

Linear Filtering as an imputation method for Singular Value Decomposition inference of host-virus associations

Marie-Andrée Ouellet¹ Gregory Albery² Dan Becker³ Colin J. Carlson^{4,5} Timothée Poisot^{1,6,*}

¹ Département de Sciences Biologiques, Université de Montréal ² Department of Biology, Georgetown University ³ Department of Biology, University of Oklahoma ⁴ Center for Global Health Science and Security, Georgetown University Medical Center ⁵ Department of Microbiology and Immunology, Georgetown University Medical Center ⁶ Québec Centre for Biodiversity Sciences

* timothee.poisot@umontreal.ca

Abstract: The current pandemic of SARS-CoV-2 is a stark reminder that we need a better understanding of the movements of viruses through novel animal hosts, and ultimately to humans. The task of predicting which virus can infect which host, and where spillovers are likely to happen, still remains difficult. Typically, anticipatory approaches can be limited by numerous difficulties (lack of suitable data, disagreement between models, etc.), and would therefore benefit from adding methods allowing imputation and producing results that could easily be added to ensemble models. In this study, we explore the potential of using the Singular Value Decomposition (SVD) technique as an imputation method to predict host-virus interactions.

1 **TK** rework this paragraph: need to predict host-virus associations

2 The current pandemic of SARS-CoV-2 is a stark reminder that movement of viruses through
3 novel animal hosts, and ultimately to human through zoonotic spillovers (Plowright et al. 2017),
4 requires that we understand the complexity of our biological surroundings. Indeed, the fact that
5 the majority of emerging infectious diseases are caused by zoonotic pathogens from wildlife
6 sources (Jones et al. 2008) gives some urgency to the task of predicting which viruses can be
7 found in which hosts, so as to provide guidance on where and what species to sample and where
8 spillovers are likely to happen (Johnson et al. 2020; Albery et al. 2020).

9 As seen with SARS-CoV and MERS-CoV epidemics, novel human infections by viruses are
10 representing a serious threat to global public health, and being able to prevent future viral emer-
11 gence now appears as a fundamental tool among our society. Zoonotic dynamics usually involve
12 three main stages: transmission within the animal reservoir, cross-species spillover and trans-
13 mission to human, and finally, transmission among humans (Lloyd-Smith et al. 2009). In the
14 past decades, substantial research effort has been put in studying and predicting dynamics at the
15 animal-human interface, but tracing back the ultimate origin of novel zoonotic viruses remains
16 a major difficulty (Becker et al. 2020). Also, the main strategy adopted so far against infectious
17 diseases consists in taking actions after the emergence by increasing the health infrastructures
18 and vigilance, as well as developing vaccines or medical treatments (Han and Drake 2016).

19 As suggested by Han and Drake (2016), a more efficient approach would be anticipatory. Yet
20 an anticipatory approach can be limited by lack of suitable data, and as Becker et al. (2020)
21 highlighted, by disagreement between models. The task of predicting possible host-virus inter-
22 actions would therefore benefit from adding methods that allow imputation, and can produce
23 results that are easily added to ensemble models. Here, we explore an approach focusing on
24 the first stage of zoonoses dynamics, by using the Singular Value Decomposition (SVD) as an
25 imputation method for identifying unobserved host-virus interactions, acting as potential inter-
26 mediate hosts in diseases transmissions.

27 **TK** SVD is a way to do link prediction in the absence of external information, but we can rely
28 on info contained in the network itself

29 **TK** main results: optimal rank, number of new associations, top 10 zoonoses

30 **Dataset**

31 **TK** this actually uses CLOVER now

32 We apply SVD imputation to the data on wildlife hosts of beta-coronaviruses collected by
33 Becker et al. (2020). This host-virus network is composed of 710 mammalian hosts (resolved
34 at the species level) and 72 viruses (resolved at the genus level). Full data are available from
35 <https://github.com/viralemergence/virionette/>. While the host-virus interaction have
36 been pulled from published sources, specific attention has been paid to betacoronaviruses, a vi-
37 ral genus at high risk of spillover, and to their potential bat hosts, a mammalian order known
38 to be evolutionary involved in the main viruses zoonotic historical epidemics (Shiple et al.
39 2019; Ren et al. 2006). Data on interactions between these groups were augmented by a Gen-
40 Bank search to retrieve the hosts associated to sequences of betacoronaviruses. Altogether, this
41 dataset represents a total of 1731 unique interactions, and 49389 host-virus pairs for which no
42 interaction were reported; these can be true negatives (the virus is unable to infect the host), or
43 false negatives (the virus can infect the host but the infection has not been reported). This type
44 of problem lends itself well to an approach using a recommender system.

45 **The model**

46 We ran all analyses in *Julia* 1.5.3 (Bezanson et al. 2017), on the *Beluga* supercomputer operated
47 by the Calcul Québec consortium.

48 **Low-rank approximation with Singular Value Decomposition**

49 Singular Value Decomposition (SVD; Gene H. Golub and Reinsch 1971; Forsythe and Moler
50 1967) is a linear algebra technique used to decompose a data matrix in a product of three matri-
51 ces:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (1)$$

52 Where \mathbf{X} is a $m \times n$ data matrix ($m \geq n$), \mathbf{U} is an unitary $m \times m$ matrix containing the left singular
53 vectors, \mathbf{V} is an unitary $n \times n$ matrix containing the right singular vectors and $\mathbf{\Sigma}$ is a diagonal
54 matrix containing the singular values ordered in decreasing order of importance, in regard of
55 the quantity of information that they present. This process allows data reduction by finding key
56 correlations among entries and then by approximating the original matrix.

57 Optimal truncation of the SVD at rank r (Eckart and Young 1936; G. H. Golub, Hoffman, and
58 Stewart 1987) of the singular values will allow data reduction while keeping enough information
59 to obtain a balance between complexity and accuracy within the model. Truncation at rank r was
60 performed by setting values $\Sigma_{(r+1)..m}$ to 0 (we note the resulting vector $^{(r)}\mathbf{\Sigma}$), and the resulting
61 low-rank approximation was obtained by

$$^{(r)}\mathbf{X} = \mathbf{U}^{(r)}\mathbf{\Sigma}\mathbf{V}^T \quad (2)$$

62 We illustrate the process on our dataset in fig. 1. Removing signal from the matrix through a low-
63 rank approximation hinges on the assumption that most data are generated by “low-rank” pro-
64 cesses, whereas the additional ranks would reflect noise or idiosyncracies acting in the dataset.
65 Under this assumption, an imputation method using a low-rank approximation would have a
66 good performance.

67 [Figure 1 about here.]

68 **Model structure**

69 For each non-interaction in the dataset, the model assigns an initial value to it and performs
70 iteratively the SVD at chosen rank, until it reaches convergence. During this step, the cells in the
71 matrix that are *not* being imputed are kept at their actual value. We capped the maximal number
72 of iterations at 50, even though the value of the imputed cells stopped changing (defined as a
73 step-wise change lower than $10 \times \epsilon$) after less than 10 steps in most cases. The initial value that
74 we first picked for this illustration is the connectance of the global host-virus interaction dataset,
75 which amounts to the probability that any pair of organisms are found to interact (0.03). Yet,

76 this can overestimate the importance of viruses with a narrow host range, or underestimate the
77 importance of generalist viruses. For this reason, the assignment of the initial value was then
78 determined based (Stock et al. 2017) work on linear filtering. This method provides a convenient
79 way to assign weights to various aspects of network structure, and has been revealed to provide a
80 good baseline estimate of how likely it is that a missing interaction actually exists, based on the
81 structure of the interaction matrix, without the need of having other side information, such as
82 traits or phylogeny. Considering our $m \times n$ data matrix \mathbf{X} , the initial value of a missing interaction
83 was fixed to the filtered value \mathbf{F}_{ij} :

$$\mathbf{F}_{i,j} = \alpha_1 \mathbf{X}_{i,j} + \alpha_2 \frac{1}{m} \sum_{k=1}^m \mathbf{X}_{k,j} + \alpha_3 \frac{1}{n} \sum_{l=1}^n \mathbf{X}_{i,l} + \alpha_4 \frac{1}{mn} \sum_{k=1}^m \sum_{l=1}^n \mathbf{X}_{k,l} \quad (3)$$

84 where $\sum_{i=1}^4 \alpha_i = 1$ and $\alpha_i \in [0, 1]$.

85 Prediction scoring

86 Using the linear filter allows to explore different hypotheses as to which parts of network struc-
87 ture are important for predictive ability. As we assume that the initial value of 0 in the matrix
88 can be a false positive, we give it no weight in the model $\alpha_1 = 0$. **TK change from here** We
89 then varied the other parameters on a regular grid of 304 points, where the values for α_4 (impact
90 of connectance), α_2 (impact of the number of hosts), and α_3 (impact of the number of viruses)
91 was varied between 0 and 1. We then applied SVD imputation for each of these parameters
92 combinations for ranks 1 to 3.

93 To rank the predictions made by the SVD-imputation, we took the value for every missing in-
94 teraction after imputation, and divided it by the initial value, then subtracted one. This gives
95 an evidence score in \mathbb{R} , which we can transform into a probability in $[0, 1]$ by taking its logistic;
96 therefore, the final probability of an interaction is defined as

$$P(x) = \frac{1}{1 + e^{-x}},$$

97 where x is the evidence for this interaction under our scoring system.

98 **Model tuning and thresholding**

99 One of the challenges associated with link prediction in this dataset is that non-interactions are
100 not necessarily true negatives; most are simply missing data. To reach the best prediction, we
101 need to answer three related questions. First, what model to assign initial values performs best?
102 Second, what rank is sufficient to give the most accurate approximation of the matrix? Finally,
103 what threshold on the interaction probability should be applied to the results of the best model
104 at the appropriate rank?

105 To answer this question, we first ran the LF-SVD imputation on a sample of 768 positive and
106 768 supposed negative interactions, at all ranks from 1 to 20, under the three initial value models
107 above (degree, hybrid, and connectance). For each of these models, we measured the AUC of the
108 ROC curve **REF**. To identify the optimal cutoff in this curve, we selected the probability score
109 that maximizes Youden’s index of informedness, which works as a “total evidence” measure of
110 model confidence, especially in datasets with severe imbalances in prevalence.

Table 1: Summary statistics of the performance for the top 5 models, ranked according to the area under the ROC curve. For the sake of completeness, the best Youden’s index (at the threshold) is reported, as well as the rates of false discovery and false omission.

	model	rank	threshold	AUC	Youden’s index	false discovery	false omission
1	connectance	12	0.846	0.849	0.64	0.09	0.23
2	connectance	11	0.908	0.846	0.62	0.08	0.25
3	connectance	17	0.929	0.844	0.62	0.08	0.24
4	connectance	8	0.705	0.842	0.59	0.13	0.24
5	hybrid	12	0.707	0.841	0.58	0.14	0.25

111 The result of hyper-parameters tuning is presented in [tbl. 1](#). The best performing model, using
112 network connectance as an initial value, and a rank 12 approximation of the matrix, had a positive
113 predictive value of 0.90, and a negative predictive value of 0.76, for an overall accuracy of 0.82.
114 All things considered, given that the prevalence in the dataset is very low (only six out of every
115 thousand species pair do have an interaction), the best model has strong predictive power. The

116 ROC curve for this model is presented in fig. 2.

117 [Figure 2 about here.]

118 **Results and Discussion**

119 First, we report the top 10 likely hosts for betacoronaviruses, using the connectance of the net-
120 work as initial values, which are ranked by their final value post imputation; larger values should
121 indicate that the interactions are more likely to be possible. We report the novel hosts (identified
122 post Becker et al. (2020), according to <https://www.viralemergence.org/betacov>). These
123 results are presented in tbl. 2 - the novel hosts are presented in **bold**. Using a rank 2 approxima-
124 tion of the dataset, we have 5 novel hosts, and 4 identified as “suspected” hosts by the Becker
125 et al. (2020) ensemble model, currently lacking empirical evidence. This suggests that rank 2
126 contains the most information about the processes generating the data, and can therefore be used
127 to infer other associations.

Table 2: Top 10 likely hosts for betacoronaviruses using the connectance of the network as initial values

Rank 1	Rank 2
Artibeus jamaicensis	Hipposideros pomona
Scotophilus kuhlii	Scotophilus kuhlii
Molossus rufus	Artibeus jamaicensis
Sturnira lilium	Carollia brevicauda
Desmodus rotundus	Chaerephon pumilus
Glossophaga soricina	Molossus rufus
Eptesicus fuscus	Glossophaga soricina
Tadarida brasiliensis	Desmodus rotundus
Myotis nigricans	Sturnira lilium
Myotis lucifugus	Hipposideros larvatus

128 Based on this information, we have also extracted the 10 highest scoring interactions across
 129 the entire matrix at rank 2 (Table 2). The results demonstrates that within the entire dataset,
 130 including all mammalian hosts and viruses' genus, 5 out of the 10 highest scoring interactions
 131 are involving bat hosts (presented in *italic*), and 8 out of the 10 interactions are involving the
 132 lyssavirus genus. This genus includes the rabies virus (RABV), and other neurotropic rabies-
 133 related viruses (Warrell and Warrell 2004).

134 [Table 2: Top 10 likely missing interactions across the entire dataset using the connectance of
 135 the network as initial values]

Hosts species	Viruses genus
Sus scrofa	Lyssavirus
<i>Hipposideros armiger</i>	Lyssavirus
Rattus norvegicus	Lyssavirus
Myodes glareolus	Lyssavirus
<i>Pipistrellus abramus</i>	Lyssavirus
Sus scrofa	Orbivirus
Capra hircus	Alphavirus
<i>Rhinolophus sinicus</i>	Lyssavirus
<i>Myotis ricketti</i>	Lyssavirus
<i>Rhinolophus affinis</i>	Lyssavirus

136 Once those results were obtain, further investigations in the form of literature surveys allowed
 137 to identify that the interaction between *Pipistrellus abramus* } and lyssaviruses has already been
 138 noted by Hu et al. (2018); Shipley et al. (2019) reported lyssavirus prevalence in the genus
 139 *Pipistrellus*, *Myotis*, and *Rhinolophus*. Other confirmed hosts of lyssaviruses are *Sus scrofa*
 140 (Sato et al. 2004), and *Rattus norvegicus* (Wang, Tang, and Liang 2014). Surveillance for
 141 novel lyssaviruses infections is of great public health interest, since the rabies virus is fatal in all
 142 cases, once the onset of clinical symptoms has started (Banyard and Fooks 2017). Although it
 143 is recognized that bats are identified as reservoir hosts for lyssaviruses, the mechanism allowing
 144 the maintenance of the virus in those populations is still poorly understood (Banyard and Fooks

145 2017), and these predictions of interactions might serve as guidance in the monitoring of new
146 infections.

147 The two non-lyssaviruses associations have been previously reported in the literature (*Sus scrofa*
148 and orbivirus by Belaganahalli et al. (2015); *Capra hircus* and the equine encephalomyelitis
149 caused by an alphavirus as early as Pursell et al. (1972)). This suggests that Singular Value
150 Decomposition of available data on host-virus associations can uncover results that have been
151 reported in the primary literature, but not incorporated in the main databases used in the field;
152 based on the fact that the majority of the top 10 overall associations were able to be validated
153 from the literature, we suggest that interactions that have no empirical evidence could be targets
154 for additional sampling.

155 The initial value to be used for the imputation was then assigned according to the linear filter,
156 as presented in the method section. The Table 3 presents the number of novel hosts predicted
157 by the model, according to the coefficients used for the filter and to the rank.

158 [Table 3: Number of novel hosts for betacoronaviruses correctly predicted by the model using
159 linear filtering for the attribution of initial values]

Alpha	Rank 1	Rank 2	Rank 3	Rank 4	Rank 5
[0, 0, 0, 1]	3	3	1	3	4
[0, $\frac{1}{2}$, $\frac{1}{2}$, 0]	3	3	1	3	3
[0, $\frac{1}{3}$, $\frac{1}{3}$, $\frac{1}{3}$]	3	3	1	4	2
[0, 1, 0, 0]	3	3	1	3	3
[0, 0, 1, 0]	3	3	1	4	3

160 From the results presented in Table 3, it is possible to see that when using linear filtering for
161 the assignment of initial values, the choice of the α parameters does not impact the accuracy
162 of the predictions for the first three rank. The fourth and fifth rank then showed a variation per
163 α values. The highest scoring interactions for every combinations was then examined and the
164 variation of its value before and after the imputation has been calculated, and the results obtained
165 are presented in Table 4.

166 [Table 4: Variation of the value pre and post imputation for the highest scoring interaction at
167 every rank]

Rank 1	Rank 2	Rank 3	Rank 4	Rank 5
0.536	0.765	0.700	0.990	1.261

168 This variation was not influenced by the α parameters, but only by the rank used. The variation
169 calculated increased as the rank got higher.

170 Being able to identify intermediate animal hosts for potential zoonotic pathogens is an important
171 step in the fight against potential threats to global public health. Using SVD as an imputation
172 method to predict those interactions has demonstrated its potential to achieve this goal by cor-
173 rectly identifying the majority of the most likely associations, as validated by literature surveys,
174 and by suggesting interactions with no empirical evidence as targets for additional sampling.
175 Host-virus associations are a challenging imputation problem, because organized datasets are
176 scarce – as a result, a lot of missing associations are reported in the literature, but not available
177 in an easily usable format. Yet this also presents an opportunity to validate the performance of
178 recommender systems that is far more interesting than cross-fold or leave-one-out validation:
179 the existence of these interactions in the literature can provide validation on data that have never
180 been used in the modeling process, and therefore provide an accurate estimate of how frequently
181 existing interactions are identified. By this measure, that most of the top 10 recommendations
182 on this dataset were validated through *de novo* sampling (for bat hosts of betacoronaviruses) or
183 by a literature survey (for the global dataset) is a strong indication that SVD is able to uncover
184 likely host-virus pairs.

185 Future work on the use of SVD for virus host associations will have to adress the question of
186 the initial value used in the imputation process in further details. As of now, we relied on the
187 average number of interactions in the matrix, and on weighted allocations for different aspects
188 of the network structure, based on Stock et al. (2017) work on linear filtering. This method
189 can provide a good baseline estimate of how likely it is that a missing interaction could actually
190 exist (and in fact was developed for this purpose). For this reason, we are confident that the

191 performance of the approach can further be improved by fine-tuning the choice of the initial
192 value used for imputation, according to the dataset used, or by relying on ensemble models that
193 would aggregate the output of the best recommenders. Combining an accurate model for the
194 initial value with the SVD imputation is likely to generate predicted interactions that are strong
195 candidates for empirical validation.

196 **Acknowledgements:** This research was enabled in part by support provided by Calcul Québec
197 (www.calculquebec.ca) and Compute Canada (www.computecanada.ca). TP and CC were funded
198 by IVADO through the rapid response to COVID special initiative.

199 **References**

200 Albery, Gregory F., Evan A. Eskew, Noam Ross, and Kevin J. Olival. 2020. “Predicting the
201 Global Mammalian Viral Sharing Network Using Phylogeography.” *Nature Communica-*
202 *tions* 11 (1): 2260. <https://doi.org/10.1038/s41467-020-16153-4>.

203 Banyard, Ashley C., and Anthony R. Fooks. 2017. “The Impact of Novel Lyssavirus Discovery.”
204 *Microbiology Australia* 38 (1): 17–21.

205 Becker, Daniel J., Gregory F. Albery, Anna R. Sjodin, Timothée Poisot, Tad A. Dallas, Evan A.
206 Eskew, Maxwell J. Farrell, et al. 2020. “Predicting Wildlife Hosts of Betacoronaviruses for
207 SARS-CoV-2 Sampling Prioritization.” *bioRxiv*, May, 2020.05.22.111344. [https://doi.](https://doi.org/10.1101/2020.05.22.111344)
208 [org/10.1101/2020.05.22.111344](https://doi.org/10.1101/2020.05.22.111344).

209 Belaganahalli, Manjunatha N., Sushila Maan, Narender S. Maan, Joe Brownlie, Robert Tesh,
210 Houssam Attoui, and Peter P. C. Mertens. 2015. “Genetic Characterization of the Tick-
211 Borne Orbiviruses.” *Viruses* 7 (5): 2185–2209. <https://doi.org/10.3390/v7052185>.

212 Bezanson, J., A. Edelman, S. Karpinski, and V. Shah. 2017. “Julia: A Fresh Approach to Numer-
213 ical Computing.” *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.

214 Eckart, Carl, and Gale Young. 1936. “The Approximation of One Matrix by Another of Lower
215 Rank.” *Psychometrika* 1 (3): 211–18. <https://doi.org/10.1007/BF02288367>.

216 Forsythe, George, and Cleve Moler. 1967. *Computer Solution of Linear Algebraic Systems*.
217 Englewood Cliffs, New Jersey: Prentice Hall.

218 Golub, G. H., Alan Hoffman, and G. W. Stewart. 1987. "A Generalization of the Eckart-
219 Young-Mirsky Matrix Approximation Theorem." *Linear Algebra and Its Applications* 88-89
220 (April): 317–27. [https://doi.org/10.1016/0024-3795\(87\)90114-5](https://doi.org/10.1016/0024-3795(87)90114-5).

221 Golub, Gene H., and Christian Reinsch. 1971. "Singular Value Decomposition and Least
222 Squares Solutions." In *Linear Algebra*, 134–51. Springer.

223 Han, Barbara A., and John M. Drake. 2016. "Future Directions in Analytics for Infectious Dis-
224 ease Intelligence: Toward an Integrated Warning System for Emerging Pathogens." *EMBO*
225 *Reports* 17 (6): 785–89. <https://doi.org/10.15252/embr.201642534>.

226 Hu, Shu-Chia, Chao-Lung Hsu, Ming-Shiuh Lee, Yang-Chang Tu, Jen-Chieh Chang, Chieh-Hao
227 Wu, Shu-Hwae Lee, et al. 2018. "Lyssavirus in Japanese Pipistrelle, Taiwan - Volume 24,
228 Number 4 April 2018 - Emerging Infectious Diseases Journal - CDC." *Emerging Infectious*
229 *Diseases*. <https://doi.org/10.3201/eid2404.171696>.

230 Johnson, Christine K., Peta L. Hitchens, Pranav S. Pandit, Julie Rushmore, Tierra Smiley Evans,
231 Cristin C. W. Young, and Megan M. Doyle. 2020. "Global Shifts in Mammalian Population
232 Trends Reveal Key Predictors of Virus Spillover Risk." *Proceedings of the Royal Society*
233 *B: Biological Sciences* 287 (1924): 20192736. [https://doi.org/10.1098/rspb.2019.](https://doi.org/10.1098/rspb.2019.2736)
234 [2736](https://doi.org/10.1098/rspb.2019.2736).

235 Jones, Kate E., Nikkita G. Patel, Marc A. Levy, Adam Storeygard, Deborah Balk, John L. Git-
236 tleman, and Peter Daszak. 2008. "Global Trends in Emerging Infectious Diseases." *Nature*
237 451 (7181): 990–93. <https://doi.org/10.1038/nature06536>.

238 Lloyd-Smith, James O., Dylan George, Kim M. Pepin, Virginia E. Pitzer, Juliet R. C. Pulliam,
239 Andrew P. Dobson, Peter J. Hudson, and Bryan T. Grenfell. 2009. "Epidemic Dynamics
240 at the Human-Animal Interface." *Science* 326 (5958): 1362–67. [https://doi.org/10.](https://doi.org/10.1126/science.1177345)
241 [1126/science.1177345](https://doi.org/10.1126/science.1177345).

242 Plowright, Raina K., Colin R. Parrish, Hamish McCallum, Peter J. Hudson, Albert I. Ko, Andrea
243 L. Graham, and James O. Lloyd-Smith. 2017. "Pathways to Zoonotic Spillover." *Nature*

244 *Reviews Microbiology* 15 (8): 502–10. <https://doi.org/10.1038/nrmicro.2017.45>.

245 Pursell, A. R., J. C. Peckham, J. R. Cole, W. C. Stewart, and F. E. Mitchell. 1972. “Natu-
246 rally Occurring and Artificially Induced Eastern Encephalomyelitis in Pigs.” *Journal of the*
247 *American Veterinary Medical Association* 161 (10): 1143–47.

248 Ren, Wuze, Wendong Li, Meng Yu, Pei Hao, Yuan Zhang, Peng Zhou, Shuyi Zhang, et al. 2006.
249 “Full-Length Genome Sequences of Two SARS-Like Coronaviruses in Horseshoe Bats and
250 Genetic Variation Analysis.” *Journal of General Virology* 87 (11): 3355–59. [https://](https://doi.org/10.1099/vir.0.82220-0)
251 doi.org/10.1099/vir.0.82220-0.

252 Sato, Go, Takuya Itou, Youko Shoji, Yasuo Miura, Takeshi Mikami, Mikako Ito, Ichiro Kurane,
253 et al. 2004. “Genetic and Phylogenetic Analysis of Glycoprotein of Rabies Virus Isolated
254 from Several Species in Brazil.” *Journal of Veterinary Medical Science* 66 (7): 747–53.
255 <https://doi.org/10.1292/jvms.66.747>.

256 Shipley, Rebecca, Edward Wright, David Selden, Guanghui Wu, James Aegerter, Anthony R
257 Fooks, and Ashley C Banyard. 2019. “Bats and Viruses: Emergence of Novel Lyssaviruses
258 and Association of Bats with Viral Zoonoses in the EU.” *Tropical Medicine and Infectious*
259 *Disease* 4 (1). <https://doi.org/10.3390/tropicalmed4010031>.

260 Stock, Michiel, Timothée Poisot, Willem Waegeman, and Bernard De Baets. 2017. “Linear
261 Filtering Reveals False Negatives in Species Interaction Data.” *Scientific Reports* 7 (April):
262 45908. <https://doi.org/10.1038/srep45908>.

263 Wang, Lihua, Qing Tang, and Guodong Liang. 2014. “Rabies and Rabies Virus in Wildlife in
264 Mainland China, 1990.” *International Journal of Infectious Diseases* 25 (August): 122–29.
265 <https://doi.org/10.1016/j.ijid.2014.04.016>.

266 Warrell, M. J., and D. A. Warrell. 2004. “Rabies and Other Lyssavirus Diseases.” *The Lancet*
267 363 (9413): 959–69. [https://doi.org/10.1016/S0140-6736\(04\)15792-9](https://doi.org/10.1016/S0140-6736(04)15792-9).

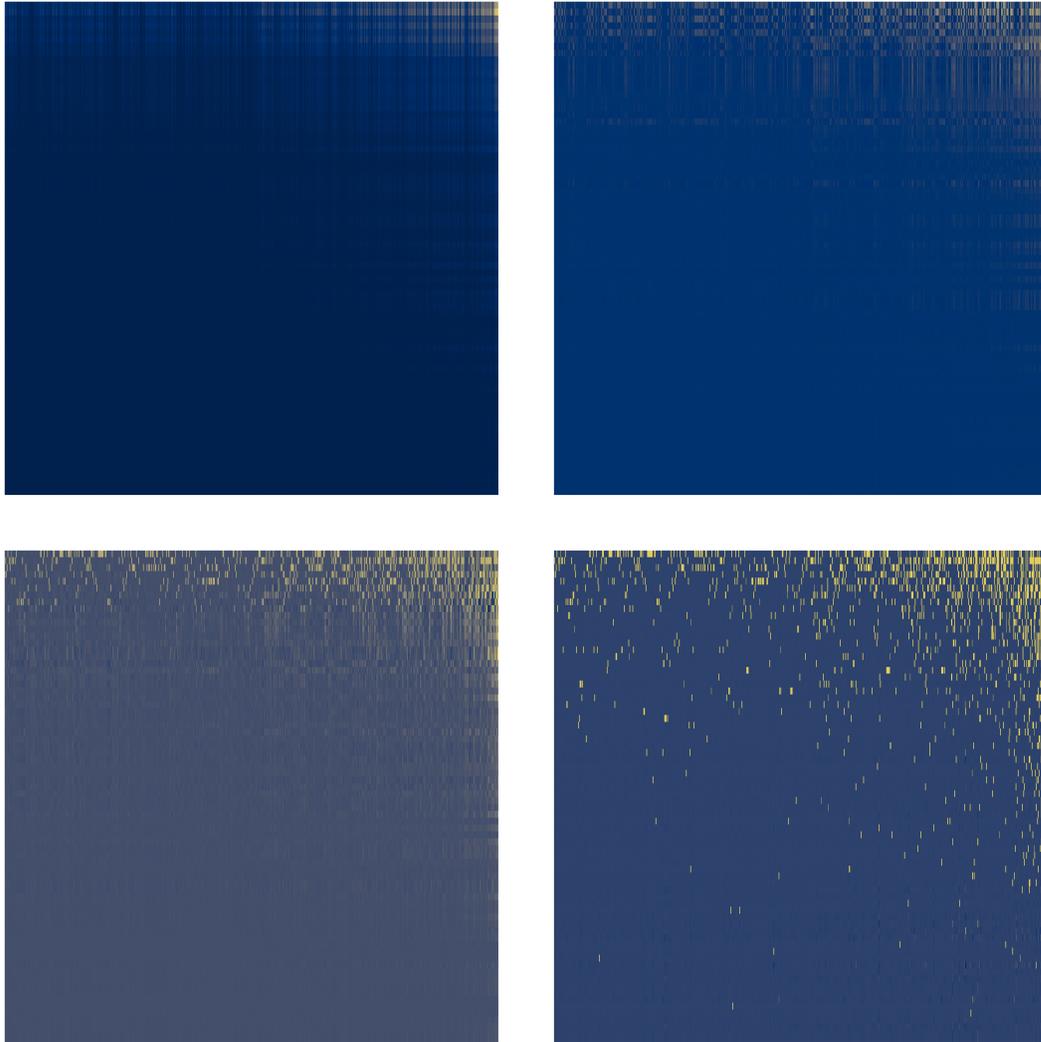


Figure 1: Overview of the dataset (yellow means interaction is more likely, blue means interactions is less likely) at different levels of approximation. At rank very low rank (top row; from left to right, $r = 1$ and $r = 3$) the matrix is mostly capturing the degree of the different species. At higher ranks (bottom row; from left to right, $r = 10$ and $r = 60$), the matrix is capturing increasing differences in species interactions.

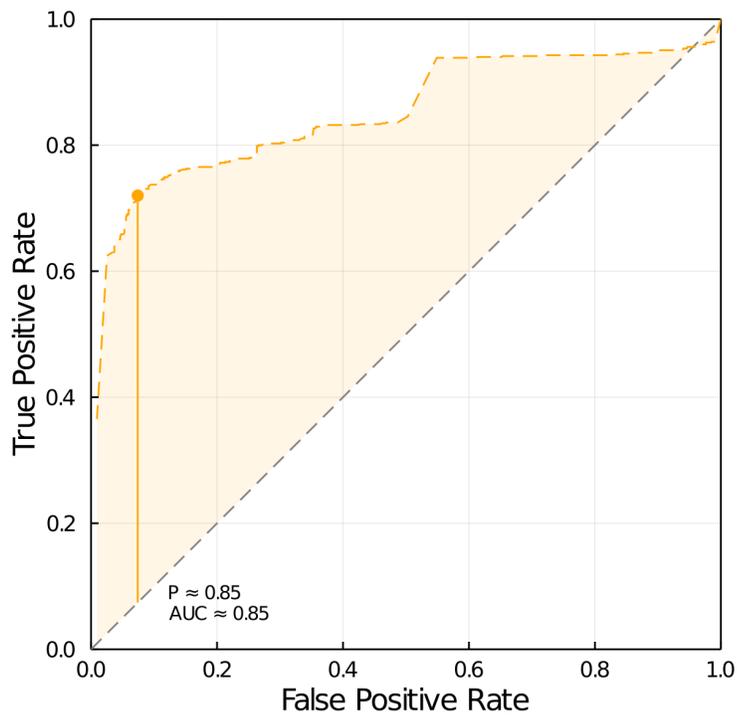


Figure 2: ROC curve for the best model, using network connectance as an initial value, and a rank 12 approximation. This model was used to run the prediction of false negatives in the entire dataset.