

# On the Friendship Paradox and Inversity: A Network Property with Applications to Privacy-sensitive Network Interventions

Vineet Kumar<sup>a,2</sup>, David Krackhardt<sup>b</sup>, and Scott Feld<sup>c</sup>

This manuscript was compiled on March 1, 2024

We provide the mathematical and empirical foundations of the friendship paradox (FP) in networks, often stated as: “Your friends have more friends than you.” We prove a set of network properties on friends of friends, and characterize the concepts of ego-based and alter-based means. We propose a network property called inversity that quantifies the imbalance in degrees across edges, and prove that the sign of inversity determines the ordering between ego-based or alter-based means for any network, with implications for interventions. Network intervention problems benefit from using highly-connected nodes, e.g., immunization. We characterize two intervention strategies based on the friendship paradox to obtain such nodes, with the novel alter-based and ego-based strategy. Both strategies provide provably guaranteed improvements for any network structure with variation in node degrees. We demonstrate that the proposed strategies obtain several-fold improvement (100-fold in some networks) in node degree relative to a random benchmark, for both generated and real networks. We evaluate how inversity informs which strategy works better based on network structure and show how network aggregation can alter inversity. We illustrate how these strategies can be used to control contagion of an epidemic spreading across a set of village networks, finding that these strategies require far fewer nodes to be immunized (less than 50%, relative to random). The interventions do not require knowledge of network structure, are privacy-sensitive, are flexible for time-sensitive action, and only require selected nodes to nominate network neighbors for intervention.

Network Intervention | Friendship Paradox | Inversity | Contagion

We examine the underlying mathematical and empirical foundations of the friendship paradox, and define a new network property called inversity, which has implications for network interventions. The friendship paradox has often been simply referred to by the maxim, “Your friends have more friends than you do.” However, we show that there are two different ways of understanding this statement, which lead to different network properties that we term as the ego-based and alter-based mean number of friends of friends. We find that both means are higher than the average degree across nodes in the network. We show that the properties are not just conceptually distinct, but they are also empirically different across a wide class of generated and real-world networks. We identify a novel network property, inversity, that connects the two means, and for any network, the sign of inversity determines whether the ego-based mean or alter-based mean is higher. Since these mathematical properties apply to any network, not just those based on friendship, we use neighbor in place of friend henceforth.

The above results have direct implications for interventions by finding highly connected nodes in a network using privacy-sensitive methods based on the friendship paradox. The two means lead to corresponding ego-based and alter-based strategies for obtaining highly connected nodes, of which the global strategy has not been used in network interventions before. The inversity of the network indicates which strategy (ego-based or alter-based) is better for obtaining highly connected nodes. We show in empirical real-world networks, and in generated networks, that the strategy used can make a meaningful difference. We illustrate, using a simplified application with real networks, how using inversity to choose the ego-based or alter-based strategy to identify inoculation candidates considerably reduces the epidemic threshold and peak infection relative to random selection.

## Friendship Paradox

The friendship paradox, which our interventions are based on, is colloquially stated as the idea that people’s neighbors

### Significance Statement

Networks across many different settings — including social, economic and natural — are powerful tools for interventions due to the cascading impact of one individual node on others. All networks with degree variation exhibit the friendship paradox phenomenon. We demonstrate its multifaceted nature, and provide its foundations mathematically and empirically. We identify a new network property — inversity — and propose novel network intervention strategies based on the friendship paradox. Inversity uniquely determines the best performing strategy. These strategies provide a privacy-sensitive approach to obtaining highly connected individuals without knowing the network, and are guaranteed to obtain a greater than average degree for almost any network. Finally, we characterize the value of these strategies theoretically and with real-world networks.

Author affiliations: <sup>a</sup>Yale School of Management, Yale University, 165 Whitney Avenue, New Haven CT 06511. ORCID: 0000-0001-8784-6858; <sup>b</sup>Heinz College, Carnegie Mellon University, 4800 Forbes Avenue, Pittsburgh, PA 15213. ORCID: 0000-0001-9487-9973; <sup>c</sup>Department of Sociology, College of Liberal Arts, Purdue University, 700 West State Street, West Lafayette, IN 47907-2059. ORCID: 0000-0003-4820-1365

VK, DK and SF designed and performed research. VK conceptualized analytical tools and theoretical results, collected and analyzed data, performed empirical analyses, and wrote the paper. DK and SF contributed to the writing and discussion.

The authors have no competing interests as it relates to this paper.

<sup>2</sup>To whom correspondence should be addressed. E-mail: vineet.kumar@yale.edu

are more popular than them (1, 2).<sup>\*</sup> The intuition for why the friendship paradox helps obtain well-connected nodes is this: since highly-connected nodes (hubs) are connected to many other nodes (by definition), obtaining a random friend (neighbor) of a random node is likely to result in hubs with greater likelihood, compared to the case of randomly selecting nodes.

We establish that the friendship paradox is not just one statement, but rather a set of distinct claims (All theorems and proofs are in Supplement §S.B). First, we find an impossibility, i.e., the individual-level friendship paradox cannot hold for all individuals in any given network (Theorem S1). In practice for real networks, it can hold for a large proportion of nodes in the network (Figs. S3 and S4 in Supplement §S.D). Second, we demonstrate that in contrast to the impossibility of the individual friendship paradox, the network level friendship paradox always holds for any non-regular network. We show how the *average* number of neighbors of neighbors across the network can be characterized in two different ways, using the ego-based and alter-based means, as defined below. Both ego-based and alter-based means are greater than the mean degree of the network, and they are related through a novel network characteristic we term **inversity**.

**Ego-based and Alter-based Means.** We formally characterize the two distinct but related network properties deriving from the friendship paradox relating to the “average number of neighbors of neighbors.” We denote an undirected network (see Table S1 for full notation) as a graph  $\mathcal{G} = (V, E)$  with  $V$  the set of (non-isolate) vertices or nodes, and  $E$  the set of edges ( $e_{ij} \in \{0, 1\}$ , denoting absence or presence of a connection between  $i$  and  $j$ ).  $D_i$  refers to the degree of node  $i$ , and  $\mathcal{N}(i)$  refers to the set of  $i$ ’s neighbors. We specify the ego-based mean as the average number of neighbors of neighbors across the nodes in the network:

$$\mu_E = \frac{1}{N} \sum_{i \in V} \left[ \frac{1}{D_i} \sum_{j \in \mathcal{N}(i)} D_j \right] \quad [1]$$

The alter-based mean is defined as the ratio of the total number of neighbors of neighbors to the total number of neighbors in the network, consistent with (1):

$$\mu_A = \frac{\sum_{i \in V} \left[ \sum_{j \in \mathcal{N}(i)} D_j \right]}{\sum_{i \in V} D_i} \quad [2]$$

The above means arise from conceptualizing the average degree across neighbors differently.<sup>†</sup> Both means above are consistent with the notion of “average number of neighbors of neighbors,” although they are distinct network properties (see Fig. S1 for an example and detailed explanation). The alter-based mean was theoretically investigated earlier and found to be greater than (or equal to) the average degree and is independent of the network topology, given node degrees (Theorem S2). Equality holds only when the network is

regular, with all nodes the same degree within and across components.<sup>‡</sup>

The ego-based mean is shown here to be greater than (or equal to) the mean degree (Theorem S3).<sup>§</sup> However, the contrast is that the ego-based mean has distinct properties that depend on network topology (i.e., who is connected to whom). Equality for the ego-based mean only holds when each component is regular, with no degree variation within components.

We identify network structures that result in a greater divergence between the ego-based and alter-based mean, and between these means and the average degree, including whether one of the means is always greater than the other, and whether they always exhibit correlated variation away from the mean degree. In Figure S2, we find that both the ego-based and alter-based means can be much greater than the mean degree, and that between these two means, either one of them can be greater than the other. In some network structures, both can be relatively high compared to the mean degree. We also see that the alter-based mean is invariant to rewiring the network while keeping the degree distribution the same, whereas the ego-based mean is impacted by rewiring (Theorem S6).

## Inversity: Connecting Ego-based and Alter-based Means

We identify and define a network property, *inversity*, that determines when the ego-based mean is greater than the alter-based mean. This property captures all local network information related to the ego-based mean. We prove that the sign of *inversity* determines whether the ego-based mean or the alter-based mean is higher for any given network.

*Inversity* is a correlation-based metric that relates the alter-based and ego-based means for any network, and is obtained as follows. First, define the following edge-based distributions to examine the relationship between the means. The *origin* degree (**O**),  $D^O(e)$ , *destination* degree (**D**),  $D^D(e)$ , and *inverse destination* degree (**ID**) distribution,  $D^{ID}(e)$ , are defined across directed edges  $e \in \hat{E}$  as:  $D^O(e_{jk}) = D_j$ ,  $D^D(e_{jk}) = D_k$ ,  $D^{ID}(e_{jk}) = \frac{1}{D_k}$ . We define the *inversity* across the edge distribution as the Pearson correlation across the origin and inverse degree distributions.

$$\rho = \text{Cor}(D^O, D^{ID}) \quad [3]$$

We show that the ego-based and alter-based means are connected by *inversity* and the degree distribution ( $\kappa_m = \sum_{i \in V} D_i^m$ ) as follows:

$$\mu_E = \mu_A + \rho \Psi(\kappa_{-1}, \kappa_1, \kappa_2, \kappa_3) \quad [4]$$

where  $\Psi$  is a positive function of moments of the degree distribution (see Theorem S4).

*Inversity* is a measure of *imbalance* between the degrees of the nodes connected by an edge. This imbalance for edge  $(i, j)$  is characterized by the ratio of degrees  $\left(\frac{D_i}{D_j}\right)$  and  $\left(\frac{D_j}{D_i}\right)$ .

<sup>\*</sup> The phenomenon has also been generalized to the idea that individual attributes and degree are correlated (3), e.g., an individual’s co-authors are more likely to be cited (4), or that neighbors are more important (5), or more socially active (6). Of specific relevance to this research is the mathematical generalization to distributions examined in (7).

<sup>†</sup> We note that a node is also its neighbor’s neighbor under both ego-based and alter-based properties.

<sup>‡</sup> There are a number of phenomena that share a similar underlying structure, e.g., disproportionately many people grow up in large families, or students experience a class size that is larger than the average size of classes. The underlying selection process here is commonly termed probability proportional to size (PPS) (8, 9). In the case of the friendship paradox, we find that the alter-based mean has a direct mathematical connection to PPS (neighbors have disproportionately more neighbors). However, the ego-based mean operates through a different mechanism.

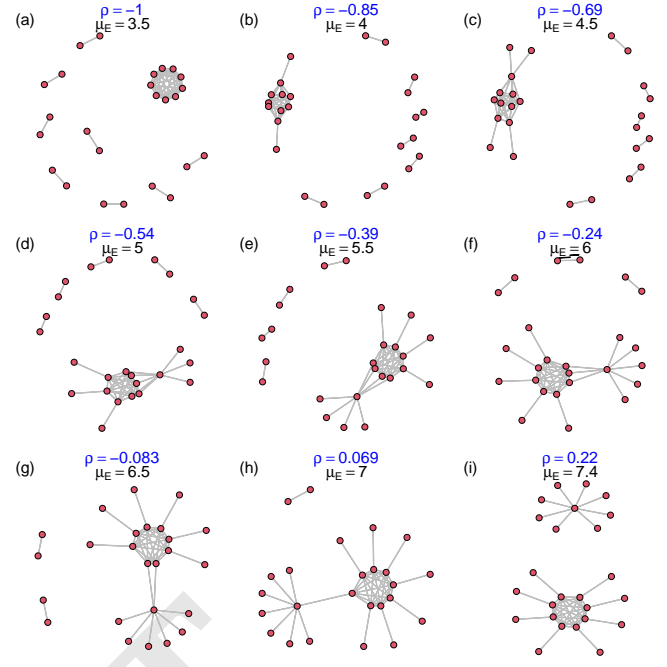
<sup>§</sup> We term these means ego-based or alter-based since **INCLUDE REASONING HERE**.

That is, more imbalanced edges tend to have nodes with high degree on one end and low degree on the other end. This imbalance plays a significant role in both inversivity and the ego-based mean. Consider how the ego-based mean is obtained: for each node, we take the mean of degrees across the node's neighbors. The expected degree of neighbors of nodes (or ego-based mean) is then  $\frac{1}{N} \sum_i \sum_{j \in \mathcal{N}(i)} \left( \frac{D_i}{D_j} \right)$ . This imbalance measure based on ratios explains why inversivity is highly sensitive to stars, whereas assortativity is more sensitive to the presence of clique-like structures, which we examine below.

When connections (edges) are mostly between nodes of similar degree, then inversivity  $\rho$  is likely to be more negative. In such a case, the global mean is greater than the ego-based mean. In contrast, when connections are more likely to be between nodes of dissimilar degree, then inversivity  $\rho$  is positive, and the ego-based strategy is likely to obtain higher degree nodes. Therefore, if inversivity is known, we don't need the entire degree distribution to obtain the ego-based mean. Rather, *four* moments of the degree distribution are sufficient for that purpose. Inversivity captures the local information on imbalances in the degrees of nodes across edges, whereas the moments of the degree distribution represent global information about the network. Inversivity  $\rho$  has a critical role in determining whether the ego-based or alter-based mean is greater for a network; specifically,  $\rho < 0$  indicates the alter-based mean is higher than the ego-based mean, whereas  $\rho > 0$  indicates the opposite. Thus, knowing inversivity can help us determine which strategy to use. Even computing inversivity is information-light, requiring only the  $2k$  distribution, which represents the degrees of nodes at the termini of each edge, rather than the entire network (10).

**How Inversivity Depends on Network Topology.** Inversivity is the only term that depends on the network topology, or structure of connections (who is connected to whom) in the relationship between ego-based and alter-based means. We examine how inversivity changes as we change the structure, while *simultaneously* preserving the degree distribution. In Figure 1, we start with a network with the minimum inversivity  $\rho = -1$ , and then use the rewiring theorem (Theorem S6) to examine how inversivity increases while the degree distribution, and consequently, the mean degree, variance of degree, minimum and maximum degree, as well as the alter-based mean, all remain fixed and identical across each of the networks (a)–(i). Specifically, each network has the following properties in common:  $N = 25$  nodes,  $\mu_D = 3.5$ ,  $\sigma_D^2 = 3.4$ ,  $D_{\min} = 1$  and  $D_{\max} = 8$ . Observe that the alter-based mean,  $\mu_A = 6.7 > \mu_D = 3.5$ , is also the same for each of the networks.

The motivation for keeping degree distribution fixed across the networks is that the alter-based mean does not change. Specifically, noting from panel (a), the degree distribution includes 16 nodes with degree 1 (the dyads), and 9 nodes with degree 8 (clique or complete subgraph). We note the rewiring patterns, beginning with (a), which displays a network with a fully connected complete component with 9 nodes, and 8 dyads. This network has the lowest possible inversivity of  $\rho = -1$ , consistent with the idea that no edge connects nodes of different degrees, which is essential for inversivity to be greater than its minimum value.



**Fig. 1. How Inversivity Changes with Rewiring.** The network is changed by rewiring, starting with the top left ( $\rho = -1$ ), to increase inversivity  $\rho$  as we traverse from panel (a) to (i). Observe that the number of nodes,  $N = 25$ , the number of edges  $|E| = 44$ , as well as the degree distribution for each of the networks in panels (a)–(i) is identical, with 16 nodes with degree 1 and 9 nodes with degree 8. We note that the ego-based mean is  $\mu_E = \mu_D = 3.5$  for network (a), but increases along with inversivity in panels (b)–(i), reaching  $\mu_E = 7.4$  for network (i). The alter-based mean,  $\mu_A = 6.7$ , remains constant across all the networks.

We use the rewiring theorem (Theorem S6) to increase inversivity; this approach *connects* low degree nodes to high degree nodes, while *removing* connections between nodes of intermediate degree. The rewiring increases the variation in the degrees of the nodes connected by an edge, as the network transforms from (a) to (b) and in each further step. We observe that the nodes in each dyad break up their edge (which connects nodes of identical degrees), and connect to nodes in the large component, which contains high degree nodes. We next observe a star-like structure form, beginning with panel (d). Finally, as the star-like structure expands, in panel (h) and (i), we find that inversivity has changed sign to become positive.

A few general observations are worth noting. First, we see that inversivity and the ego-based mean are highly sensitive to network topology, whereas the alter-based mean is impacted only by the degree distribution, specifically its mean and variance (as shown earlier). Second, we note that networks which display little or no variation among the node degrees connected by an edge have negative inversivity, like in network (a). Third, we find a wide range of possible networks with different levels of inversivity and ego-based mean for a fixed alter-based mean. These inversivity values range from negative to positive. Finally, the degree distribution can constrain the range of inversivity. We next examine the implications of these findings for network interventions.

There are many reasons why networks may take forms with high or low inversivity. For the present purpose, we provide some intuition of relevant processes. We can expect hub-based



(or star) networks to have high inversity, where most nodes have a few ties that largely go to the relatively few hubs with large numbers of ties (e.g. Twitter celebrity based networks). In contrast, we expect that networks of clusters of various sizes will have low inversity. In such clusters, members of the large clusters are tied to one another, obtaining high degree nodes within the cluster. Similarly, nodes in small clusters tend to have few ties, mostly with one another, reflecting friendship networks based upon group membership. The various causes of network structures having different levels of inversity may be the subject of extensive theoretical and empirical study in the future. We provide a discussion of this in Supplement §S.F. For details on inversity values in real-world networks and related findings, see Supplement §S.C.

**Inversity and Assortativity.** Inversity  $\rho = \text{cor}(D_i, \frac{1}{D_j})$  is related to, but distinct from, the commonly used measure of assortativity, defined as  $\rho_a = \text{cor}(D_i, D_j)$ . We observe that the formulation of both network properties appears similar, and both have values ranging from  $-1$  to  $+1$ . We might therefore expect inversity to be a reverse of assortativity, or more specifically  $\rho_a \approx -\rho$ . For details of each of these arguments and examples, please see Supplement §S.G.

*Can assortativity serve as a proxy for inversity?* We show that inversity and assortativity can both be the same sign — both positive (or both negative) — with substantial effect size. We demonstrate that this same-sign property can hold across a range of networks. This leads to a wide array of cases wherein reliance on assortativity as a proxy for inversity would lead to incorrect decisions about optimal intervention strategies. Next, these decisions could also require knowing not just the sign of inversity, but the difference in magnitudes of expected degree across the intervention strategies. Approximating  $\rho \approx -\rho_a$  to obtain expected degree differences would further magnify these errors. Broadly, we document how using assortativity as a proxy metric for inversity is not conceptually appropriate, and not required since inversity is equally easy to compute. However, it would be worthwhile to examine whether this distinction is substantial for real-world networks.

*No monotonic ordering:* We observe that there are many pairs of networks for which we do not obtain appropriate monotonic ordering across inversity and assortativity. Specifically, there are pairs of networks  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , such that we have  $\rho(\mathcal{G}_2) > \rho(\mathcal{G}_1)$  and  $\rho_a(\mathcal{G}_2) > \rho_a(\mathcal{G}_1)$ , so *both* assortativity and inversity are higher for one of the networks. Such an ordering should not be possible to obtain if assortativity and inversity are a reverse ordering of each other. The implication is that there is not a clear one-to-one mapping between assortativity and inversity.

## Network Interventions

Consider the following problems: (a) (Reducing) A new infectious disease is spreading through a large population. We want to minimize the number of infected individuals by inoculating the population using a new vaccine. However, we only have a limited number of doses to administer. (b) (Accelerating) We have a new highly effective medical device with limited samples that we would like to provide to select medical professionals, who can then share the information

through word-of-mouth. (c) (Observing) We would like to identify a viral contagion as quickly as possible by choosing individuals as observation stations (or for contact tracing). Although seemingly disparate, these problems of how to reduce, accelerate or observe dynamic contagion, represent a class of network interventions in which we benefit from identifying more central or highly connected individuals in the network (11).<sup>¶</sup>

We show how seeding interventions using the friendship paradox, based on the ego-based and alter-based strategies developed here, impact network interventions by helping to obtain highly connected seed nodes in a privacy-sensitive manner from the relevant network. Note that while the ego-based strategy has been commonly suggested and used in practice (14–16), its theoretical and empirical properties have not been examined and characterized for general networks. The alter-based strategy is novel and first proposed here, and to our knowledge, has not been suggested or used for network interventions.

Our approach stands in contrast to most existing methods of identifying seeds for interventions, which focus on taking advantage of detailed network data on social connections, and even on activity to identify influential individuals (17–20). However, privacy concerns are increasingly relevant in such settings, making it challenging to obtain network data (21–25). Users are also concerned that their data may be used in algorithms (26), and even result in discrimination against them (27).

*Relevant Network Structure:* Even, if network data is available, the challenge in many cases is that we do not have access to the *relevant* network structure. In application (a), having the Facebook (or similar) network structure might not be useful, since the relevant network would be the *physical contact* network, which might be more challenging to obtain. In contrast, for application (b), finding a high degree node using a physical contact network of everyone who interacts with a medical professional is unlikely to be informative in characterizing opinion leadership in the profession. For (c), carrying out contact tracing for all individuals can be expensive in effort and time. These factors emphasize the importance of being able to leverage the structure of the relevant network, while being sensitive to privacy concerns.

These friendship paradox based strategies have several advantages for implementation. First, despite being informationally-light, the strategies here provide provable advantages for virtually any network structure, in contrast to strategies that don't provide such guarantees for general networks. The network structure may also be expensive to collect, may not be possible to obtain in a timely manner, or may vary over time, making the proposed method more valuable. Second, the strategies are much more privacy-sensitive than mapping out social networks. Third, the strategies can be implemented quickly since they only require local network information obtained by querying individuals or interaction data. Finally, the class of interventions here can be used for both advance and consequent interventions, i.e., for both prevention and treatment interventions.

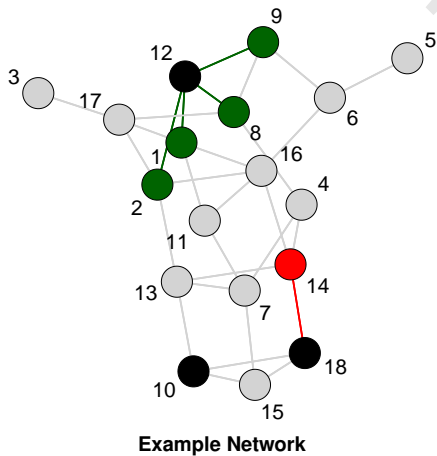
<sup>¶</sup>We focus on the class of "simple contagion" problems, which require only one exposure, rather than "complex contagion" problems, which require multiple exposures (12, 13).

## Implementing the Intervention Seeding Strategies

The above formulations of the ego-based and alter-based means suggest distinct strategies for choosing seeds for interventions, or intervention strategies. We illustrate *random*, *ego-based*, and *alter-based* strategies to choose a “seed” node in the network beginning with an initial randomly chosen node (Table 1). For example, the *ego-based strategy* would query randomly selected individual nodes with the query, “could you suggest the name of a randomly chosen neighbor?”

**Table 1. Implementation of Seeding Strategies**

Step	Details
<b>0</b>	Fix $p \in (0, 1]$ (only used for alter-based strategy in Step 2G). Repeat Steps 1-2 below until at least $k$ seeds are present in the seed set $\mathcal{S}$ .
<b>1</b>	Draw a random node $r$ uniformly from set of nodes, $V$ . If $r$ is an isolate, repeat this step. In Example Network, Nodes 10, 18, and 12 (in black) are drawn for (R), (E) and (A) strategies respectively.
<b>2</b>	Depending on the strategy Random (R), Ego-based (E) or Alter-based (A), do the following:
<b>2R (Random):</b>	Add $r$ to the seed set $\mathcal{S}$ . In Example Network, add node 10 to the seed set.
<b>2L (Ego-based):</b>	Obtain a node $s$ chosen with uniform probability from $r$ 's neighbors, i.e., $s \in \mathcal{N}_r$ . Add the neighbor $s$ to the seed set $\mathcal{S}$ . In Example Network, one of node 18's neighbors, node 14 (in red), is chosen at random. Add node 14 to the seed set.
<b>2G (Alter-based):</b>	For each of $r$ 's neighbors, $s \in \mathcal{N}_r$ : with probability $p$ ( $0 < p \leq 1$ ), add $s$ to the seed set $\mathcal{S}$ . In Example Network, each of node 12's neighbors, nodes 1, 2, 8, and 9 (in green), are added probabilistically (with probability $p$ ) to the seed set. Implementation: For each $s \in \mathcal{N}_r$ , draw from an independent uniformly distributed random variable $z_s \sim U[0, 1]$ . If $z_s < p$ , add $s$ to the seed set $\mathcal{S}$ .
Note: With Random and Ego-based strategies, we will obtain exactly $k$ nodes in the seed set $\mathcal{S}$ . With the Alter-based strategy, we might obtain more than $k$ nodes in the seed set. In such a case, we select $k$ nodes at random from the seed set $\mathcal{S}$ without replacement.	



**Example Network**

The *alter-based strategy* would ask individual nodes to choose each of their neighbors with a probability that is fixed across nodes. The probability can be set to be small

(say  $p=0.05$ ) based on how many total seeds are required for interventions, and also to balance privacy concerns. The alter-based strategy gives each neighbor of each random person an equal chance of being selected, and we prove that the expected degree of chosen nodes is equal to the alter-based mean (Theorem S5).<sup>||</sup>

We illustrate how our approach is able to obtain the *relevant network structure* in a straightforward manner. Specifically, we query nodes to select from the *relevant network*. For instance, in application (a) where the focus was on physical contagion, the relevant network is the in-person contact network. The query would then be phrased as “among the people you have interacted with in-person, choose one at random.” The idea of such queries to obtain the relevant network is general, and conditions can be added to the query (e.g. specifying a time period), depending on the desired intervention. Similar conditions can be used for applications (b) and (c). We can thus view the above as a query that provides a network that is relevant to the specific application.

The crucial distinction between the ego-based and alter-based strategies lies in whether we are choosing *one* random neighbor (ego-based) or a *fixed probability* for each neighbor (alter-based) of randomly chosen individuals. Table 1 details the algorithms to obtain  $k$  seeds in a network of size  $N \gg k$ . Next, we evaluate the relative effectiveness of these intervention strategies.

## Effectiveness of Strategies: Leverage

To evaluate how much of an improvement over the random strategy is possible, and how this varies across a variety of generated and real networks, we examine the relative effectiveness of strategies, with the random strategy as the baseline and characterize leverage as the improvement in expected degree. Leverage for strategy  $s$  on network  $\mathcal{G}$  is defined as  $\lambda_s(\mathcal{G}) = \frac{\mu_s(\mathcal{G})}{\mu_D(\mathcal{G})}$  for  $s \in \{R, E, A\}$  (since the random strategy obtains the mean degree in expectation, the leverage for  $R$  is  $\lambda_R(\mathcal{G}) = 1$  and it serves as a baseline). We examine the leverage of both generated and real networks.

**Generated Networks.** The generated networks were obtained using three generative mechanisms (29–31): (a) Random or Erdos-Renyi (ER), (b) Scale Free (SF), and (c) Small World (SW) models, as detailed in the Supplement S5.E.

The results are detailed in Fig. S5 in the Supplement. We find that for ER networks, at very low density (edge probability), the leverage is very low because most edges connect nodes that have a degree of 1. As density increases, we obtain more variation in degrees, and ego-based leverage increases. However, beyond an edge probability of  $p = 0.05$ , leverage decreases as the density of the network increases. Ego-based leverage thus forms a non-monotonic pattern with ER networks. For SF networks, rather than density or edge probability, we initially examine leverage as the network becomes more centralized (as  $\gamma$  increases above 1, very high degree nodes have a lower probability of occurring). We find

<sup>||</sup> These ego-based and alter-based strategies also have connections with respondent driven sampling (RDS), in which respondents nominate random neighbors or alters, e.g., by giving them participation tickets (28). An additional advantage of using such an approach is that the privacy risks are reduced further. The fact that these RDS based approaches have been commonly used in earlier interventions indicates that our proposed strategies are practical and knowledge about implementing them in specific contexts is likely to already exist.

that as  $\gamma$  increases from 1 to 2, the leverage increases, but then decreases beyond 2. For SW networks, unlike in the ER and SF networks, leverage is monotonically decreasing with number of neighbors (or density), and is monotonically increasing with rewiring probability.

**Real Networks.** The range of real networks is detailed in Supplement §S.C. First, observing the ego-based strategy (Fig. 2A), we find that for all networks, as expected, the friendship paradox strategies are at least as good as the random strategy. Second, for networks like Twitter (OS4) or Internet Topology (C1), the leverage can be highly substantial — on the order of 100 — implying that obtaining a connection of a random node will provide a 100-fold increase in expected degree. Third, we observe that both ego-based and alter-based leverage (Figs. 2A and 2B) are higher for nodes when average degree is intermediate, i.e., not too low or high. Some networks like the CA Roads network (I3) have very little degree variation, and ego-based and alter-based strategies are relatively less effective. Finally, we examine the conditions under which ego-based and alter-based strategies have a relative difference (Fig. 2B). We find that the highest ratio of ego-based to alter-based mean is for the Twitter network (OS4), whereas the lowest ratio (indicating that global strategy has a higher expected mean degree) is shown by Flickr (OS2), both of which belong to the same category of online social networks. Citation networks tend to have a higher global mean, whereas for infrastructure networks, both strategies seem to work just as well.

## Application: Controlling Contagion in Networks

We next illustrate the approach of using the friendship paradox strategies to obtain seeds for intervention, specifically vaccination in the face of simple contagion spreading through a network. Our goal is not for the application to directly inform immunization policy for a particular disease, but rather to serve as a proof of concept. The virus propagation model here is simple and reduced to essential components. To be more realistic, the model would be more general, e.g. richer spatial models incorporating heterogeneity, potentially continuous and discrete time, and having parameters calibrated to match epidemiological data (32, 33).

We focus on simple models of contagion that can be characterized by a single parameter termed the *epidemic threshold* to focus our analysis on the benefit provided by the ego-based and alter-based strategies. The epidemic threshold captures the idea that a contagion introduced into the network will die out if the reproductive number ( $\mathcal{R}_0$ ) is below the epidemic threshold, and will lead to an epidemic if  $\mathcal{R}_0$  is above the threshold. Thus, a network with a higher epidemic threshold would be able to better withstand or control an infection. We then examine how the epidemic threshold changes as a function of the proportion of nodes vaccinated (removed), using each strategy (random, ego-based, and alter-based).

For a wide class of *virus propagation models* (VPM), the epidemic threshold is characterized as the inverse of the greatest (first) eigenvalue of the adjacency matrix  $A$  of the network, denoted as  $\tau(E) = \frac{1}{\lambda_1(E)}$  (details in the Supplement §S.H). The above formulation applies to a range

of VPMs, including SIR, SEIR, etc., which include models commonly used for infectious diseases (34).

We use in-person contact networks for modeling contagion, with data on 75 village social networks from India (35). The social networks are captured at two different levels of aggregation, at the level of individuals and of households. The advantage of this dataset is that villages are relatively geographically isolated and can therefore be treated as separate networks. Details of the network dataset are provided in the Supplement §S.C.

We find that the village networks can have either positive or negative inversity depending on how nodes and edges are defined and aggregated. Figure 3(a) illustrates the inversity values across the 75 villages separately for individual and household networks. When nodes are defined as individuals, we find that the networks have negative inversity, whereas if the nodes are defined as households, the inversity values of the resulting networks are mostly positive. Since household-level ties are aggregated from the individual-level ties, we find that networks obtained from similar underlying relationships can result in dramatically different inversity characteristics, which can lead to different interventions. Considering interventions, the inversity values suggest that a household-based intervention might use the ego-based strategy, whereas the individual-based intervention might use the alter-based strategy.

We next evaluate how the epidemic threshold  $\tau$  changes as we immunize nodes from the network for each of the strategies (random, ego-based and alter-based). While immunizing (or removing) any node from the network is likely to increase the epidemic threshold, immunizing highly-connected nodes is likely to prove especially beneficial. In Figure 3(b), we show how the epidemic thresholds vary across strategies and proportion of nodes immunized (1% — 75%). In both household and individual networks, we find that the friendship paradox strategies obtain higher epidemic thresholds than the random strategy, for the same proportion of nodes immunized. For instance, in the household networks, to achieve a epidemic threshold  $\tau = 0.15$ , the random strategy needs to have about 50% of nodes immunized, but the ego-based and alter-based strategies require less than half of that, at around 25%. For the household networks, we find that the ego-based strategy is better than the alter-based strategy, especially at higher levels of removal. However, for individual networks, we find that the alter-based strategy obtains greater thresholds than the ego-based strategy. This broadly signifies that it is helpful to know which among the alter-based or ego-based strategies to use, as determined by the sign of inversity.

Finally, we simulate an infection process and evaluate the epidemic characteristic of peak infection using an SIR virus propagation model (details in Supplement §S.H) (36), with parameters of the simulation detailed in Table S5. We examine peak infection since it is known to be an important characteristic of epidemics (37), directly impacting the load on the healthcare system. We denote  $I_{it} \in \{0, 1\}$  as an indicator of whether an individual  $i$  is infected at time  $t$ . We evaluate the proportion of the population infected at the peak of the epidemic ( $\frac{1}{N} \max_t (\sum_i I_{it})$ ), which is a useful measure in cases where healthcare capacity is constrained. There has been much discussion about the value of interventions to avoid and minimize such a peak (38). A strategy with a density

plot to the left of another is better in terms of reducing the severity of the epidemic. Thus, for household networks, the ego-based strategy (in red) is better than the alter-based, which in turn is better than the random strategy in reducing peak infection. For the individual networks, however, the alter-based strategy is better than the ego-based strategy. Overall, we find that the friendship paradox based ego-based and alter-based strategies clearly improve upon the random strategy.

## Conclusion

We have shown fundamental mathematical properties that underlie the friendship paradox, which we find to be multifaceted. We define and characterize the properties of the ego-based mean, alter-based mean, and inversity to connect the means for any network. We show that for unknown networks, the ego-based and alter-based strategies based on these means have theoretical guarantees on obtaining better-connected individuals from the relevant network. With both generated random networks and real networks, our results show the substantial value of using the friendship paradox strategies to obtain highly connected nodes. In the vast majority of networks, we obtain at least double the average degree, and some networks show increases of close to a hundred-fold increase in node degree. We expect the advantages of these strategies, including sensitivity to privacy, speed of implementation, and generality of application areas, to prove important in using these strategies for interventions in unknown network structures.

1. S Feld, Why your friends have more friends than you do. *Am. J. Sociol.* pp. 1464–1477 (1991).
2. E Zuckerman, J Jost, What makes you think you're so popular? self-evaluation maintenance and the subjective side of the "friendship paradox". *Soc. Psychol. Q.* pp. 207–223 (2001).
3. HH Jo, YH Eom, Generalized friendship paradox in networks with tunable degree-attribute correlation. *Phys. Rev. E* **90**, 022809 (2014).
4. YH Eom, HH Jo, Generalized friendship paradox in complex networks: The case of scientific collaboration. *Sci. reports* **4**, srep04603 (2014).
5. DJ Higham, Centrality-friendship paradoxes: when our friends are more important than us. *J. Complex Networks* **7**, 515–528 (2019).
6. NO Hodas, F Kooti, K Lerman, Friendship paradox redux: Your friends are more interesting than you. *ICWSM* **13**, 8–10 (2013).
7. GT Cantwell, A Kirkley, M Newman, The friendship paradox in real and model networks. *J. Complex Networks* **9**, cnab011 (2021).
8. R Czaja, Sampling with probability proportional to size. *Encycl. Biostat.* **7** (2005).
9. A Nigam, P Kumar, V Gupta, Some methods of inclusion probability proportional to size sampling. *J. Royal Stat. Soc. Ser. B: Stat. Methodol.* **46**, 564–571 (1984).
10. C Orsini, et al., Quantifying randomness in real networks. *Nat. communications* **6**, 1–10 (2015).
11. T Valente, Network interventions. *Science* **337**, 49–53 (2012).
12. D Centola, M Macy, Complex contagions and the weakness of long ties. *Am. journal Sociol.* **113**, 702–734 (2007).
13. D Centola, The spread of behavior in an online social network experiment. *Science* **329**, 1194–1197 (2010).
14. R Cohen, S Havlin, D Ben-Avraham, Efficient immunization strategies for computer networks and populations. *Phys. review letters* **91**, 247901 (2003).
15. DA Kim, et al., Social network targeting to maximise population behaviour change: a cluster randomised controlled trial. *The Lancet* **386**, 145–153 (2015).
16. M Alexander, L Forastiere, S Gupta, NA Christakis, Algorithms for seeding social networks can enhance the adoption of a public health intervention in urban india. *Proc. Natl. Acad. Sci.* **119**, e2120742119 (2022).
17. D Kempe, J Kleinberg, É Tardos, Maximizing the spread of influence through a social network in *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*. pp. 137–146 (2003).
18. Z Katona, PP Zubcsek, M Sarvary, Network effects and personal influences: The diffusion of an online social network. *J. marketing research* **48**, 425–443 (2011).
19. C Wilson, A Sala, KP Puttaswamy, BY Zhao, Beyond social graphs: User interactions in online social networks and their implications. *ACM Transactions on Web (TWEB)* **6**, 1–31 (2012).
20. MA Al-Garadi, et al., Analysis of online social network connections for identification of influential users: Survey and open research issues. *ACM Comput. Surv. (CSUR)* **51**, 1–37 (2018).

## Materials and Methods

Our analysis combines theoretical results along with simulation and empirical analysis on generated and real-world networks, in order to characterize the fundamental properties of the friendship paradox and related constructs.

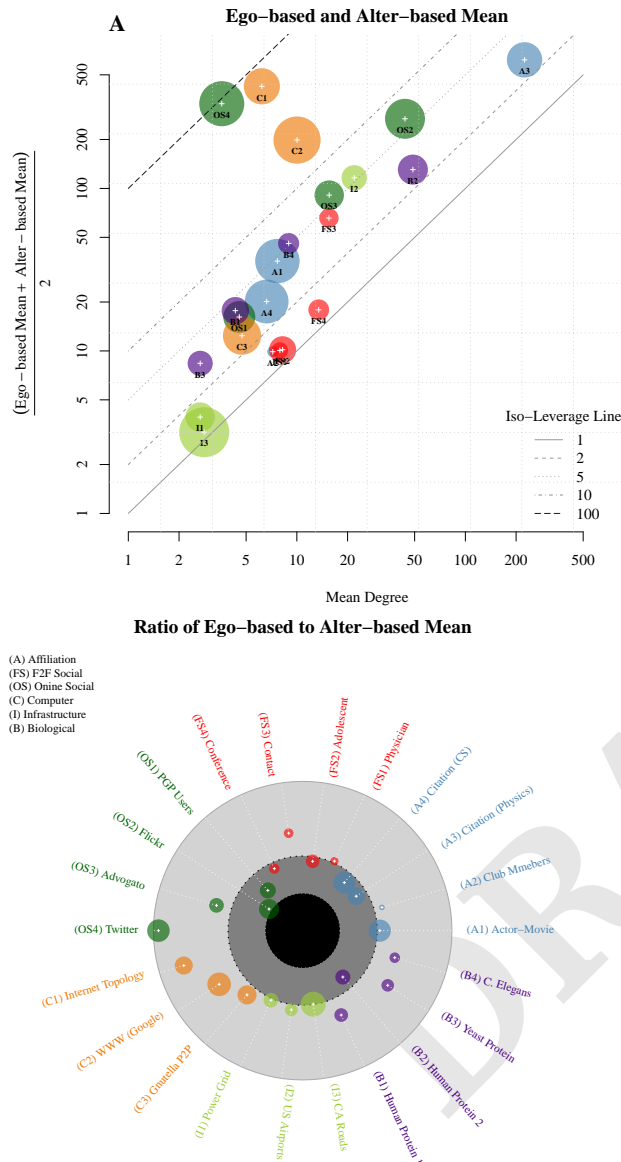
**Theoretical Properties.** The theoretical results are contained in Supplement §S.B. We prove the Theorems on the individual friendship paradox, the properties of ego-based mean and alter-based means, and inversity by using the properties of Networks (Graphs) and identifying the conditions required for these relationships to hold.

**Empirical Analysis.** For the empirical analysis, there are two separate but related parts. First, for the generated networks, in Supplement §S.B, we examine the most commonly used generative mechanisms, i.e., Random Graphs (Erdos-Renyi), Scale Free (Barabasi-Albert), and Small World (Watts-Strogatz) networks. Second, we use real-world network data for simulations, and include the empirical analysis in Supplement §S.C. Finally, we conducted a study of virus propagation under immunization carried out using the ego-based, alter-based, and random strategies. The model specification, simulation, parameterization, and values are contained in Supplement §S.H. Simulation and empirical analysis was performed in R software, using igraph and sna packages.

**ACKNOWLEDGMENTS.** We acknowledge helpful comments from Steven Strogatz, Ed Kaplan, the audiences at 2013 and 2016 Sunbelt Social Network Conferences, and seminar audiences at Massachusetts Institute of Technology, Carnegie Mellon University, Yale University, University of Texas at Austin, University of California, San Diego and Washington University at St. Louis. The authors each acknowledge financial support from their respective universities.

21. A Acquisti, L Brandimarte, G Loewenstein, Secrets and likes: The drive for privacy and the difficulty of achieving it in the digital age. *J. Consumer Psychol.* **30**, 736–758 (2020).
22. F Cerruto, S Cirillo, D Desiato, SM Gambardella, G Polese, Social network data analysis to highlight privacy threats in sharing data. *J. Big Data* **9**, 19 (2022).
23. A Praveena, S Smys, Anonymization in social networks: a survey on the issues of data privacy in social network sites. *J. Int. J. Of Eng. And Comput. Sci.* **5**, 15912–15918 (2016).
24. W Xie, K Karan, Consumers' privacy concern and privacy protection on social network sites in the era of big data: Empirical evidence from college students. *J. Interact. Advert.* **19**, 187–201 (2019).
25. S Torabi, K Beznosov, Privacy aspects of health related information sharing in online social networks in 2013 *USENIX Workshop on Health Information Technologies (HealthTech 13)*. (2013).
26. B Liu, et al., When machine learning meets privacy: A survey and outlook. *ACM Comput. Surv. (CSUR)* **54**, 1–36 (2021).
27. S Wachter, Normative challenges of identification in the internet of things: Privacy, profiling, discrimination, and the GDPR. *Comput. law & security review* **34**, 436–449 (2018).
28. SK Thompson, *Sampling*. (John Wiley & Sons) Vol. 755, (2012).
29. P Erdős, A Rényi, On random graphs, i. *Publ. Math. (Debrecen)* **6**, 290–297 (1959).
30. AL Barabási, R Albert, Emergence of scaling in random networks. *Science* **286**, 509–512 (1999).
31. D Watts, S Strogatz, Collective dynamics of "small-world" networks. *nature* **393**, 440–442 (1998).
32. LJ Thomas, et al., Spatial heterogeneity can lead to substantial local variations in covid-19 timing and severity. *Proc. Natl. Acad. Sci.* **117**, 24180–24187 (2020).
33. LJ Thomas, et al., Geographical patterns of social cohesion drive disparities in early covid infection hazard. *Proc. Natl. Acad. Sci.* **119**, e2121675119 (2022).
34. BA Prakash, D Chakrabarti, M Faloutsos, N Valler, C Faloutsos, Got the flu (or mumps)? check the eigenvalue! *arXiv preprint arXiv:1004.0060* (2010).
35. A Banerjee, AG Chandrasekhar, E Duflo, MO Jackson, The diffusion of microfinance. *Science* **341**, 1236498 (2013).
36. J Tolles, T Luong, Modeling epidemics with compartmental models. *Jama* **323**, 2515–2516 (2020).
37. A Soría, et al., The high volume of patients admitted during the SARS-Cov-2 pandemic has an independent harmful impact on in-hospital mortality from covid-19. *PloS one* **16**, e0246170 (2021).
38. RD Booton, et al., Estimating the COVID-19 epidemic trajectory and hospital capacity requirements in South West England: a mathematical modelling framework. *BMJ Open* **11**, e041536 (2021).





**Fig. 2.** Alter-based and Ego-based Leverage in Real Networks Ego-based and Alter-based Means across Networks (each circle is a network). Area of circles indicates size of networks (number of nodes) in log scale. Color of circle indicates network category. (A) The average of ego-based and alter-based mean is higher than mean degree in all real networks, with the highest differences occurring in online social networks and computer networks. Most large networks also tend to show a higher leverage ratio. For in-person or face to face networks, the pattern is more variable. The iso-leverage line indicates leverage levels of 1,2,5,10 and 100. We find that all networks have leverage greater than 1, a majority of networks have leverage greater than 5, and 2 networks have leverage close to 100. (B) **Comparison:** Ratio of Ego-based to Alter-based Mean. The ratio of ego-based to alter-based mean  $\frac{\mu_E}{\mu_A}$  is represented as follows ( $< \frac{1}{2}$  in black circle,  $\frac{1}{2} < \frac{\mu_E}{\mu_A} < 1$  in dark gray circle and  $1 < \frac{\mu_E}{\mu_A} < 2$  in light gray circle. For example, in the Twitter network, ego-based mean is almost twice the alter-based mean, whereas in the Flickr network, alter-based mean is almost twice the ego-based mean. Computer networks have higher values of the ratio, whereas Infrastructure networks have similar values of ego-based and alter-based means.

**Fig. 3.** Inversity and Epidemic Characteristics on Village Networks. Data for these networks (N=75) obtained from Indian villages is publicly available and detailed in (35). The data includes both individual-level (individual network) ties as well as connections between households (household network). (a): Inversity: Frequency plot of inversity across village networks. (b) Epidemic Thresholds with Immunization: The epidemic threshold ( $\tau = \frac{1}{\lambda_1}$ ) is computed at different levels of immunization, and for different immunization strategies (random, ego-based and alter-based). The strategies are used to select nodes, which upon immunization lose the capability to transmit contagion, and infect other nodes. (c) Epidemic Peak Infection: The proportion of the network nodes that are infected during the peak of the epidemic is represented as a density plot. Variation is obtained due to differences in outcomes across villages as well as simulation variation. See Table S5 in the Supplement for simulation parameters.

