

# Agricultural drought stress forecasting across different biomes in Brazil using machine learning

## 1. Abstract

Drought events across Brazil have become more common in recent years. The northeast, marked by the dry semi-arid Caatinga biome, experiences the most extreme drought events, which can have a large impact on economic productivity including disruptions to energy and transport infrastructure as well as agricultural losses. However, droughts are an issue across Brazil and can have significant effects on vegetation across the six biomes which make up the country. Drought stress, defined by a reduction in the vegetation health index (VHI) below the 40% threshold is indicative of significant crop yield reductions. Here we evaluate monthly forecasted reductions in VHI below the 40% threshold using machine learning. We then use the results to understand how climatological features influence drought stress in each biome. This is achieved through a combination of Empirical Orthogonal Function (EOF) analysis, Shapley plots, and correlation analysis of predicted spatial trends. We observe distinct behaviour in the Caatinga biome which is rooted in its unique regional climate patterns and physical geography. For other biomes, VHI, RZSM and precipitation are less strongly linked; however high inertia in VHI allows for high forecast performance. The findings strongly imply that drought forecasting models should treat Northeast Brazil as a distinct system, training models specifically for this region to improve the prediction accuracy of vegetation health and integrated drought indices.

## 2. Introduction

In recent years, drought episodes have become increasingly frequent across Brazil. The Northeast, marked by its semi-arid climate and the Caatinga biome, is the most severely affected by extreme droughts, which can significantly undermine economic productivity—leading to agricultural losses and disruptions in energy and transportation systems (Cunha et al., 2019). Drought stress impacts in Northeast Brazil have been reported since the sixteenth century (Marengo et al., 2017). For example, drought in 1997–1998 resulted in agricultural losses of 57% with the total economic damage estimated to be 5% of the regions GDP (Marengo et al., 2017). Agriculture in northeast Brazil is both important and vulnerable to drought impacts. The region has 7.8 M ha of agricultural lands with a share of 6–13% of the Brazilian production of soybean, maize, coffee, sugarcane, milk and beef. However, 95% of farmed land consists of rainfed agriculture, with a lack of irrigation increasing vulnerability to drought impacts (Marengo et al., 2022).

Many studies have been conducted on droughts in Northeast Brazil (Marengo et al., 2022, Cunha et al., 2019, Marengo et al., 2017, Lopes Ribeiro et al., 2021, Gallear et al., 2025). Previous drought monitoring studies estimated soil moisture trends in the region (Zeri et al., 2022). They found that soil moisture anomaly (SMA) and standardized precipitation index (SPI), showed a lagged correlation of 1 to 1.5 months in the annual scale with vegetation health index (VHI), suggesting that negative trends in SMA and SPI can be an early warning for yield losses (approximated by vegetation health) during the growing season. Gallear et al. (2025) built on this work by showing high correlations between soil moisture, SPEI and precipitation with a one-month time delay.

45 However, drought impacts are not limited to the northeast region of Brazil. Most recently,  
46 severe drought in Amazonia saw precipitation deficits of the order of 50 to 100 mm/month  
47 and temperatures 3°C greater than normal, leading to reduced evapotranspiration and soil  
48 moisture indicators (Marengo et al., 2024). In the west, drought has also impacted the  
49 Pantanal wetland biome. During the period of 2019 – 2020, a lack of precipitation caused by  
50 predominance of warm dry air masses from subtropical latitudes led to reduced river levels  
51 which subsequently affected transport systems (Marengo et al., 2021, Ferreira Barbosa et  
52 al., 2022). Droughts have also impacted South Brazil, causing US\$ 3.5 billion in soybean  
53 yield losses from 1974 to 2019 (Miyamoto, 2024).

54 Here, we build on the work of Zeri et al. (2022) and Gallea et al. (2025) to determine how  
55 variations in climate across the six biomes of Brazil affect the probability of agricultural  
56 drought stress, and the propagation from meteorological variables to agricultural drought  
57 stress. Machine learning models are trained to forecast agricultural drought for each biome  
58 and the results are analysed to extract the importance of different predictors and study  
59 spatial-temporal variations. Further contextual analyses are performed using Principal  
60 Component Analysis (PCA, also referred to as Empirical Orthogonal Function or EOF  
61 analysis) and simple linear correlations to help interpret the results.

62

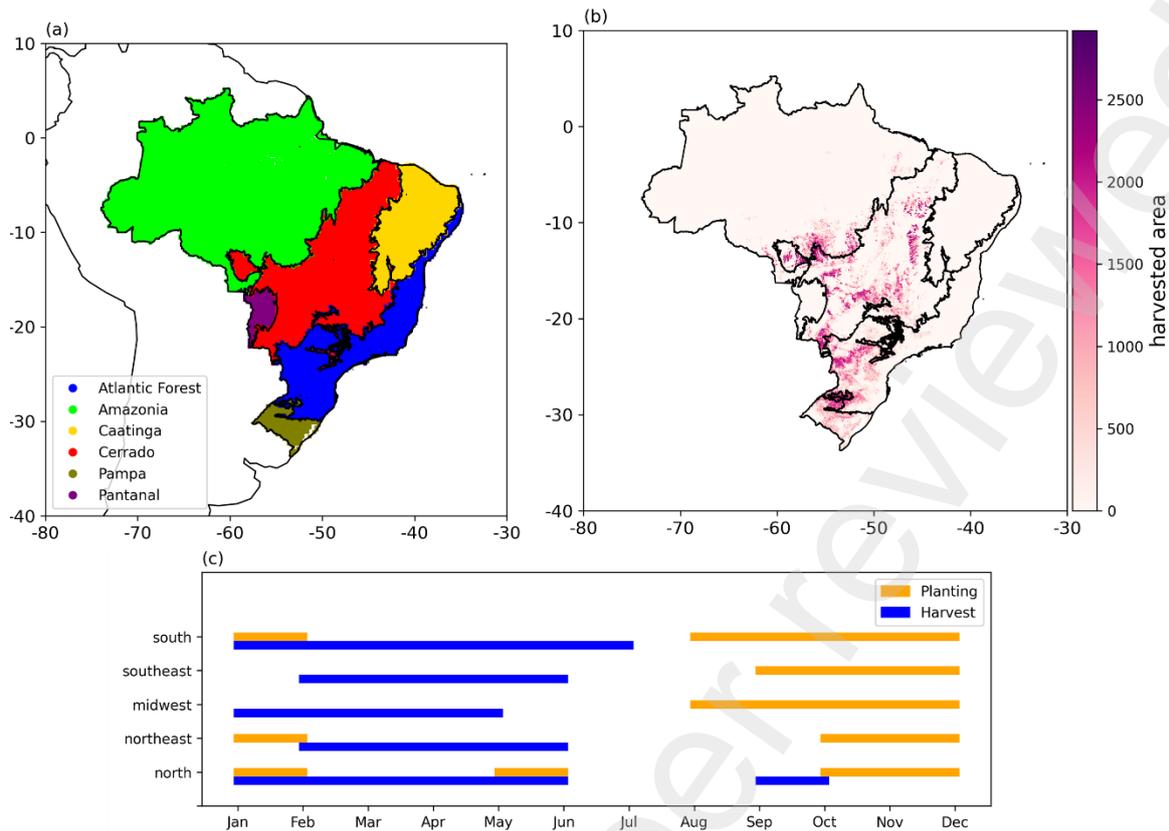
63 We seek to answer the following research questions:

- 64 1. How may monthly agricultural drought stress forecast model performance vary  
65 across Brazil's six biomes?
- 66 2. How does monthly propagation from meteorological variation to agricultural drought  
67 stress vary across Brazil
- 68 3. How does variation between biomes in Brazil affect probability of agricultural drought  
69 stress?

## 70 **3. Methodology**

### 71 **3.1. Study regions**

72 Figure 1 shows the biomes used to separate the data used in this study, the soybean  
73 harvested area in relation to biome boundaries and the cropping calendar with planting and  
74 harvesting periods of soybean for the different regions of Brazil.



**Figure 1.** (a) Each biome of Brazil: Atlantic Forest, Amazonia, Caatinga, Cerrado, Pampa and Pantanal. (b) Soybean harvested area, data from CROPGRIDS dataset (Tang et al., 2024). (c) Growing season calendar for regions of Brazil (BOPAR, 2025).

75

76 The Caatinga biome averages an annual rainfall below 800 mm per year (Alvares et al.,  
 77 2013). This, combined with high evapotranspiration rates (exceeding 2500 mm yr<sup>-1</sup>) (Lopes  
 78 Ribeiro et al., 2021) contributes to reduced water availability and a very limited storage  
 79 capacity of rivers which are mostly intermittent (Lopes Ribeiro et al., 2021). Furthermore,  
 80 Caatinga soils are generally shallow (0-50 cm) meaning limited soil moisture storage  
 81 capacity (Cirilo, 2008). Annual precipitation in the Cerrado varies between 600–2200 mm, it  
 82 is the second largest biome with wettest regions bordering the amazon and driest regions  
 83 close to the Caatinga biome. Amazonia is the largest biome, mainly characterized by  
 84 rainforest areas and an equatorial climate, with daily mean temperatures ranging from 22–28  
 85 °C. Torrential rains are distributed throughout the year (Lopes Ribeiro et al., 2021). The  
 86 Pantanal is the smallest biome in Brazil; however, it is one of the world's largest wetlands  
 87 (Marengo et al., 2021, Lopes Ribeiro et al., 2021). The Pantanal acts as a large reservoir,  
 88 causing a lag of up to 5 months between inflows and outflows; summer rains determine the  
 89 flooding seasons, being between November and March in the north and May and August in  
 90 the south (Marengo et al., 2021). The Atlantic Forest is the biome which includes the major  
 91 urban centres of São Paulo and Rio de Janeiro. The Atlantic Forest is extremely  
 92 heterogenous in composition, spanning from 4°–32° S and covers a wide range of climates  
 93 and vegetation formations. Elevation ranges from sea level to 2900 m and there can be  
 94 abrupt changes in rainfall, soil moisture, depth, and type as well as average air temperature  
 95 (Tabarelli et al., 2005). The Pampa biome located in South Brazil has a wet subtropical  
 96 climate, characterized by high precipitation throughout the year with hot summers and cold  
 97 winters, and mostly consists of grasslands with some small patches of forest (Lopes Ribeiro  
 98 et al., 2021, Overbeck et al., 2007).

99 Agricultural land is unevenly distributed across each biome with large areas of soybean  
100 farmland found in the Cerrado biome and the boundary between the Atlantic Forest and  
101 Pampa biomes. The boundary between Cerrado and the southern part of the Amazonia  
102 biome also shows large areas of expanding soybean producing land which has been shown  
103 to be of ecological concern (Song et al., 2021, Barona et al., 2010). Unlike the previous  
104 study by Gallear et al. (2025), the soybean growing area is not used to filter  
105 agrometeorological data to exclude grid cells without significant soybean farming activity.  
106 This decision was made to ensure inclusion of areas which may be used by small holder  
107 farmers in the analysis and to minimise bias of region selection towards large industrial scale  
108 farms. Although there are some changes in planting windows between the regions, October  
109 to December is the common planting window across all five regions in Figure 1 (c). For this  
110 reason, October to December are the most important months for risk to agricultural  
111 productivity, particularly because soybean exports are a significant driver of economic  
112 growth for the Brazilian economy (Cattelan and Dall’Agnol, 2018). Therefore, some analysis  
113 of model forecast performance should focus on this period of months.

### 114 3.2. Data

115 Data used in this study were either obtained from satellite products such as the Climate  
116 Hazards Center Infrared Precipitation with Station data (CHIRPS) (Funk et al., 2015) or from  
117 reanalysis data (such as ECMWF’s ERA5 (Hersbach, 2023)) (Table 1). Table 1 shows the  
118 data sources with the abbreviations used henceforth throughout the text. The variables of  
119 interest are derived from the data sources described in Table 1 at the 0.25° grid scale across  
120 all biomes shown in Figure 1. The variables were split into separate datasets according to  
121 biome with a separate model trained on each. Each grid cell in the dataset is a time series of  
122 months running from February 2003 to December 2021. 2 metre Temperature (t2m),  
123 potential evaporation (pev), and longwave radiation (longrad) data was obtained from the  
124 monthly averaged ERA5 reanalysis database. The reason for including these three variables  
125 in the model input data is to capture temperature effects which may result in drought  
126 conditions. T2m is defined as temperature two metres above the surface of the land,  
127 Potential evaporation (pev) is defined as a measure of the extent to which the near surface  
128 atmospheric conditions are conducive to evaporation, and mean surface downward  
129 longwave radiation flux (longrad) is included to capture potential heat effects. ERA5 is a  
130 reanalysis database which combines climate model data with observations using data  
131 assimilation to provide better estimates of meteorological variables at the grid scale  
132 (Hersbach, 2023).

133  
134  
135  
136  
137  
138  
139  
140  
141  
142

143  
144  
145

**Table 1.** Variables used in this study with data sources and abbreviations used throughout the text.

Variable (units)	Abbreviation	Source	Usage
Precipitation (mm/month)	Precip	CHIRPS	Feature
2 metre temperature (K)	T2M	ERA5	Feature
Potential evaporation (kg m <sup>-2</sup> )	pev	ERA5	Feature
Mean surface downward long-wave radiation flux (W m <sup>-2</sup> )	longrad	ERA5	Feature
Root Zone Soil Moisture (%)	RZSM	NASA GRACE	Feature
Vegetation health index (%)	VHI	NOAA STAR	Next month's value used to derive dependent variable
Standardized evapotranspiration-precipitation index (Unitless)	SPEI	Calculated from CHIRPS precipitation data and ERA5 data	Feature

146

### 147 3.3. Vegetation health index (VHI)

148 The vegetation health index (VHI) is a satellite-based proxy for estimating vegetation health,  
149 with values below 40% representing stress conditions. The VHI is a weighted sum of the  
150 vegetation condition index (VCI) and the temperature condition index (TCI), determined  
151 through the following formulae:

$$152 \quad VCI = \frac{100(NDVI - NDVI_{min})}{NDVI_{max} - NDVI_{min}} \#(1)$$

$$153 \quad TCI = \frac{100(BT_{max} - BT)}{BT_{max} - BT_{min}} \#(2)$$

$$154 \quad VHI = \alpha VCI + (1 - \alpha) TCI \#(3)$$

155 where BT is the brightness temperature recorded from a thermal sensor, max/min  
156 represents the maximum and minimum values of a variable over the study period and  $\alpha$  is a  
157 fixed coefficient used to determine the relative contribution of VCI and TCI to VHI. NDVI is  
158 the normalized Difference Vegetation Index. The VHI data is obtained from the NOAA STAR  
159 satellite-based vegetation health system based on the Advanced Very High Resolution  
160 Radiometer (AVHRR) found on NOAA polar orbiting satellites (Kogan, 1997). The data is  
161 upscaled from 0.036° to 0.25° spatial resolution.

### 162 **3.4. Root Zone Soil Moisture (RZSM)**

163 Root Zone Soil Moisture (RZSM) is used to determine the propagation of meteorological  
164 drought to soil moisture drought affecting plants (Gallear et al., 2025). Root zone refers to  
165 the top metre of soil, and RZSM is estimated using the NASA GRACE satellite (Li et al.,  
166 2019). The NASA GRACE satellite data are based on two satellites which record changes in  
167 the earth's gravity field caused by the redistribution of water. Similar to the previous study  
168 (Gallear et al., 2025), RZSM is averaged from weekly to monthly timescales to align with  
169 monthly vegetation health index data.

### 170 **3.5. Total monthly Precipitation (Precip)**

171 Precipitation data is obtained from the Climate Hazards Group Infrared precipitation with  
172 station data (CHIRPS) database (Funk et al., 2015). CHIRPS is a dataset which combines  
173 satellite data with in-situ measurements to provide a gridded dataset which has been used in  
174 many large-scale studies such as (Gallear et al., 2025, Lees et al., 2022, de Oliveira - Júnior  
175 et al., 2021).

### 176 **3.6. Standardized Precipitation-evapotranspiration Index (SPEI)**

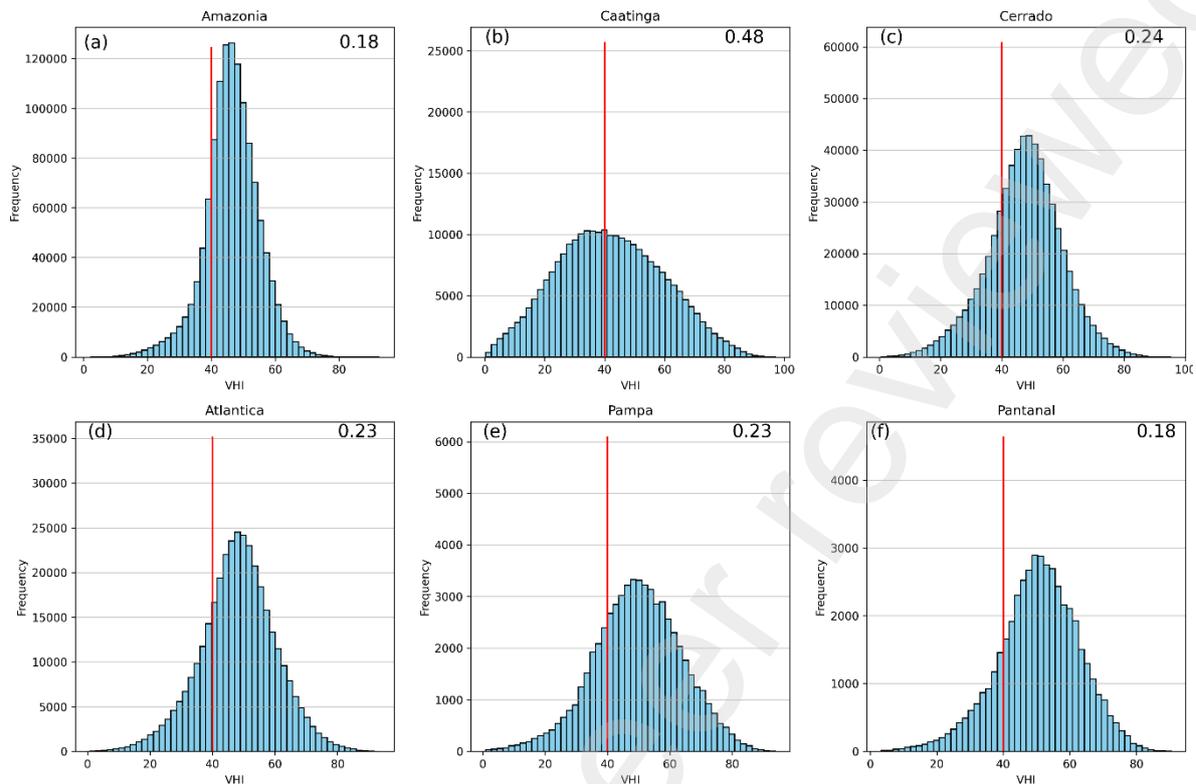
177 Standardised precipitation index (SPI) is the most widely adopted metric for monitoring  
178 meteorological drought (Svoboda and Fuchs, 2017). It's a measure of the precipitation  
179 relative to the climatology for that specific location and time of year. SPEI builds upon SPI by  
180 integrating potential evapotranspiration (PET) data to calculate a climatic water balance  
181 using the difference between precipitation and PET (Vicente-Serrano et al., 2010). In this  
182 study, daily PET is estimated using the Penman-Monteith method, as recommended by the  
183 United Nations' Food and Agriculture Organisation, and subtracted from the daily  
184 precipitation to calculate the climatic water balance (Allen et al., 1998). The climatic water  
185 balance is aggregated over running 30-, 60- and 90-day periods, with each then fitted to a  
186 log-logistic distribution to calculate daily standardised values, which were then used to  
187 calculate the SPEI (Vicente-Serrano et al., 2010). Daily SPEI values were finally averaged  
188 over monthly periods to get one, two- and three-month SPEI values used in this study. The  
189 base period used for this dataset to define the climatology is 1995–2014, reflecting the IPCC  
190 AR6 base period (IPCC, 2021). We used SPEI in this work rather than SPI due to stronger  
191 monthly correlations with VHI at the grid scale across Brazil, particularly in the Northeast  
192 Caatinga and central Midwest regions, and at shorter accumulation periods (Gallear et al.,  
193 2025).

### 194 **3.7. Definition of drought stress**

195 The accepted VHI threshold to signify serious drought stress on vegetation is 40%, this  
196 threshold has been applied globally (Kogan et al., 2013). This threshold stems from  
197 reductions in VCI and TCI being indicative of reductions in maize yield of more than 50%  
198 (Kogan, 1997). This value is therefore used to issue drought warnings and so can be  
199 interpreted as the threshold at which meteorological drought propagates to impact  
200 vegetation health to which we infer agricultural impacts such as yield losses (Kogan et al.,  
201 2013, Kogan, 1997, Gidey et al., 2018, Kloos et al., 2021, Dalezios et al., 2014).

202 This absolute definition of drought stress (rather than a relative value based on the  
203 distribution of VHI at each biome) leads to thresholds at different extremities of each dataset.  
204 Figure 2 shows the location of the threshold in relation to the monthly VHI distribution for  
205 each of the biomes. Most important to note is that the threshold of drought stress is closest  
206 to the median value of VHI for the Caatinga biome. This means that the ratio of drought to

207 non-drought cases is most equal in the Caatinga, meaning that drought occurs nearly half of  
208 all instances in this biome.



**Figure 2.** Spatio-temporal distributions of monthly VHI for each biome across both time and space showing the threshold value (40%) as a red vertical line through each histogram (a-f). The ratio of drought to non-drought cases across the spatio-temporal distribution of VHI for each biome is shown in the top right corner of each panel with drought being defined by the values which fall below the 40% threshold.

209

210 The position of the 40% threshold in relation to each biome results in differing proportions of  
211 drought cases in each biome. The drought to non-drought ratio is highest in Caatinga (0.48),  
212 with the second highest ratio (Cerrado) being half that of the Caatinga (0.24). Amazonia and  
213 Pantanal jointly have the lowest drought to non-drought ratios.

### 214 3.8. Machine Learning methods

215 Following Gallear et al. (2025), we used a random forest method to forecast the probability  
216 of drought stress for the subsequent month (defined as mean monthly VHI being below  
217 40%). Random forest is an ensemble method first developed by (Breiman, 2001). The  
218 method partitions data into subsets based on conditions at each leaf node of the tree.  
219 Random forest constructs a specified number of trees then averages the result of each  
220 individual tree (Marsland, 2011, Breiman, 2001). There are several advantages to using tree-  
221 based methods over more complex ones such as neural networks, the main advantages  
222 being better interpretability, greater robustness to outliers, and simpler calibration (Delerce et  
223 al., 2016, Gallear, 2023).

224 Spatial and per biome model validation was undertaken using a leave-one-out approach in  
225 which each model is trained and tested successively with a leave-one-year-out approach,  
226 The test scores from each set are then used to produce the distribution of test scores shown  
227 in Figure 3. This is done to provide a better overall estimate of model performance and an

228 estimate of model stability across different training and testing datasets. Hyperparameters  
229 were left at default values to reduce risk of over-tuning and provide a more general estimate  
230 of model performance. A year-based train and test split prevents spatial information leakage  
231 from training to testing set. Forecast model performance is evaluated using several  
232 classification metrics. We use accuracy, recall, precision, F1 score as well as True positive  
233 rate and false positive rate. F1 score is the harmonic mean of precision and recall; this is  
234 especially important for the grid cell level evaluation which contains many imbalanced grid  
235 cells with much fewer drought instances than non-drought instances. These metrics are  
236 defined by the following formulae:

$$237 \quad \text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \#(4)$$

$$238 \quad \text{Precision} = \frac{TP}{TP + FP} \#(5)$$

$$239 \quad \text{Recall} = \frac{TP}{TP + FN} \#(6)$$

$$240 \quad F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{precision} + \text{Recall}} \#(7)$$

241

242 Where TP is true positives, TN is true negatives, FP is false positives, and FN is false  
243 negatives.

244

245 To determine the effect of changing the drought definition threshold on model performance,  
246 for each biome we produce Receiver Operating Characteristic (ROC) curves using a range  
247 of drought thresholds. The ROC is a standard technique for summarizing classifier  
248 performance across the trade-offs between true positive and false positive rates (Chawla,  
249 2010). An ROC curve is a line plot of true positive rate versus false positive rate. The greater  
250 the area under the curve, the greater the model performance is relative to incorrectly  
251 forecasting droughts (i.e. The greater the ratio of true positives to false positives). If model  
252 results are plotted as a diagonal line across with a 1-1 ratio of true positives to false positives  
253 this would show the model performance is the same as random chance. ROC curve analysis  
254 was performed using training data from 2003 to 2018 with a testing period of 2018 to 2021.

255 SHAP (with associated Shapley values) is a feature attribution method used to determine the  
256 contribution of each feature to the model output. SHAP is a computational method which  
257 assigns proportional values to features depending on the influence of each feature on the  
258 prediction (Molnar, 2025). SHAP is computationally expensive and complex and so requires  
259 approximations and large computational resources (Chen et al., 2023). However, it is also a  
260 model agnostic method which, crucially, considers the influence of correlations between  
261 features (Rodríguez-Pérez and Bajorath, 2020). This is especially important for  
262 meteorological features such as rainfall and solar radiation which will likely be correlated.

### 263 **3.9. Empirical Orthogonal Function (EOF) analysis**

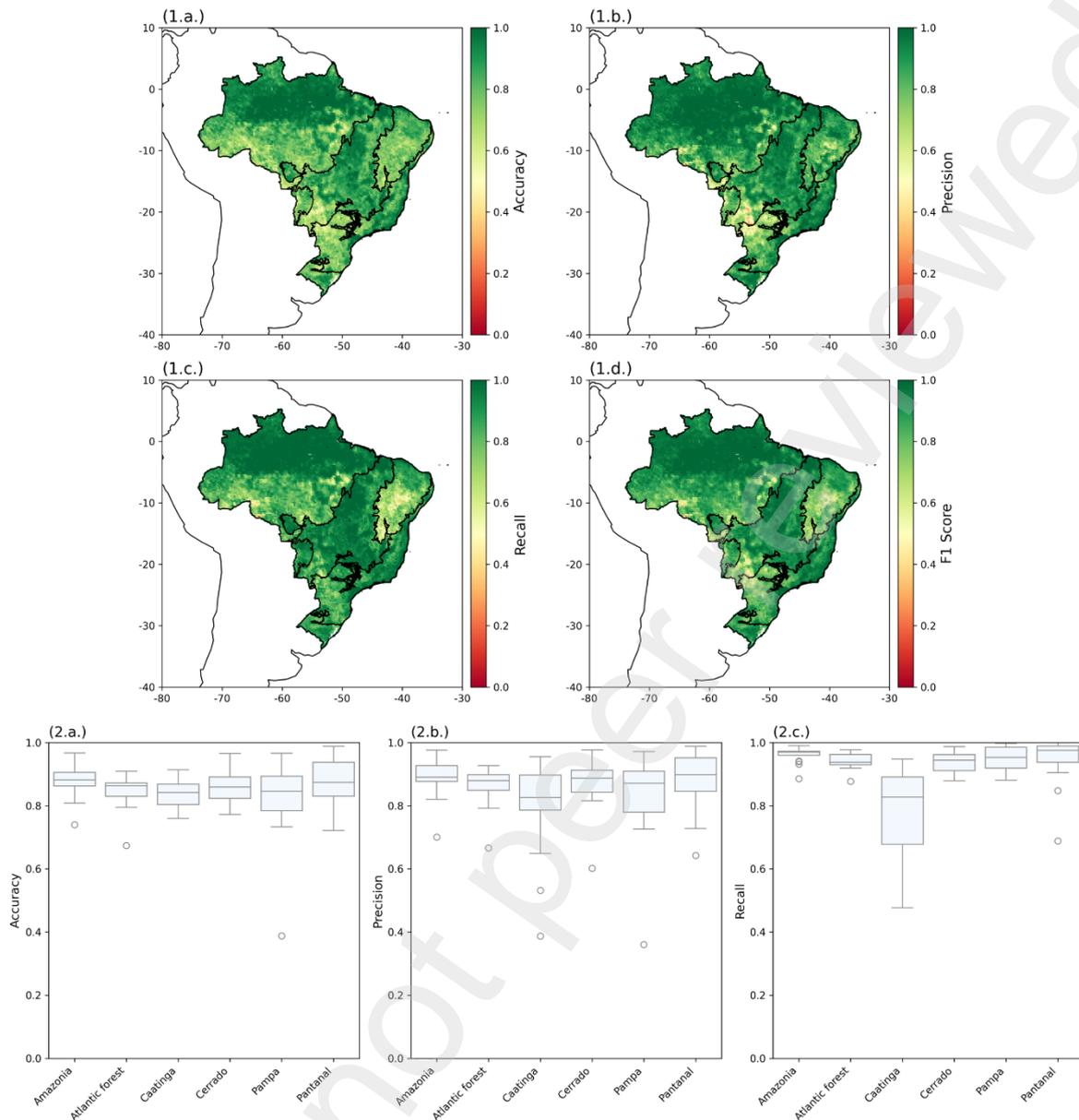
264 We use EOF analysis to describe how the spatio-temporal patterns of VHI, RZSM and  
265 precipitation may vary, providing indication for the reasons of the variation in forecast model  
266 performance. EOF analysis has been previously applied to climate and agronomic data to  
267 detect and correlate dominant patterns of variability (Challinor et al., 2003). EOF analysis  
268 (also referred to as Principal Component Analysis) reduces a larger set of variables to a new

269 smaller set of variables which are linear combinations of the original ones, chosen to  
270 represent the maximum possible fraction of the variability contained in the original data. EOF  
271 analysis has the potential for great insights into both the spatial and temporal variations in  
272 the data depending on the nature of the linear combinations which are most effective at  
273 compressing the data (Wilks, 2006).

## 274 **4. Results**

### 275 **4.1. Classification metrics and model validation**

276 The overall accuracy of the random forest models for each biome is between 0.8 and 1.  
277 Some biomes have a much higher proportion of non-drought cases than others (Figure 2).  
278 Therefore, to gain a full picture of model performance, recall and precision are also shown  
279 per biome and grid cell (Figure 3). Recall and precision are combined into the F1 score  
280 (Equation 7) and is shown with the other metrics in Figure 3. Results from Figure 3 show  
281 increased variability of model performance metrics in Caatinga, but high overall grid cell level  
282 metrics for all biomes across Brazil. High precision and recall mean that models are not  
283 significantly biased towards predicting false positives or false negatives. Greater variability in  
284 recall in the Caatinga biome indicates that in some years, more droughts are incorrectly  
285 forecasted than others. The maps also show that the central Caatinga has a lower recall  
286 score (between 0.6 and 0.4) than the rest of the biome and much of the Cerrado, this  
287 indicates that droughts are overestimated in this location. There is also a marked difference  
288 between model performance in the north of the Amazonia as compared to south Amazonia.  
289 Cross validation shows that median model performance for all metrics is between 0.8 and 1  
290 showing high consistency in model performance across test years.

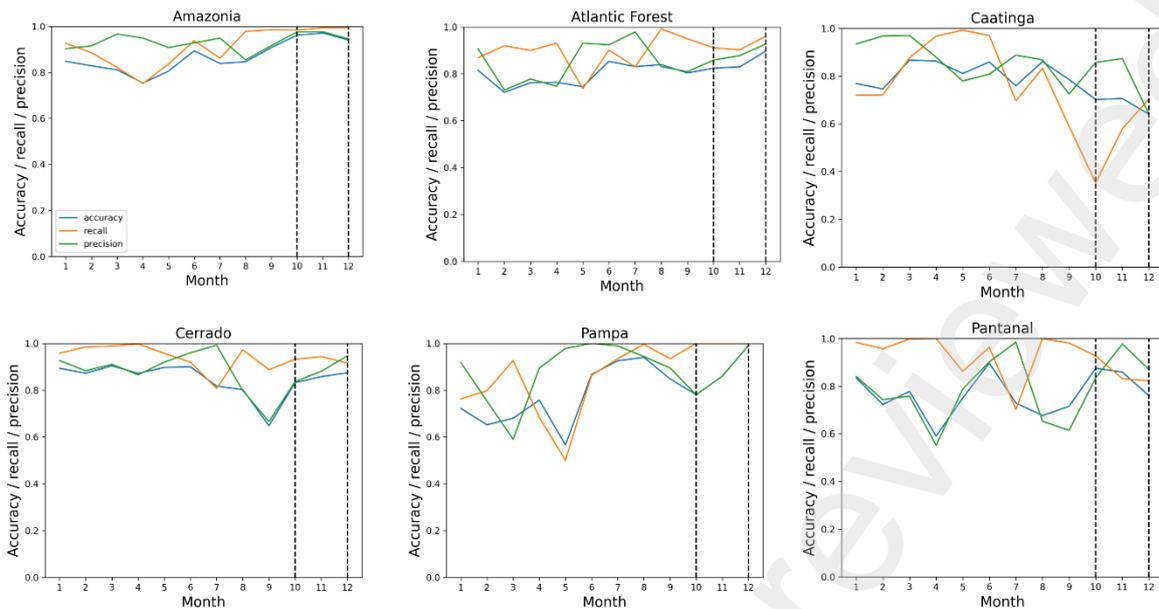


**Figure 3.** Accuracy (1.a.) precision (1.b.), recall (1.c.) and f1 score (d) of the random forest forecast model for predicting stress condition in VHI data 1 month in advance. Test data is for the last 4 years of the time series (2018 – 2021). Accuracy (2.a.), precision (2.b.) and recall (2.c.) of random forest models using leave one out cross validation for each year (2003 – 2021)

291

292 Figure 4 shows the monthly consistency of forecast accuracy, recall and precision, allowing  
 293 us to determine the usefulness of forecasts throughout the soybean growing season and  
 294 contrast this with the rest of the year. Typically, recall, accuracy and precision are consistent  
 295 throughout the year regardless of month of forecast, though there are some variability and  
 296 anomalous predictions. Dotted vertical lines which represent a typical planting window  
 297 across Brazil do not contain significantly different values for model performance than the rest  
 298 of the year in most cases (with October in Caatinga being the exception). Model  
 299 performance during this planting window tends to be in the upper end of the annual range for  
 300 Amazonia, Atlantic Forest, and Pampa, but is in the lower end for Caatinga.

301

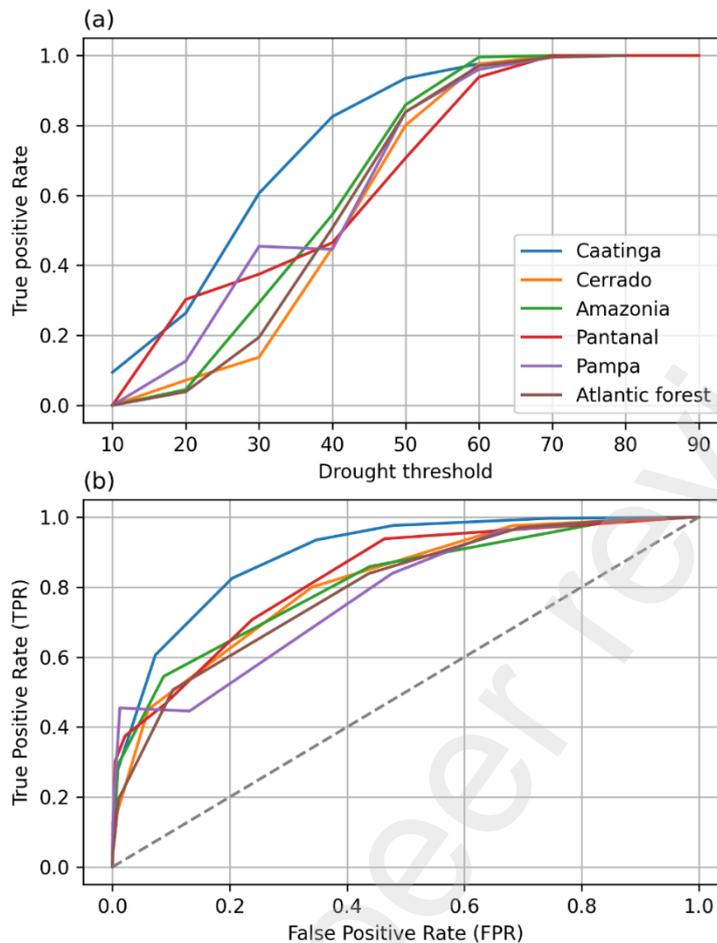


**Figure 4.** Line plots of model accuracy, recall and precision across each month for the 4-year test period for the six biomes. Dotted lines show Start of October to December period which is typically the common period of months for Soybean planting across Brazil.

302

## 303 4.2. ROC curve analysis

304 ROC curves are here used for determining the effect of changing the drought definition  
 305 threshold on the trade-off between the true positive and false positive error rates. Therefore,  
 306 we can examine to what extent more extreme definitions of drought can be forecasted using  
 307 the same modelling set up for each biome. The greater true positive rate (shown in panel (a)  
 308 of Figure 5, along with the larger area under the curve in panel (b) indicate both a higher true  
 309 positive rate at lower thresholds and a greater true positive rate in relation to the number of  
 310 false positives. panel (a) of Figure 5 shows how the drought threshold of 40% VHI has a  
 311 much higher true positive rate for the Caatinga biome than other biomes (80%). Panel (b)  
 312 shows how at the same time the ratio of true positives to false positives is also high. This  
 313 means that in Caatinga, models can forecast more severe droughts at a greater true positive  
 314 rate without incorrectly forecasting more droughts at the same time. Other biomes show very  
 315 comparable results in terms of true positive rates at different thresholds, and the ratio of true  
 316 positives to false positives at different thresholds. From this we can say that the Caatinga  
 317 biome stands out as a biome which is unique in the ability of models to forecast severe  
 318 droughts.



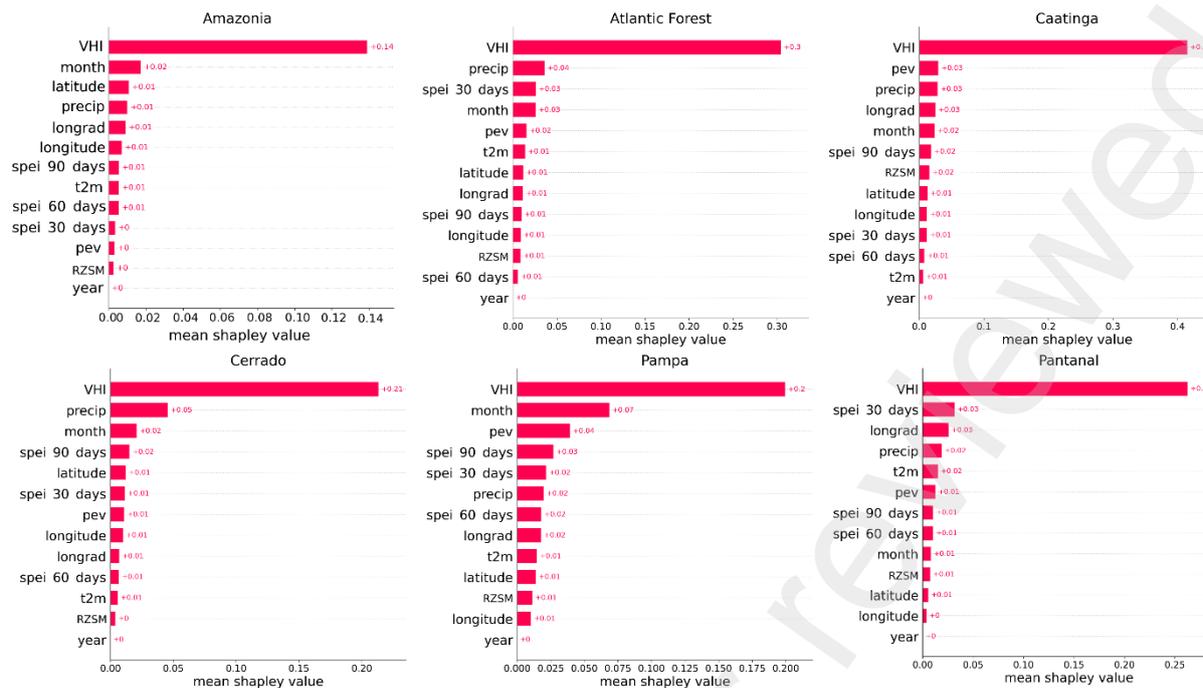
**Figure 5.** (a) Rate of true positives against threshold VHI value for each biome, (b) True positive rate against False positive rate for each drought threshold value.

319

### 320 4.3. Shapley values and spatial relationships

321 Shapley values explain how much each feature contributes to the variation in the modelled  
 322 forecasts for each biome (Figure 6). Lagged VHI dominates the contribution to model  
 323 forecasts across all biomes (ranging from 1.4 to 3 average relative contribution across  
 324 biomes). High contribution from VHI indicates the slow inertia of VHI which changes  
 325 gradually month to month. The relative contributions of other variables depend on which  
 326 biome model forecasts are trained upon. But month of the year (month) and total monthly  
 327 precipitation (precip) often make relatively high contributions to model forecasts across  
 328 multiple biomes. SPEI often appears in the top four of the features when ordered by  
 329 importance for model forecast. Which accumulation period is most important depends on the  
 330 biome, with sometimes 90 days and sometimes 30 days being a higher contributing feature.

331



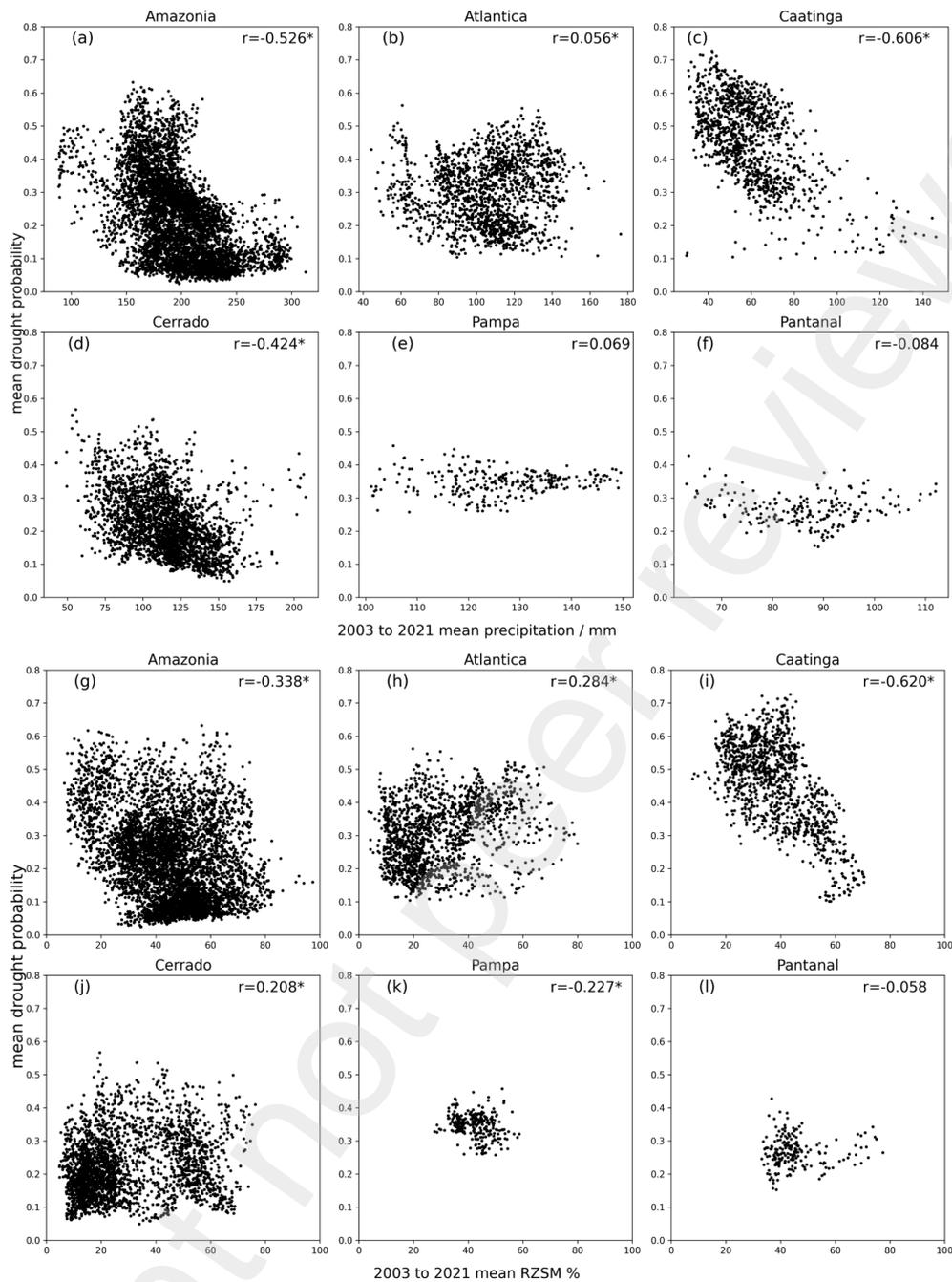
**Figure 6.** Shapley plots for each of the six biomes in Brazil.

332

333 Figure 7 shows scatter plots of mean drought probability against monthly total precipitation  
 334 (a–f) and against RZSM (g–l). Each variable has been averaged over time to produce one  
 335 value per grid cell. This allows for an understanding of the potential spatial relationships  
 336 between drought and climate for each biome. Precipitation data indicates that the strongest  
 337 spatial relationships between average rainfall and probability of drought are for the Caatinga  
 338 and Cerrado biomes. The Amazonia biome shows a medium Pearson's correlation  
 339 coefficient value (-0.526) however the relationship is not clear by eye. This may be because  
 340 the correlation value is heavily influenced by the density of data points in this biome. Other  
 341 biomes show much weaker relationships between these two variables. There is a correlation  
 342 coefficient of -0.606 between RZSM and drought probability in Caatinga, however other  
 343 biomes do not show strong correlations. There is a high degree of confidence in the r values  
 344 reported for many of the correlations, this is because the number of data points used for the  
 345 correlation is quite high.

346 The plots clearly show differences in the spatial variability in rainfall across biomes. In  
 347 Pampa and Pantanal, the range of rainfall values throughout space is small, whereas in  
 348 Cerrado, Caatinga and Amazonia, there are much larger ranges in the precipitation values.  
 349 Moreover, these latter biomes also exhibit stronger correlations between average drought  
 350 probability and average precipitation. Atlantic forest also has a larger range in precipitation  
 351 values across space but shows little relationship between average rainfall and average  
 352 drought probability. RZSM values show a similar pattern with lower spatial variability in  
 353 Pampa and Pantanal, higher spatial variability in Cerrado, Caatinga and Amazonia, but only  
 354 Caatinga shows a strong correlation between spatial RZSM values and average drought  
 355 probability.

356

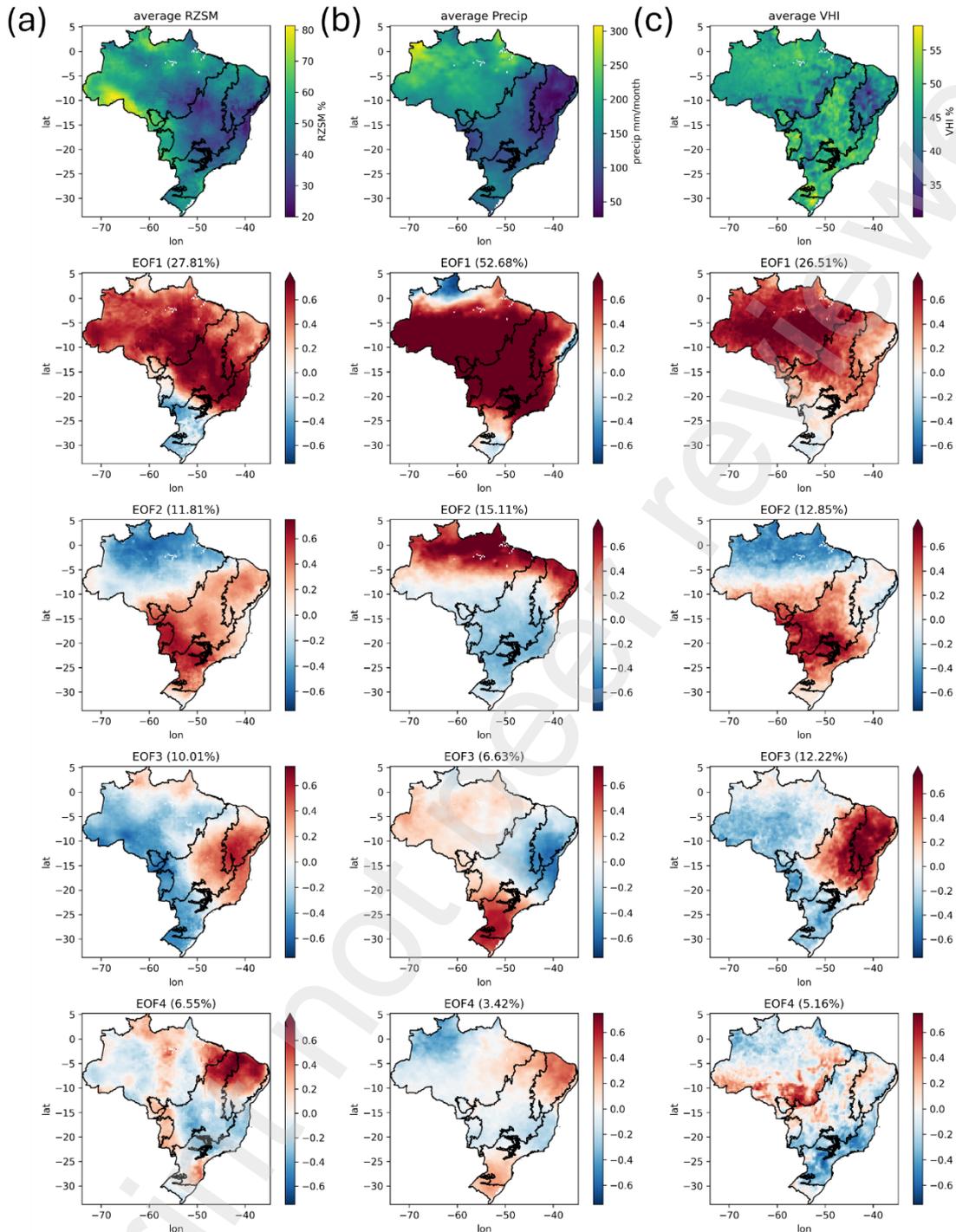


**Figure 7.** Scatter plots of average precipitation (a-f) quantities and RZSM percentages (g-l) against probability of drought threshold exceedance per grid cell. Data is shown separately for the 6 biomes of Brazil. Starred  $r$  values represent greater than a 95% degree of confidence.

357

#### 358 4.4. Empirical Orthogonal Function (EOF) analysis

359 The EOF analysis provides insights into explaining the differences in drought stress  
 360 forecasts across biomes and how propagation of drought can have different spatial  
 361 dependencies. Figure 8 shows the results of the analysis for RZSM, precipitation and VHI.



**Figure 8.** Average and 4 principal components/EOFs for RZSM (a), Precip (b) and VHI (c). The first row shows the 2003-2021 average values, followed by EOFs 1–4 for each variable in subsequent rows. Labelled percentages refer to the percentage variation of the variable explained by each EOF.

362

363 Dominant patterns of variation among the three satellite product derived variables are  
 364 captured by the first EOF modes, with decreasing percentages of variation captured with  
 365 increasing EOF mode (Figure 8). EOF 1 for each variable captures a pattern with southern  
 366 Brazil (most notably the Pampa biome) varying distinctly from the rest of the country. EOF 2  
 367 for each variable captures a clear distinction between north and south. RZSM EOF2 also

368 shows a pattern in the amazon distinct from the rest of Brazil. EOF 3 presents a distinction  
369 between northeast Brazil (where the Caatinga biome and parts of the Cerrado are located)  
370 from the rest of Brazil. RZSM, precipitation and VHI show very similar patterns for EOF 3  
371 however there are some very minor differences in the pattern of VHI EOF 3 from the other 2  
372 variables (namely a stronger positive mode in the northeast tip of the Caatinga biome). EOF  
373 4 explains less proportion of the variation for each variable and less similarities are shown  
374 between the three variables for EOF 4. EOF 4 of RZSM and precipitation show similar  
375 patterns of variation in the northeast, this contrasts with EOF 4 of VHI which is quite different  
376 and shows some patterns of variation around the edges of the amazon rainforest in the  
377 south, and east of the Amazonia biome. Further EOFs after 4 would capture less and less  
378 variation of the data however these were not plotted or analysed as the importance of these  
379 EOF modes would be less.

## 380 **5. Discussion**

381 Overall, accuracy, precision, and recall showed high values (above 0.8) for all biomes. The  
382 Caatinga biome had a lower recall median of around 0.8, but with greater spatial variability,  
383 with lower values in the central and southern parts of the biome. High accuracy and  
384 precision but slightly larger variability in recall in the Caatinga biome indicates that the model  
385 over predicts drought in some cases. Results in this biome merit closer examination as it's  
386 also more susceptible to droughts. For future forecasting work it is important to note the  
387 (Tomasella et al., 2025) study which demonstrated changes in the aridity index towards  
388 more arid conditions in the Caatinga. Future work may seek to understand if longer term  
389 trends towards increased aridity may affect the performance of agricultural drought  
390 forecasts.

391 When comparing the differences between how each of the biomes affect agricultural drought  
392 stress probability at different thresholds, ROC curve analysis shows a true positive rate of  
393 0.4 to 0.5 at the 40% threshold for all biomes excepted Caatinga, where this rate reaches  
394 0.8. When the true positive rate is 0.4 this results in a very low false positive rate of below  
395 0.1 for all biomes. For the Caatinga 0.8 true positive rate results in a false positive rate of  
396 around 0.2, and so errors are minimised. This indicates that the threshold of 40% is  
397 particularly useful for forecasting drought stress conditions for the next month in this biome.  
398 This result is most likely associated with the characteristics of the Caatinga species, which  
399 are well adapted to long periods without precipitation due to the short rainy season in most  
400 of the region (usually from February to May). Plants in this biome respond quickly to  
401 precipitation pulses and dry spells (Medeiros et al., 2022, Santos e Silva et al., 2024).  
402 Results from other biomes indicate lower drought stress probability at lower thresholds  
403 indicating that more severe droughts are not just less common but also harder to forecast.

404 Propagation of meteorological drought to agricultural drought is an important question which  
405 could be used for earlier warning of the impacts of drought on crop health depending on  
406 weather forecasts. Here, we show that the relationship between precipitation and agricultural  
407 drought stress probability varies per biome, with some biomes such as the Pantanal showing  
408 no correlation between rainfall and agricultural drought but others such as the Caatinga  
409 showing a much stronger relationship. Strong correlation between rainfall and vegetation  
410 health in the northeast is also shown in (Gallear et al., 2025). Both the results in this study  
411 and the previous, show that drought more easily propagates from meteorological to  
412 agricultural drought in the northeast, this should make it easier to forecast agricultural  
413 drought in this region at longer lead times.

414 The Shapley plots show that VHI in the current month is the strongest predictor of drought  
415 stress in the following month. This indicates that VHI does not change significantly from  
416 month to month, capturing the slow progression of plant health. This result supports the  
417 inclusion of VHI in combined drought indices, since it measures a current impact on the  
418 vegetation status and a continuous influence on drought intensity the following month. After  
419 VHI, either precipitation and/or SPEI appear among the first three features in order of  
420 importance, indicating that the water balance components are crucial to assess the  
421 likelihood of that meteorological drought conditions will propagate to agricultural drought the  
422 following month. While precipitation accounts for the input of water, SPEI considers the  
423 balance of precipitation and evapotranspiration, and soil moisture conditions as an indirect  
424 assessment, relative to the climatological average during that time of year. This results in  
425 SPEI being more able to identify longer and more severe droughts than SPI, which has been  
426 shown in multiple studies, especially for arid and semi-arid regions (Tirivarombo et al., 2018,  
427 Mwinjuma et al., 2025, Lotfirad et al., 2022).

428 Correlation analysis of climate variables and probability of drought shows how the climate of  
429 different biomes can have varying strengths of relationship with probability of drought. RZSM  
430 shows a high correlation with average probability of drought across the Caatinga biome. This  
431 is likely because Caatinga soils are generally shallow (0–50 cm) which limits water infiltration  
432 and storage capacity, hence fluctuations in soil moisture are more likely to have a significant  
433 effect on vegetation health (Lopes Ribeiro et al., 2021). Given previous research which  
434 shows decreasing soil moisture trends across the biome (Lopes Ribeiro et al., 2021), this  
435 may also affect the vegetation health to a greater extent. Precipitation is strongly correlated  
436 with vegetation health both across the Caatinga and Cerrado. Again, shallow soils likely  
437 mean that precipitation takes a more rapid effect on vegetation health in the Caatinga biome.

438 Spatio-temporal dynamics of climate are indicated by the EOF analysis. This reveals that  
439 modes of variation in VHI, RZSM and precipitation are often correlated on a coarse scale  
440 across Brazil. Most notably, the northeast exhibits specific patterns of variation in VHI,  
441 RZSM and precipitation which are distinct to the rest of Brazil. This can be used to better  
442 understand how drought may propagate in different regions from meteorological drought to  
443 agricultural drought stress. EOF 3 shows common patterns of variation between rainfall,  
444 RZSM and VHI in the northeast, therefore indicating that the chance of propagation in this  
445 region will often be more likely. This finding coincides with higher correlations between  
446 SPEI, RZSM and VHI in the northeast (Gallear et al., 2025) and therefore indicates that  
447 forecasting systems should treat northeast Brazil as a separate system when training  
448 models to forecast VHI or the integrated drought index (IDI) (which uses VHI, SPI, and  
449 RZSM).

450 We have shown that the impact of meteorological drought varies across biomes. Agricultural  
451 land is unevenly distributed across each biome with large areas of soybean farmland found  
452 in the Cerrado and Pampa regions. Therefore, forecasting would likely be most important for  
453 soybean growth in the Cerrado region. However, most of the Caatinga's rural population  
454 depends on agriculture (José Maria Cardoso da Silva, 2019). For this reason, it is important  
455 to include regions which may not be so intensely farmed as agricultural drought forecasts  
456 are still relevant and useful to smaller scale farming operations. The analysis of model  
457 performance shows that monthly performance is consistent (despite anomaly of October in  
458 the Caatinga). Therefore, model performance in most biomes can be relied upon for critical  
459 months relevant to the cropping season calendar (with the example of soybean given here in  
460 this study in Figure 1). If future work is to improve the relevance of forecasts for agriculture,  
461 improved skill may be obtained by focusing on specific growing regions where sufficient data  
462 allows for robust monthly model performance.

## 463 6. Conclusions

464 Model evaluation in this work demonstrates strong performance across all Brazilian biomes,  
465 with accuracy, precision, and recall scores generally exceeding 0.8. However, the Caatinga  
466 biome requires special attention; while median performance is high, it exhibits significant  
467 spatial variability, reflecting its inherent drought susceptibility and greater prediction  
468 uncertainty, a concern underscored by recent studies showing increasing aridity trends  
469 (Tomasella et al., 2025). The most powerful predictor for next-month drought stress is the  
470 current month's Vegetation Health Index (VHI), owing to its slow-changing nature (inertia).  
471 High inertia in VHI allows for high skill forecasting stress for the subsequent month in other  
472 biomes such as Cerrado, Amazonia and Pampa. Strong forecasts are especially important  
473 for the Cerrado where most Soybean farming occurs, an important export for the Brazilian  
474 economy. Water balance components like precipitation, root zone soil moisture (RZSM), and  
475 SPEI are also identified as crucial drivers, with a greater number of droughts forecasted at  
476 more severe thresholds in Brazil's northeast Caatinga biome.

477 The distinct behaviour observed in the Caatinga is rooted in its unique regional climate  
478 patterns and physical geography. The biome's characteristically shallow soils limit water  
479 storage capacity, making vegetation health highly and rapidly responsive to changes in  
480 RZSM and precipitation. Large-scale climate analysis confirms that Northeast Brazil  
481 functions as a separate system regarding key climate variables. Therefore, the findings  
482 strongly support that drought forecasting models should treat Northeast Brazil as a distinct  
483 system, training models specifically for this region to improve the prediction accuracy of  
484 vegetation health and integrated drought indices.

## 485 Acknowledgements

486 This work was funded by the Met Office Climate Science for Service Partnership (CSSP)  
487 Brazil project under the International Science Partnerships Fund (ISPF).

488 The authors also acknowledge support from the Growing Health (BB/X010953/1) and  
489 AgZero+ (NE/W005050/1) Institute Strategic Programmes both funded by the BBSRC.

## 490 Data availability

491 All data are available from publicly available and free-to-access data repositories. ERA 5  
492 data were obtained from <https://doi.org/10.24381/cds.f17050d7> (Hersbach et al., 2023),  
493 NASA GRACE data were obtained  
494 from <https://doi.org/10.5067/UH653SEZR9VQ> (Beaudoing et al., 2021; Li et al., 2019),  
495 CHIRPS data were obtained from <https://data.chc.ucsb.edu/products/CHIRPS-2.0/> (Funk  
496 et al., 2014, 2015) and VHI data  
497 from [https://www.star.nesdis.noaa.gov/smcd/emb/vci/VH/vh\\_ftp.php](https://www.star.nesdis.noaa.gov/smcd/emb/vci/VH/vh_ftp.php) (NOAA, 2025; Kogan, 1  
498 997), SPEI data was calculated using a combination of ERA5 data and CHIRPS data.

## 499 7. References

- 500 ALLEN, R. G., PEREIRA, L. S., RAES, D. & SMITH, M. 1998. Crop evapotranspiration-  
501 Guidelines for computing crop water requirements-FAO Irrigation and drainage paper  
502 56. *Fao, Rome*, 300, D05109.
- 503 ALVARES, C. A., STAPE, J. L., SENTELHAS, P. C., GONÇALVES, J. D. M. & SPAROVEK,  
504 G. 2013. Köppen's climate classification map for Brazil. *Meteorologische zeitschrift*,  
505 22, 711–728.

- 506 BARONA, E., RAMANKUTTY, N., HYMAN, G. & COOMES, O. T. 2010. The role of pasture  
507 and soybean in deforestation of the Brazilian Amazon. *Environmental Research*  
508 *Letters*, 5, 024002.
- 509 BOPAR. 2025. *Producers must pay attention to planting and harvesting deadlines, warns*  
510 *Mapa* [Online]. Available: [https://www.bopar.com.br/noticias-destaque/df-produtores-](https://www.bopar.com.br/noticias-destaque/df-produtores-devem-ficar-atentos-aos-prazos-de-plantio-e-colheita-das-culturas-alerta-mapa)  
511 [devem-ficar-atentos-aos-prazos-de-plantio-e-colheita-das-culturas-alerta-mapa](https://www.bopar.com.br/noticias-destaque/df-produtores-devem-ficar-atentos-aos-prazos-de-plantio-e-colheita-das-culturas-alerta-mapa)  
512 [Accessed].
- 513 BREIMAN, L. 2001. Random forests. *Machine learning*, 45, 5–32.
- 514 CATTELAN, A. J. & DALL'AGNOL, A. 2018. The rapid soybean growth in Brazil. *Ocl*, 25,  
515 D102.
- 516 CHALLINOR, A., SLINGO, J., WHEELER, T., CRAUFURD, P. & GRIMES, D. 2003. Toward  
517 a combined seasonal weather and crop productivity forecasting system:  
518 determination of the working spatial scale. *Journal of Applied Meteorology*, 42, 175–  
519 192.
- 520 CHAWLA, N. V. 2010. Data mining for imbalanced datasets: An overview. *Data mining and*  
521 *knowledge discovery handbook*, 875–886.
- 522 CHEN, H., COVERT, I. C., LUNDBERG, S. M. & LEE, S.-I. 2023. Algorithms to estimate  
523 Shapley value feature attributions. *Nature Machine Intelligence*, 5, 590–601.
- 524 CIRILO, J. A. 2008. Public water resources policy for the semi-arid region. *estudos*  
525 *avançados*, 22, 61–82.
- 526 CUNHA, A. P. M., ZERI, M., DEUSDARÁ LEAL, K., COSTA, L., CUARTAS, L. A.,  
527 MARENGO, J. A., TOMASELLA, J., VIEIRA, R. M., BARBOSA, A. A. &  
528 CUNNINGHAM, C. 2019. Extreme drought events over Brazil from 2011 to 2019.  
529 *Atmosphere*, 10, 642.
- 530 DALEZIOS, N., BLANTA, A., SPYROPOULOS, N. & TARQUIS, A. 2014. Risk identification  
531 of agricultural drought for sustainable agroecosystems. *Natural Hazards and Earth*  
532 *System Sciences*, 14, 2435–2448.
- 533 DE OLIVEIRA-JÚNIOR, J. F., DA SILVA JUNIOR, C. A., TEODORO, P. E., ROSSI, F. S.,  
534 BLANCO, C. J. C., LIMA, M., DE GOIS, G., CORREIA FILHO, W. L. F., DE BARROS  
535 SANTIAGO, D. & DOS SANTOS VANDERLEY, M. H. G. 2021. Confronting CHIRPS  
536 dataset and in situ stations in the detection of wet and drought conditions in the  
537 Brazilian Midwest. *International Journal of Climatology*, 41, 4478–4493.
- 538 DELERCE, S., DORADO, H., GRILLON, A., REBOLLEDO, M. C., PRAGER, S. D., PATIÑO,  
539 V. H., GARCÉS VARÓN, G. & JIMÉNEZ, D. 2016. Assessing Weather-Yield  
540 Relationships in Rice at Local Scale Using Data Mining Approaches. *PLOS ONE*, 11,  
541 e0161620.
- 542 FERREIRA BARBOSA, M. L., HADDAD, I., DA SILVA NASCIMENTO, A. L., MÁXIMO DA  
543 SILVA, G., MOURA DA VEIGA, R., HOFFMANN, T. B., ROSANE DE SOUZA, A.,  
544 DALAGNOL, R., SUSIN STREHER, A. & SOUZA PEREIRA, F. R. 2022. Compound  
545 impact of land use and extreme climate on the 2020 fire record of the Brazilian  
546 Pantanal. *Global Ecology and Biogeography*, 31, 1960–1975.
- 547 FUNK, C., PETERSON, P., LANDSFELD, M., PEDREROS, D., VERDIN, J., SHUKLA, S.,  
548 HUSAK, G., ROWLAND, J., HARRISON, L. & HOELL, A. 2015. The climate hazards  
549 infrared precipitation with stations—a new environmental record for monitoring  
550 extremes. *Scientific data*, 2, 1–21.
- 551 GALLEAR, J. W. 2023. *Using machine learning and process-based crop modelling for*  
552 *regional scale prediction*. University of Leeds.
- 553 GALLEAR, J. W., VALADARES GALDOS, M., ZERI, M. & HARTLEY, A. 2025. Evaluation of  
554 machine learning approaches for large-scale agricultural drought forecasts to  
555 improve monitoring and preparedness in Brazil. *Natural Hazards and Earth System*  
556 *Sciences*, 25, 1521–1541.
- 557 GIDEY, E., DIKINYA, O., SEBEGO, R., SEGOSEBE, E. & ZENEBE, A. 2018. Analysis of the  
558 long-term agricultural drought onset, cessation, duration, frequency, severity and  
559 spatial extent using Vegetation Health Index (VHI) in Raya and its environs, Northern  
560 Ethiopia. *Environmental Systems Research*, 7, 13.

561 HERSBACH, H. 2023. ERA5 monthly averaged data on single levels from 1940 to present.  
562 Copernicus Climate Change Service (C3S) Climate Data Store (CDS), Electronic  
563 resource: 10.24381/cds. f17050d7 (2023). Accessed 2023-03-01.

564 IPCC 2021. Climate Change 2021: The physical Science Basis. Contribution of working group  
565 I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change.

566 JOSÉ MARIA CARDOSO DA SILVA, I. R. L., MARCELO TABARELLI (ed.) 2019. *Caatinga:*  
567 *The largest tropical dry forest region in South America*

568 KLOOS, S., YUAN, Y., CASTELLI, M. & MENZEL, A. 2021. Agricultural drought detection  
569 with MODIS based vegetation health indices in southeast Germany. *Remote*  
570 *Sensing*, 13, 3907.

571 KOGAN, F., ADAMENKO, T. & GUO, W. 2013. Global and regional drought dynamics in the  
572 climate warming era. *Remote Sensing Letters*, 4, 364–372.

573 KOGAN, F. N. 1997. Global drought watch from space. *Bulletin of the American*  
574 *meteorological society*, 78, 621–636.

575 LEES, T., TSENG, G., ATZBERGER, C., REECE, S. & DADSON, S. 2022. Deep learning  
576 for vegetation health forecasting: a case study in Kenya. *Remote Sensing*, 14, 698.

577 LI, B., RODELL, M., KUMAR, S., BEAUDOING, H. K., GETIRANA, A., ZAITCHIK, B. F., DE  
578 GONCALVES, L. G., COSSETIN, C., BHANJA, S. & MUKHERJEE, A. 2019. Global  
579 GRACE data assimilation for groundwater and drought monitoring: Advances and  
580 challenges. *Water Resources Research*, 55, 7564–7586.

581 LOPES RIBEIRO, F., GUEVARA, M., VÁZQUEZ-LULE, A., CUNHA, A. P., ZERI, M. &  
582 VARGAS, R. 2021. The impact of drought on soil moisture trends across Brazilian  
583 biomes. *Natural Hazards and Earth System Sciences*, 21, 879–892.

584 LOTFIRAD, M., ESMAEILI-GISAVANDANI, H. & ADIB, A. 2022. Drought monitoring and  
585 prediction using SPI, SPEI, and random forest model in various climates of Iran.  
586 *Journal of Water and Climate Change*, 13, 383–406.

587 MARENGO, J. A., CUNHA, A. P., ESPINOZA, J.-C., FU, R., SCHÖNGART, J., JIMENEZ, J.,  
588 C, COSTA, M., C, RIBEIRO, J., M, WONGCHUIG, S. & ZHAO, S. 2024. The Drought  
589 of Amazonia in 2023-2024. *American Journal of Climate Change*, 13, 567–597.

590 MARENGO, J. A., CUNHA, A. P., CUARTAS, L. A., DEUSDARÁ LEAL, K. R., BROEDEL,  
591 E., SELUCHI, M. E., MICHELIN, C. M., DE PRAGA BAIÃO, C. F., CHUCHÓN  
592 ANGULO, E., ALMEIDA, E. K., KAZMIERCZAK, M. L., MATEUS, N. P. A., SILVA, R.  
593 C. & BENDER, F. 2021. Extreme Drought in the Brazilian Pantanal in 2019–2020:  
594 Characterization, Causes, and Impacts. *Frontiers in Water*, Volume 3 - 2021.

595 MARENGO, J. A., GALDOS, M. V., CHALLINOR, A., CUNHA, A. P., MARIN, F. R., VIANNA,  
596 M. D. S., ALVALA, R. C., ALVES, L. M., MORAES, O. L. & BENDER, F. 2022.  
597 Drought in Northeast Brazil: A review of agricultural and policy adaptation options for  
598 food security. *Climate Resilience and Sustainability*, 1, e17.

599 MARENGO, J. A., TORRES, R. R. & ALVES, L. M. 2017. Drought in Northeast Brazil—past,  
600 present, and future. *Theoretical and Applied Climatology*, 129, 1189–1200.

601 MARSLAND, S. 2011. *Machine learning: an algorithmic perspective*, Chapman and  
602 Hall/CRC.

603 MEDEIROS, R., ANDRADE, J., RAMOS, D., MOURA, M., PÉREZ-MARIN, A. M., DOS  
604 SANTOS, C. A., DA SILVA, B. B. & CUNHA, J. 2022. Remote sensing phenology of  
605 the Brazilian caatinga and its environmental drivers. *Remote Sensing*, 14, 2637.

606 MIYAMOTO, B. C. B. 2024. *Estimation of agricultural revenue losses due to droughts in the*  
607 *state of Rio Grande do Sul*. Statistics, Federal  
608 University of Rio Grande do Sul.

609 MOLNAR, C. 2025. *Interpretable machine learning*.

610 MWIJUMA, M., WANG, R., MTUPILI, M. & TWAHA, M. 2025. Comparisons of SPI and  
611 SPEI in capturing drought dynamics: A Global assessment across arid and humid  
612 regions. *Atmospheric Research*, 108475.

613 OVERBECK, G. E., MÜLLER, S. C., FIDELIS, A., PFADENHAUER, J., PILLAR, V. D.,  
614 BLANCO, C. C., BOLDRINI, I. I., BOTH, R. & FORNECK, E. D. 2007. Brazil's

615 neglected biome: the South Brazilian Campos. *Perspectives in Plant Ecology,*  
616 *Evolution and Systematics*, 9, 101–116.

617 RODRÍGUEZ-PÉREZ, R. & BAJORATH, J. 2020. Interpretation of machine learning models  
618 using shapley values: application to compound potency and multi-target activity  
619 predictions. *Journal of computer-aided molecular design*, 34, 1013–1026.

620 SANTOS E SILVA, C. M. S., BEZERRA, B. G., MENDES, K. R., MUTTI, P. R.,  
621 RODRIGUES, D. T., COSTA, G. B., DE OLIVEIRA, P. E. S., REIS, J., MARQUES, T.  
622 V. & FERREIRA, R. R. 2024. Rainfall and rain pulse role on energy, water vapor and  
623 CO<sub>2</sub> exchanges in a tropical semiarid environment. *Agricultural and Forest*  
624 *Meteorology*, 345, 109829.

625 SONG, X.-P., HANSEN, M. C., POTAPOV, P., ADUSEI, B., PICKERING, J., ADAMI, M.,  
626 LIMA, A., ZALLES, V., STEHMAN, S. V. & DI BELLA, C. M. 2021. Massive soybean  
627 expansion in South America since 2000 and implications for conservation. *Nature*  
628 *sustainability*, 4, 784–792.

629 TABARELLI, M., PINTO, L. P., SILVA, J. M., HIROTA, M. & BEDE, L. 2005. Challenges and  
630 opportunities for biodiversity conservation in the Brazilian Atlantic Forest.  
631 *Conservation Biology*, 19, 695–700.

632 TANG, F. H., NGUYEN, T. H., CONCHEDDA, G., CASSE, L., TUBIELLO, F. N. & MAGGI,  
633 F. 2024. CROPGRIDS: a global geo-referenced dataset of 173 crops. *Scientific Data*,  
634 11, 413.

635 TIRIVAROMBO, S., OSUPILE, D. & ELIASSON, P. 2018. Drought monitoring and analysis:  
636 standardised precipitation evapotranspiration index (SPEI) and standardised  
637 precipitation index (SPI). *Physics and Chemistry of the Earth, Parts a/b/c*, 106, 1–10.

638 TOMASELLA, J., DO AMARAL CUNHA, A. M., ZERI, M. & COSTA, L. C. 2025. Changes in  
639 the aridity index across Brazilian biomes. *Science of The Total Environment*, 989,  
640 179869.

641 VICENTE-SERRANO, S. M., BEGUERÍA, S. & LÓPEZ-MORENO, J. I. 2010. A multiscalar  
642 drought index sensitive to global warming: the standardized precipitation  
643 evapotranspiration index. *Journal of climate*, 23, 1696–1718.

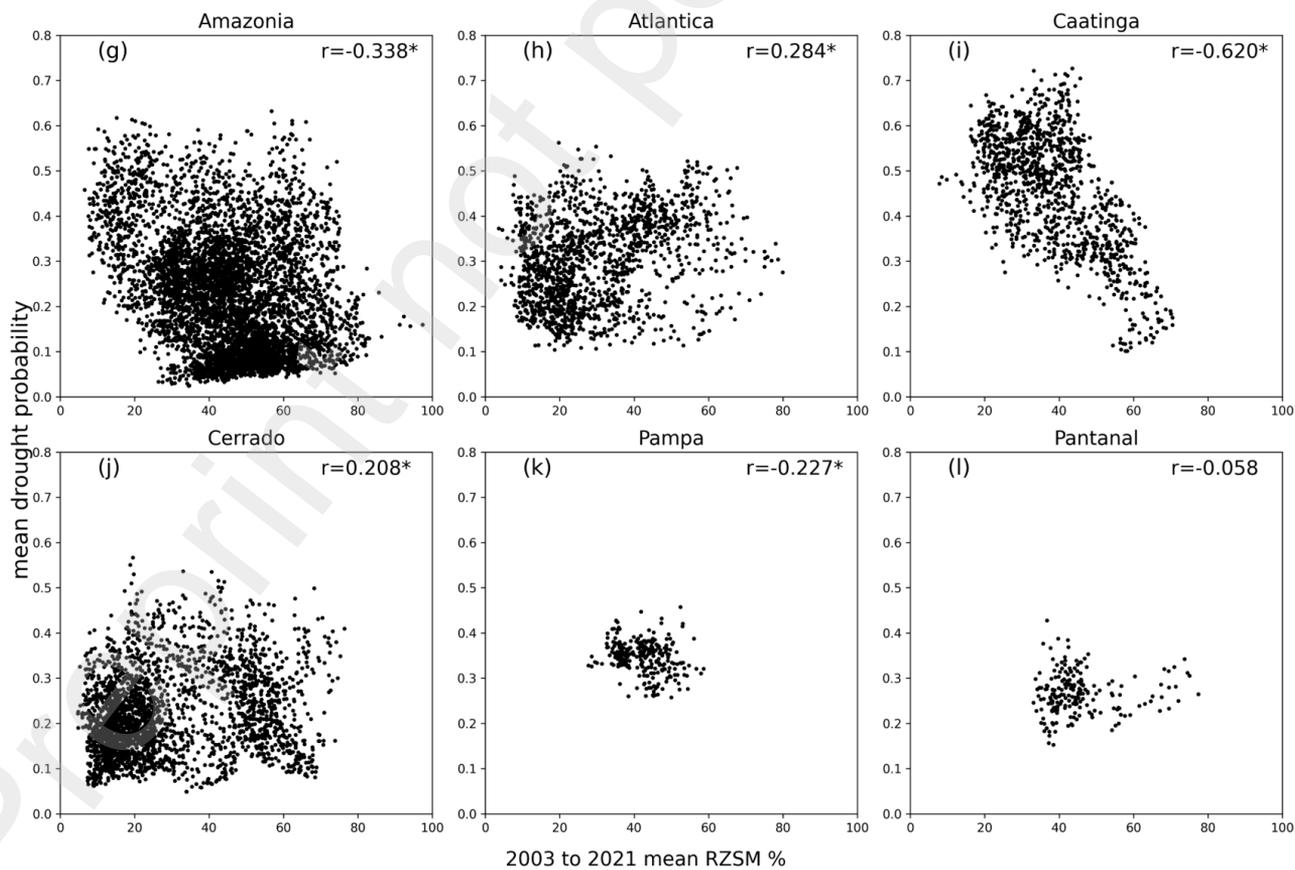
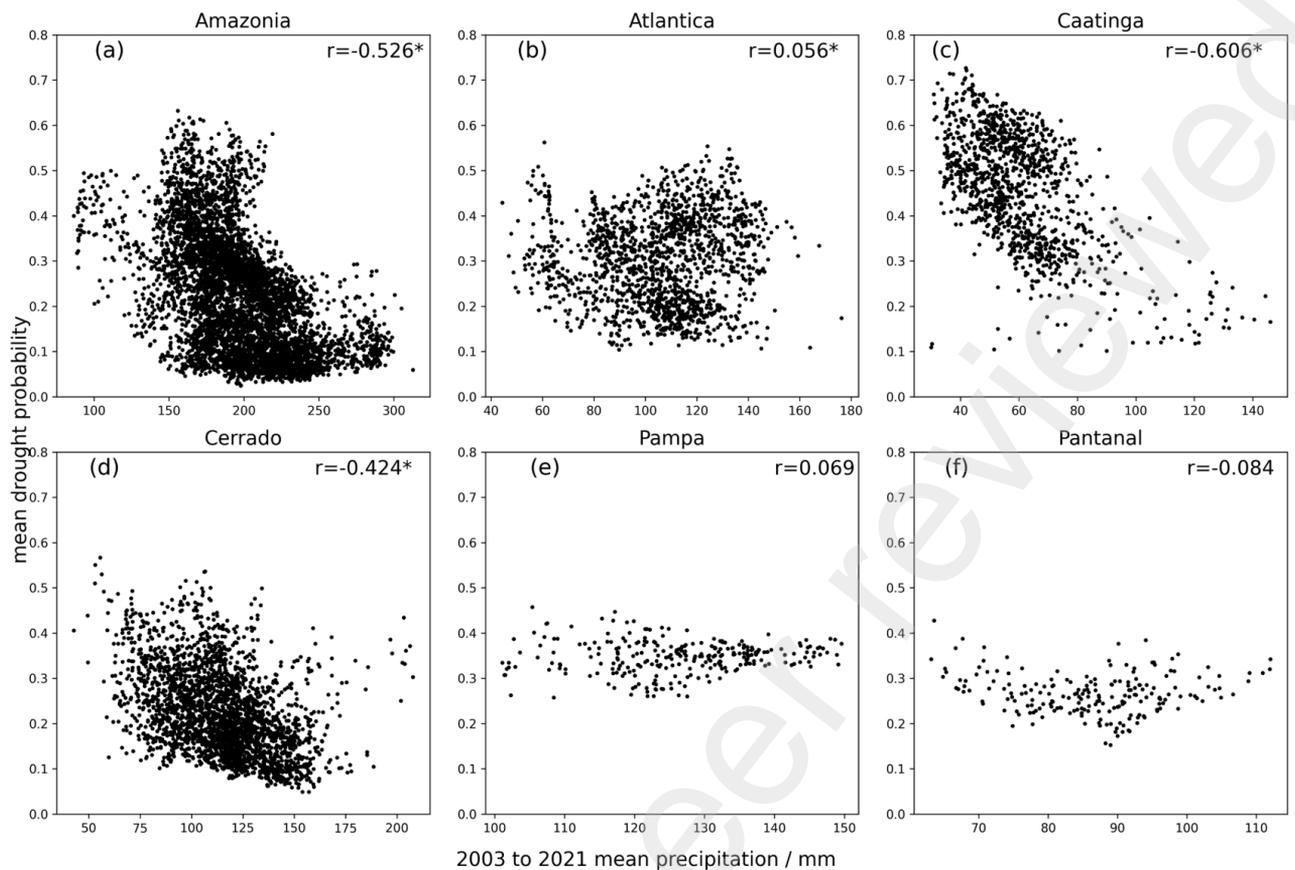
644 WILKS, D. S. 2006. *Statistical methods in the atmospheric sciences*, Academic press.

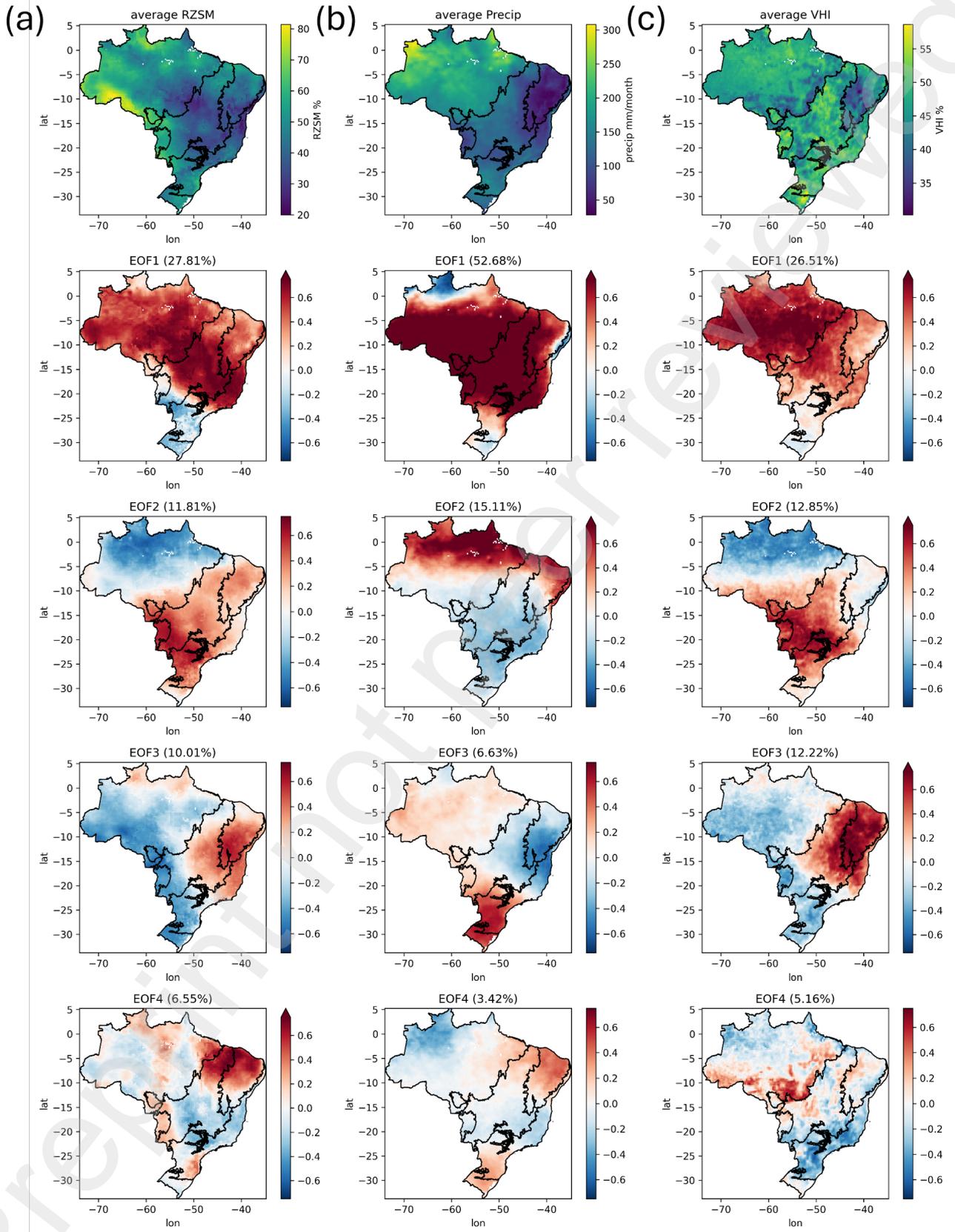
645 ZERI, M., WILLIAMS, K., CUNHA, A. P. M., CUNHA-ZERI, G., VIANNA, M. S., BLYTH, E.  
646 M., MARTHEWS, T. R., HAYMAN, G. D., COSTA, J. M. & MARENGO, J. A. 2022.  
647 Importance of including soil moisture in drought monitoring over the Brazilian  
648 semiarid region: An evaluation using the JULES model, in situ observations, and  
649 remote sensing. *Climate Resilience and Sustainability*, 1, e7.

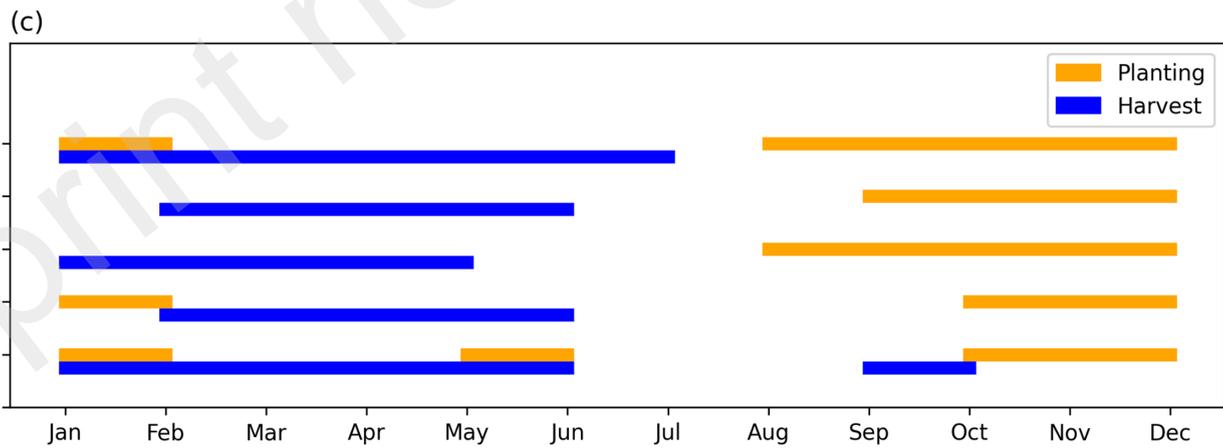
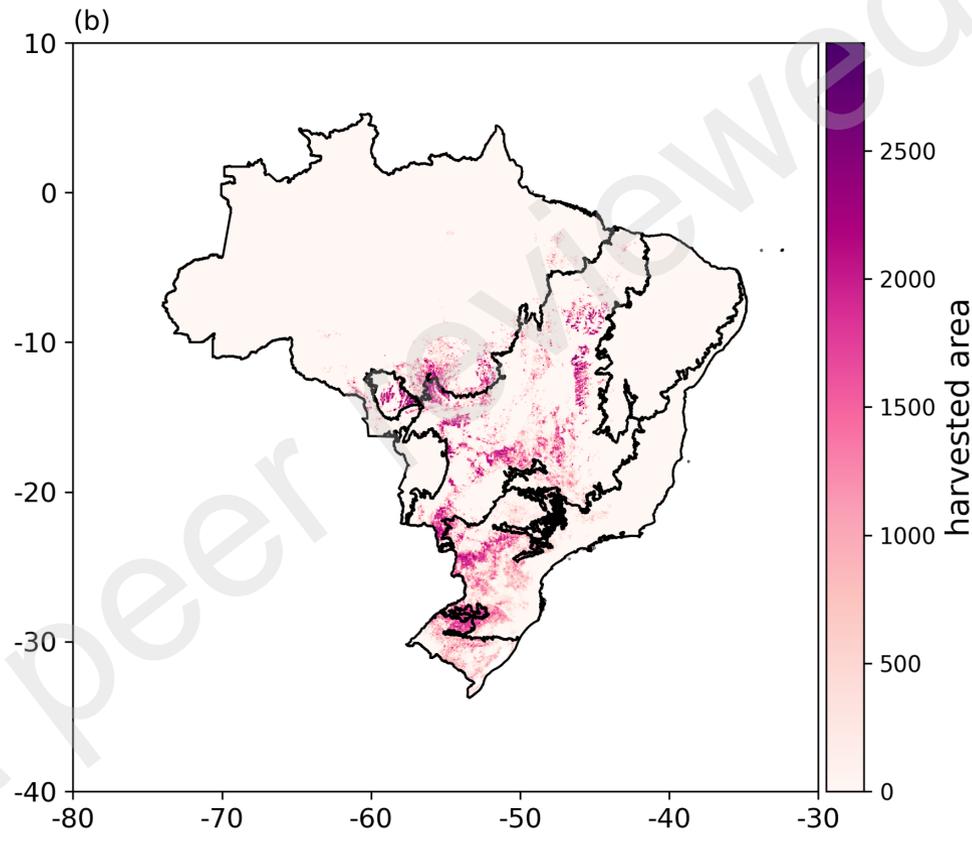
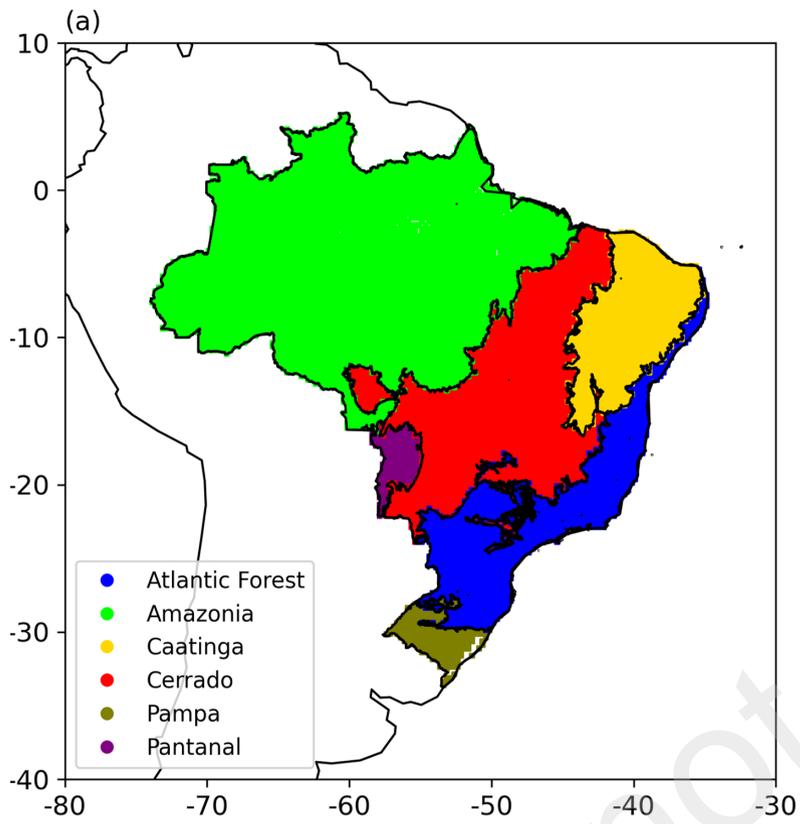
650

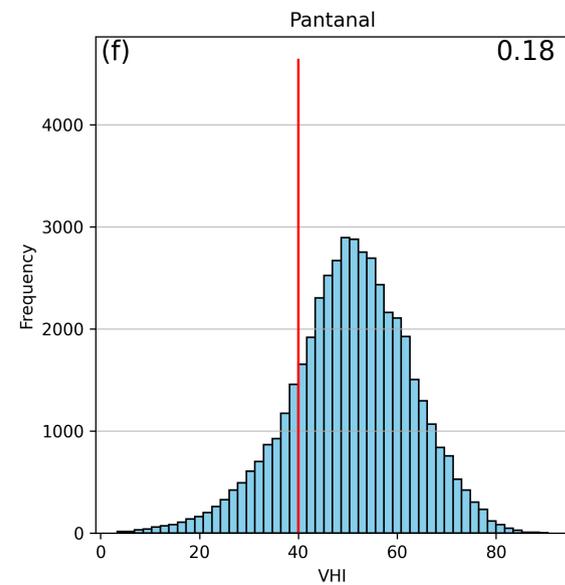
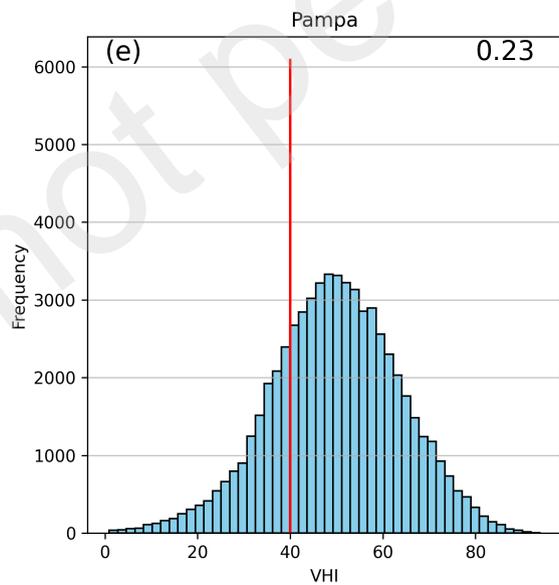
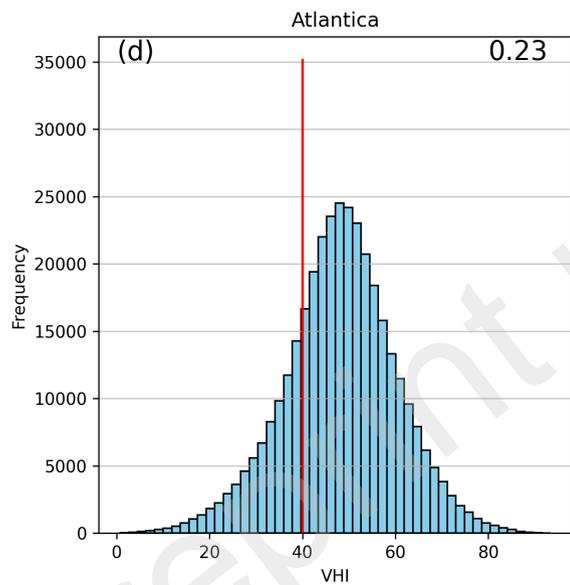
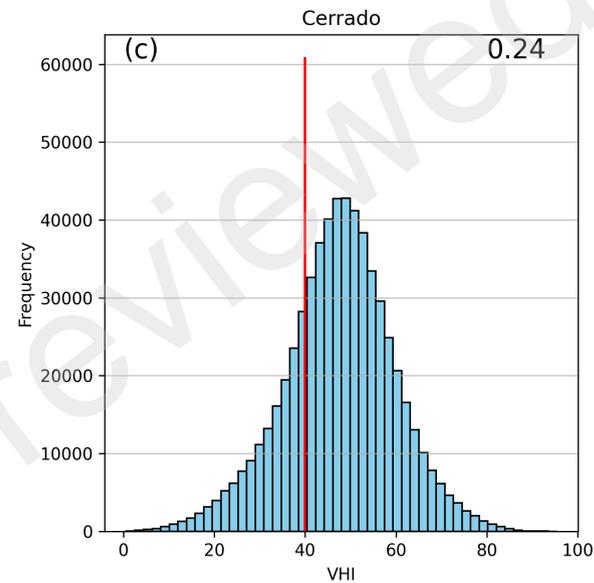
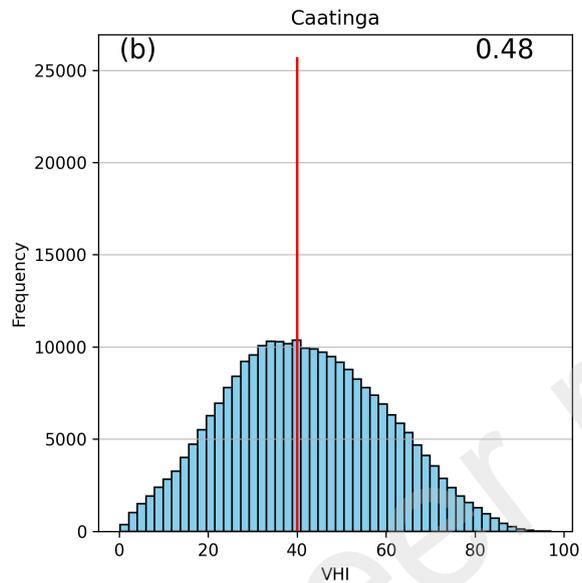
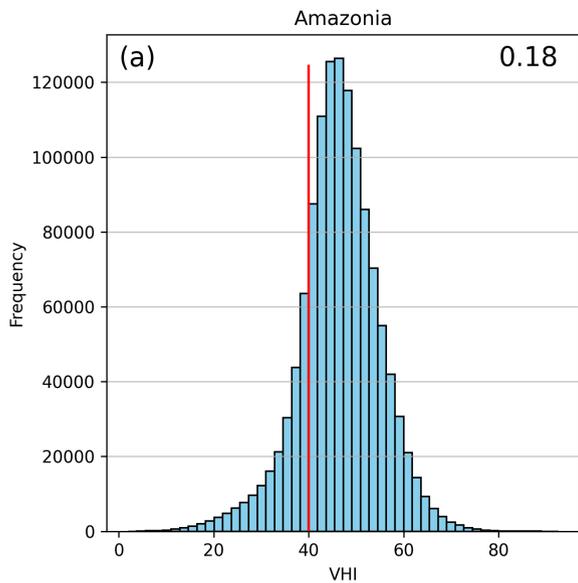
**Table 1.** Variables used in this study with data sources and abbreviations used throughout the text.

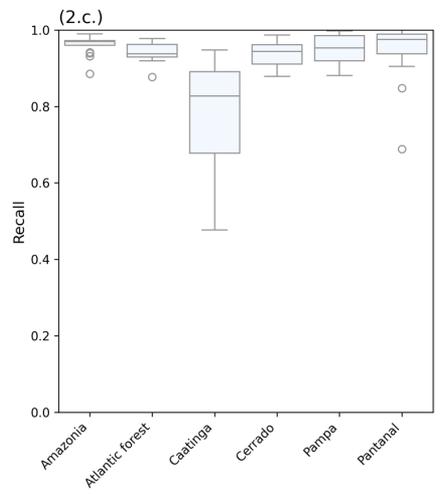
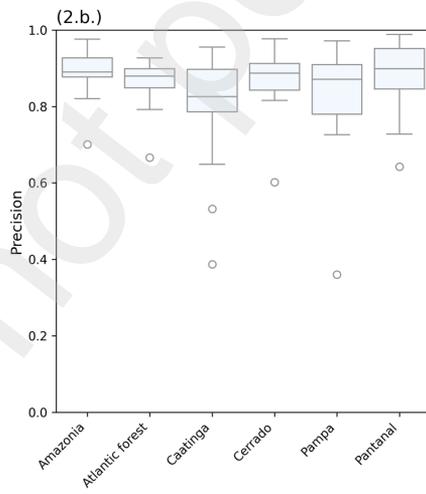
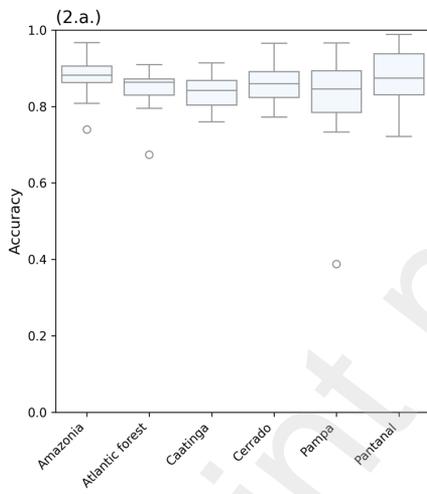
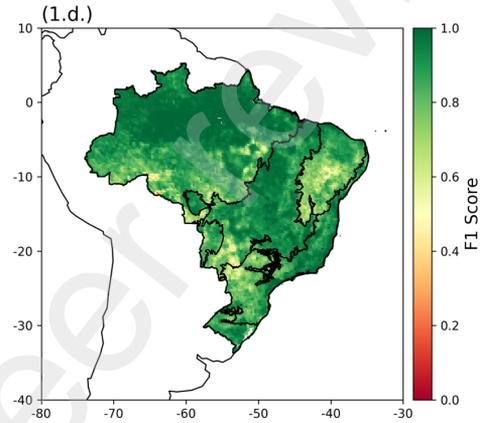
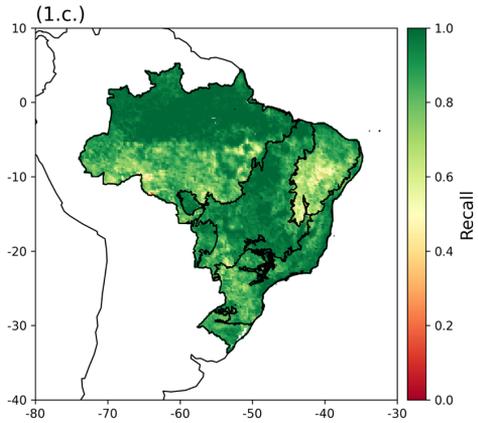
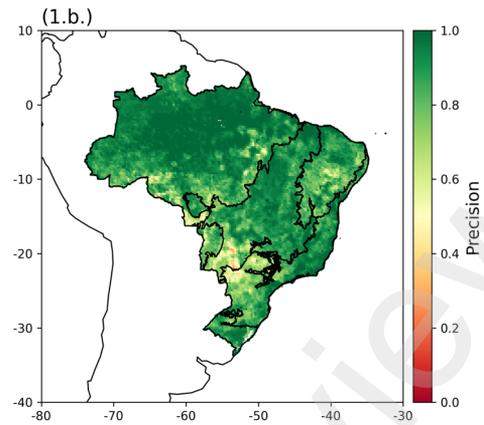
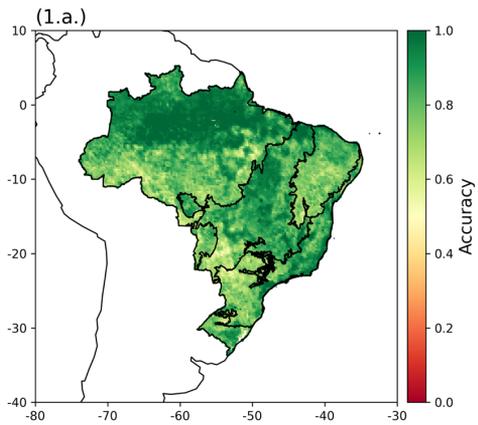
<b>Variable (units)</b>	<b>Abbreviation</b>	<b>Source</b>	<b>Usage</b>
Precipitation (mm/month)	Precip	CHIRPS	Feature
2 metre temperature (K)	T2M	ERA5	Feature
Potential evaporation (kg m <sup>-2</sup> )	pev	ERA5	Feature
Mean surface downward long-wave radiation flux (W m <sup>-2</sup> )	longrad	ERA5	Feature
Root Zone Soil Moisture (%)	RZSM	NASA GRACE	Feature
Vegetation health index (%)	VHI	NOAA STAR	Next month's value used to derive dependent variable
Standardized evapotranspiration-precipitation index (Unitless)	SPEI	Calculated from CHIRPS precipitation data and ERA5 data	Feature

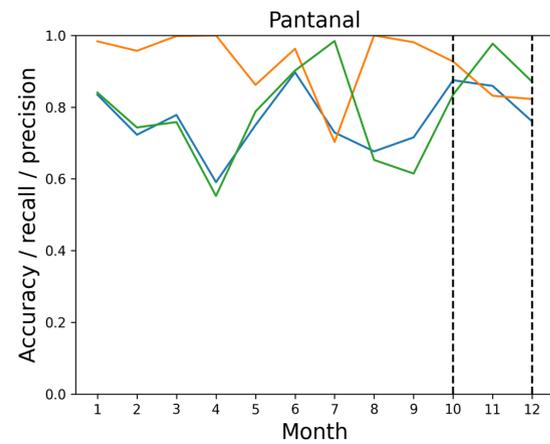
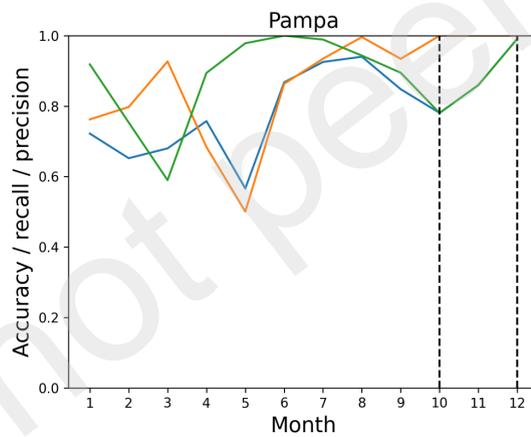
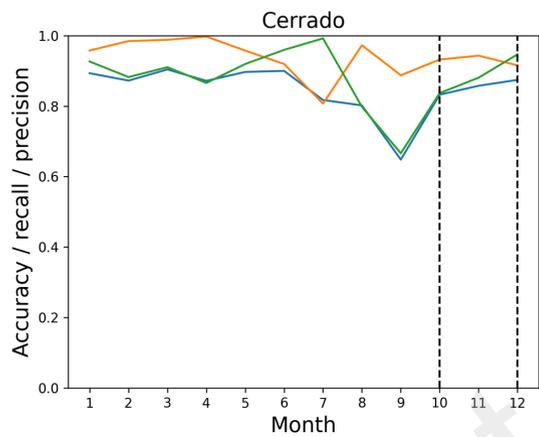
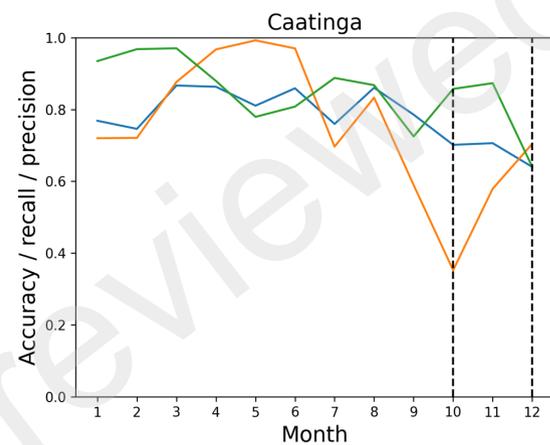
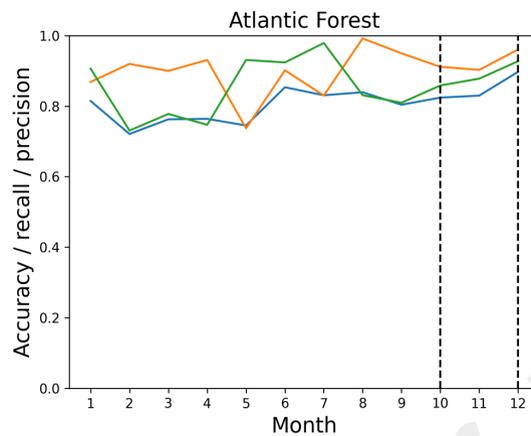
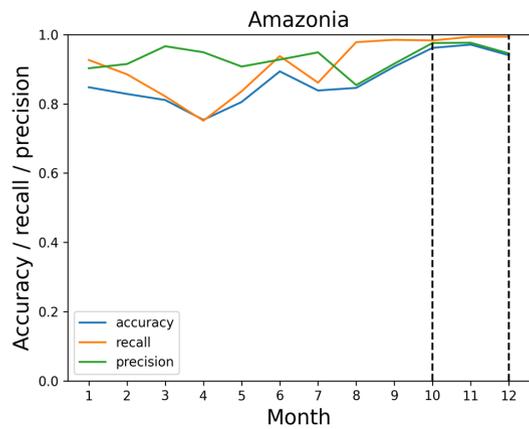


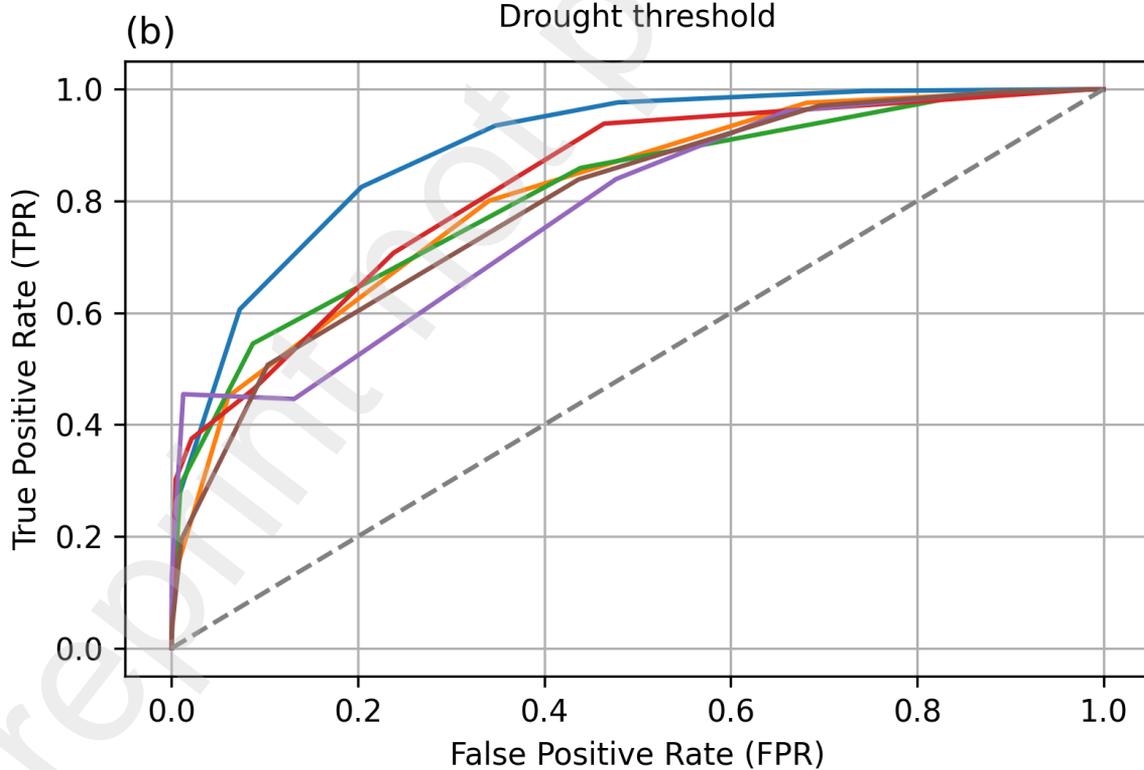
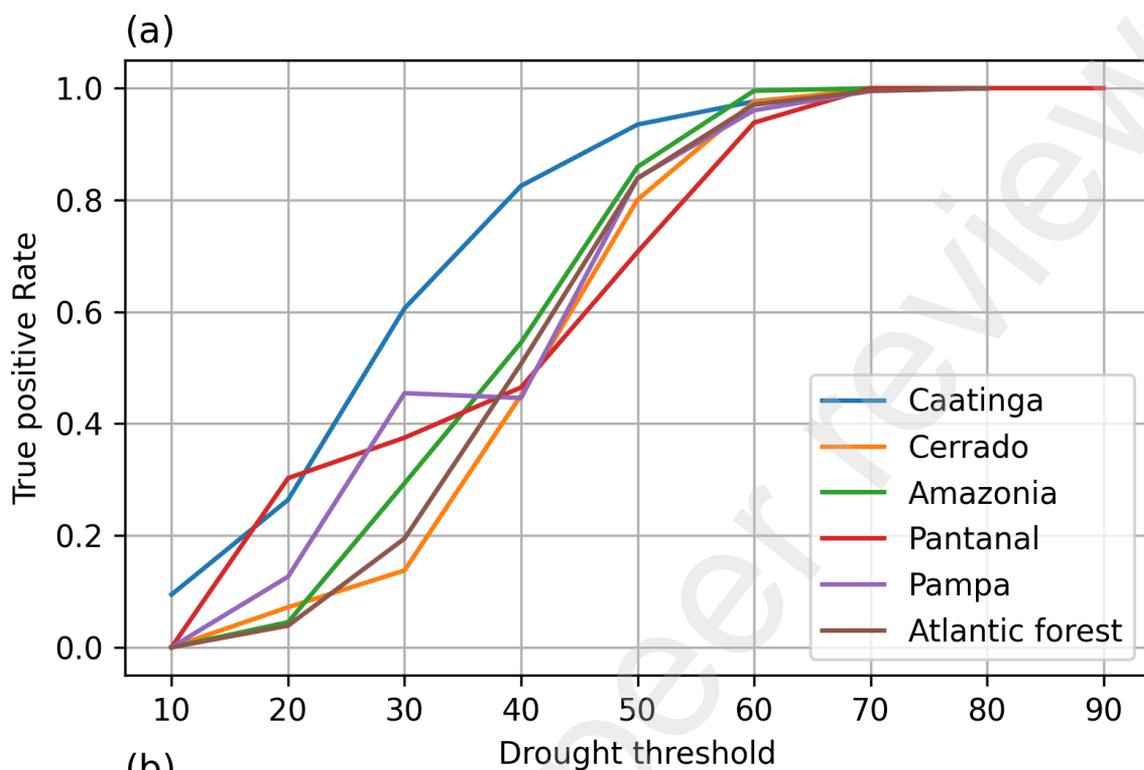


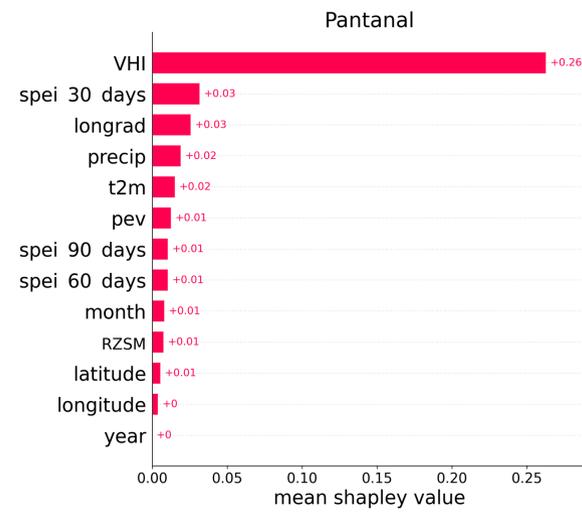
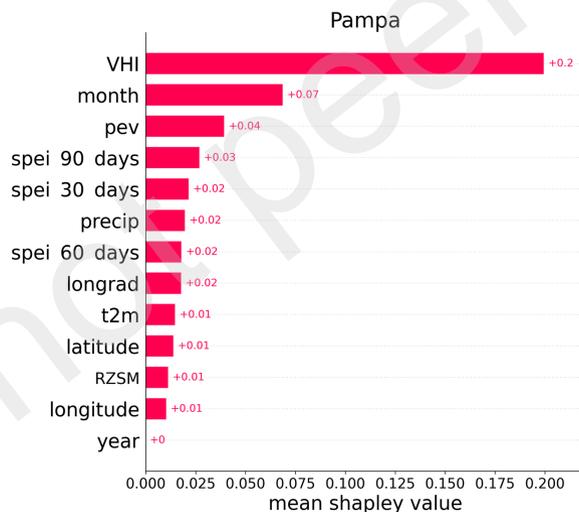
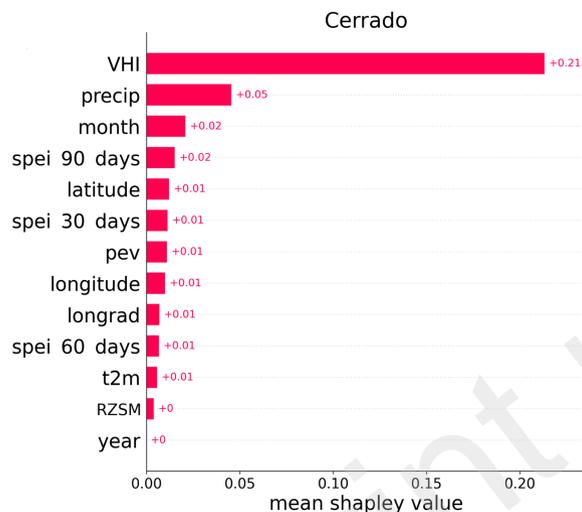
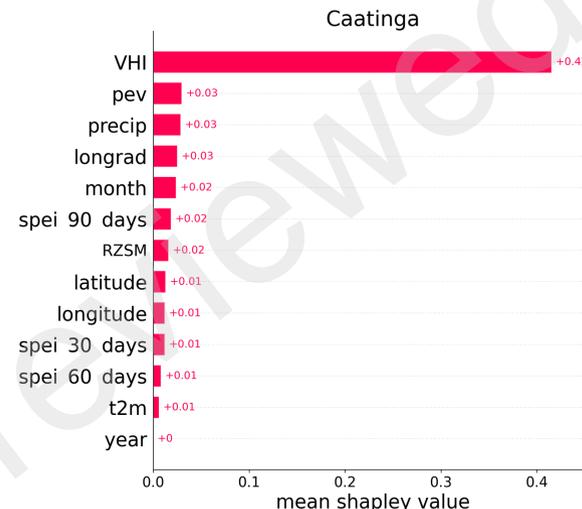
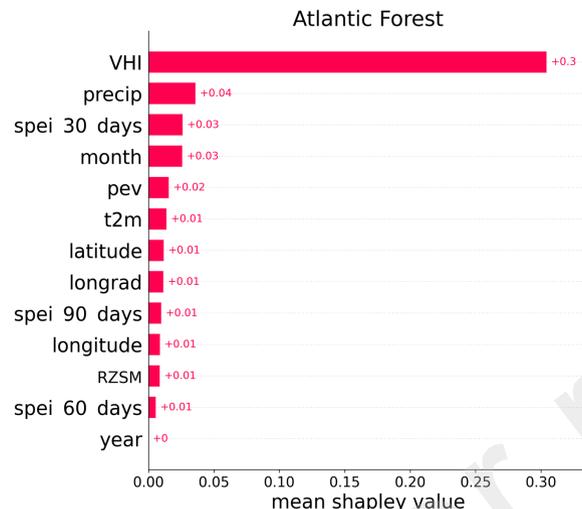
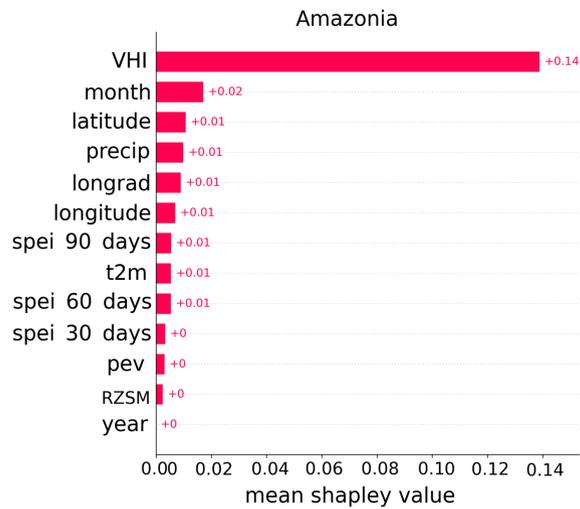












Preprint not certified by peer review