

# Influencing Network Graph Perception and Judgment: Effects of Direct Connections, Base Rates, and Visual Layout Proximity on Social Network Analysis

Rebecca E. Rhodes<sup>1</sup> (rebecca.rhodes@jhuapl.edu), Isaiah Harbison<sup>1</sup> (james.harbison@jhuapl.edu),  
Nathan Bos<sup>1</sup> (nathan.bos@jhuapl.edu), Celeste Lyn Paul<sup>2</sup> (clpaul@tycho.ncsc.mil),  
Clay Fink<sup>1</sup> (clay.fink@jhuapl.edu), Anthony Johnson<sup>1</sup> (anthony.johnson@jhuapl.edu)

1. Johns Hopkins University Applied Physics Laboratory, Laurel, MD  
2. U.S. Department of Defense, USA

## Abstract

Social network graphs are often used to help inform judgments in a variety of domains, such as public health, law enforcement, and political science. Across two studies, we examined how graph features influenced probabilistic judgments in graph-based social network analysis and identified multiple heuristics that participants used to inform these judgments. Study 1 demonstrated that participants' judgments were influenced by information about direct connections, base rates, and layout proximity, and participants' self-reported strategies also reflected use of this information. Study 2 replicated findings from Study 1 and provided additional insight into the hierarchical ordering of these strategies and the decision process underlying judgments from social network graphs.

**Keywords:** social network analysis; graph comprehension; data visualization; judgment and decision making

## Introduction

### Social Network Analysis

Social network analysis is an analytical method for understanding data that depicts relationships between entities, such as communications between people or the flow of information through a community. Social network analysis (SNA) can be used to inform judgments such as who the most connected or influential person is in a community. The applications for SNA are widespread. The CDC has used social networks to map the spread of infectious diseases (Cook, 2007), and intelligence professionals have used SNA to map potential terrorist suspects. For example, Krebs (2002) was able to map a network of suspected 9/11 terrorists after the attacks based on publicly available information in the news and provide insight into the terrorist organization.

Social networks are typically depicted using node-link diagrams, in which nodes (depicted with circles) represent objects, such as people, and edges between the nodes (depicted with lines) represent a connection, such as a friendship or a recent communication. Although there are quantitative measures of graph structures, such as network centrality, SNA often involves some degree of visual interpretation. In many applications of SNA, such as tactical decision making in military organizations, users of network graphs have limited time to make their decisions and may not have the background to supplement their visual

interpretation with more objective, quantitative measures. Thus, clearly communicating the important information in network graphs is essential. In practice, however, these network visualizations are often complex with hundreds or thousands of nodes, and relatively little is known about the cognitive strategies people use to make probabilistic judgments from these graphs. In the present study, we aimed to identify graph features and heuristics that influenced probabilistic judgments in scenarios that mirrored real-world SNA tasks.

### Graph Perception

Much research has applied perceptual and cognitive theories to improve graph comprehension and readability. Purchase, Cohen, and James (1997) first demonstrated that several aesthetic qualities could improve comprehension, such as maximizing the symmetry of the graph, minimizing the number of intersecting edges, and minimizing the number of bends in a series of segments. Purchase (2000) found that graph layouts that minimized intersecting edges significantly improved participants' abilities to answer questions about the relationships in the graph. Similarly, Ware, Purchase, Colpoys, and McGill (2002) found that cutting unnecessary edges significantly improved cognitive performance.

Other perceptual features in graphs have been studied as well. For example, McGrath, Blythe, & Krackhardt (1997) found that nodes that were positioned closer to the center of the graph were perceived as more prominent, suggesting a bias towards layout centrality. Additionally, nodes were more likely to be considered to be in the same group if they were positioned closer to each other, and nodes positioned directly between two nodes (as opposed to at an angle) made them more likely to be considered a "bridge" between the two nodes.

Gestalt principles may help explain why some aesthetic features can influence graph interpretation by exploiting people's tendency to identify patterns (Novick & Bassok, 2005). The Gestalt principle of good continuation suggests that paths are more easily recognized when they contain less jaggedness, suggesting that the fewer the bends between two nodes, the more readily the path connecting them will be observed. The principle of proximity suggests that spatial proximity between objects implies logical groupings. Other cognitive heuristics may also influence graph interpretation.

For example, people sometimes ignore base rates when making probabilistic judgments. The extent to which this occurs when making judgments from network graphs, which may depict base rate information in a visual way, is unclear (but see Micallef, Dragicevic, & Fekete, 2012).

### Pilot Study

We conducted a pilot study with 30 participants recruited from Amazon Mechanical Turk (AMT) to identify graph features people attend to when making judgments about the likelihood that a particular node belonged to a defined category. We planned to use the results of the pilot study to inform our decision about graph features to manipulate in our experiments<sup>1</sup>. Participants were asked to play the role of a data analyst at a law enforcement agency and use node-link graphs to estimate the likelihood that they would further investigate a person of interest (POI) in a community that has known drug users (Figure 1). They were also asked to describe the strategies they used to make their judgments. The POI was indicated as a black node, known drug users as blue nodes, and everyone else as yellow nodes. Lines were drawn between nodes to indicate recent communications. Two graph features were varied, 1) the number of direct connections the POI had, and 2) the number of drug users in the whole graph (base rate). The graph visual was generated using a force-directed layout, a popular way of visualizing graph data. Participant responses revealed two common strategies. The first was based on the number of drug users within the POI's direct connections (what we call the "ratio" strategy,  $N = 10$ ). The second strategy was based on the spatial proximity of drug user nodes to the POI node in the visual layout, regardless of whether they were actually closely connected to the POI ("proximity" strategy,  $N = 8$ ). We expected to see instances of the ratio strategy, but not necessarily the proximity strategy, since the position of nodes in a force-directed layout is a function of connections rather than nodes. In general, layout position in a graph can be misleading, and two nodes visually close to each other but not connected may not necessarily have a strong relationship. Three participants in the pilot study also mentioned relying on the total number of drug users in the graph. This strategy was expected to be more frequent considering that base rates were explicitly manipulated across graphs. However, participants apparently found base rates less informative than ratio and proximity features for the pilot set of graphs.

### Overview of Experiments

We conducted two experiments to understand how the graph features identified in the pilot study influence probabilistic judgments. In each study, we manipulated the three features identified during our pilot study: 1) ratio, or the number of connections the POI has to salient nodes, 2)

base rate of salient nodes, and 3) proximity of salient nodes to the POI.

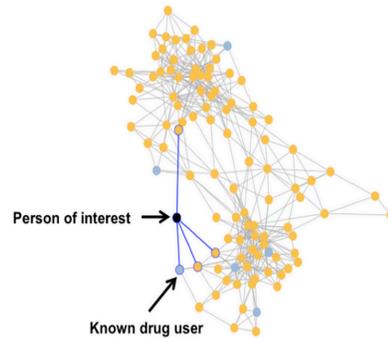


Figure 1: Example network graph provided in pilot study.

## Study 1

### Participants and Procedure

**Participants.** We recruited 30 participants from AMT. Participants were required to have a "Masters" qualification, indicating that they had repeatedly demonstrated good performance in the AMT marketplace. Most participants reported having no experience (72.41%) or only slight experience (17.24%) with SNA. One participant was excluded from analysis for failing an attention check.

**Scenario.** Similar to the pilot study, participants were asked to take on the role of a data analyst at a law enforcement agency and use graphs to identify potential drug users in a community. They were provided a set of 36 graphs that described communications in a small community over the past month. Participants were told that links between nodes represented non-family relationships with others, such as friends and coworkers. No information was provided on how the graphs were constructed. For each graph, participants made a judgment of how likely they would be to investigate the POI and described how they used the graph to make their judgment.

**Materials.** The graphs were small real-world graphs generated from the Polbooks dataset of co-purchased political books (V. Krebs, unpublished, <http://www.orgnet.com/>). Graphs were rendered with the default graph plotting parameters in the R iGraph package, which uses the Fruchterman-Rheingold force-directed algorithm. Each graph consisted of 104 nodes. As in the pilot study, the POI was indicated as a black node, known drug users as blue nodes, and everyone else as yellow nodes. Direct connections to the POI were highlighted in the graphs to aid interpretation. Emphasizing direct connections in this way may have led to more emphasis on the ratio strategy, but highlighting connections was deemed important to help participants distinguish between graph connectedness from layout proximity (i.e. visual closeness).

<sup>1</sup> All procedures in the pilot study and experiments were approved by the Johns Hopkins University Institutional Review Board.

**Conditions.** Three graph conditions were manipulated within-subjects: 1) number of drug users in the POI's  $N$  direct neighbors (0/ $N$ , 1/ $N$ , Half/ $N$ ); 2) overall base rate of known drug users in entire graph (5% or 50% of all nodes); and 3) proximity placement of known drug user nodes (near or far from POI). Note that proximity was the physical location of drug user nodes within the visual graph layout and not the graph connectedness. Graph connectedness was controlled by removing edges that connected proximal drug user nodes to the POI's direct neighbors (Figure 2).

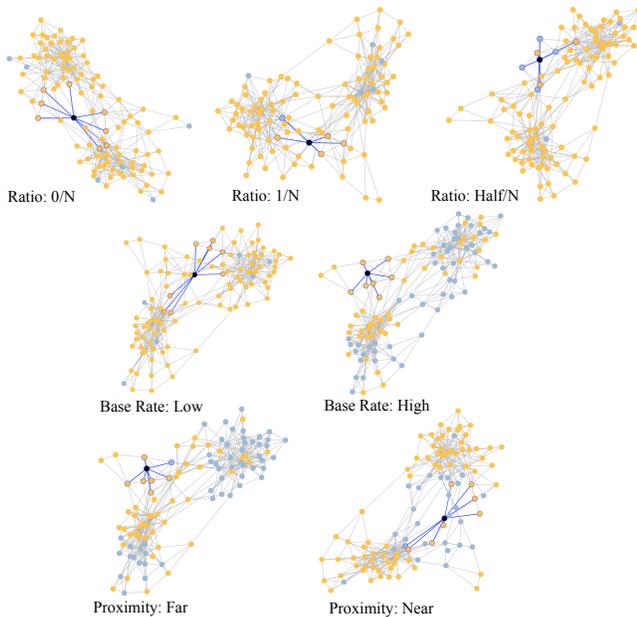


Figure 2: Graph feature conditions. Note that in the Proximity: Near graph, the proximal blue nodes are not connected to the POI or the POI's direct neighbors.

**Likelihood Judgments.** For each graph, participants were asked to rate how likely they would be to investigate the POI on a scale of 1 to 7 (1 = Extremely Unlikely, 7 = Extremely Likely).

**Strategy Use.** After giving their likelihood rating, participants were also asked to describe any specific graph information they used in making their rating. Responses were open-ended text responses and were analyzed using a structured coding approach, described in Table 1, with three coders who coded 75% of the responses independently and 25% overlapping. Inter-rater reliability on the 25% co-coded responses was high (Cohen's Kappa > .75).

Table 1: Coding scheme for open-ended responses.

Strategy	Description
Ratio	Mentions number of users in POI connections
Base Rate	Mentions number of users in graph
Proximity	Mentions visual closeness of POI to users

## Results and Discussion

**Likelihood Judgments.** We modeled likelihood judgments with a linear multilevel model with base rate, proximity, and ratio conditions treated as random effects. Graph feature conditions were nested within participant. Multilevel models are able to estimate the variance associated with each random effect; thus, the model can account for the within-individual variability for different graph feature conditions. In addition to modeling random effects, multilevel models can also simultaneously estimate fixed effects, which in this case represent the average effect for the whole sample.

We sequentially entered each graph feature into the model as a fixed effect to examine its effect on likelihood ratings. Base rate, ratio, and proximity each significantly improved the fit of the model ( $\chi^2(1) = 34.51, p < .001, \chi^2(2) = 148.62, p < .001$ , and  $\chi^2(1) = 45.07, p < .001$ , respectively). Specifically, higher base rates increased the likelihood of investigating the POI ( $b = 0.97, t(1011) = 7.54, p < .001$ ), higher ratios increased the likelihood of investigating ( $b = 1.24, t(1011) = 7.18, p < .001$  for 1/ $N$  vs. 0/ $N$  and  $b = 3.86, t(1011) = 22.34, p < .001$  for Half/ $N$  vs. 1/ $N$ ), and closer proximities also increased the likelihood of investigating ( $b = 1.02, t(1011) = 9.58, p < .001$ ). These results suggest that participants considered all three graph features when making their judgments. The observation that participants' suspicion levels of the POI increased with proximity to known drug users was significant, because it suggests that participants believed proximity to be a meaningful indicator.

We also found a significant interaction between base rate and proximity ( $\chi^2(1) = 146.00, p < .001$ ). When proximity was far, there was no difference between the base rate conditions, but when proximity was near, higher base rates significantly increased the likelihood of investigating ( $b = 1.70, t(1006) = 13.96, p < .001$ ; Figure 3). In other words, base rate information only affected judgments to the extent that it placed more users near the POI.

There were also interactions between the base rate and ratio conditions ( $\chi^2(2) = 112.90, p < .001$ ) as well as between the ratio and proximity conditions ( $\chi^2(2) = 81.89, p < .001$ ), as depicted in Figure 3. Increasing the base rate had a stronger effect in the 0/ $N$  and 1/ $N$  conditions compared to the Half/ $N$  condition ( $b = 1.70, t(1006) = 11.40, p < .001$  and  $b = 0.94, t(1006) = 6.32, p < .001$  for 0/ $N$  vs. Half/ $N$  and 1/ $N$  vs. Half/ $N$ , respectively). Likewise, proximity had a stronger effect in the 0/ $N$  and 1/ $N$  conditions as well ( $b = 1.36, t(1006) = 9.09, p < .001$  and  $b = 0.87, t(1006) = 5.85, p < .001$ , respectively). These interactions reveal that, in the highest ratio condition, base rate and proximity information had little effect. One interpretation of this result is that when ratio was high, participants did not feel a need to look at other information in the graph. In contrast, when the ratio was 0/ $N$  or 1/ $N$ , base rate and proximity information both influenced judgments, such that higher base rates and closer proximity increased likelihoods of investigating.

**Strategy Use.** Analysis of participants’ self-reported strategies validated our findings from the pilot study. Participants mentioned using the ratio strategy for the majority of graphs (85.34%). Proximity was the next most frequently mentioned strategy (22.41%), followed by base rate a small percentage of the time (12.64%).

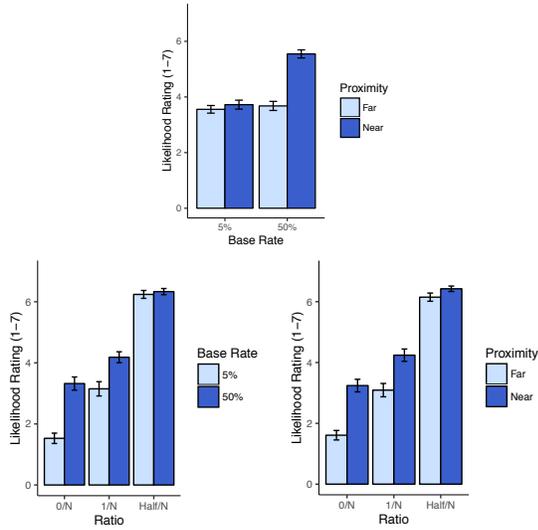


Figure 3: Interactions between graph feature conditions.

## Study 2

Study 2 sought to replicate the findings of Study 1 with a broader range of stimuli. We introduced two additional scenarios, public opinion and disease propagation. Network graphs can be used to represent relational data from a variety of domains, and assumptions about the data may change interpretations of the graph. For example, base rates may be utilized more when the activities represented in the graphs are perceived to be more mobile or contagious. We speculated that a network graph representing the spread of an infectious disease could lead people to use base rate strategies more often than a network graph representing the spread of drug use. Study 2 also sought to better understand the hierarchical ordering of cognitive strategies identified in the previous studies by constructing a decision tree that could predict the conditions under which particular strategies will be used.

### Participants and Procedure

**Participants.** We recruited 196 participants from AMT. Participants were limited to U.S. residents and had at least a 90% approval rating for previous HITs. We excluded 11 participants who either failed attention checks or had participated in a previous study. Most participants said they had no prior experience (64.86%) or only slight prior experience (28.65%) with SNA.

**Scenario and Materials.** Participants were provided the same set of 36 graphs from Study 1 and were asked to make

likelihood judgments about the POI in each graph.

**Conditions.** The same three graph conditions from Study 1 were manipulated within-subjects. We also introduced different scenarios as a between-subjects manipulation for the source of the graph data. Participants were randomly assigned to one of three data scenarios: drug use, political opinion, and infectious disease (Table 2).

Table 2: Description of data scenarios.

Scenario	Brief Description
Drug	In a community in which some percentage of people are known users, how likely will the POI become a drug user in the next six months?
Opinion	In a community in which some percentage of people are known proponents of a new proposal, how likely will the POI become a proponent in the next six months?
Disease	In a community in which some percentage of people are infected with a new disease transmitted through social contact, how likely will the POI become infected in the next six months?

**Strategy Use.** In contrast to Study 1, in which participants described their analysis strategy in an open-text form, participants in Study 2 rated the strategies from Study 1 (ratio, base rate, and proximity) as well as distractors (e.g., central or peripheral location of POI, presence of a cluster of nodes) as to how important they were to their likelihood judgment. Ratings were on a scale of 1 (Not at all important) to 5 (Extremely important). The list order of strategies was randomized at each presentation.

### Results and Discussion

**Likelihood Judgments.** We used the same model from Study 1 but added scenario to the model as a fixed effect. There was a significant main effect of scenario on likelihood judgments ( $\chi^2(2) = 26.40, p < .001$ ); specifically, participants in the disease condition gave higher likelihood ratings on average than participants in the drugs or opinion conditions ( $b = 0.58, t(182) = 4.65, p < .001$  and  $b = 0.63, t(182) = 4.71, p < .001$  for disease vs. drugs and disease vs. opinion, respectively). As in Study 1, there were also main effects of base rate ( $\chi^2(1) = 288.97, p < .001$ ), proximity ( $\chi^2(1) = 291.26, p < .001$ ), and ratio ( $\chi^2(2) = 671.77, p < .001$ ). As depicted in Figure 4, higher base rates led to higher likelihood ratings ( $b = 1.21, t(6471) = 24.06, p < .001$ ), higher ratios led to higher ratings ( $b = 0.62, t(6471) = 8.92, p < .001$  for 1/N vs. 0/N, and  $b = 2.16, t(6471) = 31.39, p < .001$  for Half/N vs. 1/N), and closer proximities also led to higher ratings ( $b = 0.99, t(6471) = 25.00, p < .001$ ).

We again found two-way interactions between base rate and proximity ( $\chi^2(1) = 645.45, p < .001$ ), base rate and ratio ( $\chi^2(2) = 424.43, p < .001$ ), and proximity and ratio ( $\chi^2(2) =$

501.50,  $p < .001$ ). Consistent with Study 1, increasing the base rate only mattered when proximity was high ( $b = 1.36$ ,  $t(6466) = 28.31$ ,  $p < .001$ ); increasing the base rate had stronger effects when ratios were 0/N or 1/N compared to Half/N ( $b = 1.28$ ,  $t(6466) = 21.70$ ,  $p < .001$  for 0/N vs. Half/N, and  $b = 0.79$ ,  $t(6466) = 13.43$  for 1/N vs. Half/N); and closer proximity also had stronger effects when ratios were 0/N or 1/N ( $b = 1.30$ ,  $t(6466) = 22.16$ ,  $p < .001$  for 0/N vs. Half/N, and  $b = 0.36$ ,  $t(6466) = 6.19$ ,  $p < .001$  for 1/N vs. Half/N). In other words, when the ratio of users in the POI's connections was very high (Half/N), other graph features had less influence on judgments.

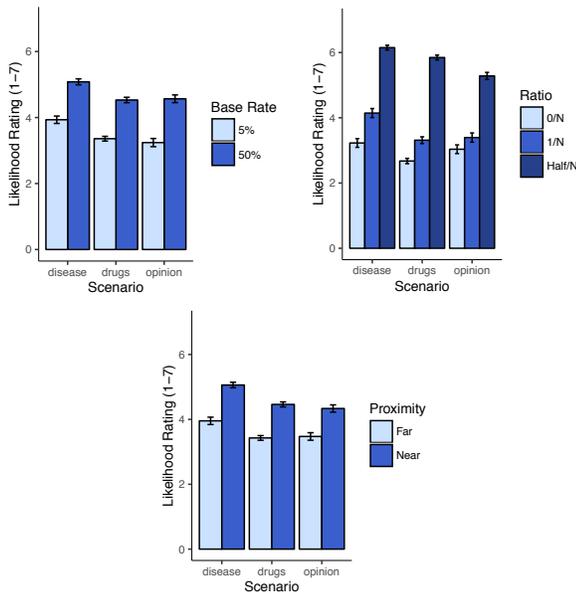


Figure 4: Main effects of data scenario and graph features.

As depicted in Figure 4, main effects of graph features were found for each scenario we tested. The main interaction we were interested in was between scenario and base rate, since we expected that base rates would be utilized more when the scenario described more contagious activities, such as the spread of disease. There was not a significant interaction between scenario and base rate, although significant interactions did emerge between scenario and the proximity and ratio graph features ( $\chi^2(2) = 6.31$ ,  $p = .04$  and  $\chi^2(4) = 40.21$ ,  $p < .001$ , respectively). However, Figure 4 suggests that differences were small and, overall, graph features affected judgments in a consistent way across scenarios.

**Strategy Use.** Consistent with Study 1, participants' ratings of strategy importance reflected a tendency to rely on the ratio strategy more than other strategies. A multilevel model with strategy treated as a random effect and nested within participant revealed that participants rated proximity and ratio information as more important to their judgment than base rate information ( $b = 0.24$ ,  $t(368) = 3.43$ ,  $p < .001$  and  $b = 0.78$ ,  $t(368) = 11.32$ ,  $p < .001$ , respectively). Participants

rated ratio information as more important than proximity information,  $b = 0.54$ ,  $t(368) = 7.89$ ,  $p < .001$ .

**Decision Tree.** Given the multiple interactions between graph features, we modeled the influence of graph features on likelihood judgments with a decision tree using the rpart package in R. Human judgment is often based on 'fast and frugal' heuristics (Gigerenzer, & Goldstein, 1996), which can be modeled by decision trees. Decision trees identify a series of binary decisions to maximize prediction accuracy of an outcome variable.

The model in Figure 5 shows the decision tree for the Study 2 data. The first decision point splits the data based on ratio. If the ratio was high (many salient nodes in the POI's direct connections), the model estimated that the likelihood rating was high (5.8), and no other variables were considered (far right side of Figure 5).

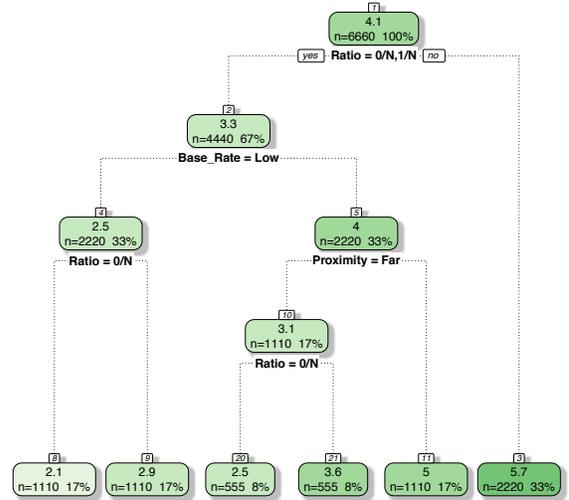


Figure 5: Decision tree predicting likelihood of investigating the POI.  $n$  represents number of judgments collected from participants.

The left side of the first decision point indicates that when there were zero or one salient node(s) in the POI's direct connections, the model next split the data based on the base rate. If the base rate was low (left side of Figure 5), the model then made another split based on ratio, and proximity was not used. If the base rate was high, the model then used proximity. If proximity was far, the model used ratio information again to make the final classification.

## General Discussion

In this research we identified three graph features that influenced judgments about a specific person of interest (POI) in a social network graph: 1) number of salient nodes within the POI's direct connections, 2) base rate of salient nodes in the graph, and 3) proximity of salient nodes to the POI. Across two experiments, we demonstrated the influence of these features on probabilistic graph-based

judgments by manipulating their presence in a series of graphs. Additionally, through our analysis of participants' self-reported strategies and strategy importance ratings, we demonstrated that participants consciously use these graph features in their judgments. Participants reported primarily using the ratio strategy to make their judgments, followed by the proximity and base rate strategies. No matter which data scenario was used (e.g., drug use, disease, public opinion), base rate, proximity, and ratio manipulations influenced judgments in similar ways, with only slight differences across scenarios.

### Strategy Use

Participants consistently used the ratio strategy to make likelihood judgments. The three-level manipulation of this variable (0/N, 1/N, Half/N) was the strongest determinant of likelihood judgments and the ratio strategy was the most frequently self-reported strategy. Furthermore, a decision tree analysis suggests that when the number of connections was high, participants made 'fast and frugal' decisions without using the other strategies.

In other cases, participants made use of base rate and layout proximity. The decision tree analysis suggests that each of these graph features was considered by participants when making likelihood judgments, and manipulating each of these variables led to significant main effects on judgments. However, the relative importance of these two strategies is less clear. Participants were more likely to mention using the proximity strategy than the base rate strategy in Study 1, and in Study 2 participants rated the proximity strategy as more important than the base rate strategy. Yet the decision tree analysis revealed that base rate was a decision point before proximity was considered, suggesting that base rate may be more important than proximity to participants' judgments. These results suggest that there may be a disconnect between participants' self-reported strategies and the strategies revealed by their actual judgments. One interpretation of this finding could be that participants were mistaken in how much they considered each graph feature. Future research could further explore the relative importance of these strategies by testing different response formats, scenarios, and graphs.

### Is Use of Proximity an Error?

An important question about these results is whether the use of proximity should be considered a reasoning error. It is true that in force-directed layouts, which was the layout algorithm used to generate the graphs in these studies, layout distance does often have a relationship with graph distance. In other words, the physical distance between two nodes is somewhat related to the number of edges that separate those nodes. However, it could be misleading to rely only on layout proximity to make the kinds of judgments in these studies for two reasons. First, the extent to which proximity provides meaningful information about the relationship between two nodes depends critically on the layout algorithm used for graph construction. Second, in

many cases, users viewing a graph will not have knowledge of the algorithm used to construct the graph (as was the case in this study), leaving the meaningfulness of proximity unclear. Thus, although use of spatial proximity as a factor in judgment is not necessarily wrong in itself, overuse of this heuristic could lead to misinterpretation in some cases. Understanding how novice audiences interpret proximity could help inform the design of layouts and the use of graphs by analysts, particularly for graphs whose layout algorithms are independent of spatial distance.

### Future Directions

The present studies were carried out with participants who lacked a background in SNA. Future studies plan to examine how novices and experts differ in their use of graph information. We also plan to further assess the validity of proximity information by testing different graph layouts and examining correlations between path and spatial distance.

### Acknowledgments

This work was funded by the Department of Defense under Contract No.: H98230-14-D-0037 through the Computer and Analytic Sciences research group.

### References

- Cook, V. J., Sun, S. J., Tapia, J., Muth, S. Q., Anguello, D. F., Lewis, B. L., Rothenberg, R. B., McElroy, P. D., & the Network Analysis Project Team (2007). Transmission network analysis in tuberculosis contact investigations. *The Journal of Infectious Diseases*, 196, 1517-1527.
- Gigerenzer, G. & Goldstein, D.G. (1996). Reasoning the Fast and Frugal Way: Models of Bounded Rationality. *Psychological Review*, 103:4, 650-669.
- Krebs, V. E. (2002). Mapping networks of terrorist cells. *Connections*, 24(3), 43-52.
- McGrath, C., Blythe, J., & Krackhardt, D. (1997). The effect of spatial arrangement on judgments and errors in interpreting graphs. *Social Networks*, 9, 223-242.
- Micallef, L., Dragicevic, P., & Fekete, J.D. (2012). Assessing the effect of visualizations on Bayesian reasoning through crowdsourcing. *IEEE Transactions on Visualization & Computer Graphics*, 18(12), 2536-2545.
- Novick, L. R., & Bassok, M. (2005). Problem solving. In K. J. Holyoak & R. G. Morrison (Eds.), *Cambridge Handbook of Thinking and Reasoning*, pp. 321-349. New York: Cambridge University Press.
- Purchase, H. C., Cohen, R. F., & James, M. I. (1997). An experimental study of the basis for graph drawing algorithms. *Journal of Experimental Algorithmics*, 2(4).
- Purchase, H. C. (2000). Effective information visualization: A study of graph drawing aesthetics and algorithms. *Interacting with Computers*, 13, 147-162.
- Ware, C., Purchase, H., Colpoys, L., & McGill, M. (2002). Cognitive measurements of graph aesthetics. *Information Visualization*, 1, 103-110.