

Contamination-Aware Experiments on Networks

Mine Su Erturk
Graduate School of Business
Stanford University
mserturk@stanford.edu

Eray Turkel
Graduate School of Business
Stanford University
eturkel@stanford.edu

1 Introduction

We study the problem of a decision maker (DM) conducting experiments over time, on a network environment. An ‘experiment’ in the context of our model is treating a node in the network. Each node responds to the treatment differently, depending on their (observable) characteristics. We also assume that treating a node creates spillovers on neighboring nodes. The DM’s goal is learning the optimal treatment allocation over the network by sequentially conducting local experiments on different nodes.

We assume that there are multiple analysts conducting experiments on the network, as is the case in many online platforms. The DM’s experiment creates negative externalities on other ongoing experiments, or the deployed treatments from previously run experiments. We model this as a cost for contaminating the treatment regime of other experimenters. Therefore, the DM has to strike a balance between learning the treatment response function efficiently, and staying within acceptable levels of contamination. We analyze the DM’s problem as a sequential decision making problem, which puts our framework closer to the literature on online learning and bandits, instead of using a randomization based inference setting, as is common in the network experimentation literature.

2 Model

There are N units connected on an undirected, unweighted graph. Denote each unit on the graph by $i \in \{1, \dots, N\}$. The units have observable characteristics, $\{x_1, \dots, x_N\}$ with $x_i \in \mathbb{R}^k$. The DM’s objective depends on the outcomes of the units, denoted $\{Y_1, \dots, Y_N\}$ with $Y \in \mathbb{R}$, which respond to a treatment regime that will be determined. The decision maker can run experiments sequentially, at time periods $t \in \{1, \dots, T\}$. Let w_{it} denote the treatment assignment for node i during the period of experimentation t , with $w_{it} \in \{0, 1\}$, and denote the vector of treatment assignments at time t by $w_t = [w_{1t}, \dots, w_{Nt}] \in \{0, 1\}^N$. We assume that the experimentation is local, in the sense that every period, only one node can be treated. Thus for all time periods t , $\sum_{i \leq N} w_{it} = 1$.

Each node’s outcome depends on their characteristics, the treatment they receive and the exposure to the treatment through their neighbors. Define the exposure for node i at period t as e_{it} , which is assumed to be:

$$e_{it} = \frac{\sum_{j \leq N} \mathbf{1}(j \in nhbd(i)) w_{jt}}{\sum_{j \leq N} \mathbf{1}(j \in nhbd(i))}$$

Y is assumed to be a linear function of the unit characteristics and the exposure to the treatment:

$$Y_{it}(w_{it}, e_{it}) = \beta^T x_i + \Gamma^T x_i e_{it} + \epsilon_i,$$

where $\beta, \Gamma \in \mathbb{R}^k$, and for all nodes i , $\epsilon_i \sim N(0, \sigma_\epsilon^2)$.

The DM’s goal is to choose a treatment assignment policy over the network based on the value of Γ . Denote the final treatment assignment for node i by w_i , dropping the time subscript. We denote a policy by $\pi(\Gamma)$, where $\pi(\Gamma) = [w_1, \dots, w_N] \in \{0, 1\}^N$. We assume that there is an upper limit on the number of nodes that can be treated in the final deployment, which we call D , i.e., $\sum_{i \leq N} w_i \leq D$.

At every period t , the DM's experimentation creates a cost, because the experiment interacts with a previous experimenter's treatment allocation, contaminating their results. We assume that the cost is proportional to the number of units exposed to the treatment, and once a unit is exposed, future exposures of that unit do not increase the cost any further. Define the total cost of experimentation at period t by:

$$c_t = c \sum_{i \leq N} \mathbf{1}(e_{it} > 0 \text{ and } e_{is} = 0, \forall s < t).$$

Throughout, we will let $c = 1$ for simplicity. Hence, the total cost incurred through experimentation by the DM can be written as $C = \sum_{t \leq T} c_t$. We assume that the DM has to stay below a 'contamination budget' of at most \bar{C} during experimentation. Then, as a function of the contamination budget \bar{C} , we define the value of a policy as follows:

$$V(\pi, \bar{C}) = \sum_{i \leq N} \Gamma^T x_i e_i(\pi),$$

where $e_i(\pi)$ denotes the exposure function determined at the last period for final deployment, induced by the treatment allocation under policy π , decided after an experimentation period with the contamination budget \bar{C} . The DM does not observe the true vector Γ but has an imperfect estimate of it learned through sequential experimentation, which we will call $\hat{\Gamma}$. The DM maximizes the empirical value function using its estimate at the end of the experimentation period:

$$\text{Max}_{\hat{\pi}} : \sum_{i \leq N} \hat{\Gamma}^T x_i e_i(\hat{\pi}) \quad \text{s.t. } \mathbf{1}^T \hat{\pi} \leq D. \quad (1)$$

Let $X_{K \times N}$ denote the matrix of covariates for all the nodes on the graph, and let X_s denote the analogous matrix only for nodes that were exposed to the treatment during the experimentation period. Let $\hat{\Gamma}_{1 \times K}$ be the vector of estimated coefficients.

Define the vector of exposure values at time t as $\bar{e}_t = [\mathbf{1}(e_{1t} > 0), \dots, \mathbf{1}(e_{Nt} > 0)]_{1 \times N}$ and let c_t denote the cost vector at time t . We can recursively write: $c_0 = \mathbf{0}_{1 \times N}$, $c_1 = \bar{e}_1$, and generally, $c_n = (c_{n-1} + \bar{e}_n) - (\text{Diag}(c_{n-1})\bar{e}_n)$. $\text{Diag}(c_n)$ denotes the diagonal matrix with the entries of c_n on its diagonal entries. This recursive formulation captures the notion that the cost of second and future exposures to the experiment are zero.

We can also define the final exposure vector e in terms of the adjacency matrix of the network, $A_{N \times N}$ (with self-edges), and the final treatment allocation vector $\pi \in \{0, 1\}^N$. Representing e as a vector, we have: $\text{Diag}(A\mathbf{1}_N)^{-1}A\pi = e$, where $\mathbf{1}_N$ is an $(N \times 1)$ vector of 1's. Therefore, we can rewrite the maximization problem (1) as a convenient linear program, where we relax π_i and allow it to take values in $[0, 1]$, interpreting the resulting policy as a probabilistic treatment assignment rule:

$$\text{Max}_{\hat{\pi}} : \left(\hat{\Gamma}^T X \text{Diag}(A\mathbf{1}_N)^{-1} A \right) \hat{\pi}, \text{ subject to: } \mathbf{1}^T \hat{\pi} \leq D, \forall i, \hat{\pi}_i \in [0, 1] \quad (2)$$

3 Regret Analysis

Let $\hat{\pi}$ denote the optimal policy maximizing the optimization problem (2) given that the DM has estimated the coefficient vector as $\hat{\Gamma}$, which achieves the value $V(\hat{\pi}, \bar{C})$. Let $V(\pi^*)$ be the optimal value achievable by the oracle who solves the optimization problem (2) with the true coefficient vector Γ . Thus, we can define the regret as follows:

$$R = V(\pi^*) - V(\hat{\pi}, \bar{C}).$$

To analyze the regret of the empirical value maximizing policy $\hat{\pi}$, let us introduce some notation. Let $\sigma_j(A)$ denote the j^{th} eigenvalue of the matrix A , with σ_1 being the largest. Suppose that the covariate matrices X, X_s are standardized, so the matrix XX^T is the symmetric correlation matrix.

Theorem 1 Fix $\tau \geq 0$. Then, regret is bounded above by K with the following probability:

$$\mathbb{P}(R \leq K) \geq 1 - \exp\left(-\frac{1}{2}(\tau - 1)^2\right), \quad (3)$$

where

$$K = \frac{\delta \max\{\sqrt{n + D^2}, \|\Gamma\| \sigma_1(XX^T)\}}{\sqrt{n + 1}} \cdot \frac{\sqrt{n + D^2}}{\sqrt{n + 1} - \delta}, \quad (4)$$

and $\delta = \tau \sigma_\epsilon \sum_{j=1}^k \sigma_j((X_s X_s^T)^{-1})$.

The proof is based on a perturbation analysis of the LP in (2) with respect to the objective function coefficients.

3.1 Proposed Algorithm and Simulations

We propose an algorithm based on information-directed sampling and the literature on linear bandits, which adaptively chooses experimental subjects that will maximize information gain. Since our setting has an additional component which is the contamination cost due to externalities on other experiments, we use a knapsack-based heuristic $c - \mu$ rule.

In simple terms, our algorithm chooses the most valuable nodes to experiment on greedily at every period, taking into account the cost of contaminating other nodes. We measure the information gain from experimenting on a given node by the change in the trace of the matrix $(X_s X_s^T)^{-1}$ after adding all the unexposed neighbors of the target node to the set of exposed nodes. We calculate the cost of experimenting on a node by calculating the number of new exposures resulting from treating the target node. We get the value of experimenting on a node by taking the ratio of the information gain to the cost, and determine the path of experimentation by sequentially choosing the nodes with the highest value until we reach our contamination budget, \bar{C} . Using this algorithm to determine which nodes to experiment on every period, the DM updates their estimate of the coefficient vector $\hat{\Gamma}$ after every experiment, solving the maximization problem in (2) in the last period to decide on a final treatment regime that will be deployed.

We run simulation experiments to compare the performance of our proposed algorithm against the policy chosen by an oracle who knows the true response function. For each simulation, we draw an Erdos-Renyi random graph with 500 nodes and $p = 0.1$, create our random covariate matrix $X_{5 \times 500}$ using draws from independent standard normal random variables, and generate our true coefficient vector Γ drawing from a 5-dimensional normal distribution with mean zero and covariance matrix I_5 . Our results show that the proposed algorithm performs reasonably well on medium-sized graphs, and the percentage gap to the oracle value falls below 20% when \bar{C} approaches half of the graph size.

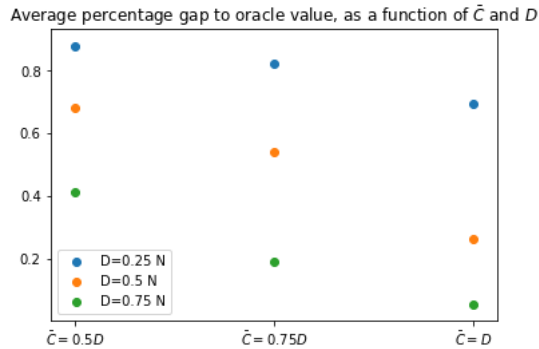


Figure 1: Percentage gap to oracle value $(\frac{V(\pi^*) - V(\hat{\pi}, \bar{C})}{V(\pi^*)})$, under different deployment and contamination budgets. $N=500$, and 100 simulations for every parameter combination.

Suppose $\|c(\Gamma) - c(\hat{\Gamma})\| \leq \delta$, then

$$R \leq K = \delta \frac{\max\{\sqrt{n+D^2}, \|\Gamma^T X\|\}}{\sqrt{n+1}} \frac{\sqrt{n+D^2}}{\sqrt{n+1}-\delta}$$

where $\|\cdot\|$ refers to the 2-norm.

$$\|(\hat{\Gamma} - \Gamma)^T X\| = \delta$$

$$R \leq \|\Delta c\| \frac{\|d\|}{\text{dist}(d, \text{Pri}\emptyset)} \frac{\|b\|}{\text{dist}(d, \text{Dual}\emptyset) - \|\Delta c\|}.$$

First, we note that the norm of the LP instance d is given by the expression:

$$\|d\| = \max\{\|A\|, \|b\|, \|c\|\} = \max\{\sqrt{n+1}, \sqrt{n+D^2}, \|\Gamma^T X\|\}.$$

Then, we refer to Corollary 2.9 of [?] for an equivalent characterization of the distance to the primal infeasibility region in terms of the condition number of the coefficient matrix A of the LP.

Lemma 1 ([?] Corollary 2.9) *For any given A , the distance to primal infeasibility, $\text{dist}(d_P, \text{Pri}\emptyset)$, is equal to*

$$\rho(d_P) = \sup\{\epsilon : \|y\| \leq \epsilon \Rightarrow y \in \{Ax : x \geq 0, \|x\| \leq 1\}\},$$

where $d_P = (A, b)$.

Hence, we compute the condition number of A as the largest singular value of the matrix A . That is, we have

$$\rho(d_P) = \|A\|_2 = \sigma_{\max}(A) = \sqrt{n+1}.$$

Similarly, we can show that the distance to the dual infeasible region is given by the largest singular value of the transpose of the coefficient matrix A^T , i.e., $\rho(d_D) = \|A^T\|_2 = \sigma_{\max}(A^T) = \sigma_{\max}(A)$.

Finally, we use the observation that the perturbation to the coefficient vector is smaller than δ . Therefore, we can conclude that regret is upper bounded by

$$K = \delta \frac{\max\{\sqrt{n+1}, \sqrt{n+D^2}, \|\Gamma^T X\|\}}{\sqrt{n+1}} \frac{\sqrt{n+D^2}}{\sqrt{n+1}-\delta}$$