

Ecological interactions and the Netflix problem

Philippe Desjardins-Proulx ^{Corresp., 1, 2}, **Idaline Laigne** ¹, **Timothée Poisot** ³, **Dominique Gravel** ¹

¹ Université de Sherbrooke, Sherbrooke, Quebec, Canada

² Université de Montréal, Montréal, Quebec, Canada

³ Université de Montréal, Montréal, Québec, Canada

Corresponding Author: Philippe Desjardins-Proulx

Email address: philippe.d.proulx@gmail.com

Species interactions are a key component of ecosystems but we generally have an incomplete picture of who-eats-who in a given community. Different techniques have been devised to predict species interactions using theoretical models or abundances. Here, we explore the K nearest neighbour approach, with a special emphasis on recommendation, along with a supervised machine learning technique. Recommenders are algorithms developed for companies like Netflix to predict whether a customer will like a product given the preferences of similar customers. These machine learning techniques are well-suited to study binary ecological interactions since they focus on positive-only data. By removing a prey from a predator, we find that recommenders can guess the missing prey around 50% of the times on the first try, with up to 881 possibilities. Traits do not improve significantly the results for the K nearest neighbour, although a simple test with a supervised learning approach (random forests) show we can predict interactions with high accuracy using only three traits per species. This result shows that binary interactions can be predicted without regard to the ecological community given only three variables: body mass and two variables for the species' phylogeny. These techniques are complementary, as recommenders can predict interactions in the absence of traits, using only information about other species' interactions, while supervised learning algorithms such as random forests base their predictions on traits only but do not exploit other species' interactions. Further work should focus on developing custom similarity measures specialized for ecology to improve the KNN algorithms and using richer data to capture indirect relationships between species.

Ecological Interactions and the Netflix Problem

Philippe Desjardins-Proulx^{0,1,2}, Idaline Laigle¹, Timothée Poisot², and
Dominique Gravel¹

⁰email: philippe.desjardins.proulx@usherbrooke.ca

¹Université de Sherbrooke, Canada.

²Université de Montréal, Canada.

June 27, 2017

0 Abstract

Species interactions are a key component of ecosystems but we generally have an incomplete picture of who-eats-whom in a given community. Different techniques have been devised to predict species interactions using theoretical models or abundances. Here, we explore the K nearest neighbour approach, with a special emphasis on recommendation, along with a supervised machine learning technique. Recommenders are algorithms developed for companies like Netflix to predict whether a customer will like a product given the preferences of similar customers. These machine learning techniques are well-suited to study binary ecological interactions since they focus on positive-only data. By removing a prey from a predator, we find that recommenders can guess the missing prey around 50% of the times on the first try, with up to 881 possibilities. Traits do not improve significantly the results for the K nearest neighbour, although a simple test with a supervised learning approach (random forests) show we can predict interactions with high accuracy using only three traits per species. This result shows that binary interactions can be predicted without regard to the ecological community given only three variables: body mass and two variables for the species' phylogeny. These techniques are complementary, as recommenders can predict interactions in the absence of traits, using only information about other species' interactions, while supervised learning algorithms such as random forests base their predictions on traits only but do not exploit other species' interactions. Further work should focus on developing custom similarity measures specialized for ecology to improve the KNN

27 algorithms and using richer data to capture indirect relationships between species.

28 1 Introduction

29 Species form complex networks of interactions and understanding these interactions is a major
 30 goal of ecology [29]. The problem of predicting whether two species will interact has been
 31 approached from various perspectives [3, 25]. Williams and Martinez [35] for instance built a
 32 simple theoretical model capable of generating binary food webs sharing important features with
 33 real food webs [17], while others have worked to predict interactions from species abundance data
 34 [1, 7] or exploiting food web topology [9, 32]. Being able to predict with high enough accuracy
 35 whether two species will interact given simply two sets of attributes, or the preferences of similar
 36 species, would be of value to conservation and invasion biology, allowing us to build food webs
 37 with partial information about interactions and help us understand cascading effects caused
 38 by perturbations. However, the problem is made difficult by the small number of interactions
 39 relative to non-interactions and relationships that involve more than two species [16].

40 In 2006, Netflix offered a prize to anyone who would improve their recommender system by
 41 more than 10%. It took three years before a team could claim the prize, and the efforts greatly
 42 helped advancing machine learning methods for recommenders [27]. Recommender systems
 43 try to predict the rating a user would give to an item, recommending them items they would
 44 like based on what similar users like [2]. Ecological interactions can also be described this
 45 way: we want to know how much a species would “like” a prey. Interactions are treated as
 46 binary variables, two species interact or they do not, but the same methods could be applied to
 47 interaction matrices with preferences. There are two different ways to see the problem of species
 48 interactions. In the positive-only case, a species has a set of preys, and we want to predict
 49 what other preys they might be interested in. This approach has the benefit of relying only on
 50 our most reliable information: positive (preferably observed) interactions. The other approach
 51 is to see binary interactions as a matrix filled with interactions (1s) and non-interactions (0s).
 52 Here, we want to predict the value of a specific missing entry (is species x_i consuming species
 53 x_j ?). For this paper, we focus on the positive-only approach, which relies on a simple machine
 54 learning approach called the K nearest neighbour.

Method	Input	Prediction
Recommender (KNN)	Set of traits & preys for each species	Recommend new preys
Supervised learning (RF)	Traits (binary and real-valued)	Interaction (1) or non-interaction (0)

Table 1: Summary of the two methods used. The recommender uses the K nearest neighbour algorithm with the Tanimoto distance measure. The Tanimoto KNN makes a recommendation, while supervised learning with random forests (RF) predict either an interaction or a non-interaction.

Statistical machine learning algorithms [27] have proven to be reliable to build effective predictive models for complex data (the “unreasonable effectiveness of data” [19]). The K nearest neighbour (KNN) algorithm is an effective and simple algorithm for recommendation, in this case finding good preys to a species with positive-only information. The technique is simple: for a given species, we find the K most similar species according to some distance measure, and use these K species to base a prediction. If all the K most similar species prey on species x , there is a good chance that our species has interest in x . In our case, similarity is simply computed using traits and known interactions, but more advanced techniques could be used with a larger set of networks. For example, it is possible to learn similarity measures instead of using a fixed scheme [4]. For this study, we use a data-set from Digel et al. [12], which contains 909 species, of which 881 are involved in predator-prey relationships and 871 have at least one prey. The data comes from soil food webs and includes invertebrates, plants, bacteria, and fungi. In total, the data-set has 34 193 interactions. The data was complemented with information on 25 binary attributes (traits) for each species, plus their body mass and information on their phylogeny. We also compare our approach to a supervised learning method, random forests, which is used to predict interactions with only the species’ traits.

A summary of the two methods used can be found in table 1. The approaches are not directly comparable. For example, the positive-only KNN recommends preys to a species. If we remove a prey from a species, ask the algorithm to recommend a prey, and check whether the prey will come up as the recommendation, there are up to 881 possibilities. On the other hand, the random forest predicts either an interaction or a non-interaction, a 50% chance of success by random. These approaches have different uses. Positive-only algorithms are interesting because we are rarely certain that two species do not interact. Also, the KNN approach uses information on what similar species do, while random forests only rely on traits.

We show the *KNN* is particularly effective at retrieving missing interactions in the positive-only case, succeeding 50% of the times at recommending the right species among 881 possibilities. With few traits, the random forests can achieve high accuracy ($\approx 98\%$ for both interactions and non-interactions) without any information about other species in the community. Random forests require only three traits to be effective: body mass and two traits based on the species' phylogeny. Our results show that, with either three traits per species or partial knowledge of the interactions, it is possible to reconstruct a food web accurately. These approaches are complementary, covering both the case where traits are readily available and when only partial knowledge of the food web is known. Both techniques can be used to reconstruct a food web with different types of information.

2 Method

2.1 Data

The first data-set was obtained from the study of Digel et al. [12], who documented the presence and absence of interactions among 881 species from 48 forest soil food webs, details of which are provided in the original publication. 34 193 unique interactions were observed across the 48 food webs, and a total of 215 418 absence of interactions. In order to improve representation of interactions involving low trophic levels species that were not identified at the species level in the first data-set, we compiled a second data-set from a review of the literature. We selected all articles involving interactions of terrestrial invertebrate species for a total of 126 studies, across these, a total of 1 439 interactions were recorded between 648 species. Only 88 absences of interactions were found. We selected traits based on to their potential role in consumption interactions (table 2). For each species or taxa, these traits were documented based on a literature review or from visual assessment of pictures. In addition to these traits, we included two proxies for hard-to-measure traits: feeding guild and taxonomy. The traits were chosen for their potential relevance for species interactions and their availability, see [22] for greater details on the data-set.

Features	Abbr.	Description	<i>n</i>
AboveGroud	<i>AG</i>	Whether the species live above the ground.	538
Annelida	<i>An</i>	For species of the annelida phylum.	34
Arthropoda	<i>Ar</i>	For species of the arthropoda phylum.	813
Bacteria	<i>Bc</i>	For species of the bacteria domain.	1
BelowGround	<i>BG</i>	For species living below the ground.	464
Carnivore	<i>Ca</i>	For species eating other animals.	481
Crawls	<i>Cr</i>	Whether the species crawls.	184
Cyanobacteria	<i>Cy</i>	Member of the cyanobacteria phylum.	1
Detritivore	<i>De</i>	For species eating detritus.	355
Detritus	<i>Ds</i>	Whether the species can be classifying as a detritus.	2
Fungivore	<i>Fg</i>	For species eating fungi.	111
Fungi	<i>Fu</i>	Member of the fungi kingdom.	2
HasShell	<i>HS</i>	Whether the species has a shell.	274
Herbivore	<i>He</i>	For species eating plants.	130
Immobile	<i>Im</i>	For immobile species.	85
IsHard	<i>IH</i>	Whether the species has a tough exterior (but not a shell).	418
Jumps	<i>Ju</i>	Whether the species can jump.	30
LongLegs	<i>LL</i>	For species with long legs.	59
Mollusca	<i>Mo</i>	Member of the mollusca phylum.	45
Nematoda	<i>Ne</i>	Member of the nematoda phylum.	5
Plantae	<i>Pl</i>	Member of the plant kingdom.	3
Protozoa	<i>Pr</i>	Member of the protozoa kingdom.	3
ShortLegs	<i>SL</i>	For species with short legs.	538
UsePoison	<i>UP</i>	Whether the species uses poison.	177
WebBuilder	<i>WB</i>	Whether the species builds webs.	89
Body mass	<i>M</i>	Natural logarithm of the body mass in grams	881
Ph_0	Ph_0	Coordinate on the first axis of a PCA of phylogenetic distances	881
Ph_1	Ph_1	Coordinate on the second axis of a PCA of phylogenetic distances	881

Table 2: The traits used. All traits are binary except for body mass, Ph_0 , and Ph_1 . We use taxonomy as a proxy of latent traits following [26]. To do so, we used the R package *ape* to obtain taxonomic distances between the species, perform classical multidimensional scaling (or principal coordinates analysis) [10] on taxonomic distances, and use the scores of each species on the first two axes (Ph_0 and Ph_1) as taxonomy-based traits. These three real-valued variables are scaled to be in the $[0, 1]$ range. For the Tanimoto similarity index, these three continuous variables have to be converted to binary features. For each, we create four binary features of equal size ($n = 881/4$).

105 2.2 *K*-nearest neighbour

106 Our recommender uses the *K*-nearest neighbour (*KNN*) algorithm [27]. The *KNN* algorithm
 107 is an instance-based method, it does not build a general internal model of the data but instead
 108 bases predictions on the *K* nearest (i.e. most similar) entries given some distance metrics. In
 109 the case of recommendation, each species is described by a set of traits and a set of preys, and
 110 the algorithm will recommend new preys to the species based on the preys of its *K* nearest
 111 neighbours. For example, if *K* = 3, we take the set of preys of the three most similar species
 112 to decide which prey to recommend. If species *A* is found twice and *B* once in the set of preys
 113 of the most similar species, we will recommend *A* first (assuming, of course, that the species
 114 does not already have this prey). See table 3 for a complete example of recommendation. In
 115 the “Netflix” problem, this is equivalent to recommending new TV series/movies to a user by
 116 searching for the users with the most similar taste and using what they liked as recommendation.
 117 It is also possible to tackle the reverse problem: Amazon uses item-based recommendations, in
 118 which case we are looking for similar items instead of similar users to base our recommendations
 119 [2].

120 Choosing the right value for *K* is tricky. Low values give high importance to the most similar
 121 entries, while high values provide a larger set of examples. Fortunately, the most computationally
 122 intensive task is to compute the distances between all pairs, a step that is independent of *K*. As
 123 a consequence, once the distances are computed, we can quickly run the algorithm with different
 124 values of *K*.

125 Different distance measures can be used. We will use the Tanimoto coefficient for recommen-
 126 dations. The Tanimoto (or Jaccard) similarity measure is defined as the size of the intersection
 127 of two sets divided by their union, or:

$$tanimoto(\mathbf{x}, \mathbf{y}) = \frac{|\mathbf{x} \cap \mathbf{y}|}{|\mathbf{x} \cup \mathbf{y}|}, \quad (1)$$

128 Since it is a similarity measure in the $[0, 1]$ range, we can transform it into a distance function
 129 with $1 - tanimoto(\mathbf{x}, \mathbf{y})$. The distance function uses two types of information: the set of traits
 130 of the species (see table 2) and their set of preys. We define the distance function with traits

131 as:

$$\text{tanimoto}_d(\mathbf{x}, \mathbf{y}, w_t) = w_t(1 - \text{tanimoto}(\mathbf{x}_t, \mathbf{y}_t)) + (1 - w_t)(1 - \text{tanimoto}(\mathbf{x}_i, \mathbf{y}_i)), \quad (2)$$

132 where w_t is the weight given to traits, \mathbf{x}_t and \mathbf{y}_t are the sets of traits for species x and y ,
133 and $\mathbf{x}_i, \mathbf{y}_i$ are their sets of preys. Thus, when $w_t = 0$, only interactions are used to compute the
134 distance, and when $w_t = 1$, only traits are used. See table 3 for an example.

135 The data is the set of preys and binary traits for each species (Table 2). To test the approach,
136 we randomly remove an interaction for each species and ask the algorithm to recommend up to
137 10 preys for the species with the missing interaction. Interactions are removed one-at-a-time and
138 similarity is computed before the interaction is removed. The code for computing similarities
139 *after* the interaction is removed is available in the code repository, but it has little effect on
140 the results while making the program much slower to run since the similarity matrix must be
141 computed for each trial. We count how many recommendations are required to retrieve the
142 missing interactions and compute the top1, top5, and top10 success rates, which are defined
143 as the probabilities to retrieve the missing interaction with 1, 5, or 10 recommendations. We
144 repeat this process 10 times for each species with at least 2 preys, totally 7200 attempts. We
145 test all odd values of K from 1 to 19, and $w_t = \{0, 0.2, 0.4, 0.6, 0.8, 1\}$. We also divided species
146 in groups according to the number of preys they have to see if it is easier to find the missing
147 interaction for species with fewer preys.

148 2.3 Supervised learning

149 We also do a simple test with random forests to see if it is possible to predict interactions in this
150 data-set using only the traits [6]. In this case, the random forests perform supervised learning:
151 we are trying to predict y (interaction) from the vector of traits \mathbf{x} by first learning a model on
152 the training set, and testing the learned model on a testing set. We keep 5% of the data for
153 testing. We perform grid search to find the optimal parameters for the random forests.

154 For our predictions, we count the number of true positives (tp), true negatives (tn), false
155 positives (fp) and false negatives (fn). The score for predicting interactions ($Score_y$), non-

Species ID	Traits	Preys	Most similar	Recommendations
0	$\{Ar, Ca\}$	$\{6, 42, 47\}$	$\{6, 28, 70\}$	$[812, 70, 72]$
6	$\{Ar, Ca\}$	$\{42, 47, 70, 72\}$		
28	$\{Ar, Ca\}$	$\{42, 47, 70, 812\}$		
70	$\{Ca\}$	$\{42, 47, 812\}$		
...		

Table 3: Fictional example to illustrate recommendations with K nearest neighbour using the Tanimoto distance measure modified to include species traits. We are trying to recommend a prey to species 0 given that the three most similar species are species 6, 28, and 70. For example, the distance from species 0 to species 70 would be $w_t 0.5 + (1 - w_t) 2/4$. To find recommendations, the set of preys found in the $K = 3$ most similar entries is computed, in this case $\{812 = 2, 70 = 2, 72 = 1\}$, leading to the list of recommendations $[812, 70, 72]$. Because they are found most often in the K most similar species, candidates 812 and 70 will be suggested before 72. To test this approach, we remove a prey from a species and check whether the algorithm recommend the missing prey. Especially with low K , it's possible that no recommendations can be found, for example if the most similar species has the exact same preys.

156 interactions ($Score_{\neg y}$) and the accuracy are defined as

$$Score_y = \frac{tp}{tp + fp}, \quad (3)$$

$$Score_{\neg y} = \frac{tn}{tn + fn}, \quad (4)$$

$$Accuracy = \frac{Score_y 34193 + Score_{\neg y} 741968}{881^2}, \quad (5)$$

157 with 34193 and 741968 being the number of observed interactions and non-interactions in
 158 the 881 by 881 matrix. We then use the True Skill Statistics (TSS) to measure how accurate
 159 the random forest is, defined a

$$TSS = \frac{(tp \times tn) - (fp \times fn)}{(tp + fn)(fp + tn)}. \quad (6)$$

160 The TSS ranges from -1 to 1.

161 2.4 Code and Data

162 Since several machine learning algorithms depends on computing distances (or similarities) for all
 163 pairs, many data structures have been designed to compute them efficiently from kd-trees discov-
 164 ered more than thirty years ago [14] to ball trees, metric skip lists, navigating nets [23], and cover

165 trees [5, 23]. We use an exact but naive approach that works well with small data-sets. Since
 166 $distance(x, y) = 0$ if $x = y$ and $distance(x, y) = distance(y, x)$, our C++ implementation stores
 167 the distances in a lower triangular matrix without the diagonal, yielding $n(n-1)/2$ distances to
 168 compute. A linear scan is then used to find the most similar species. Computing the distance
 169 matrix and testing the predictions 7000 times for a set of parameters takes less than a second.
 170 We used Scikit for random forests [28]. The C++11 code for the KNN algorithm, Python scripts
 171 for random forests, and all data-sets used are available at <https://github.com/PhDP/EcoInter>
 172 (also stored on zenodo with a DOI: [11]).

173 3 Results

174 3.1 Recommendation

175 While matrix imputation has a 50% change of success by random, the Tanimoto KNN needs to
 176 pick the right prey among up to 881 possibilities. Yet, it succeeds on its first recommendation
 177 around 50% of the times. When the first recommendation fails, the next 9 recommendations
 178 only retrieve the right species around 15% of the times so the top5 and top10 success rates are
 179 fairly close to the top1 success rate (see figure 1). The Tanimoto measure is particularly effective
 180 for species with fewer preys, achieving more than 80% success rate for species with 10 or fewer
 181 preys (Figure 2).

182 The highest first-try success rates (the probability to pick the missing interaction on the first
 183 recommendation) are found with $K = 7$ and no weights to traits, and with $K = 17$ and a small
 184 weight of 0.2 to traits (Table 4). Overall, the value of K had little effect on predictive ability.

185 3.2 Supervised learning

186 Random forests predict correctly 99.55% of the non-interactions and 96.81% of the interactions,
 187 for a TSS of 0.96. Much of this accuracy is due to the three real-valued traits (body mass, Ph_0 ,
 188 Ph_1). Without them, too many entries have the same feature vector \mathbf{x} , making it impossible
 189 for the algorithm to classify them correctly. Removing the binary traits has little effect on the
 190 model. With only body mass, Ph_0 , Ph_1 , the TSS of the random forests is 0.94.

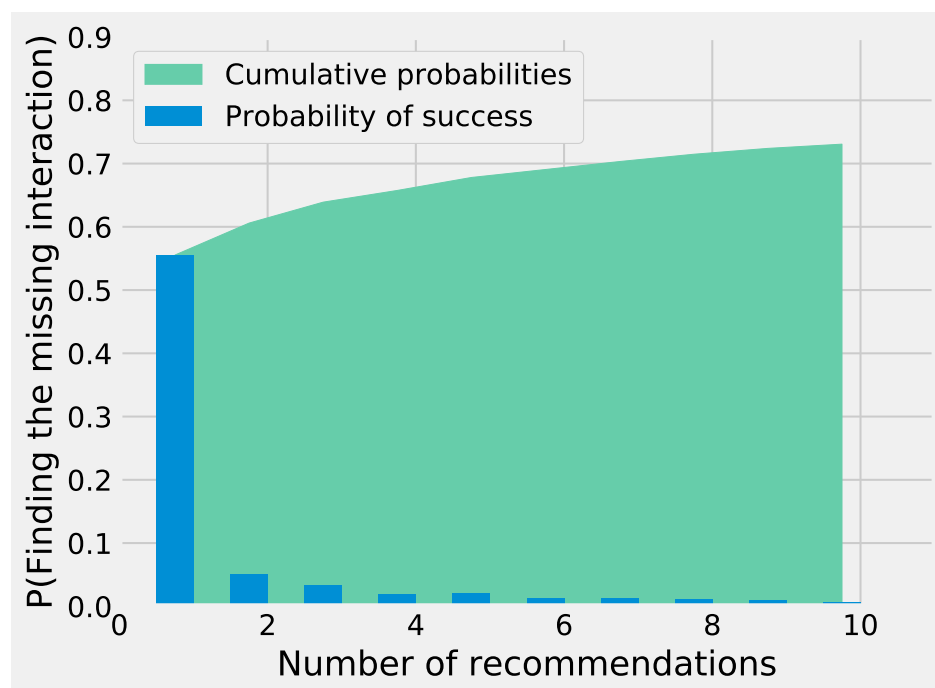


Figure 1: Finding the missing interaction with *KNN*/Tanimoto approach. After removing a prey from a predator, we ask the *KNN* algorithms with Tanimoto measure to make 10 recommendations (from best to worst). The figure shows how many recommendations are required to retrieve the missing interaction. Most retrieved interactions are found with the first attempt. This data was generated with $K = 7$ and $w_t = 0$.

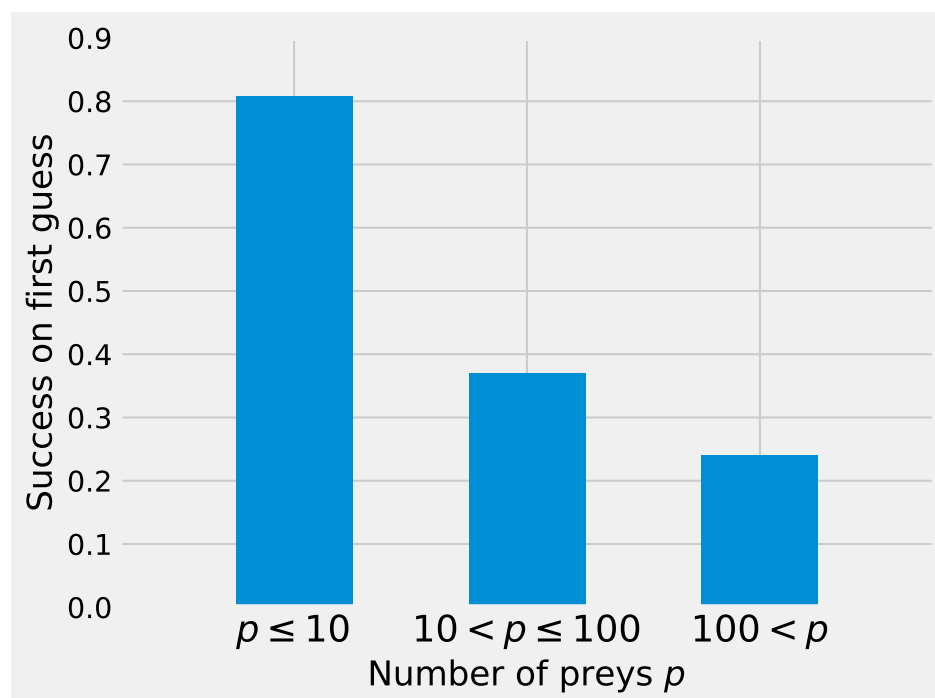


Figure 2: Success on first guess with Tanimoto similarity as a function of the number of prey. The KNN algorithm with Tanimoto similarity is more effective at predicting missing preys when the number of preys is small. This is probably in good part because there are more information available to the algorithm, since 473 species have 10 or fewer preys, 295 have between 10 and 100, 103 species have more than 100 preys.

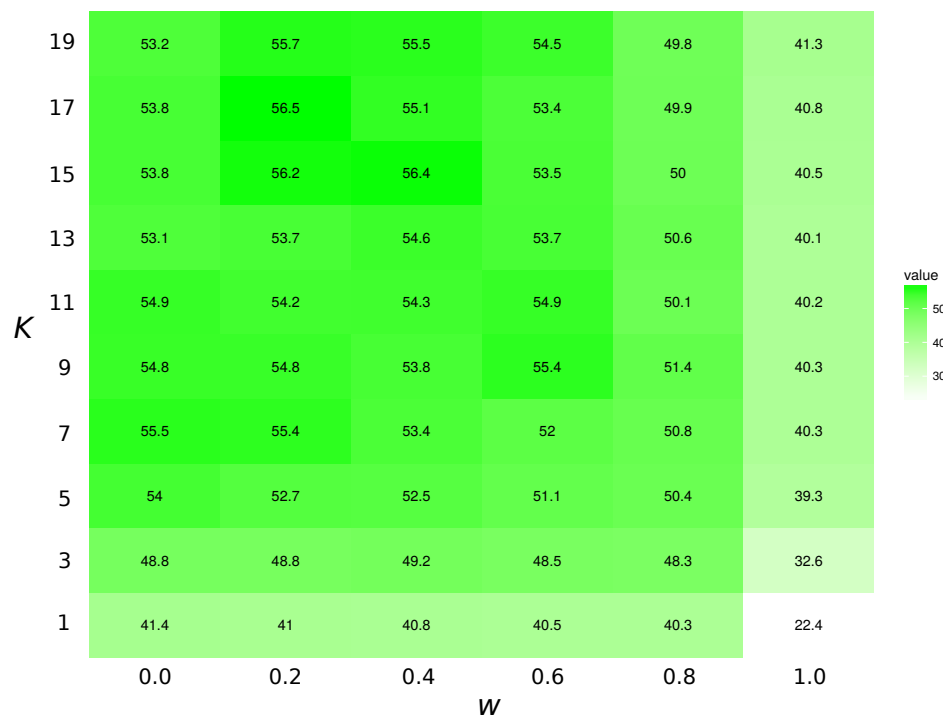


Table 4: Top1 success rates for the K NN/Tanimoto algorithm with various K and weights to traits. When $w_t = 0.0$, the algorithm will only use interactions to compute similarity between species. When $w_t = 1$, the algorithm will only consider the species' traits (see table 2). The value is the probability to retrieve the correct missing interaction with the first recommendation. For each entry, $n = 871$ (the number of species minus 10, the number of species with no preys). The best result is achieved with $K = 17$ and $w = 0.2$, although the results for most values of K and $w = [0.0, 0.2]$ are all fairly close. The success rate increases with K when only traits are considered ($w = 1$).

4 Discussion

We applied different machine learning techniques to the problem of predicting binary species interactions. Recommendation is arguably a better fit for binary species interactions, since it is essentially the same problem commercial recommenders such as Netflix face: given that a user like item i , what is the best way to select other items the user would like? In this case, users are species, and the items are their preys, but the problem is the same. In both cases, we can have solid positive evidence (observed or implied interactions), but rarely have proofs of non-interactions. The approach yields strong results, with a top1 success rate above 50% in a food web with up to 881 possibilities. The approach could be used, for example, to reconstruct entire food webs using global database of interactions [30]. The method's effectiveness rely on nestedness: how much species cluster around the same set of preys in a food web [18]. Thus, it should be less effective in food webs with more unique predators.

The KNN algorithm falls into the realm of unsupervised learning, where the goal is to find patterns in data [27]. The other class of machine learning algorithms, supervised learning, have the clearer goal of predicting a value y from a vector of features \mathbf{x} . For example, in supervised learning, we would try to predict an interaction y from the vector of traits \mathbf{x} , while a unsupervised approach can fill entries from an incomplete matrix regardless of what the entry is (interaction or trait). With a larger set of food webs, it may be possible to use an unsupervised algorithm, for example a neural network, to train a model for matrix imputation. Instead of recommending new preys, imputation would simply fill missing entries from a matrix (interaction or non-interactions).

Our random forests achieve a TSS of 0.96 using the binary traits, body mass, and the coordinates of the multidimensional scaling. This is consistent with previous research that has shown that ecological networks have relatively few dimensions [13]. A random forest can build effective predictive models by creating complex rules based on the traits, while the KNN algorithm relies on a simplistic distance metrics. However, the KNN approach has some advantages over supervised learning, namely the capacity to recommend preys using only the information from the other species' interactions. The solution to improve the KNN approach in ecology is likely to *learn* distance metrics [4] instead of using a fixed formula. This would allow complex

rules while maintaining the KNN 's ability to exploit partial food web structures. Learning distance metrics is a promising avenue to improve our results. Much efforts on the Netflix prize focused on improving similarity measures [33, 21], and custom similarity metrics can be used to improve unsupervised classification algorithms [4] by exploiting complex domain-specific rules. Maybe species with many preys, apex predators, or specialists behave differently enough to need different similarity measures. Learning distance metrics from data is a common way to improve methods based on a nearest neighbour search [37, 4], allowing the measure itself to be optimized. We only used the K nearest neighbour algorithm for unsupervised learning, but several other algorithms can be used to solve the "Netflix problem". For example: techniques based on linear programming, such as recent exact methods for matrix completion based on convex optimization [8] or low-rank matrix factorization. The latter method reduces a matrix to a multiplication between two smaller matrices, which can be used both to predict missing entries and to compress large matrices into small, more manageable matrices [34]. Given enough data, deep learning methods such as deep Boltzmann machines could also be used [38]. Deep learning revolutionized machine learning with neural networks made of layers capable of learning increasingly detailed representations of complex data [20]. Many of the most spectacular successes of machine learning use deep learning [24]. However, learning several neural layers to form a deep networks would require larger data-sets.

The low sensitivity to K in recommendations is interesting and makes the approach easier to use. This is caused by the fact that, as K grows, the set of species includes more and more unrelated species with widely different set of preys. If we increase K from k to $k + \delta$ for a recommendation, the species in δ range are not only less similar, but they are less likely to share preys among themselves. Since recommendations are based on how many times a prey is found in the K nearest species, the species in the δ range are unlikely to have as much weight as the first k species. Our KNN recommender is particularly easy to parametrize since it is neither sensible to the weight given to traits nor to the choice of K .

Our results have two limitations. It is possible that our food web was exceptionally simple, and that a food web with distinct structural properties would behave differently, especially if it has lower nestedness. The success of the KNN algorithms depends on local structure: how much can we learn from similar species. If each species has a unique set of preys, the KNN will

struggle more. Also, a deeper issue is that real food webs are not binary structures. Species, populations, and individuals have different densities, prey more strongly on some resources than others, and have preferences. In a binary matrix, we can predict if two species will interact while completely ignoring the rest of the network, but real food webs involve complex indirect relationships [36]. It is unclear how much we can learn about ecosystems and species interactions from binary matrices, and our results show that binary interactions can be predicted without direct knowledge of the community, since we are able to effectively predict if two species interact given only three traits. Species interactions are better represented with a weighted hypergraph [15], which is well-suited to model relations with an arbitrary number of participants. The hyperedge would allow for complex indirect relationships to be included. Understanding these hypergraphs is outside the scope of the *KNN* algorithm but could be understood with modern techniques such as Markov logic [31].

Recommendation (*KNN* algorithm with Tanimoto distance) and supervised learning (random forests) are complementary techniques. Supervised learning is more useful when we have traits and no information about interactions, but it is useless without the traits. On the other hand, the recommender performs well without traits but requires at least partial information about interactions, although it might be possible to use the interactions from different food webs. We suggest more research could be done on developing better distance metrics for ecological interactions or learning these metrics from data.

5 Acknowledgements

We thank Anna Eklof and one anonymous reviewer for helpful comments.

References

- [1] A Aderhold, D Husmeier, JL Lennon, CM Beale, and VA Smith. Hierarchical bayesian models in ecology: Reconstructing species interaction networks from non-homogeneous species abundance data. *Ecological Informatics*, 11:55–64, 2012.
- [2] CC Aggarwal. *Recommender Systems*. Springer, 2016.

- 276 [3] I Bartomeus, D Gravel, J Tylianakis, M Aizen, I Dickie, and M Bernard-Verdier. A common
277 framework for identifying linkage rules across different types of interactions. *Functional*
278 *Ecology*, 2016.
- 279 [4] A Bellet, A Habrard, and M Sebban. *Metric Learning*. Morgan & Claypool, 2015.
- 280 [5] A Beygelzimer, S Kakade, and J Langford. Cover trees for nearest neighbor. In *Proceedings*
281 *of the 23rd International Conference on Machine Learning*, 2006.
- 282 [6] L Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- 283 [7] EF Canard, N Mouquet, D Mouillot, M Stanko, D Miklisova, and D Gravel. Empirical
284 evaluation of neutral interactions in host-parasite networks. *American Naturalist*9, 183:468–
285 479, 2014.
- 286 [8] EJ Candès and B Recht. Exact matrix completion via convex optimization. *Foundations*
287 *of Computational mathematics*, 9(6):717–772, 2009.
- 288 [9] JE Cohen. *Food webs and niche space*. Princeton University Press, 1978.
- 289 [10] TF Cox and MAA Cox. *Multidimensional Scaling*. Chapman and Hall, 2001.
- 290 [11] P Desjardins-Proulx. <https://github.com/phdp/ecointer>. [http://doi.org/10.5281/](http://doi.org/10.5281/zenodo.161602)
291 [zenodo.161602](http://doi.org/10.5281/zenodo.161602), 2016.
- 292 [12] C Digel, A Curtsdotter, J Riede, B Klarner, and U Brose. Unravelling the complex structure
293 of forest soil food webs: higher omnivory and more trophic levels. In *Oikos*, volume 123,
294 pages 1157–1172, 2014.
- 295 [13] A Eklof, U Jacob, J Kopp, J Bosch, R Castro-Urgal, NP Chacoff, B Dalsgaard, C de Sassi,
296 M Galetti, PR Guimarães, S Beatriz Lomáscolo, AM Martín González, M Aurelio Pizo,
297 R Rader, A Rodrigo, JM Tylianakis, DP Vázquez, and S Allesina. The dimensionality of
298 ecological networks. *Ecology Letters*, 16:577–583, 2013.
- 299 [14] JH Friedman, JL Bentley, and RA Finkel. An algorithm for finding best matches in loga-
300 rithmic expected time. *Transactions on Mathematical Software*, 3(3):209–226, 1977.

- 301 [15] J Gao, Q Zhao, W Ren, A Swami, R Ramanathan, and A Bar-Noy. Dynamic shortest path
302 algorithms for hypergraphs. *Modeling and Optimization in Mobile, Ad Hoc and Wireless
303 Networks*, pages 238–245, 2012.
- 304 [16] AJ Golubski, EE Westlund, J Vandermeer, and M Pascual. Ecological networks over the
305 edge: Hypergraph trait-mediated indirect interaction (tmii) structure. *Trends in Ecology
306 and Evolution*, 31(5):1083–1090, 2016.
- 307 [17] D Gravel, T Poisot, C Albouy, L Velez, and D Mouillot. Inferring food web structure from
308 predator–prey body size relationships. *Methods in Ecology and Evolution*, 4(11):1083–1090,
309 2013.
- 310 [18] PR Guimaraes and P Guimaraes. Improving the analyses of nestedness for large sets of
311 matrices. *Environmental Modelling and Software*, 21(10):1512–1513, 2006.
- 312 [19] A Halevy, P Norvig, and F Pereira. The unreasonable effectiveness of data. *IEEE Intelligent
313 Systems*, 24:8–12, 2009.
- 314 [20] GE Hinton, S Osindero, and YW Teh. A fast learning algorithm for deep belief nets. *Neural
315 computation*, 18(7):1527–1554, 2006.
- 316 [21] T Hong and D Tsamis. Use of KNN for the Netflix Prize. 2006.
- 317 [22] Idaline I Laigle, I Aubin, C Digel, U Brose, I Boulangeat, and D Gravel. Species traits as
318 drivers of food web structure. *In preparation*, 2017.
- 319 [23] M Izbicki and CR Shelton. Faster cover trees. In *Proceedings of the 32nd International
320 Conference on Machine Learning*, 2015.
- 321 [24] V Mnih, K Kavukcuoglu, D Silver, A Graves, I Antonoglou, D Wierstra, and M Riedmiller.
322 Playing atari with deep reinforcement learning. *arXiv*, 2013.
- 323 [25] I Morales-Castilla, MG Matias, D Gravel, and MB. Araújoemail. Inferring biotic interac-
324 tions from proxies. *Ecological Informatics*, 30(6):347–356, 2015.

- 325 [26] N Mouquet, V Devictor, CN Meynard, F Munoz, LF Bersier, J Chave, P Couteron,
326 A Dalecky, C Fontaine, and D Gravel. Ecophylogenetics: advances and perspectives. *Bio-*
327 *logical reviews*, 87(4):769–785, 2012.
- 328 [27] KP Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- 329 [28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blon-
330 del, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau,
331 M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python.
332 *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- 333 [29] SL Pimm. *Food Webs*. Springer, 1982.
- 334 [30] JH Poelen, JD Simons, and CJ Mungall. Global biotic interactions: An open infrastructure
335 to share and analyze species-interaction datasets. *Ecological Informatics*, 24:148–159, 2014.
- 336 [31] M Richardson and P Domingos. Markov logic networks. *Machine Learning*, 62(1-2):107–
337 136, 2006.
- 338 [32] PPA Staniczenko, OT Lewis, NS Jones, and F Reed-Tsochas. Structural dynamics and
339 robustness of food webs. *Ecology Letters*, 13(7):891–899, 2010.
- 340 [33] A Toscher and M Jahrer. The BigChaos solution to the Netflix prize. 2008.
- 341 [34] RJ Vanderbei. *Linear programming: Foundations and extensions*. 2013.
- 342 [35] RJ Williams and ND Martinez. Simple rules yield complex food webs. *Nature*, 404:180–183,
343 2000.
- 344 [36] JT Wootton. The nature and consequences of indirect effects in ecological communities.
345 *Annual Review of Ecology and Systematics*, pages 443–466, 1994.
- 346 [37] EP Xing, AY Ng, MI Jordan, and S Russell. Distance metric learning with application
347 to clustering with side-information. *Advances in neural information processing systems*,
348 15:505–512, 2003.
- 349 [38] J Zhang. Deep transfer learning via restricted boltzmann machine for document classifica-
350 tion. In *ICMLA: Machine Learning and Applications*, volume 1, pages 323–326, 2011.