# On the Troll-Trust Model for Edge Sign Prediction in Social Networks

Géraud Le Falher[1], Nicolò Cesa-Bianchi[2], Claudio Gentile[3] and Fabio Vitale[1,4]

[1] Inria, Univ. Lille, CNRS UMR 9189 – CRIStAL, France   [2] Università degli Studi di Milano, Italy   [3] University of Insubria, Italy   [4] Department of Computer Science, Aalto University, Finland
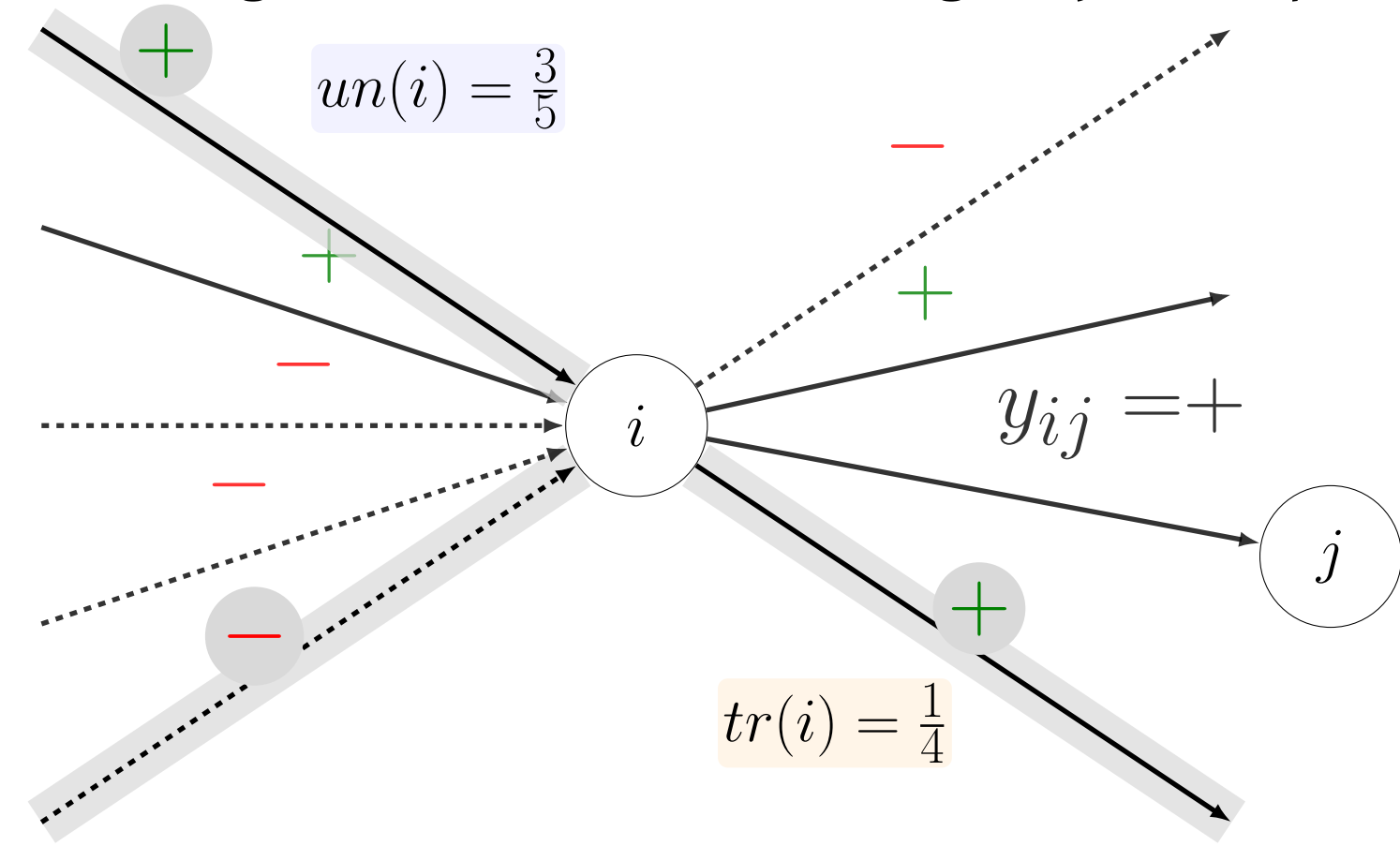
## Problem setup and notations

Given a directed graph of binary interactions between users, we want a **scalable** method to predict the sign of an interaction, using only the topology and the sign of some edges in $\mathbb{E}_0$.



$un(i) = \frac{3}{5}$

$y_{ij} = +$

$tr(i) = \frac{1}{4}$

The **trollness** $tr(i)$ of $i$ is its fraction of negative outgoing edges, its **untrustworthiness** $un(i)$ is its fraction of negative incoming edges.

## Contributions

1. A simple **generative model**, justifying existing heuristics and providing a **new principled predictor** ($\text{BLC}(tr, un)$)
2. A maximum likelihood approximation by a **label propagation algorithm** (L. Prop.), leveraging a reduction from **edge to node classification**
3. Extensive comparative **experiments on real data** against state of the art methods

## Generative model

Each node $i$ has 2 parameters, drawn from an arbitrary distribution $\mu$ over pairs in $[0,1]^2$:
- its tendency $p_i$ to send positive edges (i.e. niceness) and
- its tendency $q_i$ to receive positive edges (i.e. pleasantness)

$$i \quad \mathbb{P}(y_{ij} = +1) = \frac{1}{2}(p_i + q_j) \quad j$$

$$(p_i, q_i) \sim \mu(p, q) \qquad (p_j, q_j) \sim \mu(p, q)$$

The Bayes optimal prediction for $y_{i,j}$ is thus

$$y^*(i,j) = \text{SGN}\left(\mathbb{P}(y_{i,j} = +1) - \tfrac{1}{2}\right)$$

## Active algorithm: $\text{BLC}(tr, un)$

We use the complementary to 1 of trollness and untrustworthiness (estimated on the training set $\mathbb{E}_0$) as proxy for $p_i$ and $q_j$ and predict

$$\widehat{y}(i,j) = \text{SGN}\left(\underbrace{(1 - \widehat{tr}(i)) + (1 - \widehat{un}(j)) - \tau}_{\approx \frac{1}{2}(p_i + q_j) = \mathbb{P}(y_{ij}=+1)} - \tfrac{1}{2}\right) \qquad (1)$$

$1 - \widehat{tr}(i)$ is an indirect observation of $p_i$. Indeed, letting

$$\overline{q}_i = \frac{1}{d_{out}(i)} \sum_{j \text{ s.t.}(i,j)\in E} q_j$$

we have

$$1 - \widehat{tr}(i) \approx \tfrac{1}{2}(p_i + \overline{q}_i) \quad \text{and likewise} \quad 1 - \widehat{un}(j) \approx \tfrac{1}{2}(q_j + \overline{p}_j)$$
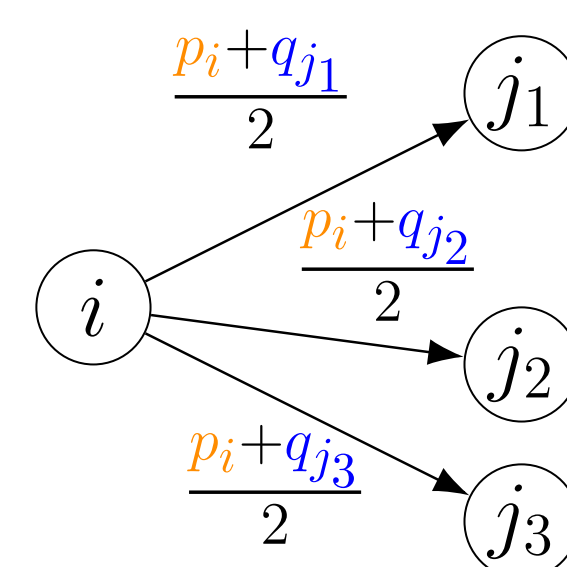
$\frac{p_i+q_{j_1}}{2}$, $\frac{p_i+q_{j_2}}{2}$, $\frac{p_i+q_{j_3}}{2}$ ($i \to j_1, j_2, j_3$)

Thus we need to subtract

$$\tau = \tfrac{1}{2}(\mu_p + \mu_q)$$

as $\overline{p}_j$ and $\overline{q}_i$ concentrate around their mean $\mu_p$ and $\mu_q$.
We sample $Q$ outgoing and incoming edges for each node, estimate empirically $\widehat{tr}(i)$, $\widehat{un}(j)$ and $\tau$ as the overall fraction of positive edges then finally predict remaining edges according to (1).
Setting $Q = \frac{1}{\varepsilon^2} \ln \frac{|V|}{\delta}$ we query $\Theta(|V| \ln |V|)$ edges.
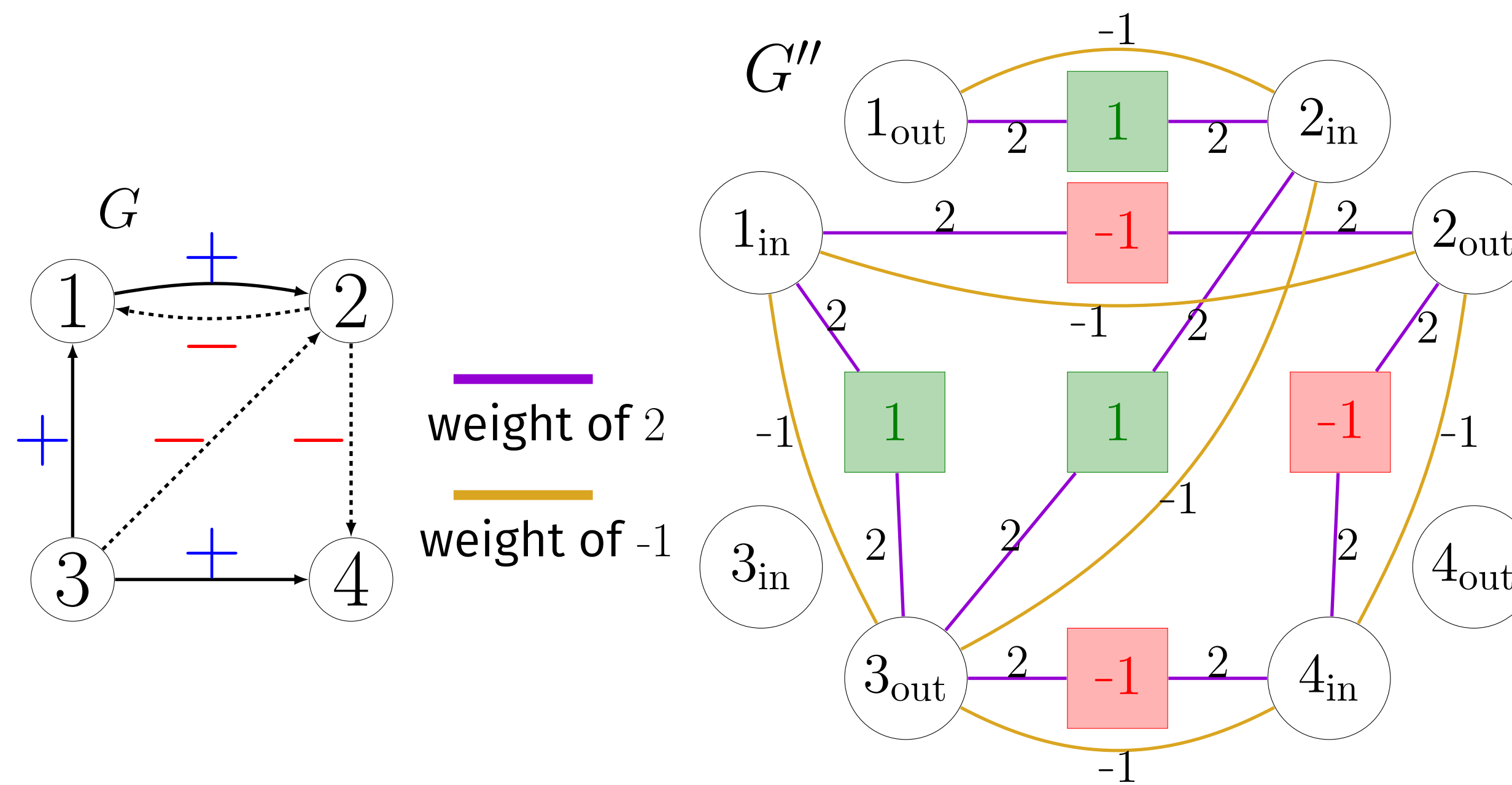
---

This is enough to guarantee that

$$\left| \left[(1 - \widehat{tr}(i)) + (1 - \widehat{un}(j)) - \widehat{\tau}\right] - \left[\frac{p_i + q_j}{2}\right] \right| \leq 8\epsilon$$

**holds with probability at least** $1 - 10\delta$ simultaneously for all non-queried edges $(i,j) \in E$ such that $d_{out}(i), d_{in}(j) \geq Q$, therefore providing correct prediction as long as $\mathbb{P}(y_{i,j} = +1)$ is bounded away from $\frac{1}{2}$.
The $\text{BLC}(tr, un)$ algorithm is inspired by this fact and also predicts according to (1), except there is no sampling but rather it uses the provided batch training set directly.

## Batch algorithm: L. Prop.

Unfortunately, a good sampling requires a rather dense graph which is not usually the case in social media. In order to address this issue, first we transform the problem from **edge to node classification**.



$G$   $G''$

weight of 2

weight of -1

Then we **approximate** $y^*(i,j)$ by resorting to a maximum likelihood estimator of the parameters $\{p_i, q_i\}_{i=1}^{|V|} \mathbb{P}\left(E_0 \mid \{p_i, q_i\}_{i=1}^{|V|}\right)$. Because the gradients of the log-likelihood function are not linear (e.g. w.r.t. $p_\ell$)

$$\sum_{\ell, j \in E_0; y_{\ell j}=+1} \frac{1}{p_\ell + q_j} - \sum_{\ell, j \in E_0; y_{\ell j}=-1} \frac{1}{2 - p_\ell - q_j}$$

we approximate them, which is equivalent to setting to zero the gradient w.r.t. $(p, q) = \{p_i, q_i\}_{i=1}^{|V|}$ of the quadratic function

$$f_{E_0}(p, q) = \sum_{(i,j)\in E_0} \left(\underbrace{\frac{p_i + q_j}{2}}_{\in[0,1]} - \underbrace{\frac{1 + y_{i,j}}{2}}_{\in[0,1]}\right)^2$$

We follow a label propagation approach by **making the test labels appear** and minimizing $f_{E_0}(p, q) + f_{E\setminus E_0}(p, q, y_{i,j})$ w.r.t. both $(p, q)$ and all $y_{i,j} \in [-1, +1]$, for $(i,j) \in E \setminus E_0$.
Setting the edge weights as in $G''$ introduces an extra regularization term.

$$\underbrace{\widehat{f}(p, q, y_{i,j(i,j)\in E\setminus E_0})}_{\text{energy function on } G''} = f_{E_0}(p, q) + f_{E\setminus E_0}(p, q, y_{i,j}) + \widetilde{\text{regul}}$$

We run $\text{diameter}(G)$ iterations of label propagation and use a binary threshold over the estimated $y_{i,j}$ to predict signs.

## Online setting

We also study an online setting where the signs are **adversarial** rather than generated by our model.

Letting $Y$ be the vector of all labels, $\Psi_{out}(i, Y)$ is the number of least used label outgoing from $i$, and $\Psi_{out}(Y) = \sum_{i\in V} \Psi_{out}(i, Y)$. Likewise for incoming edges, $\Psi_{in}(Y) = \sum_{j\in V} \Psi_{in}(j, Y)$ and finally **the regularity of a labeling** $Y$ is $\Psi_G(Y) = \min\{\Psi_{in}(Y), \Psi_{out}(Y)\}$.

We devise an algorithm consisting of a combination of Randomized Weighted Majority instances built on top of each other that on average makes $\Psi_G(Y) + O\left(\sqrt{|V|\Psi_G(Y)} + |V|\right)$ mistakes.

On the lower side, for any directed graph $G$ and any integer $K$, there exists a labeling $Y$ forcing at least $\frac{K}{2}$ mistakes to any online algorithms, while $\Psi_G(Y) \leq K$.

## Experiments

5 datasets of different size from different domains:
**Citations** $i$ the work of $j$ to praise it or criticise it.
**Wikipedia** $i$ vote for or against $j$ promotion to adminship.
**Slashdot** $i$ consider $j$ as a friend or foe.
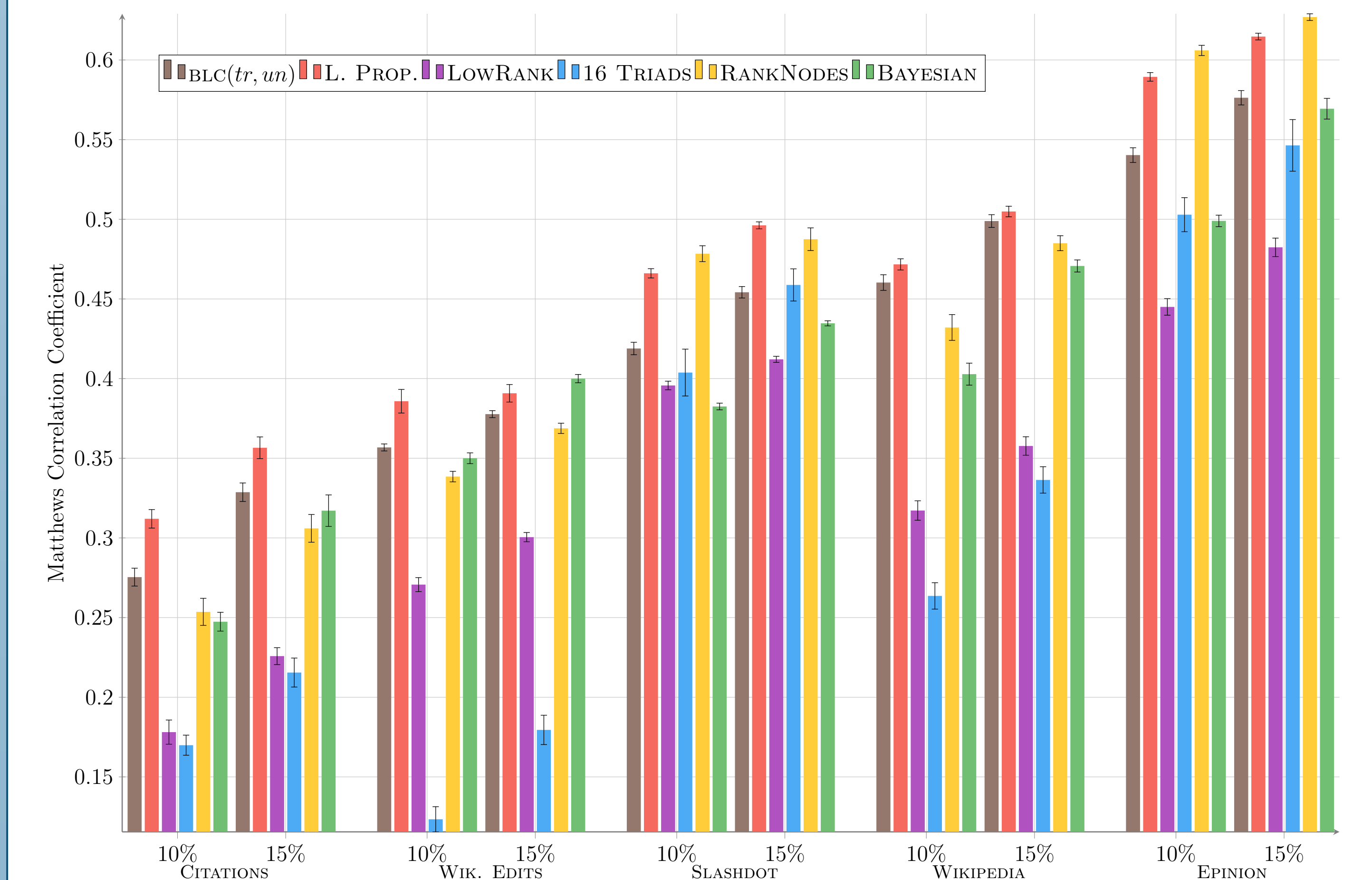**Epinion** $i$ trust or not the reviews made by $j$.
**Wik. Edits** $i$ reacted to a Wikipedia edit made by $j$, to enhance it or revert it.
Because of imbalance, we evaluate using the Matthews Correlation Coefficient (MCC):

$$\text{MCC} = \frac{tp \times tn - fp \times fn}{\sqrt{(tp + fp)(tp + fn)(tn + fp)(tn + fn)}} \begin{cases} 1 & \text{all predictions correct} \\ 0 & \text{random predictions} \\ -1 & \text{all predictions incorrect} \end{cases}$$

Competitors are representative of 2 classes of methods: LowRank **matrix completion** [1] and using **logistic regression on edge features** built:
- on 16 Triads features, as signed graphs exhibit specific triangle patterns according to the status theory [2]
- on a high number of so-called "Bayesian" features [3]
- on RankNodes scores computed with a PageRank-inspired algorithm tailored to directed graphs with negative edges [4]



| time (ms) | $\text{BLC}(tr, un)$ | L. Prop. | LowRank | 16 Triads | RankNodes | Bayesian |
|---|---|---|---|---|---|---|
| Citations | 0.6 | 19.6 | 3,279 | 6.2 | 155 | 4,813 |
| Wik. Edits | 10.1 | 1,329 | 127,654 | 209 | 3,174 | 104,305 |
| Slashdot | 8.3 | 677 | 69,988 | 131 | 2,441 | 68,085 |
| Wikipedia | 9.6 | 927 | 129,460 | 177 | 3,890 | 92,719 |
| Epinion | 1.6 | 41.9 | 8,523 | 14.8 | 249 | 12,507 |

## Discussion

We presented two batch algorithms derived from our generative model. One is local ($\text{BLC}(tr, un)$) and therefore extremely scalable while being performant both in theory and in practice. The other (L. Prop.) propagates sign information along the graph which provides more accurate results while remaining faster than previous approaches.
Later, we would like to extend our results to **weighted graphs** and **incorporate side information**.

## References

[1] K.-y. Chiang et al., "Prediction and Clustering in Signed Networks: A Local to Global Perspective," *JMLR*, vol. 15, pp. 1177–1213, 2014.

[2] J. Leskovec et al., "Predicting positive and negative links in online social networks," in *Proc. of the 19th international conference on World wide web*, 2010, p. 641.

[3] D. Song et al., "Link sign prediction and ranking in signed directed social networks," *Social Network Analysis and Mining*, vol. 5, no. 1, 2015.

[4] Z. Wu et al., "The troll-trust model for ranking in signed networks," in *WSDM*, 2016, pp. 447–456.