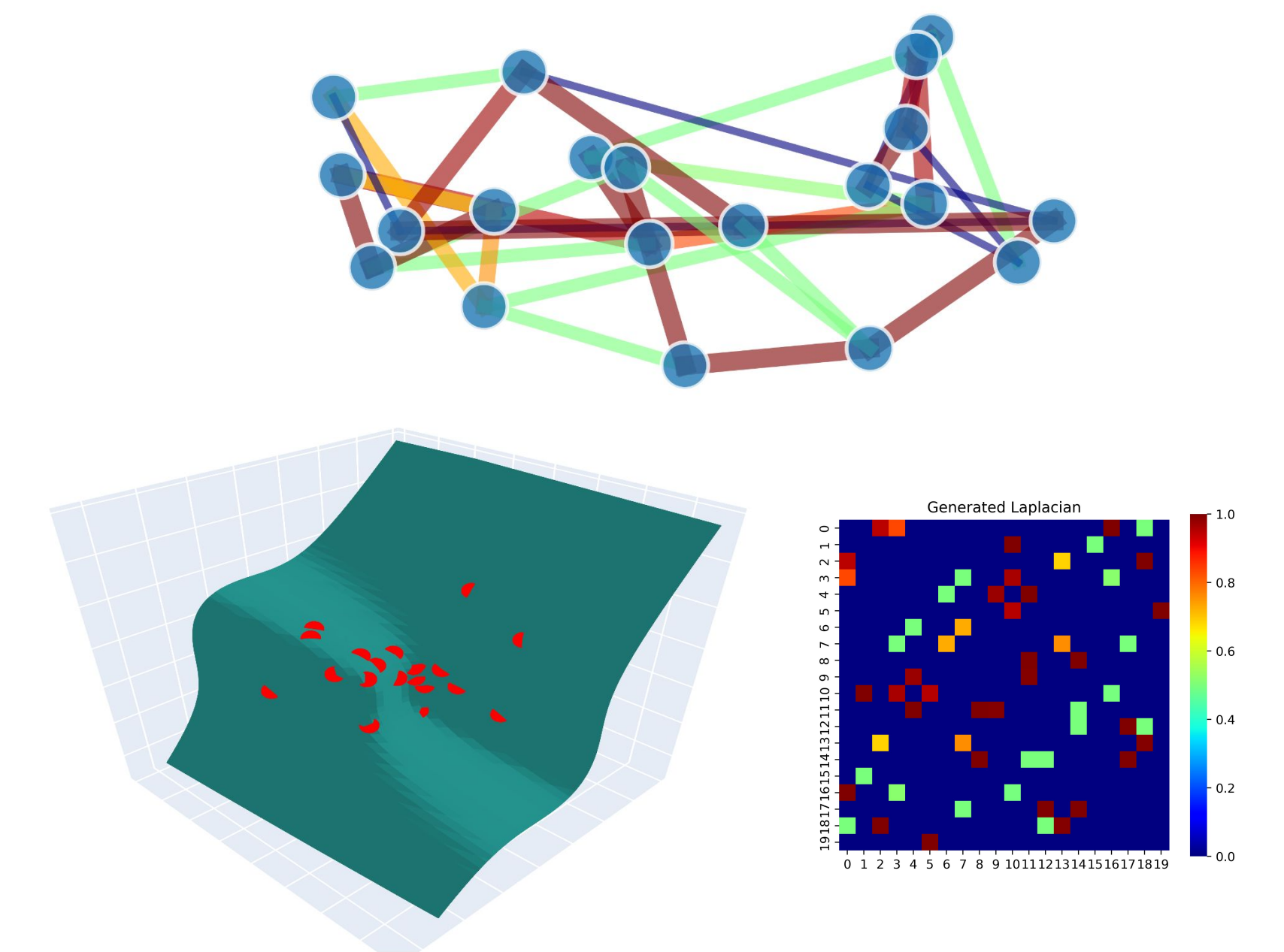
- The sparse penalty function $\rho(\cdot)$ (e.g. ℓ_1 , SCAD, MCP or LCP) prevents edge weight shrinkage.

MANIFOLD STRUCTURE

- As $n \rightarrow \infty$ and bandwidth $\epsilon \rightarrow 0$, $L \rightarrow \Delta$ (Laplace-Beltrami) under uniform sampling.
- Dirichlet energy, whose minimization solves $\Delta f = 0$:

$$\mathcal{E}(f) = \int_{\mathcal{M}} \|\nabla f\|^2 \leftrightarrow x^\top L x$$

- The Laplacian eigenvalues and eigenvectors provide a Fourier-like basis for graph signals, enabling frequency analysis.



The task shifts from learning a graph Laplacian that reflects data covariance, to identifying low-dimensional manifolds that preserve local neighborhoods.

- Diffusion maps learn a set of orthonormal basis functions—the eigenvectors of the diffusion operator H_t :

$$\frac{\partial f}{\partial t} = -L f \iff f(t) = \underbrace{e^{-tL}}_{H_t} f_0$$

FUTURE WORK

Two nodes with identical time series patterns but shifted in time, or with different overall magnitudes, will yield a larger Euclidian distance (hence smaller similarity). The question is: are these *meaningful*?

- Sampling rate: If node signals are sampled at different rates, it can artificially increase or decrease distance, since higher sampling rates in time series increase autocorrelation.
- Non-linear dependencies: For example, (constructive, or destructive) interference can lead to scale gaps that may not reflect actual dependencies.

Thus, the need to move to kernels that can capture non-linear dependencies:

- Can we generalize diffusion-based approaches to distortion-based ones?
- Can we go beyond i.i.d. time series?
- Can we have theoretical guarantees for nonlinear dynamics?

Application to Infrastructures

- Water flow in water distribution systems (WDSs)
- Water quality and propagation in sensor networks (over WDSs)
- Power grid stability