

Journal Pre-proof

Drought tolerance classification of grapevine rootstock by machine learning for the São Francisco Valley

Nina Iris Verslype , André Câmara Alves do Nascimento ,
Rosimar dos Santos Musser , Raphael Miller de Souza Caldas ,
Luiza Suely Semen Martins , Patrícia Coelho de Souza Leão

PII: S2772-3755(23)00022-9
DOI: <https://doi.org/10.1016/j.atech.2023.100192>
Reference: ATECH 100192



To appear in: *Smart Agricultural Technology*

Received date: 10 October 2022
Revised date: 23 January 2023
Accepted date: 29 January 2023

Please cite this article as: Nina Iris Verslype , André Câmara Alves do Nascimento , Rosimar dos Santos Musser , Raphael Miller de Souza Caldas , Luiza Suely Semen Martins , Patrícia Coelho de Souza Leão , Drought tolerance classification of grapevine rootstock by machine learning for the São Francisco Valley, *Smart Agricultural Technology* (2023), doi: <https://doi.org/10.1016/j.atech.2023.100192>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier B.V.
This is an open access article under the CC BY-NC-ND license
(<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

Highlights

- Comparison of six different ML algorithms to predict drought tolerance classes.
- A dataset with 42 cultivars and 30 variables was used for training and testing.
- The best model predicted drought tolerance classes of three rootstocks cultivars.
- The RF achieved 98.57% accuracy to predict rootstock drought tolerance classes.

Journal Pre-proof

Drought tolerance classification of grapevine rootstock by machine learning for the
São Francisco Valley

Nina Iris Verslype^{a,*}, André Câmara Alves do Nascimento^b, Rosimar dos Santos
Musser^a, Raphael Miller de Souza Caldas^a, Luiza Suely Semen Martins^c, Patrícia
Coelho de Souza Leão^d

^a Department of Agronomy, Federal Rural University of Pernambuco, Recife, Brazil.

^b Department of Computing, Federal Rural University of Pernambuco, Recife, Brazil.

^c Department of Biology, Federal Rural University of Pernambuco, Recife, Brazil.

^d Agricultural Research Centre for Semi-arid Tropics, Petrolina, Brazil.

Journal Pre-proof

* Corresponding author

E-mail addresses: nina.verslype@ufrpe.br (N.I. Verslype), andre.camara@ufrpe.br (A.C.A. do Nascimento), rosimar.musser@ufrpe.br (R. dos S. Musser), raphaelmillers@gmail.com (R. M. de S. Caldas), luiza.martins@ufrpe.br (L. S. S. Martins), patricia.leao@embrapa.br (P. C. de S. Leão).

ABSTRACT

Machine Learning (ML) algorithms are increasingly being used in several areas of agricultural studies, such as plant breeding. ML can assist in the recognition of relevant patterns or groups, or even in the prediction of the outcome under new settings, thus accelerating experiments and interpreting their results. The identification and selection of drought-tolerant grapevine rootstock (*Vitis* spp.) have become more relevant in late years, motivated mostly by global climate change scenarios. However, the grapevine is a perennial species, with polygenic characteristics and a complex traits inheritance by offspring, thus making it very challenging to discover new, drought tolerant cultivars. For this reason, this study's main objective was to compare the performance of six machine learning models on the prediction of drought tolerance levels of grapevine rootstock cultivars. A dataset with forty-five distinct cultivars was used to evaluate the methods, and the best performing model (AUC 0.9857) was used to predict the drought tolerance class of three cultivars (IAC 313, IAC 572, and IAC 766) whose drought tolerance level was still unknown. The results predicted a high drought tolerance for IAC 313 and IAC 766 cultivars, and a low tolerance for IAC 572.

Key words: Hydric stress, *Vitis* spp., climatic changes, artificial intelligence, supervised learning, algorithm.

1 Introduction

Grapevine is considered a crop of notable socio-economic importance (Zhang *et al.* 2016), and one of the most valuable perennial crops in the world, due to the high added value and versatility of the products (Walker *et al.*, 2019). Climate plays a crucial role in viticulture and is associated with the geography of wine (Jones *et al.*, 2012). Currently, it is singled out as one of the most critical aspects that directly interfere with the ripening and the final quality of the grapes to produce a specific style of wine (Fraga *et al.*, 2013; Jones, 2016). Climate also directly affects the choice of cultivars, planting location, vegetative potential, phytosanitary behavior, yield, and even the management and cultural practices adopted (Moura *et al.*, 2009; Jones *et al.*, 2012; Fraga *et al.*, 2013; Jones, 2016).

Climate change makes agricultural production extremely vulnerable, especially in terms of water availability, an essential component of life (Jones, 2016; Kathpalia & Bhatla, 2018). More specifically, one of the most pronounced abiotic stresses for plants is drought, which causes considerable losses in world agricultural production (Ashraf, 2010). Thus, a slight climate variation can directly affect the production and quality of grapes (Jones *et al.*, 2012; Fraga *et al.*, 2013). Besides, global warming projections indicate a high variation on rainfall patterns and intensity for the next decade (Jones & Webb, 2010; Marguerit *et al.*, 2012; Jones *et al.*, 2012; Fraga *et al.*,

2013). Therefore, mitigating the effects of climate change has become one of the main demands of the winery sector (Jones *et al.*, 2012; Fraga *et al.*, 2013; Zhang *et al.* 2016).

The scarcity of water resources emphasizes the importance of using water more efficiently in the winery sector, as most grape-growing regions around the world experience drought at some point (Chaves *et al.*, 2010; Zarrouk *et al.*, 2016; Fort *et al.*, 2017). Although several approaches can be employed to mitigate the drought problem in viticulture worldwide, using drought-tolerant rootstocks can be one of the most sustainable solutions to improve adaptability of vineyards (Marguerit *et al.*, 2012; Ollat *et al.*, 2016).

Drought tolerance in grapevines is a polygenic characteristic, controlled by many genes (Serra *et al.*, 2014; Dayer *et al.*, 2019). This hinders the identification and selection of the most promising genotypes, due to the interaction of the plant with the environment (Ashraf, 2010; Marguerit *et al.*, 2012; Serra *et al.*, 2014; Dayer *et al.*, 2019). The objectives of a rootstock breeding program must be decided effectively, given the large number of possible combinations, as well as that, the requirements to assess a wide variety of aspects, for example, drought tolerance, ease of propagation, compatibility with the graft, longevity of the scion, tolerance to environmental factors, resistance to pests and diseases in the region, as well as other plant characteristics and their influence on fruit.

Machine learning is an efficient methodology that has been increasingly used in several areas of study, such as in plant breeding (Etminan *et al.*, 2019). Since it makes it possible to identify, predict and classify genotypes according to the needs of the plant breeding program (Beiki, 2019; Etminan *et al.*, 2019). This study focused on use of machine learning algorithms to predict the drought tolerance levels of three Brazilian grapevine rootstocks cultivars (IAC 313, IAC 572, and IAC 766), in order to rank the best drought-tolerant rootstocks to be cultivated in the São Francisco Valley sub-middle (Latitude 9° S, Longitude 40° W) characterized by a semi-arid tropical climate, Brazil (Camargo & Amorim, 2004; Leão & Silva, 2014). These tropical rootstocks were selected in our study due to the lack of information about their drought tolerance degree and importance in the region of interest.

2 Materials and Methods

2.1 Plant material

Forty-five grapevine rootstock cultivars were evaluated, listed in table 1, with their drought tolerance degree and genetic origin, in order to build a dataset, which was later used to train a machine learning model. Five instances of each cultivar (i.e., plants) were used as samples, totaling 210 instances in the training dataset. Three criteria were adopted to choose these cultivars: (i) these are the most common cultivars in the São Francisco Valley region; (ii) their availability among the accessions from the Active Germplasm Bank of Embrapa (Camargo *et al.*, 2017); and (iii) the availability of information about their drought tolerance in the literature. The trained model was then used to predict the drought tolerance level of the three cultivars IAC 313 (Tropical), IAC 572 (Jales), and IAC 766 (Campinas).

Table 1 - Genetic origin and drought tolerance of the forty-five grapevine rootstock cultivars (*Vitis* spp.) considered in the training dataset.

Cultivars	Pedigree	Drought tolerance
VR 039-16	<i>Vitis vinifera</i> L. x <i>Muscadinia rotundifolia</i> Michaux ^{1;5}	Low ^{2;7}
VR 043-43	<i>Vitis vinifera</i> L. x <i>Muscadinia rotundifolia</i> Michaux ¹	Low ³
101-14 MGt	<i>Vitis riparia</i> Michaux x <i>Vitis rupestris</i> Scheele ^{1;3;5;6}	Low ^{2;3;4;5;8;10}
106-8 MGt	<i>Vitis riparia</i> Michaux x <i>Cordifolia rupestris</i> de grasset n ^o 1 ¹	High ⁶
110 R	<i>Vitis berlandieri</i> Planchon x <i>Vitis rupestris</i> Scheele ^{1;5;8}	High ^{2;4;5;6;7;8;13}
1103 P	<i>Vitis berlandieri</i> Planchon x <i>Vitis rupestris</i> Scheele ^{1;5;6;8}	High ^{2;3;5;7;8}
1202 C	<i>Vitis vinifera</i> L. x <i>Vitis rupestris</i> Scheele ¹	High ⁷
125 AA	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ¹	Low ⁴
140 Ru	<i>Vitis berlandieri</i> Planchon x <i>Vitis rupestris</i> Scheele ^{1;6}	High ^{2;4;5;7;8;12;13}
157-11 C	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ¹	Medium ⁴
1613 C	Solonis (<i>Vitis riparia</i> Michaux x <i>Vitis longii</i> Prince) x Othello ^{1;5}	Medium ⁹
161-49 C	<i>Vitis riparia</i> Michaux x <i>Vitis berlandieri</i> Planchon ^{1;6}	Low ⁵
1616 C	Solonis x <i>Vitis riparia</i> Michaux ¹	Low ²

196-17 CI	1203C x <i>Vitis riparia</i> Michaux ¹	High ^{6;7}
216-3 CI	1616C x <i>Vitis rupestris</i> Scheele ¹	Medium ^{5;6;7}
26 G	<i>Vitis vinifera</i> L. x <i>Vitis riparia</i> Michaux ¹	Medium ¹²
3306 C	<i>Vitis riparia</i> Michaux x <i>Vitis rupestris</i> Scheele ¹	Medium ⁸
3309 C	<i>Vitis riparia</i> Michaux x <i>Vitis rupestris</i> Scheele ^{1;5;6}	Low ^{4;5;7;11;13}
34 EM	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ¹	Low ^{4;5}
41 B MGt	<i>Vitis vinifera</i> L. x <i>Vitis berlandieri</i> Planchon ¹	High ^{5;10}
420 A MGt	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ^{1;3;5;6;9}	Medium ^{4;5;10;1} 2
44-53 M	<i>Vitis riparia</i> Michaux x Malegue 144 ¹	High ^{2;5;6;7}
5 BB	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ^{1;5;6}	Medium ^{2;4}
5 C	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ^{1;5}	Low ^{5;8}
62-66 C	<i>Vitis vinifera</i> L. X <i>Vitis cordifolia</i> Michaux ¹	High ¹¹
8 B	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ¹	Medium ⁵
93-5 C	<i>Vitis vinifera</i> L. x <i>Vitis rupestris</i> Scheele ¹	Low ⁶
99 R	<i>Vitis berlandieri</i> Planchon x <i>Vitis rupestris</i> Scheele ^{1;5;6;8}	High ⁸
Ganzin 1	<i>Vitis vinifera</i> L. x <i>Vitis rupestris</i> Scheele ¹	High ⁷
Dogridge	<i>Vitis rupestris</i> Scheele x <i>Vitis candidans</i> Engelmann ^{1;5}	High ⁹
Fercal	<i>Vitis berlandieri</i> Planchon x 31R ¹	Medium ^{2;5;7;8}
Freedom	Fresno 1613-59 x Dogridge ¹	Medium ^{2;8}
Golia	Castel 156-12 x <i>Vitis rupestris</i> Scheele ¹	Low ⁴
Gravesac	161-49C x 3309C ¹	Medium ^{6;7;8}
Harmony	1613C x Dogridge ¹	Medium ²
Riparia	<i>Vitis riparia</i> Michaux ^{1;5;6}	Low ^{2;5;7;10;12;13}
Gloire		
Rupestris du lot	<i>Vitis rupestris</i> Scheele ^{1;5;6;8}	Medium ^{5;7;9}
Salt Creek	<i>Vitis champinii</i> Planchon ^{1;5}	High ^{9;8}
Schwarzman n	<i>Vitis riparia</i> Michaux x <i>Vitis rupestris</i> Scheele ^{1;5}	Low ^{5;8}
SO4	<i>Vitis berlandieri</i> Planchon x <i>Vitis riparia</i> Michaux ^{1;5;6}	Low ^{3;5;8;10}

Sori	Solonis x <i>Vitis riparia</i> Michaux ¹	Medium ^{7;12}
Vitis champini	<i>Vitis champinii</i> Planchon ¹	High ¹³
IAC 313	Golia x <i>Vitis cinerea</i> Engelm ^{1;3}	Unknown*
IAC 572	<i>Vitis caribaea</i> De Candolle x 101-14MGt ^{1;3}	Unknown*
IAC 766	106-8MGt x <i>Vitis caribaea</i> De Candolle ^{1;3}	Unknown*

Information obtained in: ¹Maul *et al.* (2020); ²Sunridge Nurseries (2020); ³Embrapa Grape and Wine (2016); ⁴Vicopad (2020); ⁵Villa (2018); ⁶Storm & Krasokhina (2020); ⁷Audeguin *et al.* (2020); ⁸Wine Australia (2016); ⁹Satisha *et al.* (2006); ¹⁰ATVB (2013); ¹¹Chevalier (1925); ¹²Rebschule Mueller (2020), ¹³Carroll (2016); *No information in the literature.

2.2 Data set characteristics

The dataset holds fifty-nine variables, including physiological, biochemical, nutritional, and morpho-agronomic characteristics, manually curated by the authors from scientific articles, books, grapevine nurseries websites, thesis, conferences, and scientific research centers. The collected features are listed below:

2.2.1 Physiological characteristics: Stomatal conductance under unstressed conditions (g_s) and drought stress both in $\mu\text{mol CO}_2\cdot\text{m}^{-2}\cdot\text{s}^{-1}$; Transpiration rate under unstressed conditions (E) and drought stress both in $\text{mmol H}_2\text{O}\cdot\text{m}^{-2}\cdot\text{s}^{-1}$; Photosynthesis rate under unstressed conditions (A) and drought stress both in $\mu\text{mol m}^{-2}\cdot\text{s}^{-1}$; Intrinsic Water use efficiency under unstressed conditions (WUE) and drought stress both in $\mu\text{mol CO}_2 \text{ mol}^{-1} \text{ H}_2\text{O}$; Instantaneous water use efficiency under unstressed conditions (iWUE) and drought stress both in $\mu\text{mol CO}_2 \text{ mmol}^{-1} \text{ H}_2\text{O}$; Osmotic potential under unstressed conditions (Ψ_{os}) and drought stress both in Mpa, Hydraulic conductance (KI) in $\text{kg}/\text{MPa}\cdot\text{s}$;

2.2.2 Biochemical characteristics: abscisic acid under unstressed conditions (ABA), drought stress and rehydration condition of the plant all in $\text{ng}\cdot\text{g}^{-1}$; Proline under unstressed conditions and drought stress both in $\text{mg}\cdot\text{g}^{-1}$;

2.2.3 Nutritional characteristics: percentage of nitrogen, phosphorus, and potassium macronutrients content in the leaf petiole; and also, the uptake ability of the plant to absorb nitrogen, phosphorus, and potassium, on a rating scale assigned between 1 and 5 (high = 5, medium-high = 4, medium = 3, medium-low = 2 and low = 1);

2.2.4 Morpho-agronomic characteristics: stomata density per mm²; flower sex (male, hermaphrodite or female); grape buds, grape maturity and leaf fall in days; active limestone tolerance in percentage; wood production per plant; root distribution scaled in a range from 1 to 6 (very deep=6, deep=5, moderate-deep=4, moderate=3, shallow=2, very shallow=1); geotropic angle in degree; percentage root distributions to 60 cm and 100 cm; Chlorotic Power Index (CPI) in percentage; Drought tolerance scaled in a range from 1 to 3 (high =3, medium =2 and low = 1); Vegetative cycle (early, intermediate, late); and the variables listed below a scale of scores ranging from 1 to 5 (high=5, medium-high=4, medium=3, low-medium=2 and low=1) were assigned, for Anthracnose, Downy mildew, Fusarium, Phylloxera and Nematode resistance; lime and total limestone tolerance; ease of rooting; ease of branch-grafting; acid, sandy, wet, clay, salinity, calcareous and compactness soil tolerance; vigor; Iron chlorosis tolerance.

2.3 Data analysis methodology

All analysis were performed using the Python language on the Google Colaboratory¹ platform.

2.3.1 Pearson correlation analysis

Pearson's correlation coefficient (r) was calculated among the considered variables, using the `corr()` function available in the pandas library (McKinney, 2010), in order to examine the correlation between the variables and discard highly correlated characteristics. Pearson's correlation coefficients are calculated according to equation (01) (Ferreira *et al.*, 2011).

$$r = (\sum xy - (\sum x * \sum y/N)) / \sqrt{(\sum x^2 - ((\sum x)^2/N)) * (\sum y^2 - ((\sum y)^2/N))} \quad (01),$$

where r corresponds to the correlation coefficient that can vary between -1 and +1; N is the number of observations; x and y are the values of both variables.

Subsequently, a heat map chart, available in the seaborn library (Waskom *et al.*, 2020), was generated to improve the visualization of Pearson's correlation coefficients.

2.3.2 Dataset pre-processing

Pre-processing occurred in the same way for the training set, with 42 cultivars, and the prediction set, containing three cultivars (IAC 313, IAC 572, and IAC 766). For this, the missing values were filled in by the mean of each variable column, then

¹ <https://colab.research.google.com/>

duplicated rows were removed and finally, normalization was performed using MinMaxScaler, (Figure 1), whose transformation is given by the equation (01) and (02), available in the scikit-learn library (Pedregosa *et al.* 2011).

$$X_{std} = \frac{(x-x.min)}{(x.max - x.min)} \quad (01)$$

$$X_{scaled} = X_{std} \times (max - min) + min \quad , \quad (02)$$

where: *min*, *max* corresponds to the sample range.

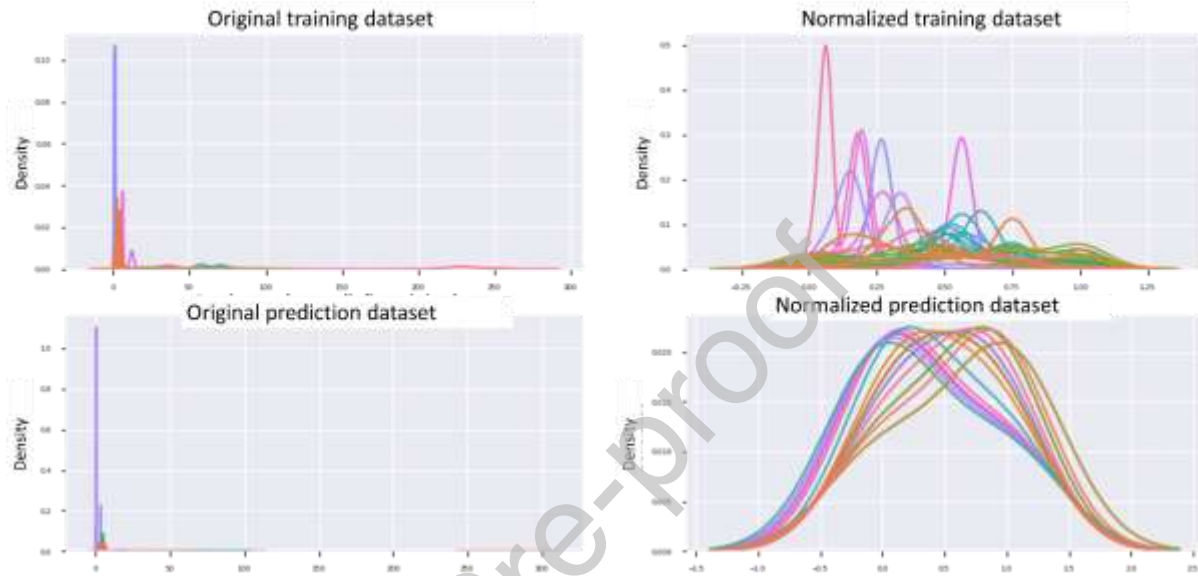


Figure 1 - Dataset before and after normalization. Source: Verslype *et al.* (2021).

2.3.3 Comparison of supervised learning algorithms

Six algorithms were considered to evaluate their ability to predict drought tolerance classes in grapevine rootstocks, namely: *Decision Tree* (DT) (Breiman *et al.*, 1984), *Random Forest* (RF) (Breiman, 2001), *K-Nearest Neighbors* (KNN) (Fix *et al.*, 1951), *XGBoost* (XGB) (Chen & Guestrin, 2016), *Support Vector Machines* (SVM) (Cortes & Vapnik, 1995) and *Linear Discriminant Analysis* (LDA) (Rao, 1948). These algorithms were select because they are easily accessible by the public, and all algorithms are available in the scikit-learn library (Pedregosa *et al.*, 2011), except for XGB available in the xgboost library (Chen & Guestrin, 2016). To determine the best algorithm in our setting, a cross-validation evaluation procedure was performed, and grid search analysis (Table 2), available in the scikit-learn library. (Pedregosa *et al.*, 2011). The precision rate (P), recall (R), accuracy, and f1-score metrics obtained by all models in the considered dataset were evaluated. Are described these evaluation metrics in equations (03), (04), (05), and (06) (Pedregosa *et al.* 2011, Shalev-Shwartz & Ben-David, 2014; Skansi, 2018).

$$P = \frac{TP}{(TP + FP)} \quad , \quad (03)$$

where: P represents the precision, TP indicates the number of true positives, and FP the number of false positives.

$$R = \frac{VP}{(TP + FN)} \quad , \quad (04)$$

where: R represents *recall*, TP is the number of true positives, and FN the number of false negatives.

$$accuracy = (TP + VN)/(TP + TN + FP + FN) \quad , \quad (05)$$

where: TP is the number of true positives, TN the number of true negatives, FP the number of false positives, and FN the number of false negatives.

$$f1 = 2 \times (precis\tilde{a}o \times recall)/(precis\tilde{a}o + recall) \quad , \quad (06)$$

where: $f1$ represents that obtained by the harmonic mean between recall and accuracy.

2.3.4 Drought tolerance prediction

A predictive drought tolerance analysis was performed on the three Brazilian grapevine rootstocks cultivars IAC 313, IAC 572, and IAC 766, using the most efficient algorithm and its respective hyperparameter configuration, provided by the grid search (Table 2).

Table 2 – Hyperparameters analyzed by grid search for each algorithm.

Algorithm	Hyperparameters
DT	Criterion: [gini, entropy]; Max_features [auto, log2, sqrt, None]
RF	N_estimators: [1, 5,10, 100, 1000]; Max_features: [1, 2, 3]
XGB	N_estimators: [10, 100, 500]; Max_features: [auto, log2, sqrt, None]
SVM	C: [0.001, 0.1, 1,10, 20, 100]; Kernel: [rbf, linear]
LDA	N_components: [10, 15, 20, 25, 30]
KNN	N_neighbors: [3, 7, 10]; Weights: [uniform, distance]

3 RESULTS AND DISCUSSION

The Pearson's correlation coefficient analysis (Figure 2) indicated 38 cases with a strong correlation between the 61 variables evaluated in this study, of which 12 had a negative correlation close to -1, and 26 had a positive correlation close to 1.

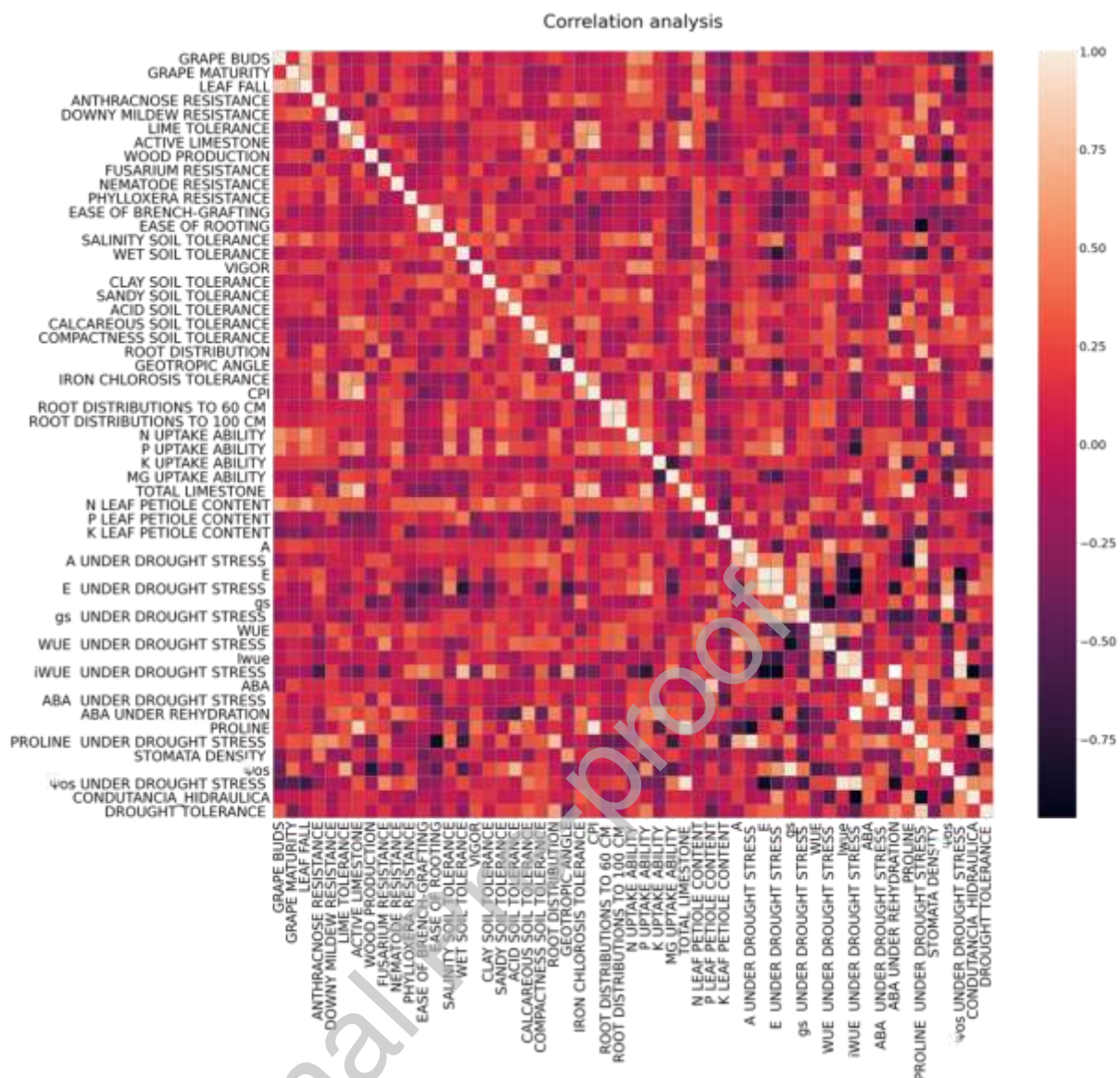


Figure 2 - Correlation coefficient analysis between the sixty-one variables evaluated.

f

In this regard, correlation analysis is interesting since it allows one to explain and determine the degree of relationship between the variables in the data set (Ferreira *et al.*, 2011). Thus, coefficients close to 1 and -1 make it possible to obtain a good prediction of one variable from the other (Forsyth, 2018). Besides, redundant features can be detected, for example, the root distributions to 100 cm and 60 cm can be highlighted among the strong correlations, due to their positive correlation ($r=0.91$). This finding indicates that there no need to evaluate both features, resulting in cost savings and a significant reduction in analytical time. Another pair of features that this analysis pointed as highly correlated were between the grape buds and grape maturity variables, as they presented, respectively ($r=0.78$) and ($r=0.73$) of positive correlation with the variable onset of leaf fall, demonstrating that only the evaluation of the variable early leaf fall would be sufficient.

The proline under drought stress conditions has a high correlation with drought tolerance ($r=0.84$). This interaction is important since proline could be used as an indirect selection criterion for the search for drought-tolerant varieties. As drought tolerance is a polygenic trait which is difficult to identify and select for superior materials in plant breeding programs (Marguerit *et al.*, 2012; Serra *et al.*, 2014; Cantu & Walker, 2019). This high correlation could be explained by the fact that proline is considered an osmotically active substance, which is accumulated in high levels in the cytoplasm when the grapevine is in water restriction periods, thus enabling it to maintain balance the water potential within the plant cell and the turgor pressure (Serra *et al.*, 2014; Keller, 2015; TAIZ *et al.*, 2017).

Thereby analyzing the variables with high correlation (Figure 3), 14 of them were discarded, namely: E under drought stress, A under drought stress, EUA under drought stress, EiUA under water stress, proline under unstressed and drought stress conditions, ABA content under unstressed conditions, osmotic potential under unstressed and drought stress conditions, IPC, total limestone tolerance, root distribution to 100 cm, grape buds and grape maturity.

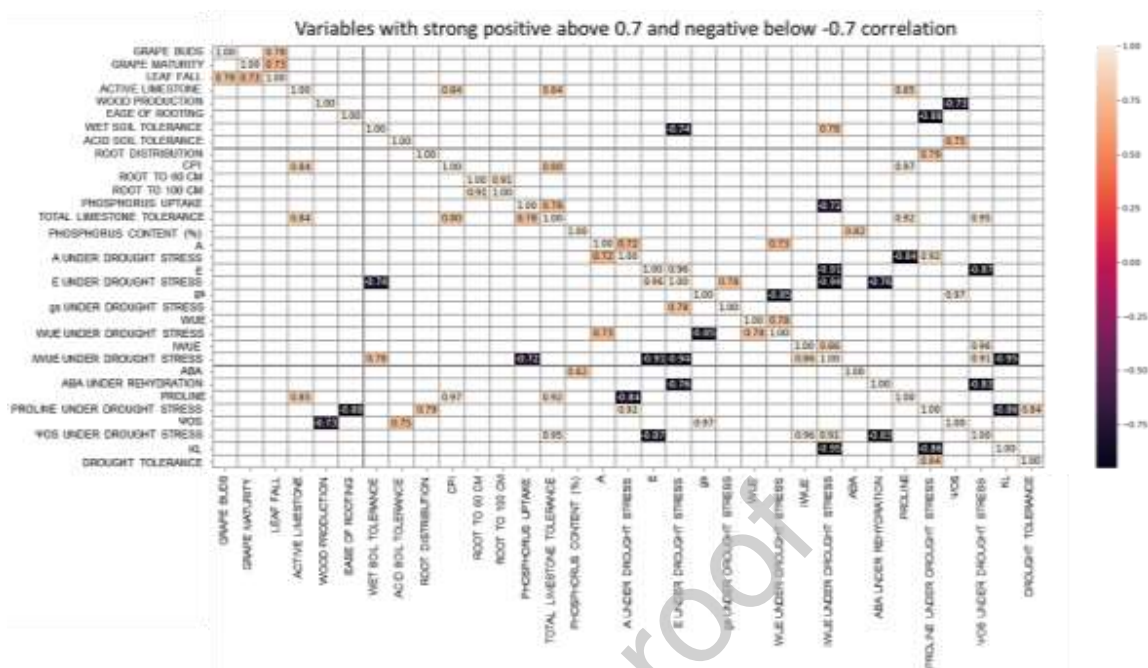


Figure 3 - Variables with strong positive and negative correlation. Source: Verslype et al. (2021).

Performance differences in the six algorithms on the training dataset were identified (Table 3). The accuracy metric range was between 98.10% and 64.29%, precision between 98% and 67%, and recall between 98% and 64%. For Skansi (2018), these evaluation metrics are an important cue to obtain a meaningful performance comparison of a set of machine learning classifiers. Since the accuracy allows identifying how good the classifier is, in average, in the task to label unseen instances, the precision determines the model's ability to avoid erroneous results like not labeling a positive sample as negative, and the recall indicates the success of the model to find all positive samples (Pedregosa *et al.*, 2011; Skansi 2018).

Table 3 - Accuracy, precision (P), recall (R), and f1-score obtained by 10-fold cross-validation evaluation on the six algorithms to predict drought tolerance classes for grapevine rootstock cultivars.

Algorithm	Accuracy (%)	P	R	f1-score
DT	82.86	0.86	0.83	0.83
RF	98.10	0.98	0.98	0.98
XGB	96.67	0.97	0.97	0.97
SVM	87.14	0.89	0.87	0.87
LDA	64.29	0.67	0.64	0.65
KNN	85.71	0.89	0.86	0.86

In this sense, the cross-validation experiments have shown that the RF algorithm as is the best classifier. Due to the highest mean values achieved for accuracy (98.10%), precision (98%), and recall (98%). Nevertheless, the LDA classifier had the lowest accuracy (64.29%), precision (67%), and recall (64%) values among all the models evaluated, indicating a low rate of correctness in the indication of the three drought tolerance classes.

Grid search analysis also was applied to the six algorithms (Table 2) to indicate the best performance algorithm to predict drought tolerance classes. The results obtained confirm that the LDA had the worst performance, reaching only 67.14% accuracy as the best possible result for all the configurations of hyperparameters tested on the dataset and consequently the lowest learning curve. Meanwhile, the RF and XGB algorithms performed the best results with 98.57% and 96.19% correctness to predict the three drought tolerance classes (Table 4). The best parameter and its respective configuration of hyperparameters and the learning curve plotted in Figure 2 indicate a higher learning rate over time.

Table 4 - Evaluation of the best hyperparameters for each model obtained by *Grid search* to predict drought tolerance classes for grapevine rootstock cultivars.

Algorithm	Best result (%)	Best parameter	Hyperparameters best configuration
SVM	83.33	{'C': 100, 'kernel': 'rbf'}	SVC (C=100, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)
RF	98.57	{'max_features': 2, 'n_estimators': 1000}	RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features=2, max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=1000, n_jobs=None, oob_score=False, random_state=None, verbose=0, warm_start=False)
KNN	90.48	{'n_neighbors': 3, 'weights': 'distance'}	KNeighborsClassifier (algorithm='auto', leaf_size=30, metric='minkowski', metric_params=None, n_jobs=None, n_neighbors=3, p=2, weights='distance')
XGB	96.19	{'max_features': 'auto', 'n_estimators': 500}	XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3, max_features='auto', min_child_weight=1, missing=None, n_estimators=500, n_jobs=1, nthread=None, objective='multi:softprob', random_state=0, reg_alpha=0, reg_lambda=1,

scale_pos_weight=1, seed=None, silent=None,
subsample=1, verbosity=1)

LDA	67.14	{'n_components': 10}	LinearDiscriminantAnalysis(n_components=10, priors=None, shrinkage=None, solver='svd', store_covariance=False, tol=0.0001)
DT	86.19	{'criterion': 'entropy', 'max_features': None}	DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='entropy', max_depth=None, max_features=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, presort='deprecated', random_state=None, splitter='best')

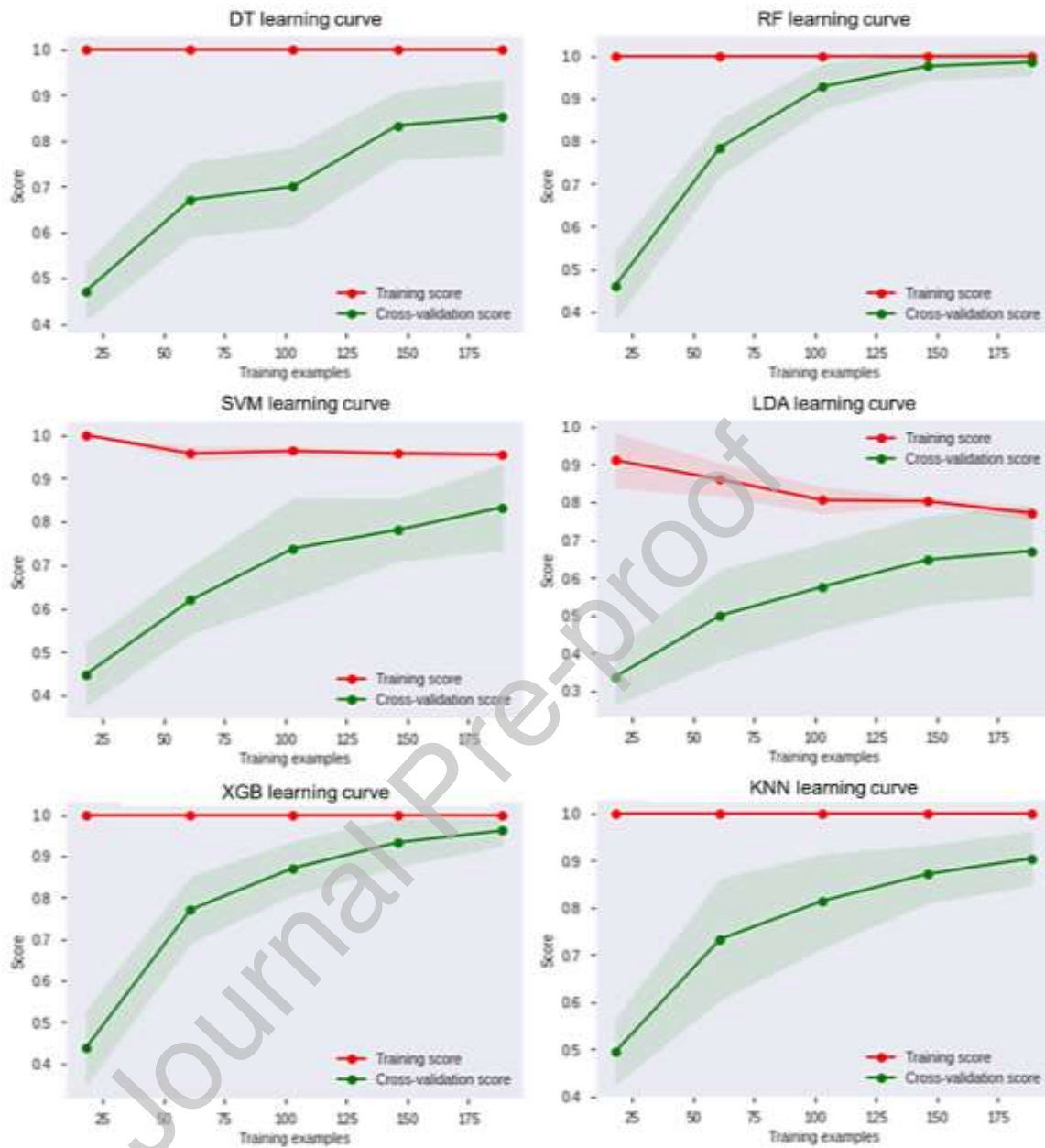


Figure 2 - Learning curve of the six algorithms on the dataset, through 10-fold cross-validation evaluation, and score metric (R^2) to predict drought tolerance classes for grapevine rootstock cultivars. Where: DT - *Decision Tree*, RF – *Random Forest*, XGB - *XGBoost*, SVM – *Support Vector Machines*, LDA - *Linear Discriminant Analysis* e KNN - *K-Nearest Neighbors*. Source: Verslype *et al.* (2021).

Table 5 presents the performance of the best model, stratified by class. The RF algorithm achieved a f1-score of 0.99 for the high tolerance class, 0.98 for low, and 0.99 for medium drought tolerance. According to Pedregosa *et al.* (2011), The f1-score values near to one indicate a better performance, while those near zero indicate the worst score in the model. In this sense, these results indicate a high success rate to classify each of the three drought tolerance classes by the RF algorithm.

Table 5 – RF algorithm classification report for 10-fold cross-validation

Class	Precision	Recall	f1-score	Support
High	1.00	0.99	0.99	70
Low	0.96	1.00	0.98	70
Medium	1.00	0.97	0.99	70
Accuracy			0.99	210
Macro average	0.99	0.99	0.99	210
Weighted average	0.99	0.99	0.99	210

A confusion matrix consists of a visualization of the number of erroneously labeled and correctly labeled results, represented in the matrix diagonal (Skansi, 2018). The confusion matrix for the best model can be seen in (Figure 3), demonstrating a low rate of erroneously classified results by the RF algorithm, which shows itself as a good classifier for our problem.

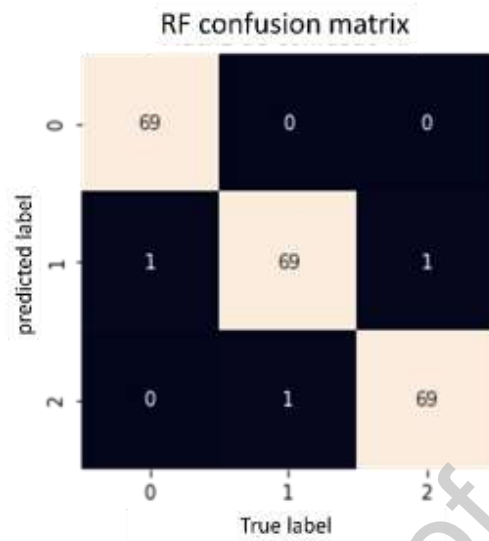


Figure 3 – RF algorithm confusion matrix. Where 0 – represents the high drought tolerance class, 1 – low, and 2 – medium. Source: Verslype *et al.* (2021).

The RF algorithm assigned a rating scale of importance to all the features evaluated in the data set. Among these, we observed that the variables cutting production, the onset of leaf fall, anthracnose resistance, geotropic angle, calcareous soil tolerance, and root distribution up to 60 cm were the first variables with the notable importance in their decision-making, as we can see it in (Figure 4).

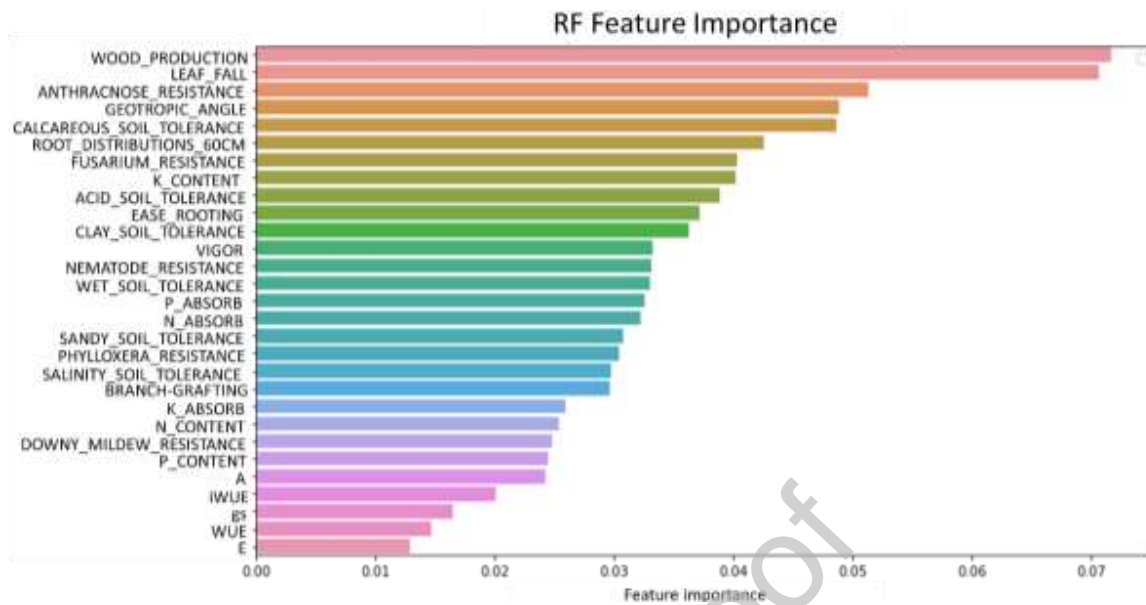


Figure 4 – Feature importance to the RF algorithm decision making to predict drought tolerance classes for grapevine rootstock cultivars. Source: Verslype *et al.* (2021).

The RF decision-making high importance attribution to wood production per plant can be explained by the fact plants with more efficient use of water produce more dry matter per gram of transpired water. Given the need for an average intake of 1 liter for each 2 grams of grapevine dry mass-produced (Ojeda *et al.*, 2004; Serra *et al.*, 2014). Consequently, these more vigorous plants can tolerate conditions with lower water availability in general (Ojeda *et al.*, 2004; Villa, 2018). While, the feature of leaf fall is related to drought tolerance, as at this stage the grapevines direct the photoassimilates to the roots. This enables the absorption area = to increase and recover from water stress (Ojeda *et al.*, 2004).

The greater importance attributed to the variables roots distribution up to 60 cm and geotropic angle in RF decision-making can be explained by the geotropic angle that determines the depth capacity of the grapevine root system (Villa, 2018). In this sense, the literature has evidenced the closed geotropic angles allow the deepening of the roots and consequently make it possible to reach deeper layers of soil remain moist, as well as to avoid the absorption of harmful elements to plants in saline soils through selective absorption (Smart *et al.*, 2006; Satisha *et al.*,

2006; Serra *et al.*, 2014; Villa, 2018). These features enable rootstocks to better adapt to water stress conditions, and consequently, a greater drought tolerance (Satisha *et al.*, 2006; Villa, 2018).

The RF model also assigned a high importance to the variable anthracnose resistance, considered a disease with a wide impact on grapevines, as it mainly affects young and tender tissues in the aerial part of the plant (Lima *et al.*, 2009). The drought stress associated with anthracnose predisposition occurrence was studied by Erbaugh *et al.* (1995) for the *Cornus florida* L. species, which observed a significant increase in disease severity in plants shaded and under drought conditions. In addition, Hsiang (2006) mentions how the *Poa annua* L. species becomes unable to respond quickly to anthracnose infection and overcome the disease when subjected to severe stresses such as cyclical water stress. In this sense, it may be a signal of some degree of interference between the two variables in grapevines.

As mentioned above, after evaluating accuracy, precision, recall, and f1-score metrics, the RF was identified as the best model produced in the scenarios and datasets test to predict drought tolerance classes in grapevine rootstocks. The prediction of drought tolerance classes indicated by the RF algorithm was high tolerance for cultivars IAC 313 and IAC 766, and a low tolerance for cultivar IAC 572. Despite being an initial study, this approach with supervised algorithms proved to be a helpful and accessible strategy for the breeder, that can help develop these cultivars by predicting, identifying and selecting promising genotypes. We would like to acknowledge some limitations of the present work, which still needs validation of the predicted classes for the IAC 313, 766 and 572 through studies in field conditions under drought conditions. Besides, the training dataset is rather a limited sample of grapevine rootstocks. In future works, we intend to increase the size of the training dataset or explore data augmentation strategies.

4 CONCLUSION

This paper proposed a pipeline to predict rootstock drought tolerance for the São Francisco Valley in Brazil, through the use of machine learning methods. A manually curated dataset was produced, which was used to evaluate several machine learning algorithms. The results indicated that the best performing classifiers were RF, followed by XGB with 98.57% and 96.19% correctness to predict drought tolerance classes of grapevine rootstocks.

The best RF hyperparameter configuration was *max_features=2* and the number of estimators equal to 1000 for the evaluated scenarios and datasets. While, the LDA classifier had the lowest efficiency, with a hit rate of 67.14%, to predict the classes of tolerance to water deficit.

The trained model was used to predict drought tolerance of 03 the grapevine rootstock cultivars that have never been experimentally evaluated. The model classified grapevine rootstock cultivars IAC 313 and IAC 766 with high drought tolerance, while IAC 572 with a low tolerance. In this regard, cultivars IAC 313 and IAC 766 could be the best option for vintners in the São Francisco Valley. It is important to note that this is a preliminary study, and these predictions, are an indicative of the real drought tolerance of the three cultivars, requiring more research to validate such results.

Machine learning algorithms demonstrated to be a helpful tool in plant breeding studies that can contribute to identifying and selecting drought-tolerant grapevine rootstock genotypes in breeding programs. As a suggestion to extend this study, we envision the validation of the predictions obtained by the ML model in field tests under drought stress, as well as an increase of the training dataset.

5 Acknowledgments

We would like to thank the National Council for Scientific and Technological Development (CNPq), for granting a master's scholarship to the first author.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

6 References

ATVB – Association Technique Viticole de Bourgogne (2013) **Caractéristiques des porte-greffes et greffons utilisés en Côte-d’Or**. Retrieved from <http://www.atvb-bourgogne.com/>. Accessed August 27, 2020.

Ashraf, M. (2010). Inducing drought tolerance in plants: recent advances. **Biotechnonology advance** **28**, 169-183.

Audeguin, L., Bécart, V., Besnard, E., Boursiquot, J.M., Ciccolini, G., Dedet, S., Delorme, G., Follet, J., Froelhy, A., Garcin, L., Gonin, N., Lacombe, T., Laffargue, F., Marchal, C., Pesato, M., Sereno, C., Uriel, G., Wentzel, & M., Yobregat, O. (2020) **PI@ntGrape**: le catalogue des vignes cultivées en France. Montpellier, France: IFV – INRA – Montpellier SupAgro. Retrieved from <http://plantgrape.plantnet-project.org/>. Accessed September 15, 2020.

Beiki, A.H. (2019). Molecular Predicting Drought Tolerance in Maize Inbred Lines by Machine Learning Approaches. **BioRxiv** **10**, 1-33.

Breiman, L., Friedman, J., Olshen, R.A., & Stone, C.J. (1984). **Classification and Regression Trees**. Wadsworth, Belmont, CA: CRC press.

Breiman, L. (2001). Random forests. **Machine learning**, 45(1), 5-32.

Camargo, U.A., & Amorim, F.D. (2004). Análise dos atuais sistemas de produção de uvas para vinho no Vale do São Francisco. In **Workshop Internacional De Pesquisa** **1**, 97-101.

Camargo, U.A., Ritschel, P.S., Maia, J.D.G., Quecini, V., Machado, C.A.E., Bosco, D.D., Sinski, I., Comachio, & V., Zilio, R. (2017). **Banco ativo de germoplasma de uva**. Embrapa Uva e Vinho, Bento Gonçalves. Retrieved from <https://www.embrapa.br/en/uva-e-vinho/banco-ativo-de-germoplasma-de-uva/pesquisa>. Accessed September 16, 2020.

Carroll, B. (2020). **Characteristics of *Vitis* grape rootstocks and adaptability**. Oklahoma State University, Oklahoma, 4p.

Chaves, M.M., Zarrouk, O., Costa, J.M., Santos, T., Regalado, A.P., Rodrigues, M.L., & Lopes, C.M. (2010). Grapevine under deficit irrigation: hints from physiological and molecular data. **Annals Appl. Biol.** **1**, 661-676.

- Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In **Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining** (pp. 785–794). New York, USA: ACM.
- Chevalier, A. (1925). L'amélioration de la vigne en France et les Travaux de G. Couderc sur l'Hybridation et le Greffage. **Revue de botanique appliquée et d'agriculture coloniale** **52**: 926-945.
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. **Machine Learning** **20**, 273–297.
- Dayer, S., Reingwartz, I., McElrone, A.J., & Gambetta, G.A. (2019). Response and Recovery of Grapevine to Water Deficit: From Genes to Physiology. In Cantu, D., Walker, M. (eds) **The Grape Genome**. Compendium of Plant Genomes. Springer, Switzerland, p. 223-245.
- Erbaugh, D.K., Windham, M.T., Stodola, A.J., & Augé, R.M. (1995). Light intensity and drought stress as predisposition factors for dogwood anthracnose. **Journal of Environmental Horticulture**, **13**(4), 186-189.
- Etminan, A., Pour-Aboughadareh, A., Mohammadi, R., Shooshtari, L., Yousefiazarkhanian, M., & Moradkhani, H. (2019). Determining the best drought tolerance indices using artificial neural network (ANN): Insight into application of intelligent agriculture in agronomy and plant breeding. **Cereal Research Communications**, **47**: 170-181.
- Ferreira, P.V. (2011). **Estatística experimental**. Maceió. Ed.UFAL. 559p.
- Fix, E., & Hodges, J.L. (1951). **Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties**. Texas: USAF School of Aviation Medicine, Randolph Field.
- Forsyth, D. (2018). **Probability and statistics for computer science** (pp. 36-42). Springer International Publishing.
- Fort, K., Fraga, H., Grossi, D., & Walker, M.A. (2017). Early measures of drought tolerance in four grape rootstocks. **J. Amer. Soc. Hort. Sci.**, **1**: 36-46.
- Fraga, H., Malheiro, A.C., Moutinho-Pereira, J., & Santos, J.A. (2013). An overview of climate change impacts on European viticulture. **Food Energy Security**, **1**: 94-110.
- Hsiang, T. (2006). **Anthracnose**: terrorizing turf grass. GCSAA: 2p. Retrieved from <https://www.gcsaa.org/>. Accessed May 14, 2020.
- Jones, G.V. (2016). Grapevines in a changing environment: a global perspective. In Gerós H, Chaves MM, Gil HM and Delrot S (eds) **Grapevine in a changing environment**: a molecular and ecophysiological perspective. John Wiley & Sons, UK, p. 23-34

Jones, G.V., Reild, R., & Vilks, A. (2012). Climate, grapes, and wine: structure and suitability in a variable and changing climate In Dougherty, Percy H (eds) **The Geography of Wine: regions, terroir and technique**. Springer, Netherlands, p. 109-36.

Jones, G.V., & Webb, L.B. (2010). Climate change, viticulture, and wine: challenges and opportunities. **Journal of Wine Research**, **21**: 103-106.

Kathpalia, R., & Bhatla, S.C. (2018). Plant Water Relations. In Bhatla SC and Lal MA (eds.) **Plant physiology, development and metabolism**. Springer, Singapore, p. 3- 36p.

Lima, M.R., Lopes, D.B., Tavares, S.C.C.H., Tessmann, D.L., & Melo, N.F. (2009). Doenças e alternativas de controle. In Soares, J.M., Leão, P.C.S. (eds.) **A vitivinicultura no Semiárido Brasileiro**. Brasília, DF, Embrapa Informação Tecnológica, Petrolina, Embrapa Semi-Árido, p.543-596.

Leão, P.C.S., & Silva, D.J. (2014). Cultivo da videira no Semiárido brasileiro. In Pio R (ed) **Cultivo de fruteiras de clima temperado em regiões subtropicais e tropicais**. UFLA, Lavras, p. 578-618.

Marguerit, E., Brendel, O., Lebon, E., Van Leeuwen, & C., Ollat, N. (2012). Rootstock control of scion transpiration and its acclimation to water deficit are controlled by different genes. **New Phytologist**, **194**: 416–429.

Maul, E., Töpfer, R., Röckel, F., Brühl, U., Hundemer, M., Mahler-Ries, A., Walk, M., Kecke, S., Wolck, A., & Ganesch, A. (2020). **Vitis International Variety Catalogue**. Federal Research Centre for Cultivated Plants, Institute for Grapevine Breeding. Retrieved from <https://www.vivc.de/>. Accessed July 30, 2020.

McKinney, W. (2010). **Data Structures for Statistical Computing in Python**. In Stéfan van der Walt and Jarrod Millman, editors, Proceedings of the 9th Python in Science Conference, pages 56 – 61.

Moura, M.S.B., Teixeira, A.H.C, & Soares, J.M. (2009). Exigências climáticas. In Soares, J.M., Leão, P.C.S. (eds.) **A vitivinicultura no Semiárido Brasileiro**. Brasília, DF, Embrapa Informação Tecnológica, Petrolina, Embrapa Semi-Árido, p.36-69.

Ojeda, H., Deloire, A., Wang, Z., & Carbonneau, A. (2004). Detereminación y control del estado hídrico de la vid: efectos morfológicosy fisiológicos de la restricción hídrica en vides. **Viticultura/Enología Profesional 90**: 27-43.

Ollat, N., Peccoux, A., Papura, D., Esmenjaud, D., Marguerit, E., Tandonnet, J.P., Bordenave, L., Cookson, S.J., Barriue, F., Rossdeutsch, L., Lecourt, J., Lauvergeat, V., Vivin, P., Bert, P.F., & Delrot, S. (2016). Rootstocks as a component of adaptation to environment. In Gerós, H., Chaves, M.M., Gil, H.M.,

- Delrot, S. (eds) **Grapevine in a changing environment: a molecular and ecophysiological perspective**. John Wiley & Sons, UK, p. 102-152
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. **Journal of Machine Learning Research** 12: 2825-2830.
- Rao, R.C. (1948). The Utilization of Multiple Measurements in Problems of Biological Classification. **Journal of the Royal Statistical Society** 10: 159-203.
- Rebschule Mueller (2020). **Unterlagen für den weinbau**. Retrieved from <http://www.rebschulemueller.de/>. Accessed August 15, 2020.
- Riaz, S., Pap, D., Uretsky, J., Laucou, V., Boursiquot, J. M., Kocsis, L., & Walker, M. A. (2019). Genetic diversity and parentage analysis of grape rootstocks. **Theoretical and Applied Genetics** 132: 1847–1860.
- Satisha, J., Prakash, G.S., Murti, G.S.R., & Upreti, K.K. (2006). Response of grape rootstocks to soil moisture stress. **Journal of Horticultural Sciences** 1: 19-23p.
- Serra, I., Strever, A., Myburgh, P.A., & Deloire, A. (2014). The interaction between rootstocks and cultivars (*Vitis vinifera* L.) to enhance drought tolerance in grapevine. **Australian Journal of Grape and Wine Research** 20: 1-14.
- Shalev-Shwartz, S., & Ben-David, S. (2014). **Understanding machine learning: from theory to algorithms**. Cambridge university press, New York, 409 p.
- Skansi, S. (2018). **Introduction to Deep Learning: from logical calculus to artificial intelligence**. Springer. 196p.
- Smart, D.R., Schwass, E., Lakso, A., & Morano, L. (2006). Grapevine rooting patterns: a comprehensive analysis and a review. **Am. J. Enol. Vitic.**, 57: 89 -104.
- Storm, B., & Krasokhina, C. (2020). **Rootstock grape varieties, gardening varieties, wild grapes, breeding varieties**. Retrieved from <https://vinograd.info/>. Accessed August 5, 2020.
- Sunridge Nurseries (2020). **Rootstock chart**. Bakersfield, California. Retrieved from <https://www.sunridgenurseries.com/>. Accessed August 25, 2020.
- Taiz, L., Zeiger, E., Møller, I. M., & Murphy, A. (2017). **Fisiologia e desenvolvimento vegetal**. Artmed Editora, Porto Alegre, 888p.
- VICOPAD - Vivai Cooperativi di Padergnone (2020). **Scelta del portinnesto**. Padergnone, Italy. Retrieved from <http://www.vicopad.it/>. Accessed August 19, 2020.
- Villa, P. (2018). **Cultivar la vid**. Editorial de Vecchi, USA, 160p.

Walker, M.A., Heinitz, C., Riaz, S., & Uretsky, J. (2019). Grape Taxonomy and Germplasm. In Cantu D and Walker MA (eds) **The grape genome**. Springer, Switzerland, p. 25-38.

Waskom, M. (2021). Seaborn: statistical data visualization. **The Open Journal** 60: 3021p. <https://doi.org/10.21105/joss.03021>.

Wine Australia (2016). **Grapevine rootstock selector tool**. Australia. Retrieved from <https://grapevine-rootstock.com/>. Accessed August 25, 2020.

Zhang, L., Marguerit, E., Rossdeutsch, L., Ollat, N., & Gambetta, G.A. (2016). The influence of grapevine rootstocks on scion growth and drought resistance. **Theoretical and Experimental Plant Physiology**, 28: 143-157.

Zarrouk, O., Costa, J.M., Francisco, R., Lopes, C., & Chaves, M.M. (2016). Drought and water management in Mediterranean vineyards In Gerós H, Chaves MM, Gil HM and Delrot S (eds) **Grapevine in a changing environment: a molecular and ecophysiological perspective**. John Wiley & Sons, UK, p. 67 – 88.

Journal Pre-proof