# Distributionally Robust Removal of Malicious Nodes from Networks

### Sixie Yu 1 Yevgeniy Vorobeychik 1

### **Abstract**

An important problem in networked systems is detection and removal of suspected malicious nodes. A crucial consideration in such settings is the uncertainty endemic in detection, coupled with considerations of network connectivity, which impose indirect costs from mistakely removing benign nodes as well as failing to remove malicious nodes. A recent approach proposed to address this problem directly tackles these considerations, but has a significant limitation: it assumes that the decision maker has accurate knowledge of the joint maliciousness probability of the nodes on the network. This is clearly not the case in practice, where such a distribution is at best an estimate from limited evidence. To address this problem, we propose a distributionally robust framework for optimal node removal. While the problem is NP-Hard, we propose a principled algorithmic technique for solving it approximately based on duality combined with Semidefinite Programming relaxation. A combination of both theoretical and empirical analysis, the latter using both synthetic and real data, provide strong evidence that our algorithmic approach is highly effective and, in particular, is significantly more robust than the state of the art.

# 1. Introduction

One of the major problems in networked settings is to identify and remove potentially malicious nodes. For example, in social networks, malicious nodes may correspond to accounts created by malicious parties which spread social spam, hate speech, fake news, and the like, with considerable deliterious effects (Allcott & Gentzkow, 2017; Cheng et al., 2015). Major social network platforms consequently devote considerable efforts to identify and remove fake or malicious accounts (Rodriguez, 2018; Scott & Isaac, 2017).

Nevertheless, evidence suggests that the problem remains pervasive (Andrade, 2018; Narayanan et al., 2018). Similarly, in cyber-physical systems (e.g., smart grid infrastructure), computing nodes compromised by malware can cause catastrophic losses, and mitigation through detection and removal of such malicious nodes is a major problem (Mo et al., 2012; Yang et al., 2017).

A central challenge faced in deciding which potentially malicious nodes to remove is to account for the combination of uncertainty about whether particular nodes are malicious, and the indirect (network) effects of the decision. This combination makes the decision about which nodes to remove fundamentally a subset selection problem—a challenging combinatorial optimization problem. Recently, Yu & Vorobeychik proposed an approach for solving it they term MINT, where the problem is captured by approximately minimizing loss which involves three terms: direct loss from removing benign nodes, indirect loss from cutting links in the benign subgraph, and indirect loss from maintaining connectivity between malicious and benign nodes. This model is il-

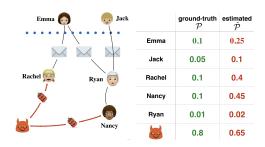


Figure 1. An illustration of a decision to remove two nodes, Jack and Emma, from the network, on our loss function.

lustrated in Fig. 1, where we consider removing Jack and Emma, two benign nodes above the dotted blue line (and failing to remove the malicious node). Suppose that we pay a penalty of  $\alpha_1$  for each benign node we remove, a penalty  $\alpha_2$  for each link we cut between benign nodes, and  $\alpha_3$  for each link between remaining malicious nodes and benign nodes. Since we remove 2 benign nodes, cut 3 links between benign nodes (one between Emma and Rachel, one between Emma and Ryan, and another between Jack and Ryan), and the malicious node is still connected to 2 nodes (Rachel and Nancy), our total loss is:  $2\alpha_1 + 3\alpha_2 + 2\alpha_3$ .

<sup>&</sup>lt;sup>1</sup>Computer Science and Engineering Department, Washington University in St. Louis. Correspondence to: yvorobey-chik@wustl.edu.

A major shortcoming of MINT is that it assumes that the distribution of node maliciousness is known. In practice, such a distribution is estimated from limited evidence, such as node behavior and other characteristics, and this estimation may be quite inaccurate (particularly if our modeling assumptions are poor, for example, if we erroneously assume that maliciousness probabilities of nodes are independent). More precisely, consider an unknown ground-truth  $\mathcal{P}$ , as illustrated in Fig. 1 in green. Whereas MINT assumes we know  $\mathcal{P}$ , in reality we only have an estimate  $\hat{\mathcal{P}}$  (shown in red in Fig. 1). To address this issue, we propose a new approach, MINT\_DRO, which is a distributionally robust framework for optimal node removal. We design an uncertainty set around the estimate  $\hat{P}$  and optimize with respect to the worst-case scenario. We propose a principled algorithmic approach for solving this problem approximately based on duality combined with Semidefinite Programming relaxation, and prove that the uncertainty set in our model contains the ground-truth distribution  $\mathcal{P}$  with high probability. This in turn implies that with high probability MINT\_DRO is robust with respect to the ground-truth distribution. Finally, we conducted extensive experiments using both synthetic and real data to show that our model is significantly more robust than MINT.

Related Work There are several prior efforts considering a related problem of *graph scan statistics* and hypothesis testing (Arias-Castro et al., 2011; Priebe et al., 2005; Sharpnack et al., 2013). These study the following problem: given a graph *G* where each node is associated with a random variable with an exogenously specified probability distribution, find a subset of nodes that maximizes a scan statistic defined over subsets of nodes (for example, this statistic may generalize log-likelihood ratio). The recent MINT approach (Yu & Vorobeychik, 2018) can be viewed through this lens as well, but as it has been shown to have state-of-the-art performance, our comparison, our experimental evaluation focuses on comparing to MINT.

Also closely related to our problem is the broader literature on distributionally robust optimization (DRO) (Scarf, 1958). In the DRO framework one defines a set of probability distributions that is assumed to contain the true stochastic model of the problem. Many solutions have been proposed to solve specific problems under the DRO framework (Xu & Mannor, 2010; Calafiore & El Ghaoui, 2006; Yue et al., 2006; Cheng et al., 2014; Wiesemann et al., 2014), although this framework has not been applied in the context of choosing which potentially malicious nodes to remove from a network.

Our design of the uncertainty set is inspired by the idea of moment-constrained uncertainty set (Delage & Ye, 2010; Popescu, 2007; Calafiore & El Ghaoui, 2006). Yet another related research strand is in using Semidefinite Programming (SDP) to approximate combinatorial optimization

problems (Goemans & Williamson, 1995; Luo et al., 2010; Bertsimas & Sethuraman, 2000), although such approaches are domain specific. Finally, our work bears some relationship to the burgeoning field of adversarial machine learning (Vorobeychik & Kantarcioglu, 2018), although we do not explicitly consider issues of adversarial response (such as evasion attacks) in our setting.

#### 2. Model

We consider a network that is represented by a graph G=(V,E), where V (|V|=N) is the set of nodes and E the set of edges connecting them. Each node  $i\in V$  represents a user and each edge (i,j) represents an edge (e.g., friendship on Facebook) between user i and user j. We focus our attention on undirected graphs. We denote the adjacency matrix of G by  $\mathbf{A}\in\mathbb{R}^{N\times N}$ . The elements of  $\mathbf{A}$  are binary if the graph is unweighted, or some non-negative real numbers if the graph is weighted. To make exposition easier we focus on unweighted graphs. Generalization to weighted graphs is straightforward.

We consider the problem of removing malicious nodes from the network G. A configuration of the network is denoted by  $\pi \in \{0,1\}^N$ , with  $\pi_i = 1$  indicating that a node i is malicious, with  $\pi_i = 0$  when i is benign. For convenience, we also let  $\bar{\pi}_i = 1 - \pi_i$  to indicate that i is benign. Consequently,  $\pi$  (and  $\bar{\pi}$ ) assigns malicious or benign label to each node. The identity of malicious and benign nodes are usually uncertain. So instead we have a probability distribution over the configurations. Formally, let  $\pi \sim \mathcal{P}$ , where  $\mathcal{P}$  captures the joint probability distribution over node configurations.

Our work builds upon the following model proposed by Yu & Vorobeychik (2018). Let S denote the set of nodes to remove. Define a vector  $\mathbf{x} \in \{-1,1\}^N$ , where  $\mathbf{x}_i = 1$  if and only if node i is removed ( $i \in S$ ), and  $\mathbf{x}_i = -1$  if node i remains in the network ( $i \in V \setminus S$ ). The goal of their model is to identify a subset of nodes S to remove so as to minimize the impact of the remaining malicious nodes on the network, while at the same time minimizing disruptions caused to the benign subnetwork. This goal is naturally captured by the loss function given in Eq. (1).

$$\alpha_{1} \underbrace{\sum_{i=1}^{N} \mathbf{x}_{i} \mathbb{E}_{\pi \sim \mathcal{P}}[\bar{\pi}_{i}]}_{\mathcal{L}_{1}} - \alpha_{2} \underbrace{\sum_{i,j}^{N} \mathbf{A}_{i,j} \mathbf{x}_{i} \mathbf{x}_{j} \mathbb{E}_{\pi \sim \mathcal{P}}[\bar{\pi}_{i} \bar{\pi}_{j}]}_{\mathcal{L}_{2}} + \alpha_{3} \underbrace{\sum_{i,j}^{N} \mathbf{x}_{i} \mathbf{x}_{j} \mathbf{A}_{i,j} \mathbb{E}_{\pi \sim \mathcal{P}}[\pi_{i} \bar{\pi}_{j}]}_{\mathcal{L}_{2}}.$$
(1)

As we can observe, the loss function is composed of three

components. The first component,  $\mathcal{L}_1$ , of the loss function is the direct loss associated with removing benign nodes. The second component,  $\mathcal{L}_2$ , penalizes cutting connections between benign nodes that are removed and benign nodes that remain; in other words, it penalizes the degradation of connectivity within the benign subgraph. The third component of the loss function,  $\mathcal{L}_3$ , captures the consequence of failing to remove malicious nodes in terms of connections from these to benign nodes. The nonnegative trade-off parameters  $\alpha_1$ ,  $\alpha_2$ , and  $\alpha_3$  satisfy  $\alpha_1 + \alpha_2 + \alpha_3 = 1$ , and weigh the relative importance of the three components of the loss function.

The configuration  $\pi$  is a random variable distributed according to  $\mathcal{P}$ . Let  $\mu \in \mathbb{R}^N$  and  $\Sigma \in \mathbb{R}^{N \times N}$  denote its mean and covariance, respectively. The loss function defined in Eq. (1) depends on both  $\mu$  and  $\Sigma$ . To make the dependency explicit we define several matrices and re-write the loss function in a matrix-vector form. We define the matrices  $\mathbf{B}(\mu)$ ,  $\mathbf{P}(\mu, \Sigma)$ ,  $\mathbf{M}(\mu, \Sigma)$  as follow.

$$\mathbf{B}(\mu) := diag(\mathbb{E}_{\pi \sim \mathcal{P}}[\bar{\pi}])$$

$$\mathbf{P}(\mu, \mathbf{\Sigma}) := \mathbf{A} \odot \mathbb{E}_{\pi \sim \mathcal{P}}[\bar{\pi}\bar{\pi}^T]$$

$$\mathbf{M}(\mu, \mathbf{\Sigma}) := \mathbf{A} \odot \mathbb{E}_{\pi \sim \mathcal{P}}[\pi\bar{\pi}^T]$$

Note that the elements of these matrices are not constant, but depend on  $\mu$  and  $\Sigma$  (see the appendix for their detailed dependency).

Slightly abusing notation, we define two additional matrices,  $\mathbf{Q}(\mu, \mathbf{\Sigma})$  and  $\mathbf{b}(\mu)$ . Note that  $\mathbf{Q} \in \mathbb{R}^{N \times N}$  is a symmetric matrix:

$$\mathbf{Q}(\mu, \mathbf{\Sigma}) := (\alpha_3/2) \left[ \mathbf{M}(\mu, \mathbf{\Sigma}) + \mathbf{M}(\mu, \mathbf{\Sigma})^T \right] - (\alpha_2/2) \left[ \mathbf{P}(\mu, \mathbf{\Sigma}) + \mathbf{P}(\mu, \mathbf{\Sigma})^T \right],$$

and  $\mathbf{b}(\mu) := (\alpha_1/2)\mathbf{B}(\mu)\mathbf{1}$ . We can now rewrite the loss function in a compact matrix-vector form as the following:

$$\mathcal{L}(\mathbf{x}; \mu, \mathbf{\Sigma}) = \mathbb{E}_{\pi \sim \mathcal{P}} \left[ \alpha_1 \mathcal{L}_1 + \alpha_2 \mathcal{L}_2 + \alpha_3 \mathcal{L}_3 \right]$$
$$= \mathbf{x}^T \mathbf{Q}(\mu, \mathbf{\Sigma}) \mathbf{x} + 2 \mathbf{x}^T \mathbf{b}(\mu)$$

Optimizing the loss function above (as done by Yu & Vorobeychik) critically assumes that the maliciousness distribution  $\mathcal P$  is known. In reality, this is typically not the case, and such a distribution is estimated from data. Let  $\hat{\mathcal P}$  denote the estimated distribution. The mean of  $\hat{\mathcal P}$  is denoted by  $\hat{\mu}$ , where  $\hat{\mu}_i$  is the estimated probability that node i is malicious given its features from past data. Similarly, the estimated covariance matrix is represented by  $\hat{\Sigma}$ . The model proposed by Yu & Vorobeychik is called MINT, which is to solve the following optimization problem:

$$\min_{\mathbf{x}} \quad \mathcal{L}(\mathbf{x}; \hat{\mu}, \hat{\mathbf{\Sigma}})$$

$$s.t. \quad \mathbf{x} \in \{-1, 1\}^{N}$$
(MINT)

Although MINT has been shown to perform well on several real-world datasets, its performance is strongly influenced by the estimation error of  $\mu$ . In fact, in Section 5 we show that even a small estimation error can severely undermine the performance of MINT.

In order to mitigate the sensitivity of MINT to estimation error, we propose a novel *Distributionally Robust Optimization* (DRO) approach for solving the problem posed above. The general idea is to design a distributional set to capture the uncertainty about the estimated mean  $\hat{\mu}$  and make decisions considering the *worst-case* scenario. Specifically, we propose a model named MINT\_DRO, which aims to solve the following optimization problem:

$$\min_{\mathbf{x}} \sup_{\mathcal{F} \sim \Pi} \mathbb{E}_{\mathcal{F}} \left[ \mathcal{L}(\mathbf{x}; \mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}}) \right]$$

$$s.t. \qquad \mathbf{x} \in \{-1, 1\}^{N},$$
(MINT\_DRO)

where the set  $\Pi$  captures uncertainty about the true mean  $\mu$ . There are several fundamental differences between MINT\_DRO and MINT. First, there is an additional inner maximization problem in MINT\_DRO. The inner maximization is optimized over a set  $\Pi$ , which contains a set of probability distributions, where  $\mathcal F$  is any distribution sampled from  $\Pi$ , and  $\mu_{\mathcal F}$  are random variables distributed according to  $\mathcal F$ . Inspired by Delage & Ye (2010) and Cheng et al. (2014), we parametrize the set  $\Pi$  by the first and second moments of the distributions in it. Specifically, let  $\mathcal F$  be any distribution in  $\Pi$ . Consider the following two constraints:

$$(\mathbb{E}[\mu_{\mathcal{F}}] - \hat{\mu})^T \hat{\mathbf{\Sigma}}^{-1} (\mathbb{E}[\mu_{\mathcal{F}}] - \hat{\mu}) \le \gamma_1$$

$$\mathbb{E}[(\mu_{\mathcal{F}} - \hat{\mu})(\mu_{\mathcal{F}} - \hat{\mu})^T] \le \gamma_2 \hat{\mathbf{\Sigma}},$$
(2)

where  $\hat{\mu}$  and  $\hat{\Sigma}$  are the mean and covariance matrix estimated from data.  $\mu_{\mathcal{F}}$  are random variables distributed according to  $\mathcal{F}$ . The first constraint defines an ellipsoid, which indicates that the expectation of  $\mathcal{F}$  lies in the ellipsoid centered at the estimate  $\hat{\mu}$ . The size of this ellipsoid is determined by  $\gamma_1$ , which provides a natural measure to quantify our uncertainty about  $\mu$  given  $\hat{\mu}$ . Note that the second constraint also defines the support of the distribution  $\mathcal{F}$ . The second constraint enforces the covariance matrix of  $\mathcal{F}$  to lie in a positive semi-definite cone. Intuitively, the second constraint captures how likely it is that the random variable  $\mu_{\mathcal{F}}$  is close to  $\hat{\mu}$ . The set  $\Pi$  is then characterized by Eq. (3):

$$\Pi(\hat{\mu}, \hat{\Sigma}, \gamma_1, \gamma_2) := \begin{cases} \mathcal{F} \middle| (\mathbb{E}[\mu_{\mathcal{F}}] - \hat{\mu})^T \hat{\Sigma}^{-1} (\mathbb{E}[\mu_{\mathcal{F}}] - \hat{\mu}) \leq \gamma_1 \\ \mathbb{E}[(\mu_{\mathcal{F}} - \hat{\mu})(\mu_{\mathcal{F}} - \hat{\mu})^T] \leq \gamma_2 \hat{\Sigma} \end{cases}$$
(3)

The set  $\Pi$  is always non-empty, since it must contain the distribution  $\hat{\mathcal{P}}$ . In Section 4 we provide probabilistic arguments to show that  $\Pi$  contains ground-truth distribution  $\mathcal{P}$ 

 $<sup>^{1}</sup>diag(\mathbf{x})$  returns a diagonal matrix with diagonal elements equal to  $\mathbf{x}$ .

with high probability, which guarantees that with high probability our model MINT\_DRO is robust with respect to the ground-truth distribution  $\mathcal{P}$ . The choice of the two parameters  $\gamma_1$  and  $\gamma_2$  is important for the robustness of MINT\_DRO. If their values are too small the benefit from the distributionally robust formulation is limited. In the extreme case where  $\gamma_1$  and  $\gamma_2$  are zeros our model MINT\_DRO reverts to MINT. On the other hand if their values are too large, our model would make excessively conservative decisions. In Section 4 we show how to make sensible choice of these values.

# 3. Solution Approach

In this section we derive the algorithm to solve our model MINT\_DRO. The optimization problem of MINT\_DRO is a binary quadratic program, which is diffcult to solve even if the loss function  $\mathcal{L}(\mathbf{x}; \mu, \Sigma)$  is convex. Additionally, in our problem the loss function is nonconvex since the matrix  $\mathbf{Q}$  is usually not positive (semi)-definite, further complicating the situation. Indeed, given that MINT, which was shown by Yu & Vorobeychik to be NP-Hard, is a special case, the following result is immediate.

**Theorem 1.** Solving MINT\_DRO is NP-Hard.

In what follows, we derive an approximation approach for solving MINT\_DRO. We first apply duality to transform the inner maximization into a minimization problem, which can be jointly minimized with the outer minimization over x. At this stage the optimization problem is still a NP-hard combinatorial optimization problem. Next, we apply Semidefinite Programming (SDP) to obtain a convex relaxation of our problem which can be solved efficiently.

The support of the distributions in  $\Pi$  is  $\mathcal{S}$ , which is defined as  $\mathcal{S} := \left\{ \mu_{\mathcal{F}} \,\middle|\, (\mu_{\mathcal{F}} - \hat{\mu})^T \hat{\Sigma}^{-1} (\mu_{\mathcal{F}} - \hat{\mu}) \leq \gamma_1 \right\}$ , where the subscript of  $\mu_{\mathcal{F}}$  indexes the distribution associated with this random variable. Note that  $\mu_{\mathcal{F}} \in \mathcal{S}$  is sufficient for the first constraint in Eq. (2) to be true, since  $\mathbb{E}[\mu_{\mathcal{F}}]$  is a convex combination of the instantiations of  $\mu_{\mathcal{F}}$  and  $\mathcal{S}$  is a convex set. We rewrite the inner maximization problem as Eq. (4):

$$\sup_{\mathcal{F} \sim \Pi} \int_{\mathcal{S}} \left[ \mathbf{x}^T \mathbf{Q}(\mu_{\mathcal{F}}, \hat{\boldsymbol{\Sigma}}) \mathbf{x} + 2 \mathbf{x}^T \mathbf{b}(\mu_{\mathcal{F}}) \right] d\mathcal{F}(\mu_{\mathcal{F}})$$
(4a)  
s.t. 
$$\int_{\mathcal{S}} d\mathcal{F}(\mu_{\mathcal{F}}) = 1$$
(4b)

$$\int_{\mathcal{S}} \left[ (\mu_{\mathcal{F}} - \hat{\mu})(\mu_{\mathcal{F}} - \hat{\mu})^T \right] d\mathcal{F}(\mu_{\mathcal{F}}) \preceq \gamma_2 \hat{\Sigma} \quad (4c)$$

$$\mu_{\mathcal{F}} \in \mathcal{S}, \forall \mu_{\mathcal{F}}.$$
 (4d)

The constraint Eq.(4b) ensures that  $\mathcal{F}$  is a valid probability distribution. The constraints Eq.(4c) guarantee that  $\mathcal{F}$  is in  $\Pi$ . The constraint Eq. (4d) ensures that any random variable  $\mu_{\mathcal{F}} \sim \mathcal{F}$  must reside in  $\mathcal{S}$ . Consequently, this constraint is actually an infinite dimensional constraint on the optimizer  $\mathcal{F}$ . Later we introduce a technique called *S-Lemma* 

to convert it to a finite dimensional constraint. We derive the lagrange function of Eq. (4), where we temporily omit constraint Eq. (4d), and pull the terms that are independent of  $\mathcal{F}$  out of the integral:

$$l(\mathcal{F}, t, \mathbf{K}) = \left[ t + Tr \left( \left[ \gamma_2 \hat{\mathbf{\Sigma}} + \hat{\mu} \hat{\mu}^T \right] \mathbf{K} \right) \right] + \int_{\mathcal{S}} \left[ \underbrace{\mathbf{x}^T \mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}}) \mathbf{x} + 2\mathbf{x}^T \mathbf{b}(\mu_{\mathcal{F}}) - t - \mu_{\mathcal{F}}^T \mathbf{K} \mu_{\mathcal{F}}}_{:= f(\mu_{\mathcal{F}})} \right],$$

where  $t \in \mathbb{R}$ , and  $\mathbf{K}$  is a real symmetric positive semidefinite matrix, and  $Tr(\mathbf{X})$  returns the trace of the matrix  $\mathbf{X}$ . where  $f(\mu_{\mathcal{F}}) \leq 0, \forall \mu_{\mathcal{F}} \in \mathcal{S}$  holds, since otherwise the solution to Eq.(4) is unbounded.

By duality (Shapiro, 2001; Delage & Ye, 2010; Cheng et al., 2014), the dual problem of Eq. (4) is formulated as the following minimization problem:

$$\min_{t,\mathbf{K}} \quad t + Tr\left(\left[\gamma_{2}\hat{\mathbf{\Sigma}} + \hat{\mu}\hat{\mu}^{T}\right]\mathbf{K}\right)$$

$$s.t. \quad f(\mu_{\mathcal{F}}) \leq 0, \forall \mu_{\mathcal{F}} \in \mathcal{S}$$

$$t \in \mathbb{R}, \mathbf{K} \in \mathbb{S}_{+}^{N}$$
(5)

where  $\mathbb{S}_+^N$  is the positive semi-definite cone. Strong duality holds between Eq. (5) and the original inner maximization problem. This is because for any  $\gamma_1>0$  and  $\gamma_2>0$ , the estimated distribution  $\hat{\mathcal{P}}$  is always in the relative interior of  $\Pi$ . Consequently, by Proposition 3.4 in Shapiro (2001) strong duality holds. Since Eq. (5) is a minimization problem, we can jointly minimize it with the outer minimization over  $\mathbf{x}$ , which results in the following:

$$\min_{\mathbf{x},t,\mathbf{K}} t + Tr\left(\left[\gamma_2 \hat{\mathbf{\Sigma}} + \hat{\mu}\hat{\mu}^T\right]\mathbf{K}\right)$$
 (6a)

s.t. 
$$(\mu_{\mathcal{F}} - \hat{\mu})^T \hat{\Sigma}^{-1} (\mu_{\mathcal{F}} - \hat{\mu}) \le \gamma_1, \forall \mu_{\mathcal{F}} \in \mathcal{S}$$
 (6b)

$$f(\mu_{\mathcal{F}}) \le 0, \forall \mu_{\mathcal{F}} \in \mathcal{S}$$
 (6c)

$$t \in \mathbb{R}, \mathbf{K} \in \mathbb{S}^N_+$$
 (6d)

where constraint Eq. (6b) is equivalent to  $\mu_{\mathcal{F}} \in \mathcal{S}, \forall \mu_{\mathcal{F}}$ . We write it this way in order to emphasize its quadratic form. Constraints Eq.(6b) and (6c) are infinite dimensional constraints. We apply a technique called *S-Lemma* to transform them to finite dimensional constraints. We first introduce the *S-Lemma*:

**Lemma 3.1** (S-Lemma (Boyd & Vandenberghe, 2004)). Let  $\mathbf{A}_1, \mathbf{A}_2 \in \mathbb{S}^n$ ,  $\mathbf{b}_1, \mathbf{b}_2 \in \mathbb{R}^n$ ,  $c_1, c_2 \in \mathbb{R}$ , where  $\mathbb{S}^n$  is the subspace of symmetrix matrices in  $\mathbb{R}^{n \times n}$ . Suppose there exists an  $\hat{\mathbf{x}} \in \mathbb{R}^n$  such that:  $\hat{\mathbf{x}}^T \mathbf{A}_1 \hat{\mathbf{x}} + 2\mathbf{b}_1^T \hat{\mathbf{x}} + c_1 < 0$ . Then the following implication holds for any  $\mathbf{x} \in \mathbb{R}^n$ :

$$\mathbf{x}^T \mathbf{A}_1 \mathbf{x} + 2 \mathbf{b}_1^T \mathbf{x} + c_1 \le 0 \implies \mathbf{x}^T \mathbf{A}_2 \mathbf{x} + 2 \mathbf{b}_2^T \mathbf{x} + c_2 \le 0$$

$$\textit{if and only if, } \exists \lambda \geq 0: \begin{bmatrix} \mathbf{A}_2 & \mathbf{b}_2 \\ \mathbf{b}_2^T & c_2 \end{bmatrix} \preceq \lambda \begin{bmatrix} \mathbf{A}_1 & \mathbf{b}_1 \\ \mathbf{b}_1^T & c_1 \end{bmatrix}.$$

Note that *S-Lemma* only requires  $A_1$  and  $A_2$  to be real symmetric matrices. In order to apply *S-Lemma* we need to have two quadratic functions. Constraint Eq. (6b) is a quadratic function in  $\mu_{\mathcal{F}}$ . Thus, what remains is to convert Eq. (6c) to a quadratic function in  $\mu_{\mathcal{F}}$ . Recall that the term,  $\mathbf{x}^T \mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}})\mathbf{x} + 2\mathbf{x}^T \mathbf{b}(\mu_{\mathcal{F}})$  in  $f(\mu_{\mathcal{F}})$ , is implicitly a quadratic function of  $\mu_{\mathcal{F}}$ . We re-formulate  $\mathbf{Q}$  and  $\mathbf{b}$  according to  $\mu_{\mathcal{F}}$ , which results in Eq.(7) (see the Appendix for details about this reformulation):

$$\mathbf{x}^{T}\mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}})\mathbf{x} + 2\mathbf{x}^{T}\mathbf{b}(\mu_{\mathcal{F}}) =$$

$$\mu_{\mathcal{F}} \left[ -(\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \left( \mathbf{x} \mathbf{x}^{T} \right) \right) \right] \mu_{\mathcal{F}}^{T} +$$

$$\mu_{\mathcal{F}}^{T} \left[ (\alpha_{3} + 2\alpha_{2}) diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} - \alpha_{1} \mathbf{x} \right] -$$

$$\left[ \alpha_{1} \mathbf{1}^{T} \mathbf{x} - \mathbf{x}^{T} \left( (\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) + \alpha_{2} \mathbf{A} \right) \mathbf{x} \right],$$
(6)

where  $diag(\mathbf{x})$  returns a diagonal matrix with diagonal elements equal to  $\mathbf{x}$ . We substitute Eq. (7) back to  $f(\mu_{\mathcal{F}})$ , which results in the following equivalence:

$$\forall \mu_{\mathcal{F}} \in \mathcal{S} : f(\mu_{\mathcal{F}}) \leq 0 \iff \mu_{\mathcal{F}}^T \mathbf{R} \mu_{\mathcal{F}} + \mu_{\mathcal{F}}^T \mathbf{r} + z \leq 0$$

where

$$\mathbf{R} = -(\alpha_2 + \alpha_3) \left( \mathbf{A} \odot \left( \mathbf{x} \mathbf{x}^T \right) \right) - \mathbf{K}$$

$$\mathbf{r} = (\alpha_3 + 2\alpha_2) diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} - \alpha_1 \mathbf{x}$$

$$z = \alpha_1 \mathbf{1}^T \mathbf{x} - \mathbf{x}^T \left( (\alpha_2 + \alpha_3) \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) + \alpha_2 \mathbf{A} \right) \mathbf{x} - t,$$

which results in a compact form of  $f(\mu_{\mathcal{F}})$ :

$$f(\mu_{\mathcal{F}}) = \mu_{\mathcal{F}}^T \mathbf{R} \mu_{\mathcal{F}} + \mu_{\mathcal{F}}^T \mathbf{r} + z$$

Note that for any  $\gamma_1 > 0$  the inequality in constraint Eq. (6b) is strict when  $\mu_{\mathcal{F}} = \hat{\mu}$ . Consequently, by *S-Lemma*, for any  $\mu_{\mathcal{F}} \in \mathcal{S}$  the implication, Eq. (6b)  $\implies \mu_{\mathcal{F}}^T \mathbf{R} \mu_{\mathcal{F}} + \mu_{\mathcal{F}}^T \mathbf{r} + z$ , is equivalent to Eq.(8):

$$\exists \lambda \geq 0 : \begin{bmatrix} \mathbf{R} & \frac{1}{2}\mathbf{r} \\ \frac{1}{2}\mathbf{r}^T & z \end{bmatrix} \leq \lambda \begin{bmatrix} \hat{\mathbf{\Sigma}}^{-1} & -\hat{\mathbf{\Sigma}}^{-1}\hat{\mu} \\ -\hat{\mu}^T\hat{\mathbf{\Sigma}}^{-1} & (\hat{\mu}^T\hat{\mathbf{\Sigma}}^{-1}\hat{\mu} - \gamma_1). \end{bmatrix}$$
(8)

The two infinite dimensional constraints Eq.(6b) and (6c) are thereby converted into a finite dimensional constraint Eq. (8). Additionally, the objective function in Eq. (6) is linear in its optimizer.

The last issue is that we still have two sources of non-convexity in Eq. (6): first,  $\mathbf{x}$  is binary, and second, the constraint represented by Eq. (8) is not convex in  $\mathbf{x}$  because of three terms involving in  $\mathbf{R}$ ,  $\mathbf{r}$  and z:

$$\mathbf{x}\mathbf{x}^{T}$$
,  $diag(\mathbf{x})\mathbf{A}\mathbf{x}$ ,  $\mathbf{x}^{T}\bigg((\alpha_{2}+\alpha_{3})(\mathbf{A}\odot\hat{\mathbf{\Sigma}})+\alpha_{2}\mathbf{A}\bigg)\mathbf{x}$ .

To deal with the first issues, we relax the feasible region of x to  $[-1,1]^N$ . To address the second, we next apply SDP relaxation to transform Eq. (6) into a convex optimization problem.

First, let us introduce a matrix  $\mathbf{X} = \mathbf{x}\mathbf{x}^T$ . Then the following three relationships hold (see the Appendix for detailed proof):

$$(r1): \left(\mathbf{A} \odot (\mathbf{x}\mathbf{x}^{T})\right) = \left(\mathbf{A} \odot \mathbf{X}\right)$$

$$(r2): diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} = diag(\mathbf{A}\mathbf{X})$$

$$(r3): \mathbf{x}^{T} \left((\alpha_{2} + \alpha_{3})(\mathbf{A} \odot \hat{\mathbf{\Sigma}}) + \alpha_{2}\mathbf{A}\right)\mathbf{x} =$$

$$(\alpha_{2} + \alpha_{3})Tr\left((\mathbf{A} \odot \hat{\mathbf{\Sigma}})\mathbf{X}\right) + \alpha_{2}Tr(\mathbf{A}\mathbf{X})$$

$$(10)$$

One problem is that the feasible regions involving  $\mathbf{X}$  and  $\mathbf{x}$  are nonconvex because of the equality  $\mathbf{X} = \mathbf{x}\mathbf{x}^T$ . In order to transform the feasible regions to be convex, we apply a two-step relaxation. The first step is to relax the equality and enforce the diagonal elements of  $\mathbf{X}$  equal to one, which results in:  $\mathbf{X} \succeq \mathbf{x}\mathbf{x}^T$  and  $\mathbf{X}_{ii} = 1, \forall i = 1, \cdots, N$ . This step transforms the feasible region of  $\mathbf{X}$  to a positive semi-definite cone, which is a convex set. However, we still have a nonconvex term  $\mathbf{x}\mathbf{x}^T$ . To handle this, in the second step we apply *Schur Complement* to transform  $\mathbf{X} \succeq \mathbf{x}\mathbf{x}^T$  to the linear matrix inequality:  $\begin{bmatrix} \mathbf{X} & \mathbf{x} \\ \mathbf{x}^T & 1 \end{bmatrix} \succeq 0$ . Combining the relationships in Eq. (10) with the results of the two-step relaxation above, the three nonconvex terms in Eq. (9) can be represented as the following convex set:

$$\hat{\mathbf{R}} = -(\alpha_2 + \alpha_3) (\mathbf{A} \odot \mathbf{X}) - \mathbf{K}$$

$$\hat{\mathbf{r}} = (2\alpha_2 + \alpha_3) \underbrace{diag(\mathbf{A}\mathbf{X})}_{(*)} - \alpha_1 \mathbf{x}$$

$$\hat{z} = \alpha_1 \mathbf{1}^T \mathbf{x} - \left( (\alpha_2 + \alpha_3) Tr((\mathbf{A} \odot \hat{\mathbf{\Sigma}}) \mathbf{X}) + \alpha_2 Tr(\mathbf{A}\mathbf{X}) \right)$$

$$- t$$

$$\begin{bmatrix} \mathbf{X} & \mathbf{x} \\ \mathbf{x}^T & 1 \end{bmatrix} \succeq 0, \mathbf{X}_{ii} = 1, \forall i = 1, \dots, N.$$

With a slight abuse of notation, the operator  $diag(\mathbf{A}\mathbf{X})$  in (\*) extracts the diagonal elements of  $\mathbf{A}\mathbf{X}$  as a column vector. Finally, by substituting  $\hat{\mathbf{R}}$ ,  $\hat{\mathbf{r}}$  and  $\hat{z}$  to the corresponding matrices in Eq. (8) we obtain the following Semidefinite Program which approximately solves MINT\_DRO (after we project the optimal solution  $\mathbf{x}$  of this problem into  $\{0,1\}^N$ , for example, by rounding):

$$\min_{\mathbf{x}, \mathbf{X}, t, \mathbf{K}, \lambda} \quad t + Tr\left(\left[\gamma_{2}\hat{\mathbf{\Sigma}} + \hat{\mu}\hat{\mu}^{T}\right]\mathbf{K}\right) 
s.t. \quad \begin{bmatrix} \hat{\mathbf{R}} & \frac{1}{2}\hat{\mathbf{r}} \\ \frac{1}{2}\hat{\mathbf{r}}^{T} & \hat{z} \end{bmatrix} \preceq \lambda \begin{bmatrix} \hat{\mathbf{\Sigma}}^{-1} & -\hat{\mathbf{\Sigma}}^{-1}\hat{\mu} \\ -\hat{\mu}^{T}\hat{\mathbf{\Sigma}}^{-1} & (\hat{\mu}^{T}\hat{\mathbf{\Sigma}}^{-1}\hat{\mu} - \gamma_{1}) \end{bmatrix} 
\begin{bmatrix} \mathbf{X} & \mathbf{x} \\ \mathbf{x}^{T} & 1 \end{bmatrix} \succeq 0, \mathbf{X}_{ii} = 1, \forall i = 1, \dots, N 
\mathbf{x} \in [-1, 1]^{N}, t \in \mathbb{R}, \mathbf{K} \in \mathbb{S}_{+}^{N}, \mathbf{X} \in \mathbb{S}_{+}^{N}, \lambda \geq 0$$
(11)

# 4. Theoretical Analysis

In this section we present a probabilistic argument that the uncertainty set  $\Pi$  defined in Eq. (3) contains the ground-truth distribution  $\mathcal{P}$  with high probability. This, in turn, implies that with high probability our model MINT\_DRO is robust with respect to the *unknown* ground-truth distribution.

We show that the ground-truth distribution  $\mathcal{P}$  belongs to  $\Pi$  with high probability in two steps, arguing first that (C1) and, subsequently, that (C2) below hold with high probability, where (C1) and (C2) are defined as follows:

$$(\mathbb{E}_{\pi \sim \mathcal{P}}[\pi] - \hat{\mu})^T \hat{\Sigma}^{-1} (\mathbb{E}_{\pi \sim \mathcal{P}}[\pi] - \hat{\mu}) \le \gamma_1$$
 (C1)

$$\mathbb{E}_{\pi \sim \mathcal{P}} \left[ (\pi - \hat{\mu})(\pi - \hat{\mu})^T \right] \leq \gamma_2 \hat{\Sigma}$$
 (C2)

The arguments in the first step are based on Lemma 4.1. For space limitation we defer its proof to the appendix.

**Lemma 4.1.** Let  $\mu$  and  $\Sigma$  denote the mean and covariance matrix of the ground-truth distribution  $\mathcal{P}$ , and suppose that  $\hat{\mu}$  is estimated from M samples,  $\hat{\mu} = \frac{1}{M} \sum_{i=1}^{M} \zeta_i$ , where  $\zeta_i$  is bounded:  $\|\mathbf{\Sigma}^{-1/2}(\zeta_i - \mu)\|_2^2 \leq R^2, \forall i$ . Then  $\hat{\mu}$  satisfies the following constraint with probability at least  $1 - \delta_1$ :

$$(\mu - \hat{\mu})^T \mathbf{\Sigma}^{-1} (\mu - \hat{\mu}) \le \beta(\delta),$$

where 
$$\beta(\delta_1) = \frac{R^2}{M} \bigg( 2 + \sqrt{2 \log \frac{1}{\delta_1}} \bigg)^2$$
.

We assume the estimated covariance matrix  $\hat{\Sigma}$  is close to  $\Sigma$ . Then, if we let  $\gamma_1 > \beta(\delta_1)$  and note that  $\mu = \mathbb{E}_{\pi \sim \mathcal{P}}[\pi]$ , a direct application of Lemma 4.1 implies that (C1) holds with probability at least  $1 - \delta_1$ .

The arguments in the second step rely on the result due to (Delage & Ye, 2010):

**Lemma 4.2** (Delage & Ye (2010)). Suppose that  $\zeta_i$  is distributed according to  $\mathcal{G}$ , and the mean  $\mu$  of the distribution is known and used to formulate the estimated covariance matrix  $\hat{\Sigma}$ , which is estimated from M samples:  $\hat{\Sigma} = (1/M) \sum_{i=1}^{M} (\zeta_i - \mu) (\zeta_i - \mu)^T$ , where  $\zeta_i$  is bounded:  $\|\Sigma^{-1/2}(\zeta_i - \mu)\|_2^2 \leq R^2, \forall i$ . Then with probability at least  $1 - \delta_2$ :

$$\Sigma \leq \frac{1}{1 - \alpha(\delta_2)} \hat{\Sigma},$$

$$t + Tr\left(\left[\gamma_2\hat{\boldsymbol{\Sigma}} + \hat{\mu}\hat{\mu}^T\right]\mathbf{K}\right) \qquad \qquad \text{where } \alpha(\delta_2) = (R^2/\sqrt{M})\left(\sqrt{1 - N/R^4} + \sqrt{\log 1/\delta_2}\right),$$

$$\begin{bmatrix} \hat{\mathbf{R}} & \frac{1}{2}\hat{\mathbf{r}} \\ \frac{1}{2}\hat{\mathbf{r}}^T & \hat{z} \end{bmatrix} \leq \lambda \begin{bmatrix} \hat{\boldsymbol{\Sigma}}^{-1} & -\hat{\boldsymbol{\Sigma}}^{-1}\hat{\mu} \\ -\hat{\mu}^T\hat{\boldsymbol{\Sigma}}^{-1} & (\hat{\mu}^T\hat{\boldsymbol{\Sigma}}^{-1}\hat{\mu} - \gamma_1) \end{bmatrix} \xrightarrow{\text{mensions of } u.} M > R^4\left(\sqrt{1 - N/R^4} + \sqrt{\log 1/\delta_2}\right)^2 \text{ and } N \text{ is the dimensions of } u.$$

In order to use Lemma 4.2 we assume that the estimated mean  $\hat{\mu}$  is close to the ground-truth  $\mu$ . Given this assumption, showing that (C2) holds with high probability is equivalent to show that the following holds with high probability:

$$\mathbb{E}_{\pi \sim \mathcal{P}}[\pi \pi^T] \preceq \gamma_2 \hat{\Sigma} + \mu \mu^T$$

by Lemma 4.2, the above is true with high probability when:  $\frac{1}{1-\alpha(\delta_2)}\hat{\Sigma} \preceq \gamma_2\hat{\Sigma} + \mu\mu^T. \text{ Consequently, by setting } \gamma_2 > \frac{1}{1-\alpha(\delta_2)}, \text{ such that the effects of } \mu\mu^T \text{ are negligible, we conclude that (C2) holds with probability at least } 1-\delta_2.$ 

Finally, by a union bound we obtain probabilistic guarantees that the uncertainty set  $\Pi$  contains  $\mathcal{P}$ .

**Theorem 2.** With probability at least  $1 - \delta$ , where  $\delta = \delta_1 + \delta_2$ , the uncertainty set  $\Pi$  defined in Eq. (3) contains the ground-truth distribution  $\mathcal{P}$ .

*Proof.* The detailed proof is deferred to the appendix.  $\Box$ 

We now demonstrate how to utilize the probabilistic arguments to make sensible choice for  $\gamma_1$ . The value of  $\gamma_2$  can be similarly obtained. Note that  $\gamma_1 > \beta(\delta_1)$  is necessary for (C1) to hold. Consider a network with N=128 nodes. Assume  $\Sigma$  is diagonal with diagonal elements equal to 0.01, which is reasonable when a single estimator is used to estimate  $\mathcal P$  and the maliciousness probabilities of nodes are independent. A reasonable estimate of R is  $\sqrt{128\times 2}$ , which is the radius of the circumcircle sphere of a hypercube with length of side equal to one. If M=5 and  $\delta_1=0.05$ , then  $\beta(0.05)=1012$ . Therefore in order for  $\Pi$  to contain  $\mathcal P$  with probability  $\geq 0.95$ , we need  $\gamma_1 \geq 1012$ . Similarly, for a network with N=500 nodes, we want  $\gamma_1 \geq 3956$ .

# 5. Experiments

In this section we present experimental results to show the effectiveness of our approach. Our experiments were conducted on both synthetic and real-world network structures, although in all cases the distribution  $\mathcal{P}$  over maliciousness of nodes was derived using real data. We considered two types of network generative models to construct synthetic networks: Barabasi-Albert (BA) (Barabási & Albert, 1999) and Watts-Strogatz networks (Small-World) (Watts & Strogatz, 1998). BA is characterized by its power-law degree distribution, where the probability that a randomly selected node has k neighbors is proportional to  $k^{-r}$ . For the BA model we experimented with three variants, BA-1, BA-2, and BA-3, which differ in the value of the exponent r of their power-law degree distributions. For Small-World networks we also experimented with three variants, SW-1, SW-2, and

SW-3, that have different local clustering coefficients. For both networks we generated instances with N=128 nodes. For real-world networks, we used a network extracted from Facebook data (Leskovec & Mcauley, 2012) which consisted of 4039 nodes and 88234 edges. We experimented with randomly sampled sub-networks with N=500 nodes. For space limitation the statistics of the networks used in our experiments are listed in the appendix.

For fair comparison with MINT (the state-of-the-art alternative), we used the same experimental setup as Yu & Vorobeychik (2018). In all of our experiments, we derived the ground-truth distribution  $\mathcal{P}$  as follow. We start with a dataset D which includes malicious and benign instances (the meaning of these designations is domain specific). The dataset **D** is partitioned into three subsets:  $\mathbf{D}_{train}$ ,  $\mathbf{D}_1$  and  $\mathbf{D}_2$ , with the ratio of 0.3:0.6:0.1. Our first step is to learn a probabilistic predictor of maliciousness as a function of a feature vector  $\mathbf{x}$ ,  $\hat{p}(\mathbf{x})$ , on  $\mathbf{D}_{train}$ . Then we randomly assign malicious and benign feature vectors from  $\mathbf{D}_2$  to the nodes on the network, assigning 10% of nodes with malicious features and 90% with benign feature vectors. For each node we use its assigned feature vector x to obtain our estimated probability of this node being malicious,  $\hat{p}(\mathbf{x})$ ; This gives us the estimated maliciousness probability distribution  $\mathcal{P}$ . This is the distribution used to solve the model MINT, and also the distribution used to construct the uncertainty set  $\Pi$  in our model. To ensure that our evaluation reasonably reflects realistic limitations of the knowledge about the ground-truth distribution  $\mathcal{P}$ , we train another predictor  $p(\mathbf{x})$  usign  $\mathbf{D}_{train} \bigcup \mathbf{D}_1$ . Applying this new predictor to the nodes and their assigned feature vectors, we obtain a distribution  $\mathcal{P}^*$  which we use to evaluate effectiveness.

We conducted three sets of experiments. In the first set of experiments we used synthetic networks and used data from the Spam (Cormack et al., 2008) dataset To simulate estimation error of  $\mathcal{P}$ , we add white Gaussian noise to the evaluation distribution  $\mathcal{P}^*$ . The standard deviation of the noise is increased from 0.1 to 0.5 to simulate different magnitudes of the estimation error.

In the second set of experiments we used real-world networks from Facebook and used Hate Speech data (Davidson et al., 2017) collected from Twitter to obtain  $\mathcal P$  as discussed above. We categorized this dataset into two classes in terms of whether a tweet represents Hate Speech. After categorization, the total number of tweets is 24783, of which 1430 are Hate Speech. We add white Gaussian noise to  $\mathcal P^*$  to simulate estimation error as discussed above. Note that in this set of experiments we used *real data* for both the networks and the maliciousness probabilities  $\mathcal P$ .

In the third set of experiments we considered the scenario that instead of being random, the location of the malicious nodes on the network is strategically determined. This scenario is not vacuous: in reality, for example, the nodes that have high degrees (e.g., celebrities with lots of followers on Twitter) may be targeted in order to maximize the influence of commercial advertisements (Kempe et al., 2003). We conducted this set of experiments on synthetic networks. A set of nodes is greedily selected from the network to maximize the number of unique neighbors connecting to them. Then we assign malicious feature vectors to these nodes.

**Experiment Results** We compared our model with a state-of-the-art approach MINT. The average losses for our first set of experiments where  $\mathcal{P}$  was simulated from Spam data are shown in Figures 2 and 3. The experimental results on BA are showed in Figure 2, with the three columns corresponding to BA-1, BA-2 and BA-3, respectively. The experimental results on Small-World are shown in Figure 3, where the three columns correspond to SW-1, SW-2, and SW-3. In both figures, each row corresponds to a combination of trade-off parameters  $(\alpha_1, \alpha_2, \alpha_3)$ ; for example, (0.2, 0.7, 0.1) corresponds to  $(\alpha_1 = 0.2, \alpha_2 = 0.7, \alpha_3 = 0.1)$ . Each bar was obtained by averaging over 30 randomly generated network topologies.

The experimental results indicate that on both BA and Small-World networks our model MINT\_DRO is significantly more robust than MINT. Additionally, when no noise is added to the evaluation distribution  $\mathcal{P}^*$  (left-most bars in all subplots), MINT\_DRO is more robust than MINT except for a few cases. this indicates that the *generalization* ability of MINT\_DRO is better than MINT.

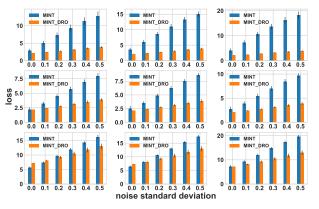


Figure 2. Experimental results on BA networks. The three columns correspond to results on BA-1, BA-2 and BA-3, respectively. **Top** row: (0.2, 0.7, 0.1); **Middle row**: (0.7, 0.2, 0.1); **Bottom row**:  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ .

The average loss on Facebook data is showed in Figure 4, with the three columns corresponding to (0.2,0.7,0.1), (0.7,0.2,0.1), and  $(\frac{1}{3},\frac{1}{3},\frac{1}{3})$ . In this experiment, both the networks and the data used to simulate maliciousness probabilities are *real* data. Each bar was averaged over 30 randomly sampled networks. Our model MITN\_DRO is significantly more robust than MINT except for the cases where no

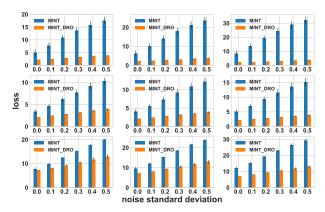


Figure 3. Experimental results on Small-World networks. The three columns correspond to results on SW-1, SW-2 and SW-3, respectively. **Top row**: (0.2, 0.7, 0.1); **Middle row**: (0.7, 0.2, 0.1); **Bottom row**:  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ .

noise is added. In this case, MINT\_DRO is only worse than MINT at the left-most bars in the middle figure, although the difference is not significant. This is actually expected since MINT\_DRO considers the *worst-case* scenario, which results in a decision that may be slightly conservative in no noise setting. One observation is that the Facebook networks used in this experiment are dramatically different from the simulated networks in terms of graph statistics (see the appendix for the detailed statistics). Particularly, the Facebook networks are disconnected, highly sparse, and have approximately 16% nodes that have zero degree. Therefore the robustness exhibited in Figure 4 provides strong evidence to the effectiveness of MINT\_DRO.

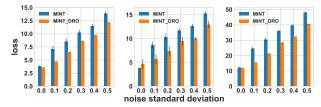


Figure 4. Experimental results on Facebook networks. The three columns correspond to (0.2, 0.7, 0.1), (0.7, 0.2, 0.1), and  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ .

The average loss on the third set of experiments are shown in Figures 5 and 6 for BA and Small-World networks, respectively. For both figures the three columns correspond to (0.2,0.7,0.1), (0.7,0.2,0.1), and  $(\frac{1}{3},\frac{1}{3},\frac{1}{3})$ . The results show that MINT\_DRO is more robust than MINT across all settings. Recall that the loss function of MINT and MINT\_DRO depends on the estimated covariance matrix  $\hat{\Sigma}$ , which encodes correlation information of the distribution  $\hat{\mathcal{P}}$ . When the actual maliciousness of nodes become correlated as we simulated in this experiment, the performance of MINT degrades since it is using the estimated distribution  $\hat{\mathcal{P}}$  which now significantly deviates from the true distribu-

tion. When  $\gamma_1$  and  $\gamma_2$  are appropriately selected,  $\Pi$  contains the distribution that characterizes the strategic correlation simulated in this experiment, resulting in significantly better robustness.

One may argue that instead of resulting from the robustness against correlation in the maliciousness distribution that comes from strategic decision about where to place the malicious nodes, the robustness exhibited in Figures 5 and 6 stems solely from the fact that MINT\_DRO is more robust than MINT when no noise is added to  $\mathcal{P}^*$ . However, consider the left-most bars in the lower-left subplot of Figure 2. In this setting MINT\_DRO performs worse than MINT. Now, consider another setting where the experimetal setup is identical except that the malicious nodes are strategically chosen. This setting corresponds to the left-most bars in the right subplot of Figure 5 where MINT\_DRO performs better than MINT. Similar observations can be found on Small-World networks. Consequently, we can see that a major advantage of MINT\_DRO is in its robustness even when the location of the malicious nodes on the graph is itself chosen strategically.

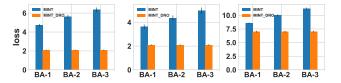


Figure 5. Experimental results on the robustness to strategic selection of malicious nodes on BA networks. Left: (0.2, 0.7, 0.1); Middle: (0.7, 0.2, 0.1); Right:  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ .

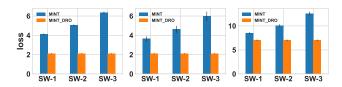


Figure 6. Experimental results on the robustness to strategic selection of malicious nodes on Small-World networks. **Left**: (0.2, 0.7, 0.1); **Middle**: (0.7, 0.2, 0.1); **Right**:  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ 

### 6. Conclusion

We considered the problem of removing malicious nodes from a network under uncertainty. We designed a model that considers the uncertainty around the estimated maliciousness probabilities, and makes decision under the *worst-case* scenario. We then proposed a principled algorithmic technique for solving it approximately based on duality combined with Semidefinite Programming relaxation. We theoretically proved that our model is robust with respect to the ground-truth, and experimentally showed that our model is more robust than the state of the art.

#### References

- Allcott, H. and Gentzkow, M. Social media and fake news in the 2016 election. *Journal of Economic Perspectives*, 31(2):211–36, 2017.
- Andrade, V. Facebook, whatsapp step up efforts in brazil's fake news battle. *Bloomberg*, 2018. URL https://www.bloomberg.com/news/articles/2018-10-23/facebook-whatsapp-step-up-efforts-in-brazil-s-fake-news-battle.
- Arias-Castro, E., Candes, E. J., and Durand, A. Detection of an anomalous cluster in a network. *The Annals of Statistics*, pp. 278–304, 2011.
- Barabási, A.-L. and Albert, R. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- Bertsimas, D. and Sethuraman, J. Moment problems and semidefinite optimization. In *Handbook of semidefinite programming*, pp. 469–509. Springer, 2000.
- Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Calafiore, G. C. and El Ghaoui, L. On distributionally robust chance-constrained linear programs. *Journal of Optimization Theory and Applications*, 130(1):1–22, 2006.
- Cheng, J., Delage, E., and Lisser, A. Distributionally robust stochastic knapsack problem. *SIAM Journal on Optimization*, 24(3):1485–1506, 2014.
- Cheng, J., Danescu-Niculescu-Mizil, C., and Leskovec, J. Antisocial behavior in online discussion communities. In *ICWSM*, pp. 61–70, 2015.
- Cormack, G. V. et al. Email spam filtering: A systematic review. *Foundations and Trends in Information Retrieval*, 1(4):335–455, 2008.
- Davidson, T., Warmsley, D., Macy, M., and Weber, I. Automated hate speech detection and the problem of offensive language. *arXiv* preprint arXiv:1703.04009, 2017.
- Delage, E. and Ye, Y. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations research*, 58(3):595–612, 2010.
- Goemans, M. X. and Williamson, D. P. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM (JACM)*, 42(6):1115–1145, 1995.
- Kempe, D., Kleinberg, J., and Tardos, É. Maximizing the spread of influence through a social network. In *Proceed*ings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 137–146. ACM, 2003.

- Leskovec, J. and Mcauley, J. J. Learning to discover social circles in ego networks. In *Advances in neural information processing systems*, pp. 539–547, 2012.
- Luo, Z.-Q., Ma, W.-K., So, A. M.-C., Ye, Y., and Zhang, S. Semidefinite relaxation of quadratic optimization problems. *IEEE Signal Processing Magazine*, 27(3):20–34, 2010.
- Mo, Y., Kim, T. H.-J., Brancik, K., Dickinson, D., Lee, H., Perrig, A., and Sinopoli, B. Cyber–physical security of a smart grid infrastructure. *Proceedings of the IEEE*, 100 (1):195–209, 2012.
- Narayanan, V., Barash, V., Kelly, J., Kollanyi, B., Neudert, L.-M., and Howard, P. N. Polarization, partisanship and junk news consumption over social media in the us. *arXiv* preprint arXiv:1803.01845, 2018.
- Popescu, I. Robust mean-covariance solutions for stochastic optimization. *Operations Research*, 55(1):98–112, 2007.
- Priebe, C. E., Conroy, J. M., Marchette, D. J., and Park, Y. Scan statistics on enron graphs. *Computational & Mathematical Organization Theory*, 11(3):229–247, 2005.
- Rodriguez, J. Facebook suspends 115 accounts for 'inauthentic behavior' as polls open, 2018. URL https://www.politico.com/story/2018/11/06/facebook-suspends-accounts-polls-2018-964325.
- Scarf, H. A min-max solution of an inventory problem. *Studies in the mathematical theory of inventory and production*, 1958.
- Scott, S. and Isaac, M. Facebook says it's policing fake accounts. but they're still easy to spot. *The New York Times*, 2017. URL https://www.nytimes.com/2017/11/03/technology/facebook-fake-accounts.html.
- Shapiro, A. On duality theory of conic linear problems. In *Semi-infinite programming*, pp. 135–165. Springer, 2001.
- Sharpnack, J. L., Krishnamurthy, A., and Singh, A. Nearoptimal anomaly detection in graphs using lovasz extended scan statistic. In *Advances in Neural Information Processing Systems*, pp. 1959–1967, 2013.
- Shawe-Taylor, J. and Cristianini, N. Estimating the moments of a random vector with applications. In *19 Colloque sur le traitement du signal et des images*, *FRA*, *2003*. GRETSI, Groupe dEtudes du Traitement du Signal et des Images, *2003*.
- Vorobeychik, Y. and Kantarcioglu, M. *Adversarial Machine Learning*. Morgan & Claypool, 2018.

- Watts, D. J. and Strogatz, S. H. Collective dynamics of small-world networks. *nature*, 393(6684):440, 1998.
- Wiesemann, W., Kuhn, D., and Sim, M. Distributionally robust convex optimization. *Operations Research*, 62(6): 1358–1376, 2014.
- Xu, H. and Mannor, S. Distributionally robust markov decision processes. In Advances in Neural Information Processing Systems, pp. 2505–2513, 2010.
- Yang, Y., Nishikawa, T., and Motter, A. E. Small vulnerable sets determine large network cascades in power grids. *Science*, 358(886), 2017.
- Yu, S. and Vorobeychik, Y. Removing malicious nodes from networks. *arXiv preprint arXiv:1812.11448*, 2018.
- Yue, J., Chen, B., and Wang, M.-C. Expected value of distribution information for the newsvendor problem. *Operations research*, 54(6):1128–1136, 2006.

# **Appendix**

# 1. Proof of Lemma 4.1

The proof is a generalization of a result proved by Shawe-Taylor & Cristianini. For completeness we list their result in Lemma 1.1.

Lemma 1.1. (Shawe-Taylor & Cristianini, 2003)

Assume  $\zeta \in \mathbb{R}^N$  is a random variable satisfying:

$$\mathbb{E}[\zeta] = \mathbf{0}$$

$$\mathbb{E}[\zeta \zeta^T] = \mathbf{I}$$

$$\|\zeta\|_2^2 \le R^2,$$

where the last inequality bounds the support of  $\zeta$ . Let  $\{\zeta_i\}_{i=1}^M$  be a set of M independently and ramdomly sampled instances of  $\zeta$ . Then with probability at least  $(1 - \delta)$ , the following inequality holds:

$$\left\| \frac{1}{M} \sum_{i=1}^{M} \zeta_i \right\|^2 \le \frac{R^2}{M} \left( 2 + \sqrt{2 \log \frac{1}{\delta}} \right)^2$$

In what follows we prove Lemma 4.1:

**Lemma 4.1.** Let  $\mu$  and  $\Sigma$  denote the mean and covariance matrix of the ground-truth distribution  $\mathcal{P}$ , and suppose that  $\hat{\mu}$  is estimated from M samples,  $\hat{\mu} = \frac{1}{M} \sum_{i=1}^{M} \zeta_i$ , where  $\zeta_i$  is bounded:  $\|\mathbf{\Sigma}^{-1/2}(\zeta_i - \mu)\|_2^2 \leq R^2$ ,  $\forall i$ . Then  $\hat{\mu}$  satisfies the following constraint with probability at least  $1 - \delta_1$ :

$$(\mu - \hat{\mu})^T \mathbf{\Sigma}^{-1} (\mu - \hat{\mu}) \le \beta(\delta),$$

where 
$$\beta(\delta_1) = \frac{R^2}{M} \left( 2 + \sqrt{2 \log \frac{1}{\delta_1}} \right)^2$$
.

*Proof.* Apply a standadization to the  $\zeta_i$ , which results in a new random variable  $\gamma_i := \mathbf{\Sigma}^{-1/2}(\zeta_i - \mu)$ . It is clear that  $\gamma_i$  satisfies Lemma 1.1. Let  $\beta(\delta_1) = \frac{R^2}{M} \left(2 + \sqrt{2\log\frac{1}{\delta_1}}\right)^2$ , then we have:

$$\mathbb{P}\left((\hat{\mu} - \mu)^T \mathbf{\Sigma}^{-1}(\hat{\mu} - \mu) \leq \beta(\delta_1)\right) = \mathbb{P}\left(\left\|\mathbf{\Sigma}^{-1/2} \left(\hat{\mu} - \mu\right)\right\|_2^2 \leq \beta(\delta_1)\right) \\
= \mathbb{P}\left(\left\|\mathbf{\Sigma}^{-1/2} \left(\frac{1}{M} \sum_{i=1}^M \zeta_i - \mu\right)\right\|_2^2 \leq \beta(\delta_1)\right) \\
= \mathbb{P}\left(\left\|\frac{1}{M} \sum_{i=1}^M \mathbf{\Sigma}^{-1/2} \left(\zeta_i - \mu\right)\right\|_2^2 \leq \beta(\delta_1)\right) \\
= \mathbb{P}\left(\frac{1}{M} \left\|\sum_{i=1}^M \gamma_i\right\|_2^2 \leq \beta(\delta_1)\right) \geq 1 - \delta_1$$

### 2. Proof of Theorem 2

**Theorem 2.** With probability at least  $1 - \delta$ , where  $\delta = \delta_1 + \delta_2$ , the uncertainty set  $\Pi$  contains the ground-truth distribution  $\mathcal{P}$ .

*Proof.* We define two events  $A_1$  and  $A_2$  as follow:

 $A_1$ : (C1) holds given  $\hat{\Sigma}$  is close to  $\Sigma$  $A_2$ : (C2) holds given  $\hat{\mu}$  is close to  $\mu$ 

Then we have:

$$\begin{split} \mathbb{P}\big(A_1 \cap A_2\big) &= \mathbb{P}\bigg((A_1^c \cup A_2^c)^c\bigg) \\ 1 - \mathbb{P}\big(A_1^c \cup A_2^c\big) \\ \text{(by union bound)} \\ &\geq 1 - \big[\mathbb{P}\big(A_1^c\big) + \mathbb{P}\big(A_1^c\big)\big] \\ &\geq 1 - \big[\delta_1 + \delta_2\big] \\ &= 1 - \delta, \end{split}$$

where  $A_1 \cap A_2$  is the event that  $\mathcal{P} \in \Pi$ . In other words,  $\mathbb{P}(\mathcal{P} \in \Pi) \geq 1 - \delta$ , which completes the proof.

# 3. Detailed dependency of $B(\mu), P(\mu, \Sigma), M(\mu, \Sigma)$ on their arguments

In the following we expand the definition of  $\mathbf{B}(\mu)$ ,  $\mathbf{P}(\mu, \Sigma)$ ,  $\mathbf{M}(\mu, \Sigma)$ , which makes their dependency on  $\mu$  and  $\Sigma$  clear:

$$\begin{split} \mathbf{B}(\mu) &:= diag \big( \mathbb{E}_{\pi \sim \mathcal{P}}[\bar{\pi}] \big) \\ &= diag \big( 1 - \mu \big) \\ \mathbf{P}(\mu, \mathbf{\Sigma}) &:= \mathbf{A} \odot \mathbb{E}_{\pi \sim \mathcal{P}}[\bar{\pi}\bar{\pi}^T] \\ &= \mathbf{A} \odot \left( \mathbf{J}(N, N) - \mathbf{J}(N, 1) \times \mu^T - \mu \times \mathbf{J}(1, N) + \mathbf{\Sigma} + \mu \times \mu^T \right) \\ \mathbf{M}(\mu, \mathbf{\Sigma}) &:= \mathbf{A} \odot \mathbb{E}_{\pi \sim \mathcal{P}}[\pi\bar{\pi}^T] \\ &= \mathbf{A} \odot \left( \mu \times \mathbf{J}(1, N) - \mathbf{\Sigma} - \mu \times \mu^T \right) \end{split}$$

### **4.** Detailed forms of the matrices $Q(\mu, \Sigma)$ and $b(\mu)$

The matrices  $\mathbf{Q}(\mu, \Sigma)$  and  $\mathbf{b}(\mu)$  defined in the paper have the following forms:

$$\begin{split} \mathbf{Q}(\mu, \mathbf{\Sigma}) &:= \frac{\alpha_3 \bigg( \mathbf{M}(\mu, \mathbf{\Sigma}) + \mathbf{M}(\mu, \mathbf{\Sigma})^T \bigg)}{2} - \frac{\alpha_2 \bigg( \mathbf{P}(\mu, \mathbf{\Sigma}) + \mathbf{P}(\mu, \mathbf{\Sigma})^T \bigg)}{2} \\ &= \alpha_3 \mathbf{A} \odot \left[ \frac{\mu \mathbf{1}^T + \mathbf{1} \mu^T}{2} - \mathbf{\Sigma} - \mu \mu^T \right] - \alpha_2 \mathbf{A} \odot \left[ J(N, N) - \mathbf{1} \mu^T - \mu \mathbf{1}^T + \mathbf{\Sigma} + \mu \mu^T \right] \\ &= \mathbf{A} \odot \left[ \bigg( \frac{\alpha_3 + 2\alpha_2}{2} \bigg) \mu \times \mathbf{1}^T + \bigg( \frac{\alpha_3 + 2\alpha_2}{2} \bigg) \mathbf{1} \times \mu^T - (\alpha_2 + \alpha_3) \mathbf{\Sigma} - (\alpha_2 + \alpha_3) \mu \times \mu^T - \alpha_2 J(N, N) \right] \\ &= \mathbf{A} \odot \left[ \bigg( \frac{\alpha_3 + 2\alpha_2}{2} \bigg) (\mu - \mathbf{1}) \times \mathbf{1}^T + \bigg( \frac{\alpha_3 + 2\alpha_2}{2} \bigg) \mathbf{1} \times (\mu - \mathbf{1})^T + (\alpha_2 + \alpha_3) J(N, N) \right. \\ &- (\alpha_2 + \alpha_3) \mathbf{\Sigma} - (\alpha_2 + \alpha_3) \mu \times \mu^T \right] \\ &\mathbf{b}(\mu) := (\alpha_1/2) \mathbf{B}(\mu) \mathbf{1} \\ &= (\alpha_1/2) (\mathbf{1} - \mu) \end{split}$$

# 5. Detailed reformulation of Eq. (7)

In the paper in order to apply the *S-Lemma* to convert the two infinite dimensional constraints, Eq. (6b) and Eq. (6c), to a finite dimensional constraint, we need two functions in quadratic forms. Notice that Eq. (6b) is already a quadratic function in  $\mu_{\mathcal{F}}$ . So what remains is to convert Eq. (6c) to a quadratic function in  $\mu_{\mathcal{F}}$ . We first convert the following to a quadratic function in  $\mu_{\mathcal{F}}$ :

$$\mathbf{x}^T \mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}}) \mathbf{x} + 2 \mathbf{x}^T \mathbf{b}(\mu_{\mathcal{F}})$$

From last section we know:

$$\mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}}) = \left(\frac{\alpha_{3} + 2\alpha_{2}}{2}\right) \underbrace{\left[\mathbf{A} \odot \left((\mu_{\mathcal{F}} - \mathbf{1}) \times \mathbf{1}^{T}\right)\right]}_{\widehat{\mathbf{I}}} + \left(\frac{\alpha_{3} + 2\alpha_{2}}{2}\right) \underbrace{\left[\mathbf{A} \odot \left(\mathbf{1} \times (\mu_{\mathcal{F}} - \mathbf{1})^{T}\right)\right]}_{\widehat{\mathbf{I}}} + \underbrace{\left(\alpha_{2} + \alpha_{3}\right)\mathbf{A} - (\alpha_{2} + \alpha_{3})\left(\mathbf{A} \odot \hat{\mathbf{\Sigma}}\right) - (\alpha_{2} + \alpha_{3})\mathbf{A} \odot \left(\mu_{\mathcal{F}} \times \mu_{\mathcal{F}}^{T}\right)}_{\widehat{\mathbf{3}}}$$

$$(12)$$

The three terms (1), (2) and (3), together with  $\mathbf{x}$ , form three quadratic functions in  $\mathbf{x}$ . In what follows, we convert them to quadratic functions in  $\mu_{\mathcal{F}}$ . Note that the operator  $diag(\mathbf{x})$  returns a diagonal matrix with diagonal elements equal to  $\mathbf{x}$ :

$$\mathbf{x}^{T} \underbrace{\left[ \mathbf{A} \odot \left( \mathbf{1} \times (\mu_{\mathcal{F}} - \mathbf{1})^{T} \right) \right]}_{\boxed{2}} \mathbf{x} \stackrel{(*)}{=} Tr \left[ diag(\mathbf{x}) \cdot \mathbf{A} \cdot diag(\mathbf{x}) \cdot (\mu_{\mathcal{F}} - \mathbf{1}) \times \mathbf{1}^{T} \right] =$$

(the trace operator is invariant under cyclic permutations)

$$\begin{split} &= Tr \left[ \mathbf{A} \cdot diag(\mathbf{x}) \cdot (\mu_{\mathcal{F}} - \mathbf{1}) \times \mathbf{1}^{T} \cdot diag(\mathbf{x}) \right] \\ &= Tr \left[ \mathbf{A} \cdot diag(\mathbf{x}) \cdot (\mu_{\mathcal{F}} - \mathbf{1}) \cdot \mathbf{x}^{T} \right] \\ &\left( Tr \left[ \mathbf{A} \cdot diag(\mathbf{x}) \cdot \mathbf{1} \cdot \mathbf{x}^{T} \right] = Tr \left[ \mathbf{A} \mathbf{x} \mathbf{x}^{T} \right] = \mathbf{x}^{T} \mathbf{A} \mathbf{x} \right) \\ &= Tr \left[ \mathbf{A} \cdot diag(\mathbf{x}) \cdot \mu_{\mathcal{F}} \cdot \mathbf{x}^{T} \right] - \mathbf{x}^{T} \mathbf{A} \mathbf{x} \\ &\left( diag(\mathbf{x}) \mu_{\mathcal{F}} = diag(\mu_{\mathcal{F}}) \mathbf{x} \right) \\ &= Tr \left[ \mathbf{A} \cdot diag(\mu_{\mathcal{F}}) \cdot \mathbf{x} \mathbf{x}^{T} \right] - \mathbf{x}^{T} \mathbf{A} \mathbf{x} \\ &= \mathbf{x}^{T} \left[ \mathbf{A} \cdot diag(\mathbf{x}) \right] \mu_{\mathcal{F}} - \mathbf{x}^{T} \mathbf{A} \mathbf{x} \end{split}$$

where (\*) comes from the fact that  $\mathbf{x}^T[\mathbf{A} \odot \mathbf{B}]\mathbf{x} = Tr[diag(\mathbf{x}) \cdot \mathbf{A} \cdot diag(\mathbf{x}) \cdot \mathbf{B}^T]$ . Similarly we have:

$$\mathbf{x}^{T} \underbrace{\left[ \mathbf{A} \odot \left( (\mu_{\mathcal{F}} - \mathbf{1}) \times \mathbf{1}^{T} \right) \right]}_{\mathbf{I}} \mathbf{x} = Tr \left[ diag(\mathbf{x}) \cdot \mathbf{A} \cdot diag(\mathbf{x}) \times \mathbf{1} \times (\mu_{\mathcal{F}} - \mathbf{1})^{T} \right]$$

$$= Tr \left[ diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} \cdot (\mu_{\mathcal{F}} - \mathbf{1})^{T} \right]$$

$$= Tr \left[ diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} \cdot \mu_{\mathcal{F}} \right] - \mathbf{x}^{T} \mathbf{A} \mathbf{x}$$

$$= \mathbf{x}^{T} \left[ \mathbf{A} \cdot diag(\mathbf{x}) \right] \mu_{\mathcal{F}} - \mathbf{x}^{T} \mathbf{A} \mathbf{x}$$

and:

$$\mathbf{x}^{T} \underbrace{\left[ (\alpha_{2} + \alpha_{3}) \mathbf{A} - (\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) - (\alpha_{2} + \alpha_{3}) \mathbf{A} \odot \left( \mu_{\mathcal{F}} \mu_{\mathcal{F}}^{T} \right) \right]}_{\mathbf{3}} \mathbf{x} = \underbrace{\left[ (\alpha_{2} + \alpha_{3}) \mathbf{x}^{T} \mathbf{A} \mathbf{x} - (\alpha_{2} + \alpha_{3}) \mathbf{x}^{T} \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) \mathbf{x} - \underbrace{\left( \alpha_{2} + \alpha_{3} \right) \mu_{\mathcal{F}}^{T} \left( \mathbf{A} \odot \left( \mathbf{x} \mathbf{x}^{T} \right) \right) \mu_{\mathcal{F}}}_{(\diamond)},$$

where (\$\dappa\$) comes from the following:

$$\begin{aligned} &-(\alpha_{2}+\alpha_{3})\mathbf{x}^{T}\left[\mathbf{A}\odot\left(\mu_{\mathcal{F}}\mu_{\mathcal{F}}^{T}\right)\right]\mathbf{x} \\ &=-(\alpha_{2}+\alpha_{3})Tr\left[diag(\mathbf{x})\cdot\mathbf{A}\cdot diag(\mathbf{x})\cdot\left(\mu_{\mathcal{F}}\mu_{\mathcal{F}}^{T}\right)\right] \\ &=-(\alpha_{2}+\alpha_{3})Tr\left[diag(\mathbf{x})\cdot\mathbf{A}\cdot diag(\mu_{\mathcal{F}})\cdot\mathbf{x}\cdot\mu_{\mathcal{F}}^{T}\right] \\ &=-(\alpha_{2}+\alpha_{3})Tr\left[\mathbf{A}\cdot diag(\mu_{\mathcal{F}})\cdot\mathbf{x}\cdot\mu_{\mathcal{F}}^{T}\cdot diag(\mathbf{x})\right] \\ &=-(\alpha_{2}+\alpha_{3})Tr\left[\mathbf{A}\cdot diag(\mu_{\mathcal{F}})\cdot\mathbf{x}\mathbf{x}^{T}\cdot diag(\mu_{\mathcal{F}})\right] \\ &=-(\alpha_{2}+\alpha_{3})Tr\left[diag(\mu_{\mathcal{F}})\cdot\mathbf{x}\mathbf{x}^{T}\cdot diag(\mu_{\mathcal{F}})\cdot\mathbf{A}\right] \\ &=-(\alpha_{2}+\alpha_{3})\mu_{\mathcal{F}}^{T}\left(\mathbf{A}\odot\left(\mathbf{x}\mathbf{x}^{T}\right)\right)\mu_{\mathcal{F}} \end{aligned}$$

Putting the above derivation together we obtain:

$$\mathbf{x}^{T}\mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}})\mathbf{x} + 2\mathbf{x}^{T}\mathbf{b}(\mu_{\mathcal{F}}) = \\ \mu_{\mathcal{F}}^{T} \left[ -(\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \left( \mathbf{x} \mathbf{x}^{T} \right) \right) \right] \mu_{\mathcal{F}} + \mu_{\mathcal{F}}^{T} \left[ (\alpha_{3} + 2\alpha_{2}) diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} - \alpha_{1} \mathbf{x} \right] - \\ \left[ \alpha_{1} \mathbf{1}^{T} \mathbf{x} - \mathbf{x}^{T} \left( (\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) + \alpha_{2} \mathbf{A} \right) \mathbf{x} \right]$$

So the function  $f(\mu_{\mathcal{F}})$  becomes:

$$\mathbf{x}^{T}\mathbf{Q}(\mu_{\mathcal{F}}, \hat{\mathbf{\Sigma}})\mathbf{x} + 2\mathbf{x}^{T}\mathbf{b}(\mu_{\mathcal{F}}) - t - \mu_{\mathcal{F}}^{T}\mathbf{K}\mu_{\mathcal{F}} =$$

$$\mu_{\mathcal{F}}^{T} \left[ -(\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \left( \mathbf{x} \mathbf{x}^{T} \right) \right) - \mathbf{K} \right] \mu_{\mathcal{F}} + \mu_{\mathcal{F}}^{T} \left[ (\alpha_{3} + 2\alpha_{2}) diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} - \alpha_{1} \mathbf{x} \right] +$$

$$\left[ \alpha_{1} \mathbf{1}^{T} \mathbf{x} - \mathbf{x}^{T} \left( (\alpha_{2} + \alpha_{3}) \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) + \alpha_{2} \mathbf{A} \right) \mathbf{x} - t \right],$$

which is a quadratic function in  $\mu_{\mathcal{F}}$ . Define **R**, **r** and z as the following:

$$\mathbf{R} = -(\alpha_2 + \alpha_3) \left( \mathbf{A} \odot \left( \mathbf{x} \mathbf{x}^T \right) \right) - \mathbf{K}$$

$$\mathbf{r} = (\alpha_3 + 2\alpha_2) diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x} - \alpha_1 \mathbf{x}$$

$$z = \alpha_1 \mathbf{1}^T \mathbf{x} - \mathbf{x}^T \left( (\alpha_2 + \alpha_3) \left( \mathbf{A} \odot \hat{\mathbf{\Sigma}} \right) + \alpha_2 \mathbf{A} \right) \mathbf{x} - t,$$

which results in a compact form of  $f(\mu_{\mathcal{F}})$ :

$$f(\mu_{\mathcal{F}}) = \mu_{\mathcal{F}}^T \mathbf{R} \mu_{\mathcal{F}} + \mu_{\mathcal{F}}^T \mathbf{r} + z$$

# 6. Proof of Eq.(10) in the paper

The relation (r1) is direct. To see why (r2) holds, note that the *i*-th element of  $diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x}$  is:

$$\left[diag(\mathbf{x}) \cdot \mathbf{A} \cdot \mathbf{x}\right]_i = \mathbf{x}_i \sum_{j=1}^N \mathbf{A}_{ij} \mathbf{x}_j,$$

which is equal to the *i*-th element of  $diag(\mathbf{AX})$ :

$$\left[diag(\mathbf{A}\mathbf{X})\right]_i = \left[diag(\mathbf{A})\mathbf{x}\mathbf{x}^T\right]_i = \mathbf{x}_i \sum_{j=1}^N \mathbf{A}_{ij}\mathbf{x}_j.$$

The relation (r3) holds because:

$$\mathbf{x}^{T} \Big( (\alpha_{2} + \alpha_{3})(\mathbf{A} \odot \hat{\mathbf{\Sigma}}) + \alpha_{2} \mathbf{A} \Big) \mathbf{x} = (\alpha_{2} + \alpha_{3}) \mathbf{x}^{T} (\mathbf{A} \odot \hat{\mathbf{\Sigma}}) \mathbf{x} + \alpha_{2} \mathbf{x}^{T} \mathbf{A} \mathbf{x}$$

$$= (\alpha_{2} + \alpha_{3}) Tr \Big[ (\mathbf{A} \odot \hat{\mathbf{\Sigma}}) \mathbf{x} \mathbf{x}^{T} \Big] + \alpha_{2} Tr \Big[ \mathbf{A} \mathbf{x} \mathbf{x}^{T} \Big]$$

$$= (\alpha_{2} + \alpha_{3}) Tr \Big[ (\mathbf{A} \odot \hat{\mathbf{\Sigma}}) \mathbf{X} \Big] + \alpha_{2} Tr \Big[ \mathbf{A} \mathbf{X} \Big]$$

# 7. Statistics of the networks used in experiments

	r	density	#edges	clustering coeff.
BA-1	2.7167	0.0461	375	0.1340
BA-2	2.2789	0.0610	496	0.1504
BA-3	2.0374	0.0757	615	0.1646
SW-1		0.0787	640	0.3664
SW-2		0.1102	896	0.3875
SW-3		0.1575	1280	0.4059
Facebook		0.0106	1325	0.3930

Table 1. Statistics of networks used in our experiments. r is the exponent of the power-law degree distribution.