

TreeSp2Vec: a method for developing tree species embeddings using deep neural networks

Miguel A. González-Rodríguez (✉ miguel.an.gon.ro@gmail.com)

Universidade de Santiago de Compostela <https://orcid.org/0000-0002-0072-3124>

Ulises Diéguez-Aranda

Universidade de Santiago de Compostela

Research Article

Keywords: species classification, national forest inventory, representation learning, multi-layer perceptron, artificial intelligence

Posted Date: March 9th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-957638/v3>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

TreeSp2Vec: a method for developing tree species embeddings using deep neural networks

González-Rodríguez, M.A.^{a,*}, Diéguez-Aranda, U.^b

^a*Guangzhou Institute of Geography, Guangdong Academy of Sciences, Guangzhou 510070, China*

^b*Unidade de Xestión Ambiental e Forestal Sostible, Departamento de Enxeñaría Agroforestal, Universidade de Santiago de Compostela. Escola Politécnica Superior de Enxeñaría, R/ Benigno Ledo, Campus Terra, 27002 Lugo, Spain*

Abstract

In recent years, Representation Learning (RL), a subdiscipline of artificial intelligence, has proved a useful resource in many research fields for effectively mapping abstract categories into numeric scales. In this study, we explored a novel application of RL to forest ecology, labeled TreeSp2Vec, for developing tree species numeric representations using deep neural networks and data from a National Forest Inventory. Our approach consisted of a supervised species classification of individual trees using as input a set phytocentric and geocentric variables derived from forest inventory data and environmental cartography. Among the tested neural network architectures, a multi-layer perceptron with two hidden layers of 1024 units and an embedding layer of 16 units provided the best apparent and test performances (Matthew's Correlation Coefficient = 0.89). The developed latent representations (\mathbf{W}), or embeddings, were evaluated intrinsically by estimating their correlations with supplementary species descriptors that were not included in the training dataset. The evaluation analysis revealed some significant associations that proved the generality of the embedding model. Some latent dimensions (e.g., W_6 and W_{16}) were useful for differentiating species general features, such as conifers vs broad-leaved species, while other dimensions (e.g., W_2 and W_5) were related to forest ecosystem characteristics such as competition intensity (relative spacing index) and biodiversity (Simpson index). We concluded that the developed embeddings provided accurate and generalizable numeric representations of the considered tree species, which can be used as a ground for further cutting-edge forest modelling approaches

28 and open a new range of artificial intelligence applications in ecology.

29 *Keywords:* species classification, national forest inventory, representation learning,
30 multi-layer perceptron, artificial intelligence

31 1. Introduction

32 Over the last decade, deep neural networks have revolutionized the areas of predictive
33 and descriptive modelling (Sejnowski, 2018). Recent innovations in network architectures
34 have paved the path for new and fruitful scientific disciplines (e.g., Computer Vision), fre-
35 quently included in the broader scope of artificial intelligence. One of the most notable and
36 recent milestones in the field of deep learning is the emergence of Representation Learning
37 discipline (RL, Goodfellow et al., 2013; Huang et al., 2014). The latter focuses primarily on
38 the mapping of qualitative objects into numeric scales, so abstract and hard to quantify dif-
39 ferences between objects (i.e., “semantic” distances) can be expressed as algebraic distances.
40 The resulting quantitative dimensions from this mapping procedure are commonly referred
41 to as “embeddings” and the vector space defined by them is called “latent” space (Yu et al.,
42 2013). From a practical modelling perspective, the most relevant implication of RL is the
43 possibility of transforming or encoding categorical variables into ensembles of continuous
44 variables. One of the first and most notorious applications of RL was the embedding of En-
45 glish words into a 100D latent space with the Word2Vec methodology (Mikolov et al., 2013),
46 which is the backbone of current Natural Language Processing research. A diverse range
47 of applications of RL have been developed in the recent years, such as sentence embedding
48 (Tang et al., 2013), biomedical notes embedding (Wang et al., 2018), semantic-based image
49 embedding (Irtaza et al., 2014) and embedding of states in dynamic systems (Lesort et al.,
50 2018; Gelada et al., 2019).

51 Despite the recent successful uses of RL in several disciplines, applications in ecology and
52 environmental sciences are still scarce or non-existent. The potential benefits of applying

*Corresponding author

Email addresses: miguelangel.gonzalez.rodriguez@rai.usc.es (González-Rodríguez, M.A.),
ulises.dieguez@usc.es (Diéguez-Aranda, U.)

53

54 RL in ecology are varied and have to do with the ability of embedding models to transform
55 usual ecological categories, such as species, populations, and communities, into quantitative
56 variables. On the one hand, this transformation might allow to overcome frequent infor-
57 mation bottlenecks in ecological research derived from the combination of categorical and
58 numeric data. For instance, the projection of individuals belonging to a community into
59 a latent space might allow to develop complex quantitative biodiversity indexes that inte-
60 grate both categorical information (e.g., species and sociological classes) and quantitative
61 variables (e.g., individual morphology and environmental conditions), which could be a cru-
62 cial upgrade regarding the study of biodiversity. Another example could be to streamline
63 the study of dissimilarity between populations, species and communities by turning it into
64 a continuous problem where differences between classes could be represented by algebraic
65 distances in the latent space. On the other hand, transforming ecological categories into
66 quantitative variables might open a new range of high generality ecological modelling appli-
67 cations. For instance, the representation of different species in a latent space might allow to
68 develop multi-species modelling approaches, able to describe ecological dynamics of many
69 species simultaneously while keeping high predictive performance.

70 In the current study, we introduce a novel application of RL on forest ecology: an
71 embedding model of tree species, labeled “TreeSp2Vec”, based on deep neural networks. As
72 a case study, we applied our modelling approach to the most relevant forest species in Spain
73 using national forest inventory data and available environmental cartography.

74 **2. Methods**

75 *2.1. Data sources and preprocessing*

76 We used tree-level data of the 40 most frequent species in the Third Spanish National
77 Forest Inventory, accounting to a total of 50K inventory plots and approximately 850K
78 trees. The dataset encompassed a relatively wide range of forest types, including native
79 secondary forests, afforestations with native species and commercial plantations of exotic
80 species, existing in different mediterranean and eurosiberian ecoregions within the country.

81 For details about characteristics and methodology of the Spanish national forest inventory
82 see the explanation by Alberdi et al. (2010). From this dataset, 44 “phytcentric” descriptors
83 were computed, including: i) the diameter at breast height (cm) and total height (m) of
84 each tree, ii) the arithmetic means of diameters and heights of all the trees in each plot, and
85 iii) the count of the rest of trees in each plot for every considered species (i.e. 40 variables
86 resulting from the sum of expanded one-hot encoded vectors). Additionally, we derived 26
87 quantitative “geocentric” descriptors from available cartography, including: i) physical and
88 soil attributes obtained from the European Soil Data Centre (Ballabio et al., 2016, 2019), ii)
89 physiography variables derived from a Digital Elevation Model of the National Geographical
90 Institute of Spain (MDT200 <2009> CC-BY 4.0 ign.es), and iii) climate variables from the
91 Wordlclim 2.1 dataset (Fick and Hijmans, 2017).

92 2.2. Embedding approach and model training

93 A common approach for learning representations is to transform the problem into a
94 supervised learning task (Yang et al., 2015). Under this setup, the latent space is defined
95 by a layer of intermediate transformations necessary for obtaining the final target. In the
96 present study, the proposed encoding strategy was to perform a multi-label classification
97 for identifying the species of every tree in the dataset as a function of the phytocentric
98 and geocentric descriptors. Being u the number of descriptors and v the number of species
99 considered (i.e., the input and output dimension, respectively), and w the dimension of
100 the latent space, the embedding procedure consisted of an “encoder” map $F : \mathbb{R}^u \rightarrow \mathbb{R}^w$,
101 that transforms an input sample into an embedded representation, and a “decoder” map
102 $G : \mathbb{R}^w \rightarrow [0, 1]^v$, that translates the representation into a boolean target. Two types of
103 embedding could be derived from this approach. On the one hand, the output of the encoder
104 function F (hereafter, \mathbf{E}) is a sample-specific or *tree-level* embedding with shape $n \times w$, being
105 n the number of samples in the dataset. On the other hand, the set of weight parameters
106 in G with shape $v \times w$ (hereafter, \mathbf{W}) represents a class-specific or *species-level* embedding.

107 Both transformations, F and G , were implemented using deep neural networks (Figure 1).
108 The F map comprised a series of fully connected layers including the embedding layer (\mathbb{R}^w) as

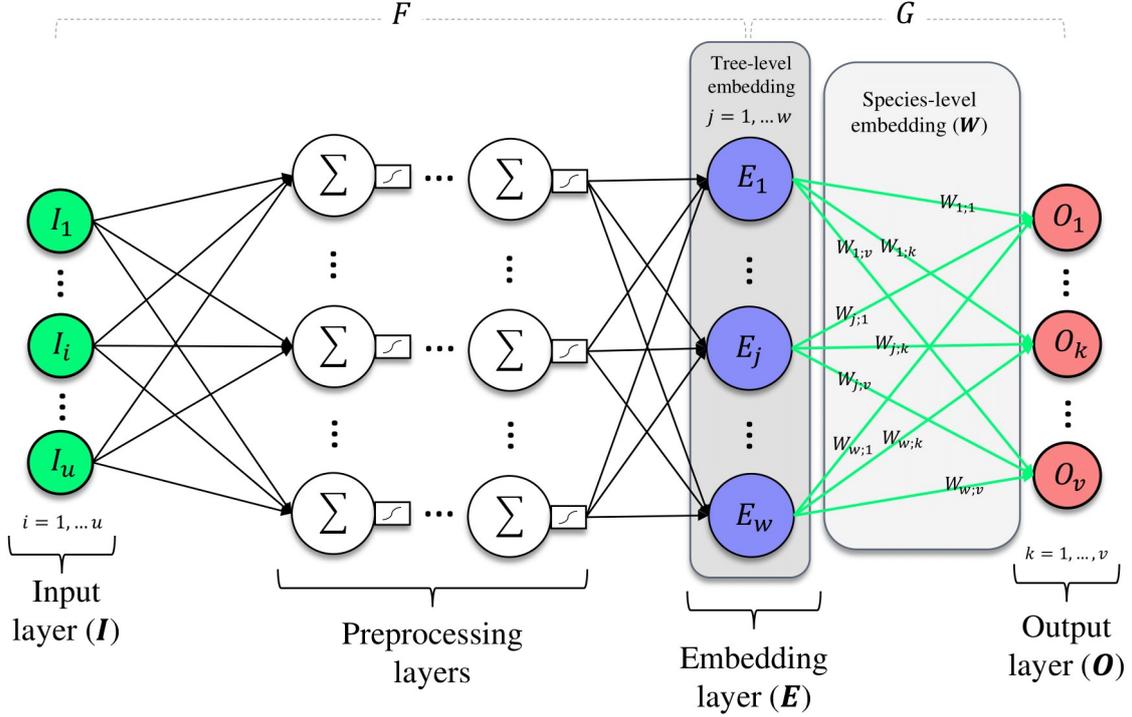


Figure 1: Diagram of the embedding approach using a deep neural network.

109 well as preprocessing layers with hyperbolic tangent activations for feature extraction. With
 110 regard to G , the embedding layer was fully connected to a v -dimensional output layer with
 111 softmax activation. We used the Keras API (Chollet et al., 2015) for Python 3 and GPU
 112 acceleration for optimizing different architectural variants with one to three preprocessing
 113 layers of 64-1024 units and an embedding layer with four to 32 units (scaling in powers of
 114 two). The different model variants were fitted using an Adam optimizer for minimizing the
 115 categorical cross-entropy loss over a maximum of 5000 epochs. Finally, the trained models
 116 were evaluated using a test subset of 20% of the initial dataset basing on the Matthew's
 117 Correlation Coefficient (MCC) implemented in the *scikit-learn* API (Pedregosa et al., 2011).

118 2.3. Evaluating the embeddings

119 To test the generality of the latent spaces developed, the ability of the model for inferring
 120 species and forest unseen characteristics was evaluated. We did so by checking potential as-
 121 sociations between embeddings and supplementary species attributes that were not explicitly

122 included as inputs. Specifically, point-biserial and Spearman correlations with latent dimen-
123 sions were estimated for three different groups of variables (hereafter, “test descriptors”):

- 124 • a set of qualitative descriptors encoded as binary variables representing whether a
125 species is a conifer (conifer = 1, broad-leaved = 0), deciduous (deciduous = 1, evergreen
126 = 0), native (native = 1, non-native = 0), invader (invader = 1, non-invader = 0), or
127 of “commercial interest” (according to recent timber production statistics in Spain);
- 128 • the relative frequencies of occurrence of each species in the eight ecoregions of Spain,
129 according to the TEOW regionalization (Olson et al., 2001); and
- 130 • the by-species mean values of a set of forest stand variables calculated for each inven-
131 tory plot (e.g., basal area).

132 In addition, we also evaluated the embeddings subjectively through visualization. To
133 further evaluate and illustrate the usability of the embeddings, an example of the method
134 for developing multi-species predictive models is presented in Supplemental Material.

135 **3. Results**

136 *3.1. Model building*

137 The trained models achieved high classification performances, being the minimum test
138 MCC score above 0.80. The more robust architecture in terms of classification performance
139 (apparent MCC = 0.89 and test MCC = 0.87) had two hidden preprocessing layers of 1024
140 units and 16 embeddings (1024:1024:16). The models including a high number of parame-
141 ters, such as the architectures with three preprocessing layers or 32 embedding units, yielded
142 very similar yet slightly lower accuracy (see Table 1). Consequently, the model with two
143 hidden layers of 1024 units and 16 embeddings (1.14×10^6 parameters), hereafter labeled
144 “TreeSp2Vec”, was selected as the best alternative for species representation learning. Some
145 variations in performance across the set of species for this model were noticed (see the con-
146 fusion matrix in Figure 2), being the minimum MCC = 0.66 for *Populus tremula* L. and the

147 maximum MCC = 0.98 for *Populus x canadensis* Moench. However, no consistent relation-
 148 ship was found between classification performance and species proportions in the dataset,
 149 since some very infrequent species, such as *Abies pinsapo* Boiss. or *Chamaecyparis lawsoni-*
 150 *ana* (A. Murray bis) Parl., were accurately classified (MCC = 0.96 and 0.94, respectively),
 151 while some relatively frequent species, such as *Quercus faginea* Lam., showed poorer per-
 152 formance (MCC = 0.78). The apparent MCC for every species is shown in Appendix A
 153 (Table A.1).

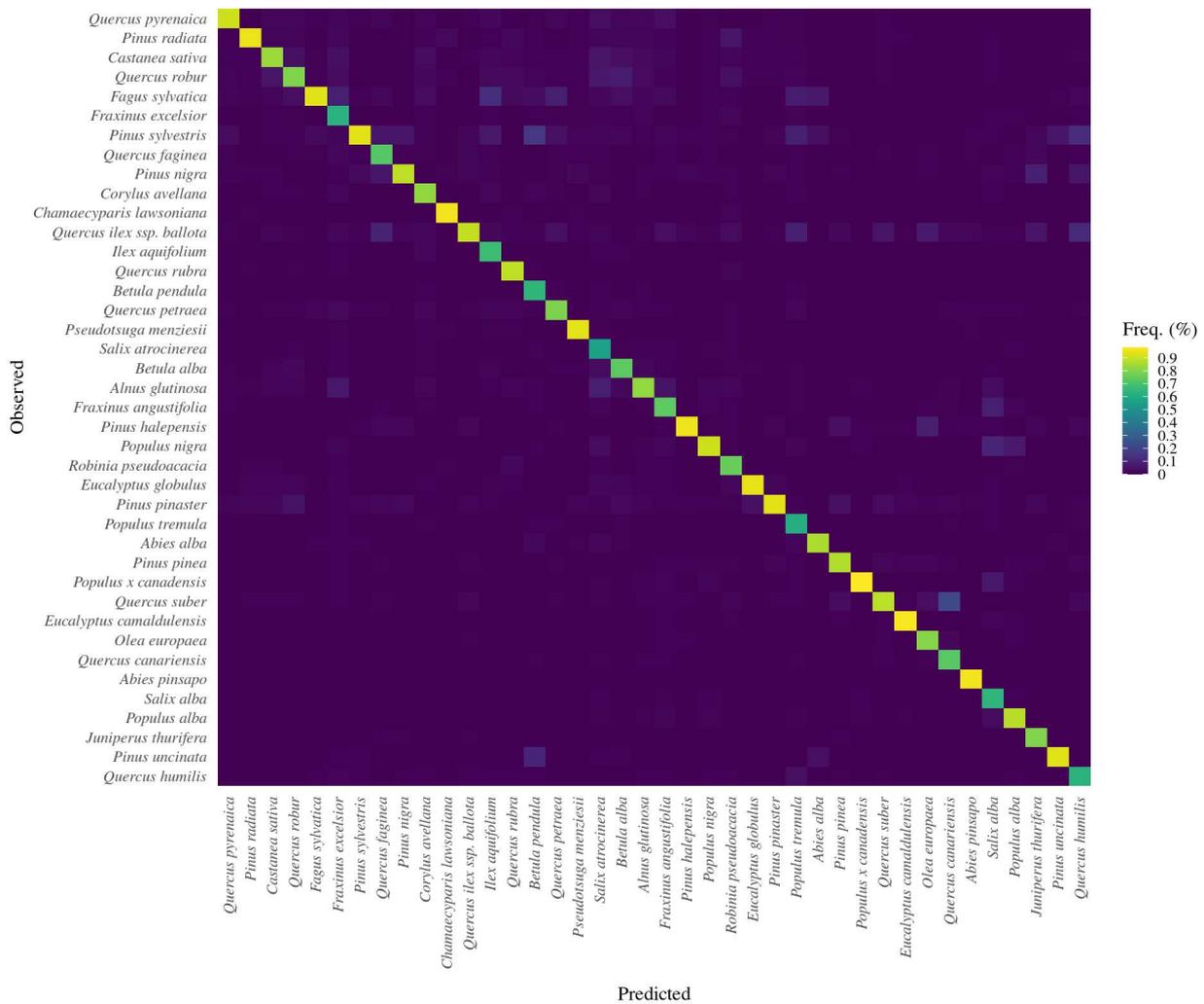


Figure 2: Confusion matrix (normalized frequencies) of the multi-label classification task using the TreeSp2Vec model for the 40 species considered.

Architecture	n° par.	MCC	MCC _{test}
1024:1024:16	1.14×10^6	0.892	0.871
1024:1024:32	1.16×10^6	0.889	0.870
1024:1024:8	1.13×10^6	0.890	0.869
1024:1024:1024:32	2.21×10^6	0.887	0.868
1024:1024:1024:16	2.19×10^6	0.885	0.866
512:512:32	3.17×10^5	0.887	0.866
512:512:512:32	5.79×10^5	0.885	0.865
1024:32	1.07×10^5	0.881	0.864
512:512:512:16	5.71×10^5	0.883	0.863
512:512:16	3.08×10^5	0.883	0.863

Table 1: Summary of MCC statistics and number of parameters in the ten best model architectures tested.

154 3.2. Embedding evaluation

155 The association analysis between the resulting species-level embeddings and the test
156 descriptors (Figure 3) showed some noticeable correlations, for a confidence level of 90%,
157 with the exceptions of W_1 and W_9 . Concerning the species qualitative attributes, the latent
158 dimensions W_6 and W_{16} showed intense correlations both with *Conifer* and *Deciduous*
159 features, while the dimensions W_4 and W_5 were respectively correlated with *Native* and
160 *Invader* attributes. W_2 and W_4 were strongly related to the commercial vs non-commercial
161 dichotomy. Regarding the species presence in the TEOW ecoregions, correlations with latent
162 dimensions were found for six out of the eight regions considered. Some of the strongest
163 correlations were found between W_5 and the presence in the “Pyrenees conifer and mixed
164 forests”, between W_4 and W_{10} and the “Iberian sclerophyllous and semi-deciduous forests”
165 and between W_{10} and the “Iberian conifer forests”. The six stand variables were significantly
166 associated with species embeddings. This was especially noticeable for the stand density,
167 which was correlated with four latent dimensions. The two variables linked to α -diversity
168 (species richness and the Simpson index), were similarly correlated with the dimensions W_2 ,

169 and W_6 . W_{15} was strongly correlated with the dominant height (m) and the relative spacing
 170 (also named Hart-Becking index).

171 The visualization of latent dimensions also allowed to detect some patterns that made
 172 species easily distinguishable, at least, for some dimensions and species. A plot of tree-level
 173 embeddings of latent dimensions E_6 and E_{15} is shown in (Figure 4).

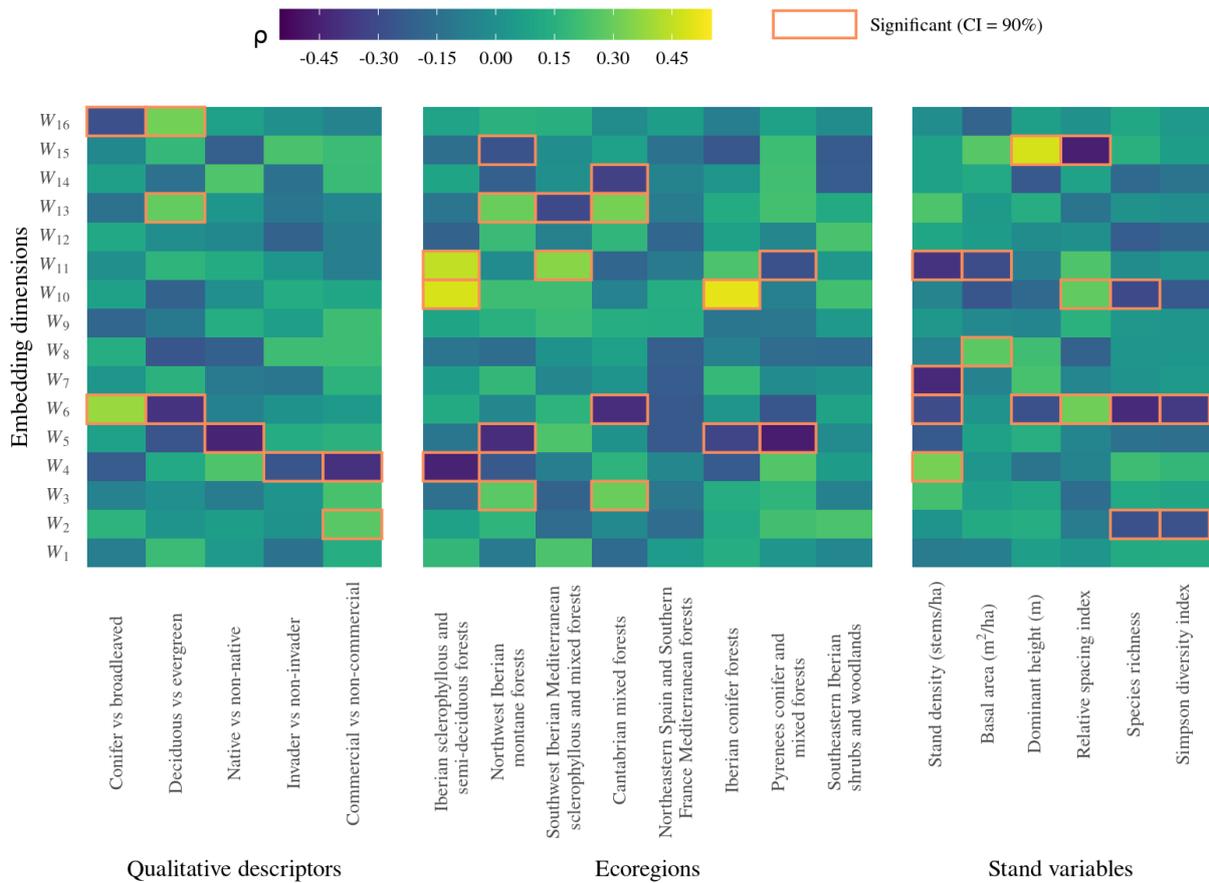


Figure 3: Heatmaps of Spearman and point-biserial correlations (ρ) between species-level embeddings and: 1) qualitative descriptors, 2) presence frequencies in the eight TEOW ecoregions, and 3) mean stand variables. Significant correlations for a Confidence Interval (CI) of 90% are highlighted.

174 4. Discussion

175 In the field of Representation Learning, the embeddings are considered to efficiently en-
176 capsule the most discriminative portion of information hidden in the data for effectively
177 representing the underlying factors that define abstract relationships between objects (Ben-
178 gio et al., 2013). In the current study, we applied a Representation Learning approach for
179 developing a latent space that allowed to quantify abstract characteristics of tree species
180 present in Spanish forests. As input data, we considered a set of phyto- and geocentric
181 descriptors that represented tree species attributes regarding three major axes: i) tree size
182 (diameters and heights), ii) specific composition in the vicinity, and iii) habitat character-
183 istics. The high classification performance yielded by the developed neural network models
184 suggests that the proposed architectures were able to extract useful information retained in
185 the inputs and efficiently translate it into a latent space. The latent dimensions resulting
186 from the TreeSp2Vec model seemed to be semantically meaningful, as they showed signif-
187 icant correlations with unseen species features (see Figure 3). The geocentric information
188 provided as input seemed to be effectively encoded by the model, as some latent dimensions
189 proved to be good identifiers of the habitat variability encompassed by the considered ecore-
190 gions. In addition, the phytocentric information was effectively compressed by the encoding
191 approach into synthetic representations related to stand variables and specific composition,
192 allowing the latent dimensions to identify species and forest characteristics coherently. For
193 instance, W_6 and W_{16} dimensions seemed to be useful for differentiating between conifers
194 and broad-leaved species, as well as between deciduous and evergreens. W_2 proved to be
195 an effective identifier of commercial plantations with low biodiversity (negative correlations
196 with species richness and Simpson index). W_8 seemed to be an exclusive identifier of tree
197 species basing on forest basal area. The strong relationship between stand density and some
198 dimensions, such as W_4 and W_{11} , provided a useful criterion for identifying species associated
199 with sparse forests, such as elm oaks and cork oaks, frequent in the Iberian sclerophyllous
200 and semi-deciduous forests. W_3 was an exclusive identifier of the Northwest Iberian montane
201 forests and the Cantabrian mixed forests, characterised by deciduous native species (W_{13})

202 and non-native plantations of evergreens (W_5 and W_6), which is consistent with the actual
203 characteristics of forests existing in those regions.

204 In the current study, the embedding evaluation procedure was based on correlations with
205 supplementary species descriptors, and could be classified as a “concept categorization”
206 intrinsic evaluator, according to the classification proposed by Wang et al. (2019). This
207 method differed notably from the most frequently used approaches in the family of intrinsic
208 evaluators in Representation Learning studies (similarity and analogy evaluators), which
209 compare statistical distances between pairs or groups of objects with subjective similarity
210 rankings previously made by domain experts (Faruqui et al., 2016). In forestry, we lack
211 of a dataset of human perceived similarity between species and, consequently, alternative
212 methods for embedding interpretation are necessary. However, even though the evaluation
213 method proposed in the present study proved to be effective for revealing the generality of the
214 encoding approach, the coherency concerning species proximity in the latent space remains
215 unclear. In this regard, visualization proved to be useful for assessing semantic coherency in
216 some cases. For instance, in Figure 4 the dimensions E_6 and E_{15} allow to visually distinguish
217 between species belonging to the same genus. In this Figure, species that never occur in
218 the same inventory plot, such as the mountain pine and eucalypts, are represented far away
219 in the subspace and do not overlap. In contrast, occasional neighbours, such as the Scots
220 pines and cork oaks happening in forests in central and northern Spain or the elm oaks and
221 the river red gums that can be present together in southwestern forests, overlap partially.
222 Regardless of the evaluation approach, the coherency of the latent space could be enhanced
223 by implementing regularization approaches, such as variational architectures, which have
224 demonstrated very good properties in previous studies for generative-oriented encoding (Pu
225 et al., 2016; Kingma and Welling, 2014). Derived from this, the development of a variational
226 embedding model for tree species might be a key improvement line in future research.

227 Concerning model performance, the variations in MCC observed across the set of species
228 might have implications regarding the uncertainty of the developed latent representations.
229 It is reasonable to think that latent vectors corresponding to species with low MCC values
230 will have more uncertainty regarding the meaningfulness of their semantic representations.

231 Though there is not an objectively definable threshold for accepting or rejecting the clas-
232 sification metrics yielded by the model, additional simulations carried out with a random
233 naive classifier yielded very low MCC values (maximum observed MCC = 0.0007), which
234 suggests that the species considered in this study were, overall, adequately classified. Even
235 so, complementary methods for analysing the potential interactions between performance
236 and semantic coherency might be necessary for clearing this matter in future research.

237 Considering the adequate results of model performance and semantic coherency, we be-
238 lieve that the tree species embeddings developed in the current study might be a valuable
239 resource for other forest research areas. The species-level embeddings can be used as in-
240 put for multi-species predictive modelling tasks by representing the different target species
241 by their corresponding latent vector instead of a categorical variable. In this case, testing
242 the predictive ability of an embedding-based multi-species model in comparison to single-
243 species models can provide a criterion for semantic coherency evaluation from an extrin-
244 sic perspective (as shown in Supplemental Material). Apart from multi-species predictive
245 modelling, the use of species-level embeddings could enhance previous methodologies for
246 biological similarity estimation between forest types (Hao et al., 2019a), by providing a
247 supervised framework for integrating non-linear relationships and interactions of a broad
248 spectrum of tree and habitat descriptors. In this regard, two major improvement lines could
249 be proposed: i) the development of new metrics for assessing the heterogeneity of the latent
250 space in terms of biodiversity, similarly to traditional indexes (e.g., the Avalanche Index,
251 Ganeshaiah et al., 1997), and, ii) the addition of new species descriptors for representing
252 morphological diversity and physiological traits (Hao et al., 2019b). Moreover, the inclusion
253 of spatial (e.g. spatial distribution indexes) and temporal (e.g. growth estimations derived
254 from plot remeasurements) descriptors could be crucial improvements for the development
255 of numeric representations of tree species considering forest structure and dynamics. These
256 developments might provide relevant advantages for forest research in the future.

257 **5. Conclusions**

258 In this study, an approach for learning abstract representations of tree species basing
259 on tree classification using deep neural networks was presented. The developed models
260 yielded high classification accuracy, reaching the best performance (Matthew's Correlation
261 Coefficient = 0.89) with a latent space of 16 dimensions. The evaluation of these embeddings
262 revealed that the resulting model, labeled TreeSp2Vec, was able to capture and compress
263 abstract characteristics of tree species in a semantically meaningful latent space, whose
264 generality allowed to predict unseen forest characteristics. We conclude that the developed
265 latent space of representations might be a useful resource for further research on forest multi-
266 species modelling and opens a new range of artificial intelligence applications in ecology.

267 **6. Acknowledgements**

268 We thank the *Banco de Datos de la Naturaleza* (BDN) for the Third National Forest
269 Inventory data provided. We also thank the IGN for the digital elevation model provided,
270 the ESDAC for the soil attributes maps and the Wordclim 2.1 project for the climate dataset
271 (Fick and Hijmans, 2017).

272 **7. Conflict of interest statement**

273 The authors declare no conflict of interest.

274 **8. Data availability**

275 The input data used in this article is available at the *Banco de Datos de la Naturaleza*
276 website (miteco.gob.es).

277 **References**

278 Alberdi, I., Condés, S., Martínez, J., Martínez, S., de Toda, S., Sánchez, G., Pérez, F., Villanueva, J.,
279 and Vallejo, R. (2010). Spanish national forest inventory. Martínez, S.; de Toda, S.; Sánchez, G.;
280 Pérez, F.; Villanueva, J.A.; Vallejo, R. Alberdi, I.; Condés, S.; Martínez, J.; Martínez, S.; de Toda, S.;

281 Sánchez, G.; Pérez, F.; Villanueva, J.A.; Vallejo, R.; Spanish national forest inventory. In National
282 Forest Inventories. Pathways for Common Reporting; Tomppo, E.; Gschwantner, T., Lawrence, M.,
283 McRoberts, R.E., Eds.; Springer: Berlin, Germany, 2010; pp. 527–54.

284 Ballabio, C., Lugato, E., Fernández-Ugalde, O., Orgiazzi, A., Jones, A., Borrelli, P., Montanarella, L., and
285 Panagos, P. (2019). Mapping lucas topsoil chemical properties at european scale using gaussian process
286 regression. *Geoderma*, 355:113912.

287 Ballabio, C., Panagos, P., and Monatanarella, L. (2016). Mapping topsoil physical properties at european
288 scale using the lucas database. *Geoderma*, 261:110–123.

289 Bengio, Y., Courville, A., and Vincent, P. (2013). Representation learning: A review and new perspectives.
290 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828.

291 Chollet, F. et al. (2015). Keras.

292 Faruqui, M., Tsvetkov, Y., Rastogi, P., and Dyer, C. (2016). Problems with evaluation of word embeddings
293 using word similarity tasks. In *Proceedings of the 1st Workshop on Evaluating Vector-Space Representa-*
294 *tions for NLP*, pages 30–35, Berlin, Germany. Association for Computational Linguistics.

295 Fick, S. E. and Hijmans, R. J. (2017). WorldClim 2: new 1-km spatial resolution climate surfaces for global
296 land areas. *International Journal of Climatology*.

297 Ganeshaiah, K. N., Chandrashekara, K., and Kumar, A. R. V. (1997). Avalanche index: A new measure of
298 biodiversity based on biological heterogeneity of the communities. *Current Science*, 73(2):128–133.

299 Gelada, C., Kumar, S., Buckman, J., Nachum, O., and Bellemare, M. G. (2019). DeepMDP: Learning
300 continuous latent space models for representation learning. In Chaudhuri, K. and Salakhutdinov, R.,
301 editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings*
302 *of Machine Learning Research*, pages 2170–2179. PMLR.

303 Goodfellow, I. J., Erhan, D., Carrier, P. L., Courville, A., Mirza, M., Hamner, B., Cukierski, W., Tang,
304 Y., Thaler, D., Lee, D.-H., Zhou, Y., Ramaiah, C., Feng, F., Li, R., Wang, X., Athanasakis, D., Shawe-
305 Taylor, J., Milakov, M., Park, J., Ionescu, R., Popescu, M., Grozea, C., Bergstra, J., Xie, J., Romaszko,
306 L., Xu, B., Chuang, Z., and Bengio, Y. (2013). Challenges in representation learning: A report on three
307 machine learning contests.

308 Hao, M., Corral-Rivas, J. J., González-Elizondo, M. S., Ganeshaiah, K. N., Nava-Miranda, M. G., Zhang, C.,
309 Zhao, X., and von Gadow, K. (2019a). Assessing biological dissimilarities between five forest communities.
310 *Forest Ecosystems*, 6.

311 Hao, M., Ganeshaiah, K. N., Zhang, C., Zhao, X., and von Gadow, K. (2019b). Discriminating among forest
312 communities based on taxonomic, phylogenetic and trait distances. *Forest Ecology and Management*, 440.

313 Huang, P., Huang, Y., Wang, W., and Wang, L. (2014). Deep embedding network for clustering. In *2014*
314 *22nd International Conference on Pattern Recognition*, pages 1532–1537.

315 Irtaza, A., Jaffar, M. A., Aleisa, E., and Choi, T.-S. (2014). Embedding neural networks for semantic
316 association in content based image retrieval. *Multimedia Tools and Applications*, 72.

317 Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes.

318 Lesort, T., Díaz-Rodríguez, N., Goudou, J.-F., and Filliat, D. (2018). State representation learning for
319 control: An overview. *Neural Networks*, 108:379–392.

320 Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of
321 words and phrases and their compositionality. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani,
322 Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 26.
323 Curran Associates, Inc.

324 Olson, D. M., Dinerstein, E., Wikramanayake, E. D., Burgess, N. D., Powell, G. V. N., Underwood, E. C.,
325 D’amico, J. A., Itoua, I., Strand, H. E., Morrison, J. C., Loucks, C. J., Allnutt, T. F., Ricketts, T. H.,
326 Kura, Y., Lamoreux, J. F., Wettengel, W. W., Hedao, P., and Kassem, K. R. (2001). Terrestrial Ecoregions
327 of the World: A New Map of Life on Earth: A new global map of terrestrial ecoregions provides an
328 innovative tool for conserving biodiversity. *BioScience*, 51(11):933–938.

329 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer,
330 P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and
331 Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
332 12:2825–2830.

333 Pu, Y., Gan, Z., Heno, R., Yuan, X., Li, C., Stevens, A., and Carin, L. (2016). Variational autoencoder
334 for deep learning of images, labels and captions. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and
335 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates,
336 Inc.

337 Richards, F. J. (1959). A flexible growth function for empirical use. *Journal of Experimental Botany*,
338 10:290–300.

339 Sejnowski, T. (2018). *The Deep Learning Revolution*. The MIT Press. MIT Press.

340 Tang, D., Qin, B., Liu, T., and Li, Z. (2013). Learning sentence representation for emotion classification
341 on microblogs. In Zhou, G., Li, J., Zhao, D., and Feng, Y., editors, *Natural Language Processing and*
342 *Chinese Computing*, pages 212–223, Berlin, Heidelberg. Springer Berlin Heidelberg.

343 Wang, B., Wang, A., Chen, F., Wang, Y., and Kuo, C.-C. J. (2019). Evaluating word embedding models:
344 methods and experimental results. *APSIPA Transactions on Signal and Information Processing*, 8:e19.

345 Wang, Y., Liu, S., Afzal, N., Rastegar-Mojarad, M., Wang, L., Shen, F., Kingsbury, P., and Liu, H. (2018).
346 A comparison of word embeddings for the biomedical natural language processing. *Journal of Biomedical*
347 *Informatics*, 87.

348 Yang, S., Luo, P., Loy, C. C., Shum, K. W., and Tang, X. (2015). Deep representation learning with target

349 coding. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1).
350 Yu, W., Zeng, G., Luo, P., Zhuang, F., He, Q., and Shi, Z. (2013). Embedding with autoencoder regulariza-
351 tion. In Blockeel, H., Kersting, K., Nijssen, S., and Železný, F., editors, *Machine Learning and Knowledge*
352 *Discovery in Databases*, pages 208–223, Berlin, Heidelberg. Springer Berlin Heidelberg.

353 **Appendix A. Classification performance by species**

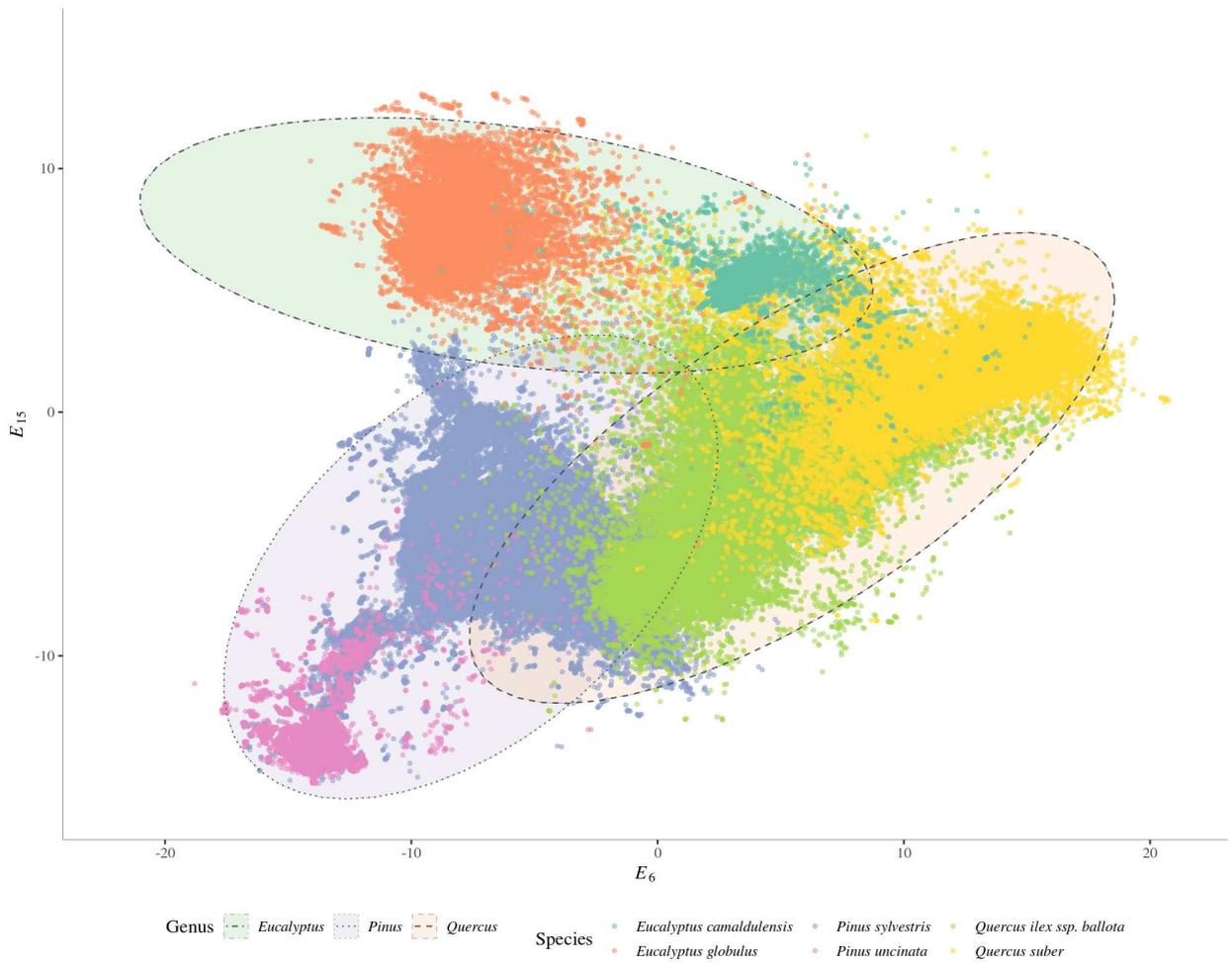


Figure 4: Scatterplot of 6th vs 15th dimensions of tree-level embeddings for *Pinus uncinata* Ramond ex A.DC., *Pinus sylvestris* L., *Quercus ilex* L. ssp. *ballota*, *Quercus suber* L., *Eucalyptus globulus* Labill., and *Eucalyptus camaldulensis* Dehnh. Shaded ellipses correspond to confidence intervals = 99% of by-genus multivariate normal distributions.

Species	MCC
<i>Chamaecyparis lawsoniana</i>	0.938
<i>Pinus sylvestris</i>	0.903
<i>Pinus uncinata</i>	0.931
<i>Pinus pinea</i>	0.875
<i>Pinus halepensis</i>	0.945
<i>Pinus nigra</i>	0.875
<i>Pinus pinaster</i>	0.918
<i>Pinus radiata</i>	0.943
<i>Abies alba</i>	0.853
<i>Abies pinsapo</i>	0.955
<i>Pseudotsuga menziesii</i>	0.927
<i>Juniperus thurifera</i>	0.829
<i>Quercus robur</i>	0.794
<i>Quercus petraea</i>	0.819
<i>Quercus pyrenaica</i>	0.914
<i>Quercus faginea</i>	0.787
<i>Quercus ilex ssp. ballota</i>	0.878
<i>Quercus suber</i>	0.857
<i>Quercus canariensis</i>	0.755
<i>Quercus rubra</i>	0.879
<i>Populus alba</i>	0.892
<i>Populus tremula</i>	0.658
<i>Alnus glutinosa</i>	0.798
<i>Fraxinus angustifolia</i>	0.785
<i>Populus nigra</i>	0.913
<i>Eucalyptus globulus</i>	0.931
<i>Eucalyptus camaldulensis</i>	0.960
<i>Ilex aquifolium</i>	0.734
<i>Olea europaea</i>	0.814
<i>Fagus sylvatica</i>	0.926
<i>Castanea sativa</i>	0.841
<i>Corylus avellana</i>	0.786
<i>Robinia pseudoacacia</i>	0.750
<i>Quercus pubescens</i>	0.718
<i>Fraxinus excelsior</i>	0.662
<i>Salix alba</i>	0.738
<i>Populus x canadensis</i>	0.975
<i>Betula alba</i>	0.747
<i>Salix atrocinerea</i>	0.677
<i>Betula pendula</i>	0.718

Table A.1: Mathew’s Correlation Coefficient of the classification task performed by the TreeSp2Vec model for each of the 40 species considered.

Supplemental Material: TreeSp2Vec: a method for developing tree species embeddings using deep neural networks

S-I. Example of embedding-based multi-species modelling

We present a simple illustrative example of how to use species-level embeddings for developing empirical multi-species models. The task is to develop a multi-species generalized height-diameter (h - d) equation for predicting the total height of individual trees as a function of their diameter at breast height.

S-I.1. Example methods

The proposed h - d model for this example is the Chapman-Richards equation (Richards, 1959):

$$h = 1.3 + a \left(1 - \exp(-bd) \right)^c, \quad (\text{S1})$$

where h is the tree total height (m), d is the tree diameter at breast height (cm), and a , b and c are parameters. We fit the Chapman-Richards equation for the 30 most frequent species in the dataset by developing 1) one global generalized h - d model (i.e., with the same a , b and c values for all the species) and 2) 30 different single-species models. Then, the parameters (a , b and c) in each single-species model are expanded using the species-level embeddings. We do so by predicting a , b and c as a function of the previously developed 16 latent dimensions (\mathbf{W}) using a machine learning model. Specifically, we use multi-layer perceptrons (MLPs) with three hidden layers of 32 units and ReLU activations, calibrated through Monte Carlo validation. Finally, we apply the expanded a , b and c parameters for each species for re-predicting the total tree height using a multi-species Chapman-Richards model:

$$h = 1.3 + MLP_a(\mathbf{W}) \left(1 - \exp \left(- MLP_b(\mathbf{W}) d \right) \right)^{MLP_c(\mathbf{W})}, \quad (\text{S2})$$

where MLP_x are the outputs of the multilayer perceptrons trained for predicting a , b and c as a function of \mathbf{W} .

S-I.2. Example results

As expected, considering the variety of species in the dataset, the global $h-d$ model performs very poorly (see Table S1), yielding an $R^2=0.49$, being the unweighted mean over the set of species $R^2=0.05$. In comparison, the single-species models perform much better, although they also show strong variations in R^2 between different species (ranging from $R^2=0.32$ to 0.73, see Figure S1). The three parameters are successfully predicted from the latent dimensions (see Figure S2), which accounts for the existence of meaningful associations with \mathbf{W} (see Figure S3). Finally, the re-prediction of tree height using the multi-species Chapman-Richards model reveals, on average, only in a slight drop in performance with respect to the single-species models (see Figure S4). This result confirms the usefulness of the approach for using only one model for all the species while maintaining good performance.

Model	R^2	R^2_{min}	R^2_{mean}	R^2_{max}
Global	0.485	-0.502	0.0566	0.501
Single-species	0.736	0.322	0.504	0.7308
Multi-species	0.662	0.0929	0.455	0.7208

Table S1: Predictive performance of the developed $h-d$ approaches.

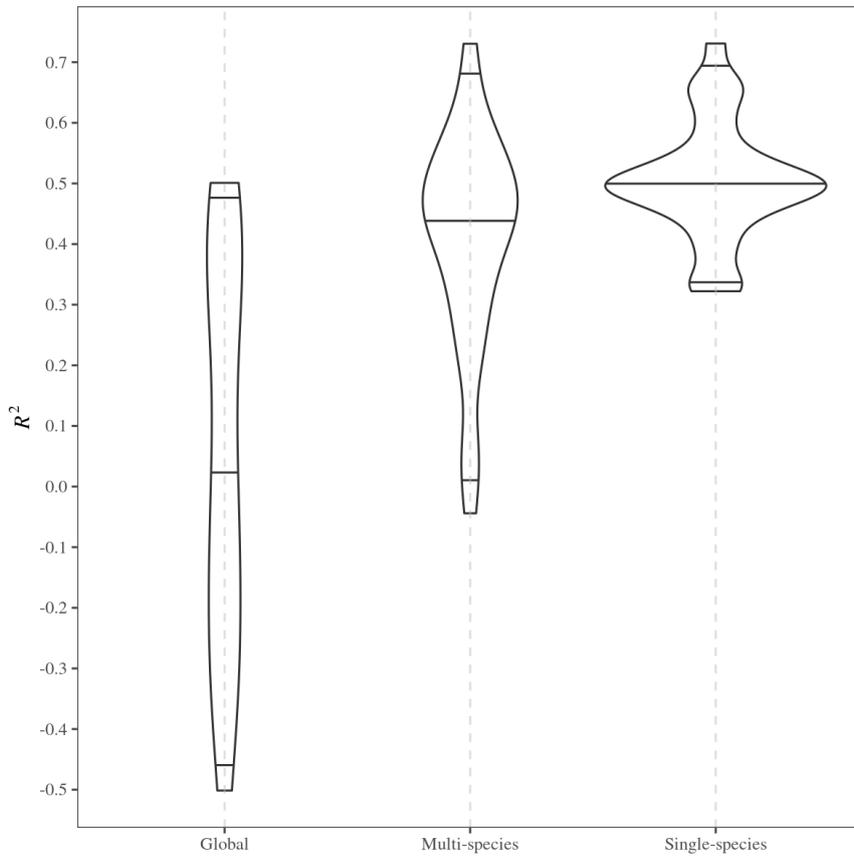


Figure S1: Violin plots representing the distribution of R^2 across species of the three $h-d$ modeling approaches.

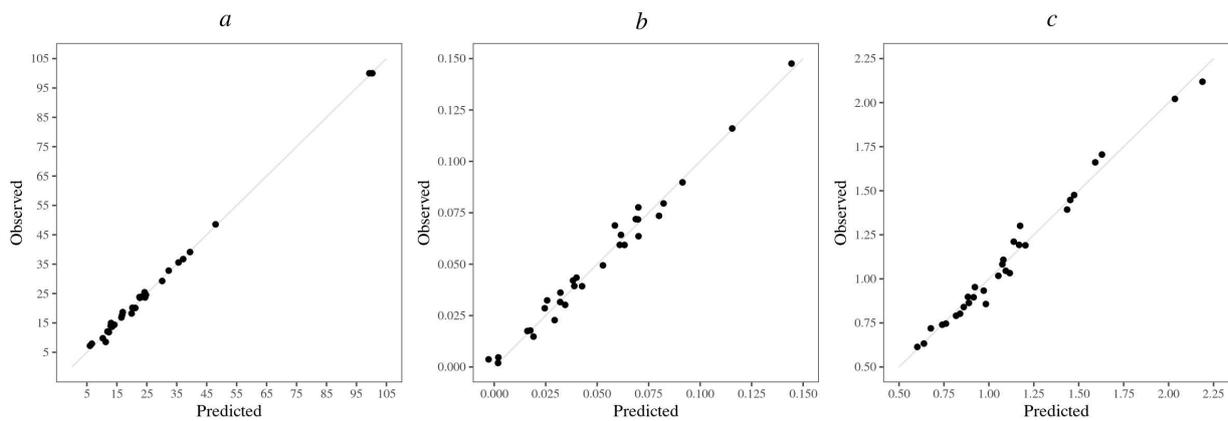


Figure S2: Observed vs predicted values of the parameters using the three MLP models.

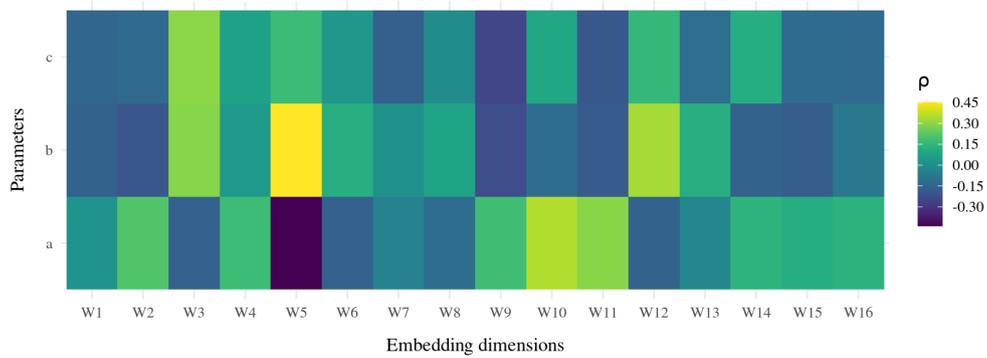


Figure S3: Heatmap of Spearman correlations (ρ) between species-level embeddings and the three parameters of each single-species *h-d* model.

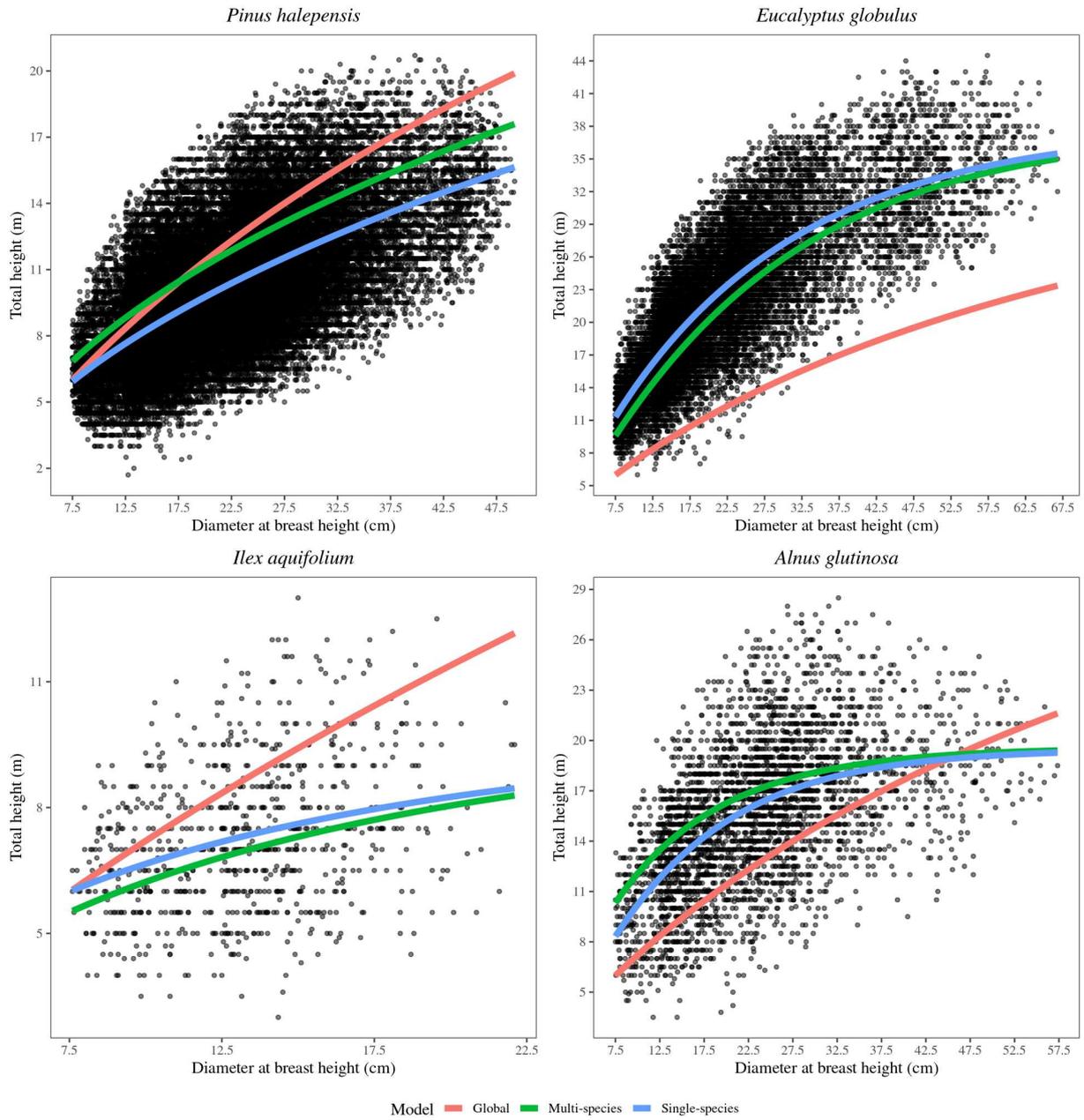


Figure S4: h - d scatter plots of four tree species with predicted trends for the three generalized h - d model approaches.