

1 THINKING OUTSIDE THE BOX -
2 PREDICTING BIOTIC INTERACTIONS IN
3 DATA-POOR ENVIRONMENTS

4 *DAVID BEAUCHESNE*^{1*}, *PHILIPPE DESJARDINS-PROULX*²,
5 *PHILIPPE ARCHAMBAULT*³, and *DOMINIQUE GRAVEL*²

6 * email: david.beauchesne@uqar.ca

7 ¹ *Université du Québec à Rimouski*

8 ² *Université de Sherbrooke*

9 ³ *Université Laval*

10 September 18, 2016

1 Abstract

Large networks of ecological interactions, such as food webs, are complex to characterize, be it empirically or theoretically. The former requires exhaustive observations, while the latter generally requires ample data to be validated. We therefore wondered whether readily available data, namely empirically described interactions in a variety of ecosystems, could be combined to predict species interactions in data deficient ecosystems. To test this, we built a biotic interactions catalogue from a collection of 94 empirical food webs, detailed predator-prey interaction databases and interactions from the Global Biotic Interactions (GloBI) database. We used an unsupervised machine learning method to predict interactions between any given set of taxa, given pairwise taxonomic proximity and known consumer and resource sets found in the interaction catalogue. Initial results suggest that pairwise interactions can be predicted with high accuracy. Although results are seemingly dependent on the comprehensiveness of the catalogue knowledge of taxonomy was found to complement well the catalogue and improve predictions, especially as empirical information available diminished. Given it's high accuracy, this methodology could democratize the use of food webs and network level descriptors in remote location where empirical data is hard to gather. Network characteristics could then be efficiently evaluated and correlated to levels of environmental stressors in order to improve vulnerability assessments of ecosystems to global changes, opening promising avenues for further research and for management initiatives.

2 Introduction

Large networks of ecological interactions, such as food webs, are complex to characterize (Martinez, 1992; Pascual and Dunne, 2006). Empirical descriptions require exhaustive observations, while theoretical inference generally requires ample data to be validated. For this reason, studies focusing on communities of interacting species remain understudied, even though we acknowledge the importance of considering the reticulated nature of complex networks (Ings et al. 2007; Tylianakis et al. 2008). When time is of the essence, the long term studies required quickly become impractical and the use of network level approaches is relegated to the sideline.

Alternatively, a currently evolving approach is to predict interactions using proxies such as functional traits, phylogenies and spatial distributions (e.g. Gravel et al., 2013; Morales-Castilla et al., 2015; Bartomeus et al. 2016). For example, multiple traits can play a significant role in community dynamics and influence the presence and intensity of biotic interactions, like the influence of body size on predator-prey interactions, a literal take on *big fish eats small fish* (Cohen et al., 2003; Brose et al., 2006; Gravel et al. 2013). However, the time required to gather the necessary data to apply those methods may still be restrictive, or the data be unavailable altogether, so much so that other methods have been developed to fill the gaps in knowledge (e.g. Schrodte et al. 2015).

We therefore wondered whether more readily available data could be used to infer interactions in data deficient ecosystems. There is an increasing amount of data describing worldwide species interactions, some freely available through the Global Biotic Interactions (GloBI) database (Poelen et al. 2014). Another readily available piece of information on species is their taxonomy, through initiatives like the World Register of Marine Species (WoRMS; Bailly et al. 2016). More than simple nomenclature, evolutionary processes are thought to influence consumer-resource relationships (Mouquet et al. 2012; Rohr and Bascompte, 2014) so that taxonomically related species would be more likely to share similar types of both consumers and resources (Eklof et al. 2012; Morales-Castilla et al. 2015; Gray et al. 2015). Based on that assumption, taxonomy might be useful in predicting interactions for species lacking detailed information on their biology, but which have a taxonomically related species for which such information is available. The objective of this work is thus to combine empirical biotic interactions originating from a variety of ecosystems with taxonomic relatedness to predict interactions in data deficient ecosystems. As an example, we compare the observed interactions in the Southern Gulf of St Lawrence (Savenkoff et al. 2004) with predictions made using our approach.

3 Methods

The objective of our methodology is to predict the interactions between all pairs of taxa within an arbitrary set N_1 , using a set of taxa N_0 with empirically described interactions from which we can extract pairs of consumers and resources and their taxonomy. We couple the use of empirical data with an unsupervised machine learning method to achieve this.

3.1 Biotic interaction catalogue

We built a biotic interaction catalogue to serve as a set of taxa N_0 for training the algorithm with empirically described interactions. The empirical data used to construct the interaction catalogue was gathered in two successive steps. The first consisted of gathering data from a collection of 94 empirical food webs in marine and coastal ecosystems from which we extracted pairwise taxa interactions (see Brose et al. 2005; Kortsch et al. 2015; GlobalWeb database for more information). We also used a detailed predator-prey interaction database describing trophic relationships between XX predators and their prey (Barnes et al. 2008). From these datasets, only interactions between taxa at the taxonomic scale of the family or higher were selected for inclusion in the catalogue.

As empirical food webs are vastly dominated by non-interactions, these datasets yielded a highly skewed distribution of interactions vs non-interactions. To counterbalance this, the second step of data compilation consisted of extracting observed interactions from the Global Biotic Interaction (GloBI) database (Poelen et al. 2014), which describes binary interactions for a wide range of taxa worldwide. We extracted all interactions available on GloBI for species

94 belonging to the families of taxa identified through step 1. Interactions were
 95 extracted using the rGloBI package in R (Poelen et al. 2015). As per step 1,
 96 only interactions between taxa at the taxonomic scale of the family or higher
 97 were retained

98 The nomenclature used between datasets and food webs varied substantially.
 99 Taxa names thus had to be verified, modified according to the scientific nomen-
 100 clature and validated. This process was performed using the Taxize package
 101 in R (Chamberlain and Szöcs 2013; Chamberlain et al. 2014) and manually
 102 verified for errors. The same package was used to extract the taxonomy of all
 103 taxa for which interactions were obtained in previous steps. The complete R
 104 code and data used for the catalogue is available at https://github.com/david-beauchesne/Interaction_catalog.
 105

106 3.2 Unsupervised machine learning

107 We use the K -nearest neighbor (KNN) algorithm (**ref**) to predict pairwise in-
 108 teractions for a set of taxa S . The KNN algorithm predicts missing entries or
 109 proposes additional entries by a majority vote based on the K nearest (*i.e.* most
 110 similar) entries (see Box 1 for an example). In this case, taxa are described by
 111 a set of resources when considered as a consumer, a set of consumers when
 112 considered as a resource and their taxonomy (*i.e.* kingdom, phylum, class, or-
 113 der, family, genus, species). Similarity between taxa was evaluated using the
 114 Tanimoto similarity measure (**ref**), which compares two vectors with i elements
 115 based on the number of elements they share and contain:

$$tanimoto(\mathbf{x}, \mathbf{y}) = \frac{\sum_i x_i \wedge y_i}{\sum_i x_i \vee y_i}, \quad (1)$$

116 where \wedge is bitwise *and*, while \vee is the bitwise *or* operators. Adding a weigh-
 117 ing scheme, we can measure the similarity using two different sets of vectors
 118 with i and j elements, respectively.

$$tanimoto_t(\mathbf{x}, \mathbf{y}, w_t) = w_t tanimoto(\mathbf{x}_i, \mathbf{y}_i) + (1 - w_t) tanimoto(\mathbf{x}_j, \mathbf{y}_j), \quad (2)$$

119 where w_t is the weight given to vector i , $\mathbf{x}_i, \mathbf{y}_i$ are the resource or consumer
 120 sets of the two taxa and \mathbf{x}_j and \mathbf{y}_j are the vectors for the taxonomy of two taxa.
 121 When $w_t = 0$ only resource or consumer sets are used to compute similarity,
 122 while $w_t = 1$ solely uses taxonomy.

123 3.3 Predicting interactions, Biotic predictor algorithm, Two- 124 way Tanimoto algorithm, Feng shui name algorithm, 125 Find a name for the algorithm

126 The XXX algorithm is built on a series of logical steps that ultimately predicts
 127 a candidate resources list C_R for each taxon in N_1 (Figure 1). For all consumer
 128 taxa T_C in N_1 , the algorithm first verify whether it has empirical resources T_R

129 listed in the catalogue (Step S1, Figure 1). When it does, if T_R are also in N_1 ,
 130 they are added as predicted resources for T_C (S2, S3). This corresponds to what
 131 we refer to as the catalogue contribution to resource predictions. Two taxa in
 132 N_1 that are known to interact through the catalogue are automatically assumed
 133 to interact in N_1 .

134 Otherwise, the algorithm passes to what we refer to as the predictive con-
 135 tribution to resource predictions (S4 to S16), with candidate resources for T_C
 136 identified with the KNN algorithm. If T_R are absent from N_1 , K most similar
 137 resource $T_{R'}$ are identified in N_1 to add to C_R (S4 to S7). Then for all T_C in N_1 ,
 138 the algorithm identifies K most similar consumer $T_{C'}$ in N_0 and extracts their
 139 resource sets (S8). As before, if those resources are found in N_1 (S9) they are
 140 added to C_R (S10 to S12), otherwise K most similar resources $T_{R'}$ are identified
 141 in N_1 (S13) to add to C_R (S14 to S16). A simple working example is presented at
 142 Box 1. Note that other parameters are used in the algorithm, but not presented
 143 here for the sake of message clarity. A more comprehensive mathematical de-
 144 scription of the algorithm and the parameters used is however available through
 145 Figure 1 and the complete R code and data used for the algorithm is available
 146 at https://github.com/david-beauchesne/Predict_interactions.

147 3.4 Algorithm prediction accuracy

148 We used the most extensive and taxonomically detailed datasets included in
 149 the catalogue (**ref**) to assess the prediction accuracy of the algorithm. Testing
 150 accuracy of a particular dataset was done by first removing from the catalogue all
 151 pairwise interacting taxa originating from that dataset. Accuracy was evaluated
 152 using three different statistics:

- 153 1. $Score_y$ is the fraction of interactions correctly predicted:

$$Score_y = \frac{a}{a + c} \quad (3)$$

- 154 2. $Score_{-y}$ is the fraction of non-interactions correctly predicted:

$$Score_{-y} = \frac{d}{b + d} \quad (4)$$

- 155 3. TSS, The True Skilled Statistics (TSS) evaluated prediction success by
 156 considering both true and false predictions, returning a value ranging from
 157 1 (perfect predictions) to -1 (inverted predictions; [Allouche2006]):

$$TSS = \frac{(ad - bc)}{(a + c)(b + d)} \quad (5)$$

158 where a is the number of links predicted and observed, b is the number
 159 predicted but not observed, c is the number of non-interaction predicted but
 160 interactions observed and d is the number of non-interaction predicted absent

and observed. These three statistics give a different perspective on prediction accuracy, focusing in turn on true interactions and non-interactions, and on both true and false predictions.

We evaluated the three statistics for the complete algorithm and for the catalogue and the predictions individually to evaluate their respective contribution to the algorithm predictive accuracy. Multiple w_t values were also tested to evaluate whether taxa similarity measured as a function of resource/consumer sets or taxonomy contributed more significantly towards increased predictive accuracy. The same was done with multiple K values.

Finally, we evaluated the influence of the comprehensiveness of the catalogue on prediction accuracy. We selected the arctic food web from Kortsch et al. (2015) as a test. This food web was selected as it is highly detailed taxonomically and because empirical data remains available for most of its taxa after its exclusion from the catalogue. We iteratively and randomly ($n = 50$ randomizations) removed a percentage of empirical data describing the food web taxa from the catalogue before generating new predictions with the algorithm. We also tested w_t values of 0.5 and 1 to evaluate whether taxonomic similarity could support predictive accuracy in cases when empirical data for species in N_1 in the catalogue is unavailable.

4 Results

4.1 Biotic interaction catalogue

The data compilation process allowed us to build an interaction catalogue composed of 276708 pairwise interactions (interactions = 72110; non-interactions = 204598). A total of 9712 taxa (Superfamily = 15; Family = 591; Subfamily = 29; Tribe = 8; Genus = 1972; Species = 7097) are included in the catalogue, 4159 of which have data as consumers and 4375 as resources.

4.2 Algorithm predictive accuracy

The overall predictive accuracy of the algorithm ranges between 80% to almost 100% in certain cases (Figure 2). Both interactions and non-interactions are well predicted by the algorithm. TSS scores are lower than $Score_y$ and $Score_{-y}$ due to misclassified interactions and non-interactions. This can also be observed through the effect of varying K values, which increases the number of potential candidate resources for each taxa in the predictive portion of the algorithm. Prediction accuracy increases for interactions, while it decreases for non-interactions, as K values increase.

Similarity being predominantly measured with resource/consumer sets (w_t closer to 0) yielded better predictions than when measured with taxonomy (w_t closer to 1; Figure 2). Resource/consumer sets therefore appears to serve as a better predictor of similarity between taxa for interactions predictions. It is nonetheless interesting to note that although the predictive contribution of the

algorithm decreases as w_t increases, an increased mean and decreased variability values for the TSS and $Score_y$ statistics is also observed (Figure 2)). This suggests that while using taxonomy for similarity measurements yields lower predictive accuracy, it may also complement the catalogue contribution by predicting interactions not captured through empirical data, effectively increasing the predictive accuracy of the complete algorithm.

The partitioning of the catalogue and predictive portions of the algorithm shows that it is dependent on the comprehensiveness of the catalogue for high prediction accuracy (Figures 2, 3). As the amount of empirical data available in the catalogue decreases so does the overall accuracy of the algorithm (Figures 3). The predictive contribution of the algorithm however slows down the decrease in the prediction efficiency of the algorithm. Prediction accuracy still remains around 75% with only 40% of N_1 taxa found in the catalogue (Figures 3). Furthermore, the use of taxonomy for similarity measurements is more efficient as empirical data becomes scarcer and no different than resource/consumer sets for the complete algorithm when ample data is available (Figures 3).

4.3 Southern Gulf of Saint Lawrence

As an example, we used the XXX algorithm to predict interactions in the Southern Gulf of Saint Lawrence (SGSL) in eastern Canada. The empirical data and taxa list come from Savenkoff et al. (2004). They present a list of 29 functional groups for a total of 80 taxa presented at least at taxonomical scale of the family. Other coarser taxa families were not used for this example (see Table S1 in Supplementary information (SI) and Savenkoff et al. 2004 for a complete description of functional groups). As their analysis was performed on the functional groups rather than the taxa themselves, we used the algorithm to predict interactions between all 80 taxa selected. We then aggregated them back to their original functional groups to compare with interactions presented in Savenkoff et al. (2004). In total, there were empirical data available in the catalogue for 78% of SGSL taxa (62/80). The algorithm correctly predicted close to 80% of interactions ($a = 135/170$) and non-interactions ($d = 354/455$) extracted from Savenkoff et al. (2004). It also predicted an additional 101 interactions (c) that were not noted in Savenkoff et al. (2004) and failed to predict 36 observed interactions that were (c), resulting in a TSS score of 0.57. A visual comparison of results obtained from the algorithm with interactions noted in Savenkoff et al. (2004) is available at Figure 4. The network presented is centered on the observed and predicted interactions of the capelin (*Mallotus villosus*) and piscivorous small pelagic feeders (*e.g.* *Scomber scombrus* and *Illex illecebrosus*).

239 5 Discussion

240 5.1 Algorithm accuracy

241 We show that out of the box interaction inference for a set of taxa with incom-
242 plete or unavailable preexisting information can be achieved with high accuracy
243 using a combination of empirical data describing biotic interactions and tax-
244 onomic relatedness. Although the efficiency of the algorithm is dependent on
245 the comprehensiveness of the interactions catalogue, taxonomic proximity acts
246 as a complement to increase the number of observed interactions correctly pre-
247 dicted. Taxonomic proximity also supports the efficiency of the algorithm when
248 catalogue comprehensiveness decreases.

249 5.2 Usefulness of taxonomic relatedness

250 While we found that taxonomy could be useful as a complement to predictions
251 made using empirical data, the accuracy of predictions made using the KNN al-
252 gorithm could be improved. Other uses of this machine learning approach have
253 achieved much higher prediction rates (*e.g.* **refs**), which suggests that taxonomy
254 may not be the optimal proxy for predicting interactions. While evolutionary
255 history plays a significant role in influencing consumer-resource trait matching
256 and food web structure (Mouquet et al. 2012; Rohr and Bascompte, 2014), phy-
257 logenetic constraints do not account efficiently for certain traits such as body
258 size (Eklöf and Stouffer 2016). Including traits like body size and metabolism
259 as an additional component of this algorithm could thus help increasing overall
260 prediction accuracy, especially in cases where the catalogue lacks data on taxa
261 for which interactions have to be predicted. Although promising, such an ap-
262 proach would undermine the premise under which this method was built and
263 which constitutes its main strength, *i.e.* predicting interactions in data deficient
264 environments using readily available data.

265 5.3 Interactions classification

266 That $Score_y$ and $Score_{\neg y}$ are inversely proportional means that non-interactions
267 are misclassified as interactions in the process of increasing $Score_y$, consequently
268 decreasing $Score_{\neg y}$. This could either stem from the algorithm poorly predict-
269 ing non-interactions or from the empirical data itself. Accuracy evaluation as-
270 sumes that non-interactions from empirical food web are observed data, yet it
271 is usually not the case. Most empirical webs have a strong focus attributed to
272 higher order consumer species and very little attention given to other taxa (**ref**).
273 Furthermore, the methodologies used to obtain consumer-resource data usually
274 relies on gut content analyses, which is efficient at observing interactions, but
275 not so for absence of interactions (**ref**). Misclassified interactions could thus be
276 real, albeit unobserved through empirical data available.

277 5.4 Southern Gulf of Saint Lawrence

278 The St Lawrence example (Figure 4 and SI) provides great material to discuss
279 predictions in greater detail. The algorithm fails to predict 20% of interactions
280 presented in Savenkoff et al. (2004). Interactions that failed to be predicted
281 were mainly centered on invertebrate species (*e.g.* polychaetes and mollusks)
282 and large functional groups described by coarse taxonomic categories (*e.g.* di-
283 atoms) alongside few species in Savenkoff et al (2004) (*e.g.* piscivorous small
284 pelagic feeders; Table S3). As we focused on the taxa at least at the scale of
285 family, it is likely that their functional groups had a broader range of possible
286 interactions included than what the algorithm could predict using only a few
287 taxa. Furthermore, the efficiency of the algorithm greatly depends on the un-
288 derlying empirical data that defines the catalogue. If the empirical data used
289 to build the catalogue focuses on higher order consumers, it should come as no
290 surprise that the algorithm would be afflicted by the same limitations.

291 The algorithm also predicts substantially more interactions than those pre-
292 sented in Savenkoff et al. (2004) (Figure 4; Table S2). The catalogue is not
293 currently built to take into account life stages of species. Considering life stages
294 and the fact that they are not explicitly considered in the catalogue could
295 explain additional interactions that seem suspicious at first, like the surprise
296 amount of additional interactions predicted for small piscivorous pelagic feeders
297 as consumers (Figure 4). Due to the aggregated nature of the SGSL web, we
298 believe the TSS score to be an underestimate of the efficiency of the algorithm.

299 5.5 Perspectives

300 Overall, we believe that the methods performs well and offers promising av-
301 enues for further applied research and management initiatives. Interaction
302 strength and species co-occurrence are major attributes affecting the probability
303 of observing interactions. Interaction strength is instrumental to understand-
304 ing community dynamics, stability and robustness (Laska and Wootton 1998;
305 Morales-Castilla et al. 2015), while the co-occurrence of species affects com-
306 munity assembly and is a pre-requisite for any given interaction to be observed
307 (Cazelles et al. 2016). Considering them in our methodology would be highly
308 valuable to correctly assess interactions in a given ecosystem and predict the
309 spatial distribution of interaction networks. Given its high efficiency and sim-
310 plicity, our methodology could broaden the use and the accessibility of food
311 webs and network level descriptors for integrative management initiatives such
312 as cumulative impacts assessments and systematic planning (Giakoumi et al.
313 2016; Beauchesne et al. 2016), especially for remote locations where empirical
314 data is hard to gather. Network characteristics could be efficiently evaluated
315 and correlated to levels of multiple environmental stressors to assess the vul-
316 nerability of ecosystems to global changes. We believe that the development of
317 such predictive approaches could represent the first much needed steps towards
318 the use of ecological networks in systematic impacts assessments.

319 **6 Acknowledgements**

320 We thank the Fond de Recherche Québécois Nature et Technologie (FRQNT)
321 and the Natural Science and Engineering Council of Canada (CRSNG) for fi-
322 nancial support. This project is also supported by Québec Océan, the Quebec
323 Centre for Biodiversity Science (QCBS), and the Notre Golfe and CHONeII net-
324 works. We also wish to thank K. Cazelles for the help, constructive comments
325 and suggestions.

DRAFT

6.1 Box 1

The XXX algorithm follows a series of logical steps to predict resources for all taxa in an arbitrary set of taxa N_1 using a set of taxa N_0 with empirically described interactions from which we can extract sets of consumers and resources and their taxonomy. In this example, we are predicting interactions for a fictitious $N_1 = \{T_1, T_9, T_{10}, T_{11}, T_{12}\}$ using N_0 with information on 12 taxa. This catalogue holds information on consumer or resource for 10 taxa and the taxonomy for all 12 taxa in the list.

N_0 taxa ID	taxonomy	resource	consumer
T_1	$\{a, b, c\}$	$\{T_2, T_3, T_{12}\}$	$\{T_4\}$
T_2	$\{e, f, g\}$		$\{T_1, T_5\}$
T_3	$\{i, j, k\}$		$\{T_5\}$
T_4	$\{m, n, o\}$	$\{T_1, T_5\}$	
T_5	$\{a, b, d\}$	$\{T_8, T_9\}$	$\{T_4\}$
T_6	$\{i, q, r\}$	$\{T_2, T_8\}$	$\{T_4\}$
T_7	$\{e, f, h\}$		$\{T_1, T_6\}$
T_8	$\{s, t, u\}$		$\{T_5, T_6\}$
T_9	$\{s, t, v\}$		$\{T_5\}$
T_{10}	$\{i, j, l\}$		
T_{11}	$\{m, n, p\}$		
T_{12}	$\{q, r, s\}$		$\{T_1\}$

Similarity between all pairs of taxa in N_0 is measured for consumer, resource and taxonomic proximity using equation 1. The upper triangular matrix represents similarity measured with taxa sets of resources/consumers, while the lower triangular represents taxonomic similarities. For consumer/resource set similarities, values of 0 mean that similarity equals 0 for both similarity measurements.

$$\text{tanimoto}(T_Cx, T_Cy) / \text{tanimoto}(T_Rx, T_Ry)$$

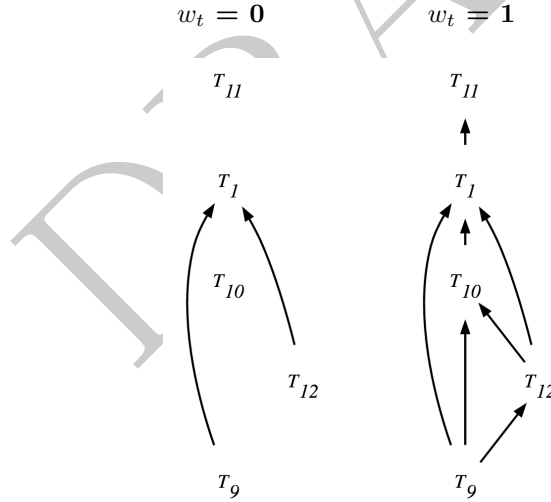
	T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8	T_9	T_{10}	T_{11}	T_{12}
T_1	-	0	0	0	0/1	0.3/1	0	0	0	0	0	0
T_2	0	-	0/0.5	0	0	0	0/0.3	0/0.3	0/0.5	0	0	0/0.5
T_3	0	0	-	0	0	0	0	0/0.5	0/1	0	0	0
T_4	0	0	0	-	0	0	0	0	0	0	0	0
T_5	0.5	0	0	0	-	0.3/1	0	0	0	0	0	0
T_6	0	0	0.2	0	0	-	0	0	0	0	0	0
T_7	0	0.5	0	0	0	0	-	0/0.3	0	0	0	0/0.5
T_8	0	0	0	0	0	0	0	-	0	0	0	0
T_9	0	0	0	0	0	0	0	0.5	-	0	0	0
T_{10}	0	0	0.5	0	0	0.2	0	0	0	-	0	0
T_{11}	0	0	0	0.5	0	0	0	0	0	0	-	0
T_{12}	0	0	0	0	0	0.5	0	0.2	0.2	0	0	-

$$\text{tanimoto}(T_Tx, T_Ty)$$

340 From these, the algorithm goes through logical steps (Figure 1) to identify
 341 a candidate resource list C_R for each taxon in N_1 using either empirical data
 342 directly or K most similar taxa with equation 2. Going through the process for
 343 T_1 , using $K = 1$ and $w_t = 1$:

Steps		Catalogue	Prediction
1	$I(T_1, T_R)$ in N_0 ?		
2	T_R in N_1 ?		
4-7	$T_2 = \text{no} \rightarrow t(T_2, T_{R'}, w_t) = \text{NA}$	$\{\}$	$\{\}$
4-7	$T_3 = \text{no} \rightarrow t(T_2, T_{R'}, w_t) = T_{10} = 0.5$	$\{\}$	$\{T_{10}\}$
3	$T_{12} = \text{yes}$	$\{T_{12}\}$	$\{T_{10}\}$
8	$t(T_1, T_{C'}, w_t) = T_5 = 0.5$		
9	$I(T_5, T_R)$ in N_1 ?		
13-16	$T_8 = \text{no} \rightarrow t(T_8, T_{R'}, w_t) = T_9 = 0.5$	$\{T_{12}\}$	$\{T_9, T_{10}\}$
10-12	$T_9 = \text{yes}$	$\{T_9, T_{12}\}$	$\{T_9, T_{10}\}$

344 The logical steps allow us to predict a set of resources for $T_1 = \{T_9, T_{10},$
 345 $T_{12}\}$. Doing it for all taxa in N_1 with $w_t = 0$ and 1 predicts the following
 346 networks:



6.2 Figures

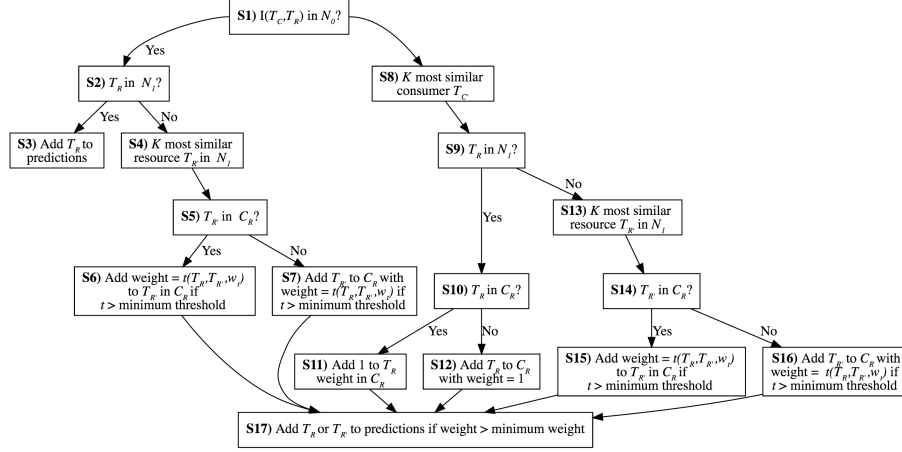


Figure 1: Description of the logical steps used by the algorithm to suggest a list of candidate resources (C_R) for each consumer tax (T_C) in an arbitrary set of N_1 for which interactions are predicted, using a set of taxa N_0 with empirically described interactions. Interactions between consumer and resource taxa are denoted as $I(T_C, T_R)$. K is the number of most similar neighbours selected for the KNN algorithm, t stands for tanimoto in equation 1, w_t is the weight given to sets of resources and consumers in equation 2, the minimum threshold is an arbitrary value setting the minimal similarity value accepted for taxa to be considered as close neighbours in the KNN algorithm, the weight is the value added to a candidate resource each time it is added to C_R and the minimum weight is the minimal weight value accepted for candidate resources to be selected as predicted sources in the algorithm.

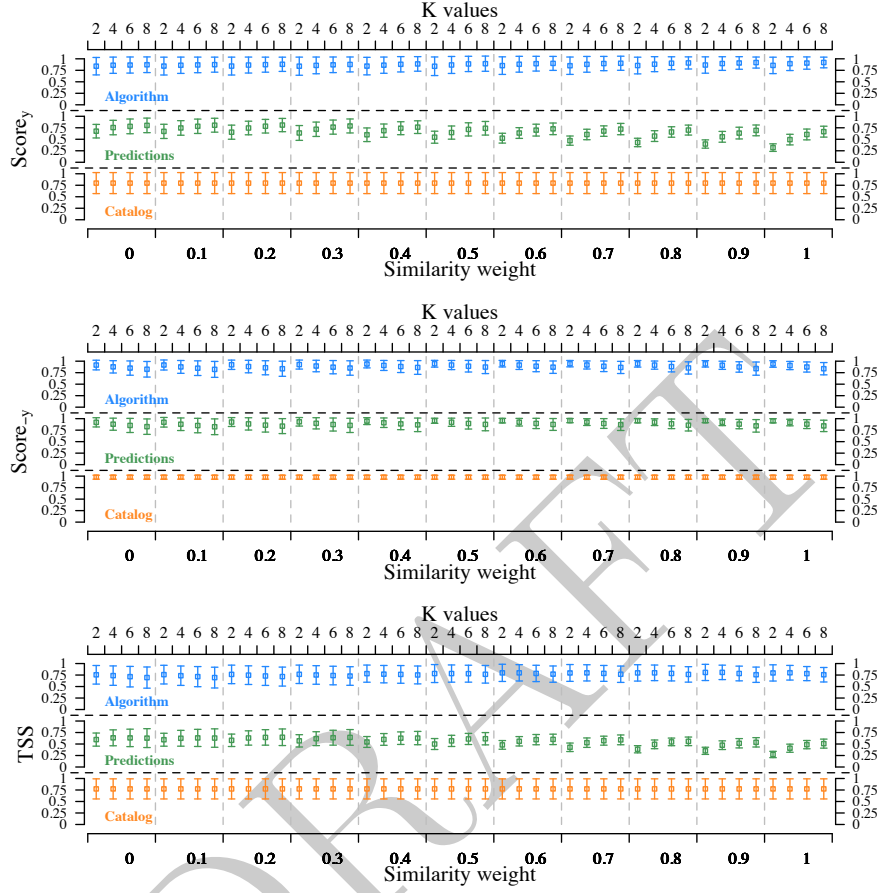


Figure 2: The graph presents the three statistics as a function of trait weight, which varies between 0 and 1. A weight of 0 means that similarity is measured only using set of resources for each taxa, while a weight equal to 1 means that similarity is based solely on taxonomy. We present 6 food webs with over 50 taxa each and the Barnes et al. (2008) dataset.

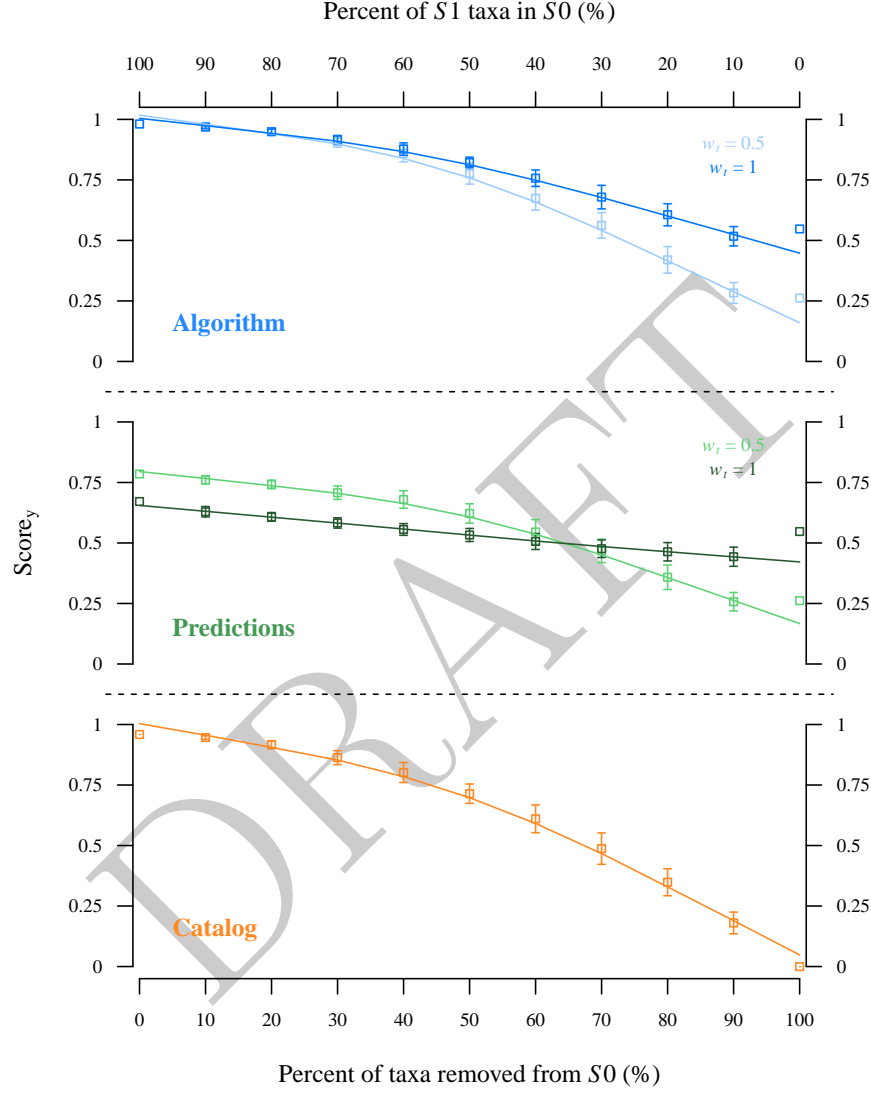


Figure 3: Graph presenting predictive accuracy as a function of the amount of information available in the catalogue. The arctic food web from Kortsch et al. (2015) was used for this, as it is highly detailed and because almost all taxa found in it had information in the catalogue even when not included in the catalogue. A random percentage of taxa in the web was iteratively removed from the catalogue ($n = 50$) before predicting interactions with the XXX algorithm.

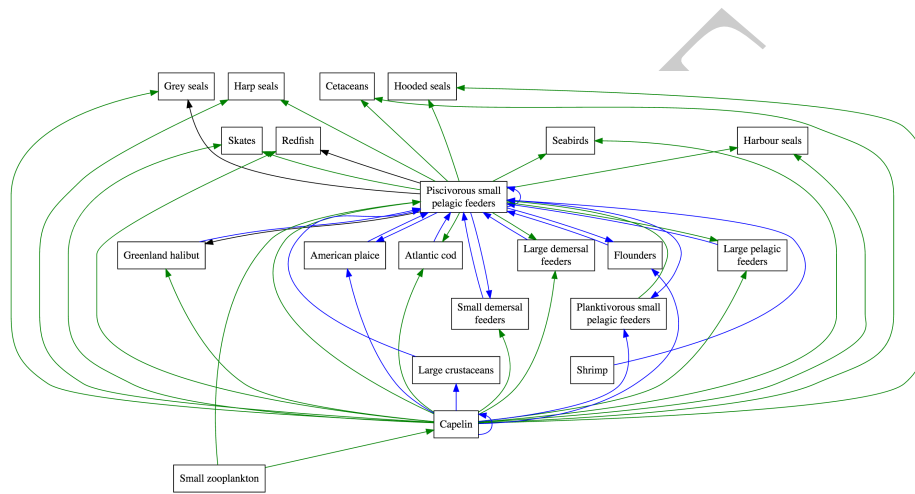


Figure 4: Example of results from the algorithm with the Network of the Southern Gulf of Saint Lawrence (Savenkoff et al. 2004) centered on interactions of the capelin (*Mallotus villosus*) and piscivorous small pelagic feeders (e.g. *Scomber scombrus* and *Illex illecebrosus*). Edge with colors green were both predicted and observed (26), black were observed only (3) and blue were predicted only (19).