
Semi-Supervised Learning with Competitive Infection Models: Supplementary Material

Nir Rosenfeld
Harvard University

Amir Globerson
Tel Aviv University

1 Proof of Proposition 2 in Main Text

In this section we prove the correctness of our algorithm 2. The proof considers the more general CTIC infection dynamics and allows for node features or priors (via a penalty function).

In the infection dynamics presented in the paper, once a node’s label is set, it remains fixed. In contrast, during the course of the algorithm’s run, a node’s label may change with each distance update. It therefore remains to show that the algorithm outputs the desired labels. For basing our claim it will be easier to assume that instead of initially inserting all seed nodes into Q , we add a dummy root node r to V , with edges of length $w_{rs} = 0$ to all $s \in S$, and initialize Q to include only r . It is easy to see that after extracting r from Q , we return to our original algorithm.

Recall that the standard single-source Dijkstra algorithm offers three important guarantees: (1) the estimated distances of extracted nodes is correct (and remains unchanged), (2) nodes are extracted in increasing order of their true distance, and (3) the distance estimates always upper-bound the true distances.

Let $v \in U$ be a node that has just been extracted, and assume by induction that the labels of all previously extracted nodes (which include all seed nodes) are correct.¹ The above guarantees tell us that the distance from r to v is correct, and all nodes on the shortest path from r to v have already been extracted. This is true even when a penalty is incurred, as it can only increase the distance estimate. As these nodes are assumed to be correctly labeled, v inherits the correct label as well, as by construction its shortest path from r goes through exactly one seed node. The correctness of the labels of the seed nodes gives the induction basis, which concludes the proof.

¹Correct in the sense of the algorithm, not in the sense of their true labels.

2 Extensions

In this section we give an in-depth description of several useful extensions of our method that were briefly discussed in the main text. These include applying our method to the Linear Threshold model, incorporating node features and priors into the infection dynamics, and a framework for using our method in an active SSL setting.

2.1 The Linear Threshold model

In this section we show how InfProp can be applied with the Linear Threshold (LT) dynamics, rather than the IC or CTIC dynamics discussed in the main text. This includes adapting the algorithm for computing expected labels to the LT model, as well as supporting node features and priors.

The input to the LT model is a weighted graph $G = (V, E, W)$ and an initial set of infected seed nodes S . We assume that weights are positive, and that for each $v \in V$, the sum of incoming weights $\sum_u W_{uv}$ is at most 1 (though it can be strictly less than 1). Before the process begins, each node u is assigned a threshold η_u sampled uniformly at random from the interval $[0, 1]$. The dynamics then progress in discrete time steps, where at time t , a susceptible node v becomes infected if the weighted sum of its infected neighbors exceeds its threshold. Denoting by $I_u(t)$ an indicator of whether u is infected at time t , v is infected at time t if:

$$\sum_u W_{uv} I_u(t-1) \geq \eta_v \quad (1)$$

Note that the randomness in this model comes from the threshold η ; given η , the dynamics are deterministic.

The authors of [6] show that the LT model can also be equivalently expressed via a graphical perspective using active edge sets. Here, however, edges are no longer sampled independently. Instead, for each node v , only at most one incoming edge will become active in each instance. Specifically, for each node v , each

incoming edge $(u, v) \in E$ is selected to be the (only) edge with probability $p_{uv} = W_{uv}$, and with probability $W_{-u} = 1 - \sum_v W_{uv}$ no incoming edges are activated. Then, for a given instance, v is infected if and only if there is an active path to v from some seed node in S .

An interesting interpretation of the above is that the chosen active edge (u, v) can be thought of as corresponding to the node u whose infection triggered the infection of v by crossing the threshold. Under this view, a label-dependent specification of the above model is one where v inherits its label from its triggering neighbor u , which we refer to as his infector. This allows the model to be applied to the competitive setting which we consider. In terms of implementation, the only necessary modification to the algorithm is the way in which active edges are sampled.

The competitive LT model can also incorporate node priors using penalty terms. Specifically, the node-label prior $\rho_{v\ell}$ will induce a multiplicative penalty $q_v(\ell) \in [0, 1]$ on the original weights W_{uv} when u tries to infect v with label ℓ . Thus, given that u has label ℓ , the penalty reduces the probability that it will be the infector of v . To implement this, when u is expanded, the edge (u, v) is sampled to be active with the penalized probability, and all other incoming weights (including the complementing W_{-u}) are re-normalized.

2.2 Incorporating node features and priors

In addition to the graph, many network-based datasets include node or edge features. These can be used to generate node-specific class priors. In this section we describe a novel generalization of the competitive infection models introduced above which incorporates class priors into the dynamics. In this setting, our approach is to first train a probabilistic classifier (e.g., logistic regression) on the labeled seed set, and then use its predictions on the unlabeled nodes as a prior for our model.

Our method utilizes node priors by transforming them into penalties on incubation times. Consider a single instance of an infection process. Assume node u has just been infected with label $\ell \in \mathcal{Y}$, and succeeded in its attempt to infect node v with an incubation time of δ_{uv} . If δ_{uv} is small, then it is very likely that v will get infected with ℓ as well. On the other hand, if δ_{uv} is large, then other nodes might have a chance to infect v with other labels. This motivates the idea of further *penalizing* the infection time of a node according to its prior. We do this by adding a label-dependent penalty $q_v(\ell)$ to δ_{uv} , as a function of the prior $\rho_{v\ell}$. We use the link function $q_v(\ell) = -\log(\rho_{v\ell})$, which maps low priors into large penalties, and high priors into low penalties, where $\rho_{v\ell} = 1$ entails no penalty. Hence,

setting $\rho_{v\ell} = 1$ for all v, ℓ recovers the original model.

Note that while the priors are deployed locally, their effect is in fact global, as penalizing a node’s infection time delays the potential propagation of its acquired label throughout the graph. This increases the significance of nodes which are central to the infection process, and reduces the significance of those which play a small role in it, a property captured by our notion of confidence. The strength of the above formulation lies in its ability to introduce non-linear label dependencies to the actual infection dynamics. To see this, we can write the original predictions as:

$$f = \mathbb{E}_{\delta \sim D} [\mathcal{A}^\delta \mathcal{S}] = \mathbb{E}_{\delta \sim D} [\mathcal{A}^\delta] \mathcal{S} \quad (2)$$

where $\mathcal{A}_{vs}^\delta = \mathbb{1}_{\{s=\alpha(v)\}}$ indicates ancestors in G^δ , and $\mathcal{S}_{s\ell} = \mathbb{1}_{\{y_s=\ell\}}$ indicates the seed nodes’ true labels. This shows that predictions are non-linear in the propagation of the seed nodes, but linear in the labels. In the prior-dependent model, the above no longer holds, as activation times are now label-dependent.

2.3 Confidence and Active Learning

Recall that a node v has a probability f_{v0} of not being infected by *any* label. This suggests a very natural measure of *confidence* in our prediction, namely:

$$\sigma_v(S) = 1 - f_{v0} = \sum_{\ell=1}^L f_{v\ell}, \quad \sigma(S) = \sum_v \sigma_v(S) \quad (3)$$

The function σ quantifies the confidence *in the labeling*. This is conceptually different from confidence *in a label*. Our model supports both concepts distinctly. The former is controlled by the activations p , as they determine reachability in the active graph and are agnostic to labels. The latter is controlled by θ , as it affects the speed of propagation of the labels.

The notion of confidence allows us to apply our method to an active SSL setting. Instead of assuming the seed is given as input, in this setting we are allowed to *choose* the seed set, often under a cardinality constraint. The goal is then to choose the seed set which leads to a good labeling. Various graph-based notions have been suggested as objectives for active seed selection, such as those based on graph cuts [5], graph signals [2], and generalization error [4]. Such methods however either optimize an adversarial objective, or simply offer a heuristic solution. In contrast, using σ as a seed-selection criterion offers an optimistic alternative, as summing over all classes makes it indifferent to the actual (latent) labels.

The confidence term σ coincides with the well-studied notion of *influence*, defined as the expected number of

nodes a seed will infect. In [6] it is shown that for various settings, influence is submodular, and therefore admits to a greedy $(1 - 1/\epsilon)$ -approximation scheme. Any algorithm for maximizing influence efficiently (e.g., [1, 3]), can therefore be adopted for our setting.

3 Non-Homogeneous Laplacian

Here we prove that:

$$\mathcal{L}(S)f = \mathbf{b}(S)$$

where:

$$\mathbf{b}_{ul}(S) = \sum_v b_{vul}^{(S)}, \quad b_{vul}^{(S)} = \text{cov}[T_{vu}(S), Y_{u\ell}]$$

For clarity we drop the notational dependence on S . We begin by expanding $f_{u\ell}$ using Y and T :

$$\begin{aligned} f_{u\ell} &= \mathbb{E}[Y_{u\ell}] = \mathbb{E}[T_{u\cdot} Y_{\cdot\ell}] \\ &= \mathbb{E}\left[\sum_v T_{uv} Y_{v\ell}\right] \\ &= \sum_v \mathbb{E}[T_{uv} Y_{v\ell}] \\ &= \sum_v (\mathbb{E}[T_{uv}] \mathbb{E}[Y_{v\ell}] + \text{cov}[T_{uv}, Y_{v\ell}]) \\ &= \sum_v (\mathbf{T}_{uv} f_{v\ell} + b_{uv\ell}) \end{aligned}$$

where the final step is true for the product of general random variables. Rewriting in matrix form gives: $f = \mathbf{T}f + \mathbf{b}$. Rearranging we get: $(I - \mathbf{T})f = \mathcal{L}f = \mathbf{b}$, as required.

The objective function can then be expressed as:

$$\begin{aligned} \|\mathcal{L}f' - \mathbf{b}\|_2^2 &= \sum_u \sum_\ell (\mathcal{L}_{u\cdot} f'_{\cdot\ell} - b_{u\ell})^2 \\ &= \sum_u \sum_\ell \left(\sum_v \mathcal{L}_{uv} f'_{v\ell} - b_{uv\ell} \right)^2 \\ &= \sum_u \sum_\ell \left(\sum_v (\mathbb{1}_{\{u=v\}} - \mathbf{T}_{uv}) f'_{v\ell} - b_{uv\ell} \right)^2 \\ &= \sum_u \sum_\ell \left(f'_{u\ell} - \sum_v (w_{uv} f'_{v\ell} + b_{uv\ell}) \right)^2 \end{aligned}$$

where $w_{uv} = \mathbf{T}_{uv}$.

4 Details for the Illustrative Synthetic Experiment

Our general hypothesis in this work is that infection dynamics are a good candidate for propagating

label information over real networks. To illustrate this, we designed a synthetic experimental setup in which our goal was to capture the structure of real world networks. One well-known property of such networks is that they often have a community-like structure, with many intra-community edges, but few inter-community edges. In many cases, only a few specific nodes within a community are also connected to other communities. Hence, we randomly created small networks with the above properties.

Specifically, each network was set to have 3 (possibly overlapping) communities, each with 64 nodes. Nodes were randomly assigned into one community, and in each community, 8 nodes were randomly assigned to an additional community. To account for some noise, all other edges were added with probability 0.05. The seed set included one randomly chosen node from each community, giving $|S| = 3$. The figure in the main text displays a random instance of the above setting, providing both the instance specific accuracies, as well as the average accuracy over 1,000 random instances.

Recall that InfProp can be interpreted both as the expected result of a dynamic infection process and as a stochastic ensemble of shortest paths. We therefore compared our method to two baselines. To compare the dynamics, we used Label Propagation (LabelProp) which is based on the more standard random-walk dynamics. As we argue in the text, these dynamics are prone to getting stuck in dense clusters. As can be seen, while InfProp provides almost exact predictions, the predictive values of LabelProp are almost uniform and hence extremely error-prone. This demonstrates the inability of label information to propagate efficiently over the network.

To demonstrate the power of using a stochastic ensemble of paths, we compared to simply setting labels according to the deterministic shortest paths given by the original graph. While correctly classifying most labels, shortest paths can be very sensitive to cross-community or noisy edges. In contrast, InfProp mitigates this noise by considering a distribution over shortest-paths.

5 Datasets

We evaluated our method on various learning tasks over three collections of benchmark datasets, which include network based datasets for multi-class learning with features² [9], multi-class learning without features³ [8], and multi-label learning⁴ [7]. The following

² <http://linqs.umiacs.umd.edu/projects//projects/lbc/>

³ <http://cs.gmu.edu/~tsaha/Homepage/Projects.html>

⁴ <http://github.com/sharadnandanwar/snbc>

table provides some statistics.

	Dataset	Nodes	Edges	Classes	Features	Avg. $ y $
Multiclass ³	CoRA	2,708	5,278	7	-	1
	DBLP	5,329	21,880	6	-	1
	Flickr	7,971	478,980	7	-	1
	IMDb	2,411	12,255	22	-	1
	Industry	2,189	11,666	12	-	1
Features ²	CiteSeer	3,132	4,713	6	3,703	1
	CoRA	2,708	5,278	7	1,433	1
	PubMed	19,717	44,324	3	500	1
Multilabel ⁴	Amazon	83,742	190,097	30	-	1.546
	CoRA	24,519	92,207	10	-	1.004
	IMDb	19,359	362,079	21	-	2.300
	PubMed	19,717	44,324	3	-	1
	Wikipedia	35,633	49,538	16	-	1.312
	YouTube	22,693	96,361	47	-	1.707

References

- [1] COHEN, E., DELLING, D., PAJOR, T., AND WERNECK, R. F. Sketch-based influence maximization and computation: Scaling up with guarantees. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management* (2014), ACM, pp. 629–638.
- [2] GADDE, A., ANIS, A., AND ORTEGA, A. Active semi-supervised learning using sampling theory for graph signals. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining* (2014), ACM, pp. 492–501.
- [3] GOMEZ-RODRIGUEZ, M., SONG, L., DU, N., ZHA, H., AND SCHÖLKOPF, B. Influence estimation and maximization in continuous-time diffusion networks. *ACM Trans. Inf. Syst.* *34*, 2 (Feb. 2016), 9:1–9:33.
- [4] GU, Q., AND HAN, J. Towards active learning on graphs: An error bound minimization approach. In *Data Mining (ICDM), 2012 IEEE 12th International Conference on* (2012), IEEE, pp. 882–887.
- [5] GUILLORY, A., AND BILMES, J. A. Active semi-supervised learning using submodular functions. In *UAI 2011, Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence* (2011), pp. 274–282.
- [6] KEMPE, D., KLEINBERG, J., AND TARDOS, É. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining* (2003), ACM, pp. 137–146.
- [7] NANDANWAR, S., AND MURTY, M. N. Structural neighborhood based classification of nodes in a network. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016), KDD '16, pp. 1085–1094.
- [8] SAHA, T., RANGWALA, H., AND DOMENICONI, C. Flip: active learning for relational network classification. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases* (2014), Springer, pp. 1–18.
- [9] SEN, P., NAMATA, G. M., BILGIC, M., GETOOR, L., GALLAGHER, B., AND ELIASSI-RAD, T. Collective classification in network data. *AI Magazine* *29*, 3 (2008), 93–106.