

Graph-based Semi-Supervised & Active Learning for Edge Flows

KDD2019

Junteng Jia Cornell University jj585@cornell.edu

Santiago Segarra Rice University segarra@rice.edu Michael T. Schaub Massachusetts Institute of Technology University of Oxford mschaub@mit.edu

> Austin R. Benson Cornell University arb@cs.cornell.edu

> > 2019.11.6

Given the labels of a subset of vertices, and our goal is to find a label assignment of the unlabeled vertices such that the labels vary smoothly across neighboring vertices.

 $\mathcal{E} = \{\mathcal{E}_1, \dots, \mathcal{E}_r, \dots, \mathcal{E}_m\}$ The incidence matrix B is defined as $B_{kr} = \begin{cases} 1, & \mathcal{E}_r \equiv (i,j), \ k = i, \ i < j \\ -1, & \mathcal{E}_r \equiv (i,j), \ k = j, \ i < j \\ 0, & \text{otherwise.} \end{cases}$

 \bigcirc labeled -1.00.0 +1.0

||B|This notion of smoothness can be defined via a loss function of the form

$$\mathbf{S}^{\mathsf{T}}\mathbf{y}\|^2 = \sum_{(i,j)\in\mathcal{E}} (y_i - y_j)^2$$

$$\|\mathbf{B}^{\mathsf{T}}\mathbf{y}\|^{2} = \mathbf{y}^{\mathsf{T}}\mathbf{L}\mathbf{y} \qquad \qquad \mathbf{L} = \mathbf{B}\mathbf{B}^{\mathsf{T}}$$
$$\mathbf{y}^{*} = \arg\min_{\mathbf{y}} \|\mathbf{B}^{\mathsf{T}}\mathbf{y}\|^{2} \qquad \text{s.t.} \quad y_{i} = \hat{y}_{i}, \ \forall \mathcal{V}_{i} \in \mathcal{V}^{\mathsf{L}}$$

Graph-Based SSL for Edge Flows

problem

The edge flows over a network can be represented with a vector f, where $f_r > 0$ if the flow orientation on edge r aligns with its reference orientation and $f_r < 0$ otherwise. In this sense, we are only accounting for the **net** flow along an edge.





To impose a **flow conservation assumption** for edge flows, we consider the divergence at each vertex, which is the sum of outgoing flows minus the sum of incoming flows at a vertex.

$$(\mathbf{Bf})_i = \sum_{\mathcal{E}_r \in \mathcal{E}} \sum_{: \mathcal{E}_r \equiv (i,j), i < j} \mathbf{f}_r - \sum_{\mathcal{E}_r \in \mathcal{E}} \sum_{: \mathcal{E}_r \equiv (j,i), j < i} \mathbf{f}_r.$$

loss

To create a loss function for edge flows that enforces a notion of flow-conservation, we use the sum-of-squares vertex divergence:

$$|\mathbf{B}\mathbf{f}||^2 = \mathbf{f}^\mathsf{T}\mathbf{B}^\mathsf{T}\mathbf{B}\mathbf{f} = \mathbf{f}^\mathsf{T}\mathbf{L}_e\mathbf{f}.$$

However, unlike the case for smooth vertex labels, requiring $\mathbf{f}^{\mathsf{T}}\mathbf{L}_{e}\mathbf{f} = 0$ is actually under-constrained, i.e., even when more than one edge is labeled, many different divergence-free edge-flow assignments may exist that induce zero loss.

$$\mathbf{f}^* = \arg\min_{\mathbf{f}} \|\mathbf{B}\mathbf{f}\|^2 + \lambda^2 \cdot \|\mathbf{f}\|^2$$

s.t.
$$f_r = \hat{f}_r, \forall \mathcal{E}_r \in \mathcal{E}^L$$
.



Computation

Let \mathbf{f}^0 be a trival feasible point where $\mathbf{f}_r^0 = \hat{\mathbf{f}}_r$ if $r \in \varepsilon^L$ and $\mathbf{f}^r = 0$ otherwise.

Denote the set of indices for unlabeled edges as $\mathcal{E}^{U} = \{\mathcal{E}_{1}^{U}, \mathcal{E}_{2}^{U}, \dots, \mathcal{E}_{m^{U}}^{U}\}$

We define the expansion operator Φ as a linear map from \mathbb{R}^{m^U} to \mathbb{R}^m given by $\Phi_{rs} = 1$ if $\varepsilon_r = \varepsilon_s^U$ and 0 otherwise. Let $\mathbf{f}^U \in \mathbb{R}^{m^U}$ be the edge flows on the unlabeled edges.

The original problem can be converted to a linear least-squares problem:

$$\mathbf{f}^* = \arg\min_{\mathbf{f}} \|\mathbf{B}\mathbf{f}\|^2 + \lambda^2 \cdot \|\mathbf{f}\|^2$$

s.t.
$$\mathbf{f}_r = \hat{\mathbf{f}}_r, \forall \mathcal{E}_r \in \mathcal{E}^{\mathrm{L}}.$$

$$\mathbf{f}^{\mathrm{U}*} = \arg\min_{\mathbf{f}^{\mathrm{U}}} \left\| \begin{bmatrix} \mathbf{B}\Phi\\ \lambda \cdot \mathbf{I} \end{bmatrix} \mathbf{f}^{\mathrm{U}} - \begin{bmatrix} -\mathbf{B}\mathbf{f}^{0}\\ 0 \end{bmatrix} \right\|^{2}$$

Any feasible point can be written as $\mathbf{f}^0 + \mathbf{\Phi} \mathbf{f}^U$

The least-squares problem can be solved with iterative methods such as LSQR or LSMR.

LSMR: An Iterative Algorithm for Sparse Least-Squares Problems.(SIAM2017) LSQR: An Algorithm for Sparse Linear Equations and Sparse Least Squares. ACM TOMS (1982). The spectral decomposition of the graph Laplacian matrix is $L = U \Lambda U^{T}$. Because $L = BB^{T}$ the orthonormal basis $U \in \mathbb{R}^{n \times n}$ for vertex labels is formed by the left singular vectors of the incidence matrix

 $\Sigma \in \mathbb{R}^{n \times m}$ is the diagonal matrix of ordered singular values with m – n columns of zero-padding on the right the right singular vectors $\mathbf{V} \in \mathbb{R}^{m \times m}$ is an orthonormal basis for edge flows.

The divergence-minimizing objective can be rewritten in terms of the right singular vectors of **B** Let $\mathbf{p} = \mathbf{V}^{\mathsf{T}} \mathbf{f} \in \mathbb{R}^m$ represent the spectral coefficients of \mathbf{f} expressed in terms of the basis **V**.

$$\mathbf{f}^* = \arg\min_{\mathbf{f}} \|\mathbf{B}\mathbf{f}\|^2 + \lambda^2 \cdot \|\mathbf{f}\|^2 \qquad \text{s.t.} \quad \mathbf{f}_r = \hat{\mathbf{f}}_r, \, \forall \mathcal{E}_r \in \mathcal{E}^{\mathrm{L}}.$$

$$\mathbf{f}^{*} = \mathbf{V} \cdot \arg\min_{\mathbf{p}} (\mathbf{V}\mathbf{p})^{\mathsf{T}} \mathbf{B}^{\mathsf{T}} \mathbf{B} (\mathbf{V}\mathbf{p}) + \lambda^{2} \cdot (\mathbf{V}\mathbf{p})^{\mathsf{T}} (\mathbf{V}\mathbf{p})$$

$$= \mathbf{V} \cdot \arg\min_{\mathbf{p}} \mathbf{p}^{\mathsf{T}} \left(\Sigma^{\mathsf{T}} \Sigma + \lambda^{2} \cdot \mathbf{I} \right) \mathbf{p}$$

$$= \mathbf{V} \cdot \arg\min_{\mathbf{p}} \lambda^{2} \cdot \sum_{\alpha} \frac{\sigma_{\alpha}^{2} + \lambda^{2}}{\lambda^{2}} p_{\alpha}^{2}$$
s.t. $(\mathbf{V}\mathbf{p})_{r} = \hat{\mathbf{f}}_{r}, \forall \mathcal{E}_{r} \in \mathcal{E}^{\mathsf{L}},$

$$\mathsf{thickness: flow magnitude}$$

By construction V is a complete orthonormal basis for the space of edge flows. This space can be decomposed into two orthogonal subspaces.

The first subspace is the **cut-space** $\mathcal{R} = im(B^T)$ spanned by the singular vectors associated with nonzero singular values.

The space \mathcal{R} is also called the space of gradient flows, since any vector may be written as $B^{T}y$,

Y is a vector of vertex scalar potentials that induce a gradient flow.

The second subspace is the cycle-space $C = \ker(B)$ spanned by the remaining right singular vectors V_C associated with zero singular values. Note that any vector $f \in C$ corresponds to a circulation of flow, and

will induce zero cost in the loss function.

Let $\mathbf{u}_{\alpha}, \sigma_{\alpha}, \mathbf{v}_{\alpha}$ a triple of a left singular vector, singular value, and right singular vector.

 $\mathbf{u}_{\alpha}^{\mathsf{T}} \mathbf{L} \mathbf{u}_{\alpha} = \sigma_{\alpha}^{2}$ provide a notion of σ_{α} "unsmoothness" of basis vector \mathbf{u}_{α} representing vertex labels,

 $\mathbf{v}_{\alpha}^{\mathsf{T}} \mathbf{L}_{e} \mathbf{v}_{\alpha} = \sigma_{\alpha}^{2}$ gives the sum-of-squares divergence of basis vector \mathbf{v}_{α} representing

edge flows.



LEMMA 2.1. Assume the ground truth flows are divergence-free. Then as $\lambda \to 0$, the solution of Eq. (8) can exactly recover the ground truth from some labeled edge set \mathcal{E}^{L} with cardinality c = m - n + 1.

$$\begin{aligned} \mathbf{f}^* &= \mathbf{V} \cdot \arg\min_{\mathbf{p}} (\mathbf{V}\mathbf{p})^\mathsf{T} \mathbf{B}^\mathsf{T} \mathbf{B} (\mathbf{V}\mathbf{p}) + \lambda^2 \cdot (\mathbf{V}\mathbf{p})^\mathsf{T} (\mathbf{V}\mathbf{p}) \\ &= \mathbf{V} \cdot \arg\min_{\mathbf{p}} \mathbf{p}^\mathsf{T} \left(\Sigma^\mathsf{T} \Sigma + \lambda^2 \cdot \mathbf{I} \right) \mathbf{p} \\ &= \mathbf{V} \cdot \arg\min_{\mathbf{p}} \lambda^2 \cdot \sum_{\alpha} \frac{\sigma_{\alpha}^2 + \lambda^2}{\lambda^2} p_{\alpha}^2 \end{aligned} \quad \text{s.t.} \quad (\mathbf{V}\mathbf{p})_r = \hat{\mathbf{f}}_r, \ \forall \mathcal{E}_r \in \mathcal{E}^\mathsf{L}. \end{aligned}$$

Recall that $p_C \in \mathbb{R}^c$ are the spectral coefficients of a basis V_C in the cycle-space, then the ground truth edge flows can be written as $\hat{f} = V_C p_C$

On the other hand, in the limit $\lambda \to 0$ the spectral coefficients of basis vectors with non-zero singular values have infinite weights and are forced to zero.

Therefore, by choosing the set of labeled edges corresponding to c = m - n + 1 linearly independent rows from V_C , the ground truth \hat{f} is the unique optimal solution.

THEOREM 2.2. Let $\mathbf{V}_C^{\mathrm{L}}$ denote c linearly independent rows of the \mathbf{V}_C that correspond to labeled edges. If the divergence-free edge flows **f** are perturbed by δ , then as $\lambda \to 0$, the reconstruction error of the proposed algorithm is bounded by $[\sigma_{\min}^{-1}(\mathbf{V}_C^{\mathrm{L}}) + 1] \cdot \|\delta\|$.

The ground truth edge flows can be written as

$$\hat{\mathbf{f}} = \mathbf{f} + \delta = \begin{bmatrix} \mathbf{f}^{\mathrm{L}} \\ \mathbf{f}^{\mathrm{U}} \end{bmatrix} + \begin{bmatrix} \delta^{\mathrm{L}} \\ \delta^{\mathrm{U}} \end{bmatrix}$$

Further, the reconstructed edge flows are given by

 $\mathbf{V}_{\mathcal{C}}(\mathbf{V}_{\mathcal{C}}^{\mathrm{L}})^{-1}(\mathbf{f}^{\mathrm{L}}+\delta^{\mathrm{L}})$

Therefore, we can bound the norm of the reconstruction error as follows:

$$\|\mathbf{V}_{\mathcal{C}}(\mathbf{V}_{\mathcal{C}}^{\mathrm{L}})^{-1}(\mathbf{f}^{\mathrm{L}} + \delta^{\mathrm{L}}) - (\mathbf{f} + \delta)\|$$

SEMI-SUPERVISED LEARNING RESULTS

Minnesota road network n = 2642, m = 3303

The water irrigation network of Balerma city

US power grid network from **KONECT** n = 4941, m = 6593

An autonomous system network n = 520, m = 1280

n = 447,m = 454

Create synthetic flows with spectral coefficients for each basis vector that are inversely proportional to the associated singular values

$$p_{\alpha} = \frac{b}{\sigma_{\alpha} + \epsilon}$$
 $b = 0.02, \epsilon = 0.1$

Performance Measurement and Baselines

Pearson correlation coefficient
$$\rho(X,Y) = \frac{E[(X-\mu_X)(Y-\mu_Y)]}{\sigma_X \sigma_Y} = \frac{E[(X-\mu_X)(Y-\mu_Y)]}{\sqrt{\sum_{i=1}^n (X_i - \mu_X)^2} \sqrt{\sum_{i=1}^n (Y_i - \mu_Y)^2}}$$

First, the **ZeroFill** baseline simply assigns 0 edge flows to all unlabeled edges.

Second, the **LineGraph** baseline uses a line-graph transformation of the network and then applies standard vertex-based SSL on the resulting graph.

The LineGraph approach performs no better than ZeroFill. This should not be surprising, since the LineGraph approach does not interpret the sign of an edge flow as an orientation but simply as part of a numerical label

In our first set of experiments, the network topology comes from real data, but we use synthetic flows to demonstrate our method.



Learning Real-World Traffic Flows



Information Flow Networks



Even though the flows are not physical, FlowSSL method still outperforms the baselines.

Two active learning algorithms

$$= [\sigma_{\min}^{-1}(\mathbf{V}_C^{\mathrm{L}}) + 1] \cdot \|\delta\|_{\cdot}$$

one strategy for selecting \mathcal{E}^{L} is to choose m^{L} rows from \mathcal{V}_{0} that maximize the smallest singular value of the resulting submatrix. This problem is known as optimal column subset selection (maximum submatrix volume) and is NP-hard.

However, a good heuristic is the rank revealing QR decomposition (RRQR) [5], which computes

$$\mathbf{V}_C^\mathsf{T} \Pi = Q \begin{bmatrix} R_1 & R_2 \end{bmatrix}$$

 Π is a permutation matrix that keeps R_1 well-conditioned.

This approach is mathematically similar to graph clustering algorithms that use RRQR to select representative vertices for cluster centers.

Recursive Bisection (RB)

The intuition behind this heuristic is that edge flows on bottleneckedges, which partition a network, are able to capture global trends in the networks' flow pattern.

We start with an empty labeled set \mathcal{E}^L , a target number of labeled edges m^L , Next, we recursively partition the largest cluster in the graph with spectral clustering and add every edge that connects the two resulting clusters into \mathcal{E}^L







"Similar methods have been shown to be effective in semi-supervised active learning for vertex labels, in these cases, the graph is first clustered, and then one vertex is selected from each cluster."

Label Selection on Graphs(NIPS2009)

While any other graph partitioning algorithm could be used and greedy recursive bisection approaches can be sub-optimal.

How Good is Recursive Bisection?(SIAM1998)

