
Active Learning for Hidden Attributes in Networks

Xiao Ran Yan

Department of Computer Science
University of New Mexico
Albuquerque, NM
everyxt@gmail.com

Yao Jia Zhu

Department of Computer Science
University of New Mexico
Albuquerque, NM
zyaojia@cs.unm.edu

Jean-Baptiste Rouquier

Complex Systems Institute Rhône-Alpes
ENS Lyon, France
jrouquie@gmail.com

Cristopher Moore

Department of Computer Science
University of New Mexico
Albuquerque, NM
and the Santa Fe Institute
moore@santafe.edu

Abstract

In many networks, vertices have hidden attributes that are correlated with the network’s topology. For instance, in social networks, people are more likely to be friends if they are demographically similar. In food webs, predators typically eat prey of lower body mass.

We explore a setting in which the network’s topology is known, but these attributes are not. If each vertex can be queried, learning the value of its hidden attributes—but only at some cost—then we need an algorithm which chooses which vertex to query next, in order to learn as much as possible about the attributes of the remaining vertices. We assume that the network is generated by a probabilistic model, but we make no assumptions about the assortativity or disassortativity of the network. We then query the vertex with the largest mutual information between its type and that of the others (a well-known approach in active learning) or with the largest average agreement between two independent samples of the Gibbs distribution which agree on its type.

We test these approaches on two networks with known attributes, the Karate Club network and a food web of species in the Weddell Sea. In several cases, we found that the average agreement algorithm performs better than mutual information, and both perform better than simpler heuristics. The algorithms appear to explore the network intelligently, first querying vertices at the centers of communities, and then vertices along the boundaries between communities.

1 Introduction

Suppose we have a network, represented by a graph $G = (V, E)$ with n vertices. Suppose further that each vertex v has a type $t(v) \in \{1, \dots, k\}$, representing the value of some hidden attribute that takes k different values. We are given the graph G , and our goal is to learn the types $t(v)$. One way we might do this is to assume that G is generated by some probabilistic model, in which its topology is correlated with these types.

The simplest such model, although by no means the only one to which our methods could be applied, is a stochastic block model. Here we assume that each pair of vertices u, v have an edge between them with a probability $p_{t(u), t(v)}$, and that these events are independent. Given an assignment $t :$

$V \rightarrow \{1, \dots, k\}$ of types to vertices, and a $k \times k$ matrix of probabilities p_{ij} , the likelihood of generating G in this ensemble is

$$\begin{aligned} \mathcal{L}(G | t, p) &= \left(\prod_{(u,v) \in E} p_{t(u), t(v)} \right) \left(\prod_{(u,v) \notin E} (1 - p_{t(u), t(v)}) \right) \\ &= \prod_{i,j=1}^k p_{ij}^{e_{ij}} (1 - p_{ij})^{n_i n_j - e_{ij}}, \end{aligned} \quad (1)$$

where $n_i = |\{v \in V : t(v) = i\}|$ is the number of vertices of type i , and $e_{ij} = |\{(u, v) \in E : t(u) = i, t(v) = j\}|$ is the number of edges from vertices of type i to vertices of type j . Note that (1) assumes that edges are directed, and allows self-loops. We can disallow self-loops, or make the edges undirected, by replacing n_i^2 with $n_i(n_i - 1)$ or $\binom{n_i}{2}$ respectively, and/or taking the product over pairs of types i, j with $i \leq j$.

We do not assume that p_{ij} takes one value when $i = j$ and a smaller value when $i \neq j$. In other words, we do not assume an assortative community structure, where vertices are more likely to be connected to other vertices of the same type. Nor do we require that $p_{ij} = p_{ji}$, since the directed nature of the edges may be important. For example, herbivores eat plants, but the reverse is usually not the case.

This kind of stochastic block model is well-known in the sociology and machine learning communities, and has also been used in ecology (e.g. [1, 2, 3, 4, 5]). If we assume a uniform prior on the probabilities p_{ij} , in which they are chosen independently and uniformly from the unit interval, then integrating over them gives a likelihood of G given t ,

$$\mathcal{L}(G | t) = \int \int \int_0^1 dp_{ij} \mathcal{L}(G | t, p) = \prod_{i,j=1}^k \frac{1}{(n_i n_j + 1) \binom{n_i n_j}{e_{ij}}}. \quad (2)$$

Of course, we could easily assume some other prior for the p_{ij} s, such as a beta distribution, and then optimize its parameters, but here we will stick to (2) for its simplicity. If we assume a uniform prior over the assignments t , then Bayes' rule gives them a Gibbs distribution

$$P(t) = P(t | G) \propto \mathcal{L}(G | t). \quad (3)$$

Note that (2) is maximized when, for each pair of types, e_{ij} is close to 0 or to $n_i n_j$. In other words, the most likely assignments are those where, for each i, j , pairs of vertices of types i and j are either mostly connected or mostly unconnected.

An alternate approach is to assume that the p_{ij} take their maximum likelihood values

$$\hat{p}_{ij} = \operatorname{argmax}_p \mathcal{L}(G | t, p) = e_{ij} / n_i n_j,$$

and set $\mathcal{L}(G | t) = \mathcal{L}(G | t, \hat{p})$. This approach was used, for instance, for a hierarchical block model in [6]. When k is fixed and n is large, this will give results similar to (2), since the integral over p is tightly peaked around \hat{p} . However, for finite n it makes more sense, at least to a Bayesian, to integrate over the p_{ij} , since they obey a posterior distribution rather than taking a fixed value. Moreover, averaging over the parameters as in (2) discourages overfitting, since the average likelihood goes down when we increase k and hence the volume of the parameter space. This should allow us to determine k automatically, although in this paper we set k by hand.

We emphasize, however, that the approaches to active learning we discuss below are not tied to this particular type of block model. They can be adapted to a wide range of other probabilistic models in which topology is correlated with hidden attributes of the vertices.

2 Active Learning

In the active learning setting, the algorithm can learn the type of any given vertex, but at a cost—say, by devoting resources in the laboratory or the field. Since these resources are limited, it has to decide which vertex to query. Its goal is to query a small set of vertices, and use their types to make good

guesses about the types of the remaining vertices. We note that active learning has been applied to the problem of learning network topology by Hanneke and Xing [7]. In contrast, here we assume that the network topology is known, but that the types of the vertices are not.

One natural approach (see, e.g., MacKay [8] or Guo and Greiner [9]) is to query the vertex v with the largest mutual information (MI) between its type $t(v)$ and the types $t(G \setminus v)$ of the other vertices. We can write this as the difference between the entropy of $t(G \setminus v)$ and its conditional entropy given $t(v)$,

$$\text{MI}(v) = I(v; G \setminus v) = H(G \setminus v) - H(G \setminus v | v).$$

Here $H(G \setminus v | v)$ is the entropy, averaged over $t(v)$, of the joint distribution of $t(G \setminus v)$ conditioned on $t(v)$. In other words, $\text{MI}(v)$ is the average decrease in the entropy of $t(G \setminus v)$ that will result from learning $t(v)$. Since the mutual information is symmetric, we also have

$$\text{MI}(v) = I(v; G \setminus v) = H(v) - H(v | G \setminus v),$$

where $H(v)$ is the entropy of the marginal distribution of $t(v)$, and $H(v | G \setminus v)$ is the entropy, on average, of the distribution of $t(v)$ conditioned on the types of the other vertices. Thus a good vertex to query is one about which we are quite uncertain, so that $H(v)$ is large—but which is strongly correlated with other vertices, so that $H(v | G \setminus v)$ is small.

We estimate these entropies by sampling from the space of assignments according to the Gibbs distribution. Specifically, we use a single-site heat-bath Markov chain. At each step, it chooses a vertex v uniformly from among the unqueried vertices, and chooses $t(v)$ according to the conditional distribution proportional to $\mathcal{L}(G | t)$, assuming that the types of all other vertices stay fixed. In addition to exploring the space, this allows us to collect a sample of the conditional distribution of the chosen vertex v and its entropy. Since $H(v | G \setminus v)$ is the average of the conditional entropy, and since $H(v)$ is the entropy of the average conditional distribution, we can write

$$I(v; G \setminus v) = - \sum_{i=1}^k \bar{P}_i \ln \bar{P}_i + \overline{\sum_{i=1}^k P_i \ln P_i},$$

where P_i is the probability that $t(v) = i$.

We offer no theoretical guarantees about the mixing time of the heat-bath Markov chain, and it is easy to see that there are families of graphs and values of k for which it grows exponentially with n . For instance, if G is an Erdős-Rényi graph $G(n, 1/2)$ and $k = 2$, it takes $2^{\Omega(n^2)}$ time to switch from a state where most vertices are of type 1 to one where most are of type 2, since the “bottleneck” states where half the vertices are of each type have total probability $2^{-\Omega(n^2)}$. However, for the real-world networks we have tried so far, the Markov chain appears to converge to equilibrium, and give good estimates of $\text{MI}(v)$, in a reasonable amount of time. We also improve our estimates by averaging over many runs, each one starting from an independently random initial state.

To complete the description of the MI active learning algorithm, we say that it is in *stage* j if it has already queried j vertices. In that stage, it estimates $\text{MI}(v)$ for each v , using the Markov chain to sample from the Gibbs distribution conditioned on the types of the vertices queried so far. It then queries the vertex v with the largest MI. We provide it with $t(v)$, and it moves on to the next stage.

Another strategy is to query the vertex that maximizes another quantity, which we call the *average agreement* (AA). Given two type assignments t_1, t_2 , define their *overlap* as the number of vertices on which they agree,

$$|t_1 \cap t_2| = |\{v : t_1(v) = t_2(v)\}|.$$

Since our goal is to label as many vertices correctly as possible, what we would really like to maximize is the overlap between an assignment t_1 , drawn from the Gibbs distribution, and the correct assignment t_2 . But since we don’t know t_2 , the best we can do is assume that it is drawn from the Gibbs distribution as well. If we think of (t_1, t_2) as having a joint distribution, then querying v would project onto the part of this distribution where $t_1(v) = t_2(v)$. So, we define $\text{AA}(v)$ as the expected overlap between two assignments t_1, t_2 drawn independently from the Gibbs distribution, conditioned on the event that they agree at v . This gives us the following Jaccard-like quantity:

$$\text{AA}(v) = \frac{\sum_{t_1, t_2: t_1(v)=t_2(v)} P(t_1)P(t_2) |t_1 \cap t_2|}{\sum_{t_1, t_2: t_1(v)=t_2(v)} P(t_1)P(t_2)}. \quad (4)$$

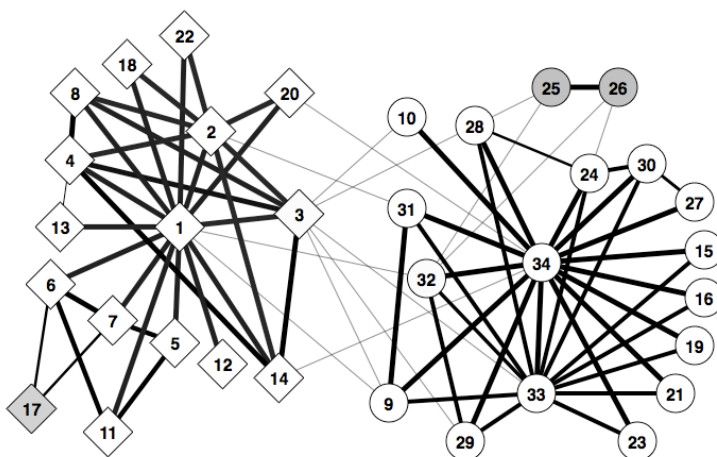


Figure 1: Zachary’s Karate Club. Vertices 1 and 34 are the instructor and president, and their communities are indicated by diamonds and circles. Shaded vertices are more peripheral, and have weaker ties to their communities.

For instance, imagine that $n = 6$ and $k = 2$, and that the types of the vertices are chosen uniformly and independently from $\{1, 2\}$, except that vertices 1, 2, and 3 always have the same type. This gives 16 possible assignments, each of which appears with probability $1/16$. If t_1 and t_2 agree at 1, then they also agree at 2 and 3, and 4, 5, and 6 are each in $t_1 \cap t_2$ with probability $1/2$. So, $AA(1) = 3 + 3/2 = 9/2$. On the other hand, if $t_1 \cap t_2$ agree at 6, then each of the other 5 vertices are in $t_1 \cap t_2$ with probability $1/2$, so $AA(6) = 1 + 5/2 = 7/2$. Thus we should query one of the first three vertices, because doing so will tell us the types of two other vertices as well.

We estimate $AA(v)$ using the same heat-bath Gibbs sampler as for $MI(v)$, except that we draw pairs of assignments (t_1, t_2) independently, by starting the Markov chain at two independently chosen initial states. We then estimate the numerator and denominator of (4) by averaging over these pairs, giving the estimate

$$AA(v)_{\text{est}} = \frac{\sum_{(t_1, t_2)} \delta(t_1(v), t_2(v)) |t_1 \cap t_2|}{\sum_{(t_1, t_2)} \delta(t_1(v), t_2(v))} \quad (5)$$

where $\delta(i, j) = 1$ if $i = j$ and 0 otherwise. We can keep track of these averages for each vertex v as follows: each time we draw a pair (t_1, t_2) , for each $v \in t_1 \cap t_2$, we increment the numerator and the denominator of (5) by $|t_1 \cap t_2|$ and 1 respectively, and for $v \notin t_1 \cap t_2$ we leave the numerator and denominator unchanged. This gives an alternate algorithm for active learning, where in each stage we query the vertex with the largest estimated AA .

We judge the performance of these algorithms by asking, at each stage and for each vertex, with what probability the Gibbs distribution assigns it the correct type. We can then plot, as a function of j , what fraction of the unqueried vertices are assigned the correct type with probability at least q , for various thresholds q , during stage j .

3 Results

We tested the MI and AA algorithms on Zachary’s Karate Club [10], shown in Fig. 1. This is a social network consisting of 34 members of a karate club, where edges represent friendships. The club split into two factions, one centered around the instructor (vertex 1) and the other around the club president (vertex 34), each of which formed their own club. The network is highly associative, with a high density of edges within each faction and a low density of edges between them.

It is no surprise that, after querying just four or five vertices, both algorithms succeed in correctly identifying the types of most of the remaining vertices—i.e., to which faction they belong—with high accuracy. The AA algorithm performs slightly better, achieving an accuracy close to 100%

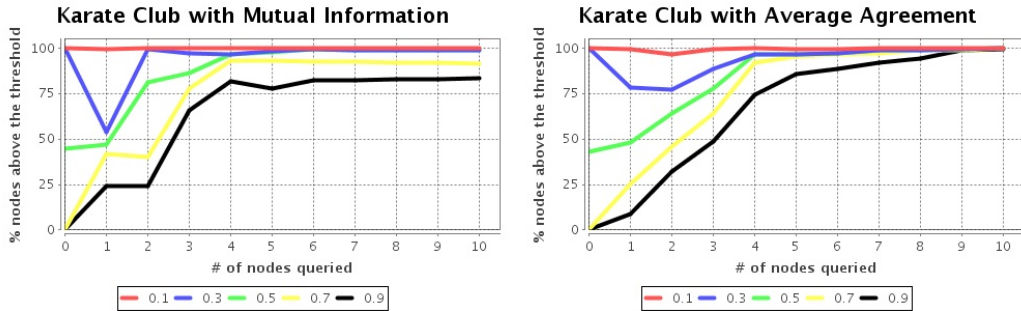


Figure 2: Results of the active learning algorithms on Zachary’s Karate Club network. In each stage we sample the Gibbs distribution using 100 independently chosen initial conditions, doing 2×10^4 steps of the heat-bath Markov chain for each one, and computing averages using the last 10^4 steps. The y axis shows the fraction of vertices, other than those queried so far, which are labeled correctly by the conditional Gibbs distribution with various probabilities. Red, blue, green, yellow, and black lines correspond to probabilities 0.1, 0.3, 0.5, 0.7, and 0.9 respectively. Left, we query the vertex with the largest mutual information (MI) between it and the rest of the network. Right, we query the vertex with the largest average agreement (AA) as defined in the text. After querying 4 or 5 vertices, both methods assign the correct label to about 80% of the remaining vertices with probability 0.9 or greater. The AA algorithm performs somewhat better, with the accuracy quickly converging to 100% as it queries more vertices.

after nine queries. These results are shown in Fig. 2. Of course, this network is quite small, and there are many community structure algorithms that identify the two factions with perfect or near-perfect accuracy (e.g. [11, 12]).

Perhaps more interesting is the order in which these algorithms choose to query the vertices. In Fig. 3, we sort the vertices in order of the average stage at which they are queried. Both algorithms start by querying the two central vertices, the instructor and the president. They then query vertices such as 3, 9, and 10, which lie at the boundary between the two communities. At that point, the algorithms “understand” that the network consists of two assortative communities, and the boundary between the communities is clear. The last vertices to be queried are those such as 2, 4, and 24, which lie deep inside their communities, so that their types are not in doubt. It is not clear why the AA algorithm performs better, but from this small experiment, it seems that it places a lower priority on peripheral vertices, such as 25, than the MI algorithm does.

We also examined a food web of 488 species in the Weddell Sea, in the Antarctic [13, 14]. This data set is very rich, but we focus on two particular attributes—the part of the environment, or habitat, in which the species lives, and its feeding type. In Fig. 4, we show results for the habitat attribute. It gives $k = 5$ types of species, namely pelagic, benthic, benthopelagic, demersal, and land-based. As for the karate club, the AA algorithm performs noticeably better than the MI algorithm. Note that the accuracy is measured as a fraction of the un-queried vertices. It can decrease, for instance, if we query “easy” vertices early on, so that “hard” vertices form a larger fraction of the remaining ones.

Fig. 4 also shows that both algorithms arrive at a stage at which they are either right most of the time, or wrong most of the time, about each of the remaining vertices. For instance, after the AA algorithm has queried 380 species, it labels 70% of the remaining vertices correctly with probability 90%—but it labels the other 30% of the remaining vertices correctly with probability less than 10%. In other words, it has a high degree of certainty about all the vertices, but it is wrong about many of them. Its accuracy improves as it continues to query the vertices, but it doesn’t achieve high accuracy on all the unqueried vertices until there are only 30 of them left. The MI algorithm gets a small fraction of the unqueried vertices wrong up until the very end of the algorithm.

Thus both these algorithms find some vertices easy to classify, but others very hard. Delving into the data, we found that, to a large extent, the blame lies not with the learning algorithms themselves, but with the stochastic block model, and its ability to model the data given this particular attribute. Pelagic, demersal, and land-based species are mostly labeled correctly. But the benthic category, which is the largest and most diverse, includes species with many feeding types and trophic levels.

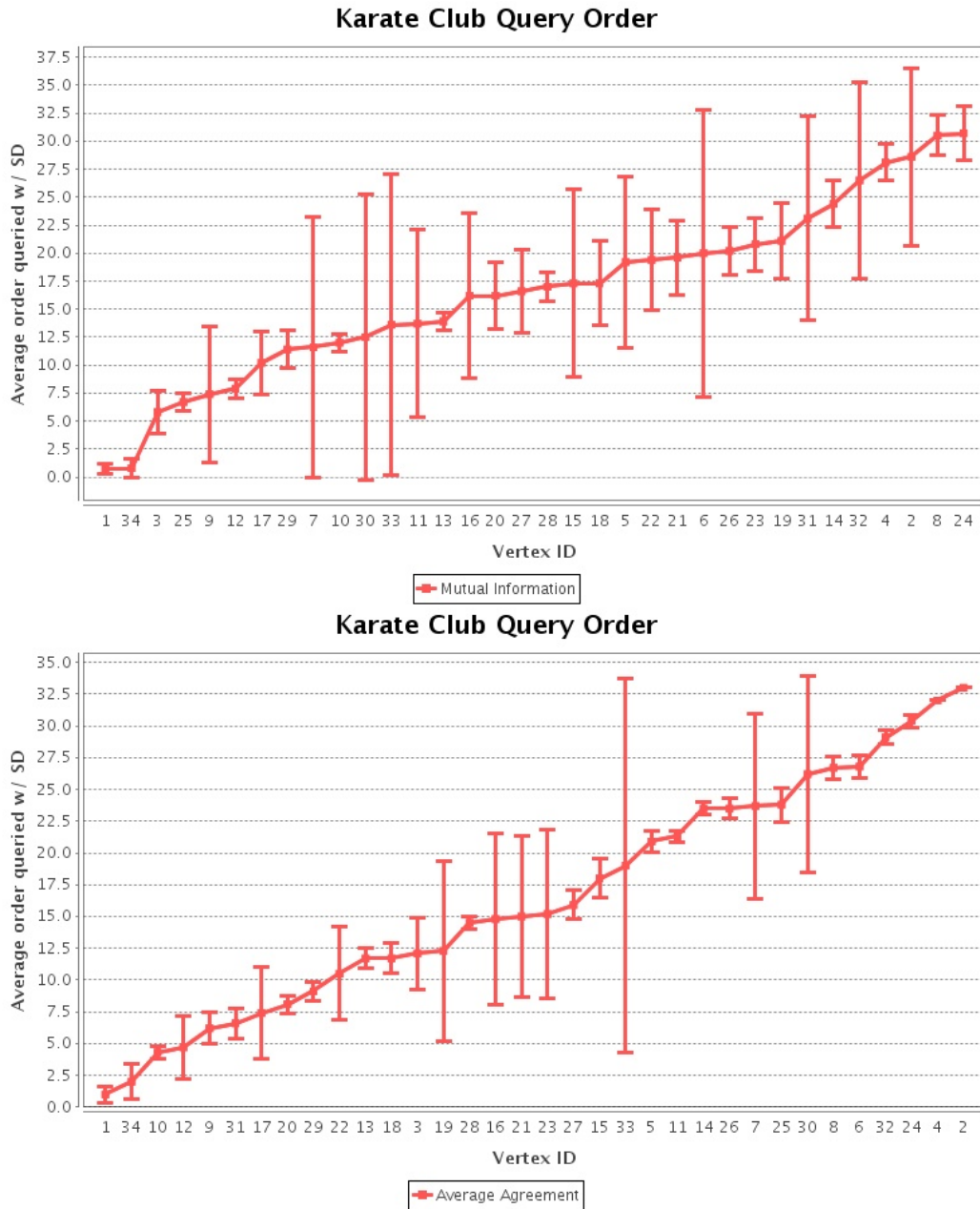


Figure 3: The order in which the active learning algorithms query vertices in Zachary’s Karate Club network, averaged over 10 independent runs of each algorithm. Error bars show the standard deviation. Both algorithms start by querying vertices 1 and 34, which are central to their respective communities, and then query vertices at the boundary between the two communities.

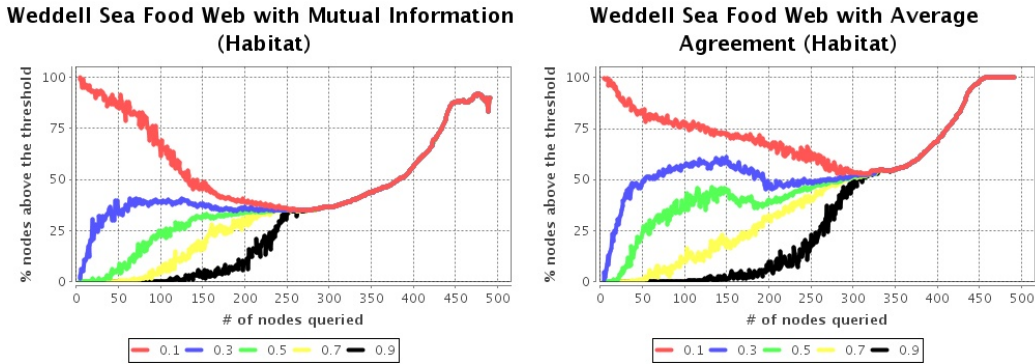


Figure 4: Results for the Weddell Sea food web, averaged over 10 runs of each algorithm. In each stage we sample the Gibbs distribution using 100 independently chosen initial conditions, doing 10^5 steps of the heat-bath Markov chain for each one, and computing averages using the last 5×10^4 steps. The y axis shows the fraction of vertices, other than those queried so far, which are labeled correctly by the conditional Gibbs distribution with various probabilities. Red, blue, green, yellow, and black lines correspond to probabilities 0.1, 0.3, 0.5, 0.7, and 0.9 respectively. As in the karate club, the AA algorithm does better than the MI algorithm. However, after querying about half the species, both algorithms get to a stage where every species is either labeled correctly with high probability, or incorrectly with high probability. In other words, the algorithm is confident, but wrong, about a significant fraction of the species. Most of these are species which are poorly modeled by the stochastic block model using habitat as the only attribute—that is, those which would be misclassified even if you knew the habitats of all the other species.

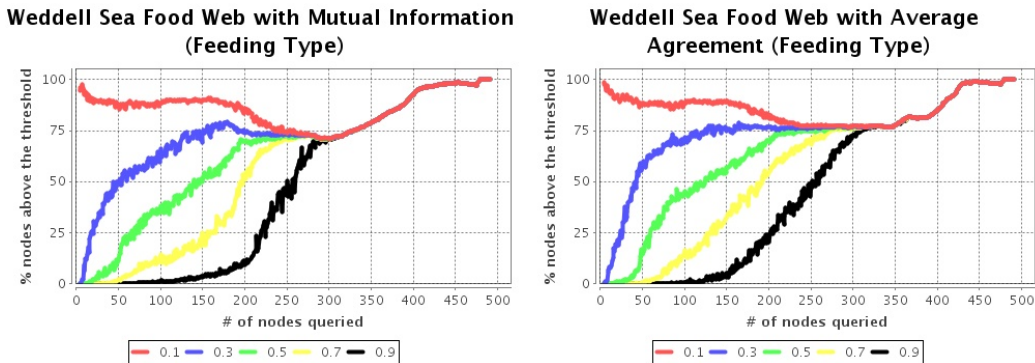


Figure 5: Performance of the learning algorithms using the feeding type of the species, rather than habitat. This attribute is more strongly correlated with the network topology, and the algorithms achieve a higher accuracy with fewer queries than they did for the habitat attribute.

These additional attributes have a large effect on the topology, but they are hidden from the block model in our experiments. As a result, more than half the benthic species are misclassified by the block model, in the sense that their most likely type is pelagic, benthopelagic, demersal, or land-based—even if the types of *all* the other vertices are known.

To draw an analogy, if a member of the karate club was good friends with the instructor, but joined the president’s club because it was close to her favorite café—and if the block model did not have access to this information—we could not expect the learning algorithm to classify her correctly until it got around to querying her. Of course, we can also regard our algorithms’ mistakes as evidence that these habitat types are not cut and dried. Biologists are well aware that there are “connector species” that connect one habitat to another, and belong to some extent to both.

When we use the feeding type attribute instead, both learning algorithms perform significantly better, presumably because this attribute is more strongly correlated with network topology. As Fig. 5

shows, after querying 350 vertices, the MI algorithm correctly labels about 80% of the remaining vertices as herbivorous, carnivorous, omnivorous, detritivorous, or a primary producer, after querying about half the vertices—without, of course, being given the definitions of these terms. The AA algorithm does about as well.

Finally, we note that the type of block model we use here does not deal well with sparse networks, or with heterogeneous degree distributions. In particular, it tends to label high-degree and low-degree vertices as belonging to different types, with higher and lower values of p_{ij} . In future work, we will test the MI and AA algorithms on a degree-corrected block model, such as in [15], where the degrees of the vertices are part of the input, as opposed to data that the model is obliged to explain.

Nevertheless, we feel that our learning algorithms, especially the AA algorithm, offer an interesting approach to dealing with networks where knowledge of vertex attributes is incomplete and costly. We have found that both AA and MI greatly outperform simple heuristics such as querying vertices in random order, or vertices with high degree in the subgraph of unqueried vertices. In future work, we will compare them to other heuristics based on various notions of betweenness and centrality.

Acknowledgments

We are grateful to Joel Bader, Aaron Clauset, Jennifer Dunne, Nathan Eagle, Brian Karrer, and Mark Newman for helpful conversations, and to Ute Jacob for the Weddell Sea food web data. C.M. and J. B. R. are also grateful to the Santa Fe Institute for their hospitality. This work was supported by the McDonnell Foundation.

References

- [1] Yuchung J. Wang and George Y. Wong, “Stochastic Blockmodels for Directed Graphs.” *Journal of the American Statistical Association* **82** (397) 8–19 (1987).
- [2] Martin Rosvall and Carl T. Bergstrom, “An information-theoretic framework for resolving community structure in complex networks.” *Proc. Natl. Acad. Sci. USA* **104** 7327–7331 (2007).
- [3] Edoardo M. Airolidi, David M. Blei, Stephen E. Fienberg, and Eric P. Xing, “Mixed membership stochastic blockmodels.” *Journal of Machine Learning Research* **9** 1981–2014 (2008).
- [4] J. Hofman and C. H. Wiggins, “Bayesian approach to network modularity.” *Physical Review Letters* **100** 258701 (2008).
- [5] Stefano Allesina and Mercedes Pascual, “Food web models: a plea for groups.” *Ecology Letters* **12** 652–662 (2009).
- [6] A. Clauset, C. Moore, and M. E. J. Newman, “Hierarchical structure and the prediction of missing links in networks.” *Nature* **453** 98–101 (2008).
- [7] Steve Hanneke and Eric P. Xing, “Network Completion and Survey Sampling.” Proc. 12th Intl. Conf. on Artificial Intelligence and Statistics (2009).
- [8] David J. C. MacKay, “Information-based objective functions for active data selection.” *Neural Computation* **4** 589603 (1992).
- [9] Yuhong Guo and Russ Greiner, “Optimistic Active Learning using Mutual Information.” Proc. IJCAI 2007.
- [10] Wayne W. Zachary, “An information flow model for conflict and fission in small groups.” *Journal of Anthropological Research* **33** 4520–473 (1977).
- [11] M. A. Porter, J.-P. Onnela, and P. J. Mucha, “Communities in Networks.” *Notices of the American Mathematical Society* **59** (9) 1082–1097 (2009).
- [12] S. Fortunato, “Community detection in graphs.” *Physics Reports*, to appear, [arxiv:0906.0612](https://arxiv.org/abs/0906.0612) (2009).
- [13] Ute Jacob, “Trophic Dynamics of Antarctic Shelf Ecosystems—Food Webs and Energy Flow Budgets.” University of Bremen, Thesis, 2005.
- [14] U. Brose et al., “Empirical consumer-resource body size ratios.” *Ecology* **86** 2545 (2005).
- [15] A. Scott Patterson, Yongjin Park, and Joel S. Bader, “Degree-corrected block models.” Manuscript.