

1 Highlights

2 **???Grape Quality and its Link to Regional Differences in the Aus-**
3 **tralian Winegrowing Industry**

4 Author

5 • ???

6 • ???

7 • ????

8 • ????

29 2. Methods

30 2.1. Data

31 Data used in this analysis were obtained from Sustainable Winegrowing
32 Australia. Australia’s national wine industry sustainability program, which
33 aims to facilitate grape-growers and winemakers in demonstrating and im-
34 proving their sustainability (SWA, 2022). Data recorded by the SWA is
35 entered manually by winegrowers using a web based interface tool. A total
36 of 6091 observations were collected from 2012 to 2022. Each observation
37 contained 23 variables reflecting a vineyards account for the given year (see
38 Table 2.1).

39 The data originally contained only two multiclass variables: year and re-
40 gion. Variables that measured the same metric from different sources (such
41 as water collected from rivers versus water from dams) were converted into
42 multiclass variables representing the source. The total amount used from
43 these variables was retained as a separate variable. Occurrences of multiple
44 sources were defined as separate classes. As harvest does not run by calendar
45 year, years are in financial years. Region represents one of the 65 Geographi-
46 cal Indicator Regions (GI Region) used to describe different unique localised
47 traits of vineyards across Australia (Halliday, 2009; Oliver et al., 2013; SOAR
48 et al., 2008). Each region is explicitly defined under the Wine Australia Cor-
49 poration Act of 1980 (Attorney-General’s Department, 2010). Profit was also
50 used as a binary variable, depicting whether a vineyard was profitable or not.

Table 1: Summary of variables used in the analysis. The recorded column indicate values that were either greater than zero or that were not missing.

Variable	Units	Recorded	Number of Classes
Water Used	Mega Litres	5846	
Diesel	Litres	5585	
Biodiesel	Litres	25	
LPG	Litres	958	
Herbicide Spray	Times per year	2026	
Year	Class	6091	10
Disease	Class	6091	2
Region	Class	6091	58
Solar	Kilowatt Hours	622	
Irrigation Type	Class	6091	20
Petrol	Litres	4309	
Slashing	Times per year	2290	
Yield	Tonnes	5935	
Irrigation Energy	Class	6091	16
Area Harvested	Hectares	6091	
Electricity	Kilowatt Hours	1015	
Insecticide Spray	Times per year	1092	
Fertiliser	Kilograms of Nitrogen	795	
Fungicide Spray	Times per year	2260	
Cover Crop	Class	6091	32
Water Type	Class	6091	39
Profit	AUD	³ 853	
Operating Costs	AUD	853	

51 2.2. XGBoosted Trees

52 XGBoosted (eXtreme Gradient Boosting) trees were created using the
 53 XGBoost library (Chen and Guestrin, 2016) in the Python Programming
 54 language (G. van Rossum, 1995). They were chosen for this analysis as they
 55 provide both a high predictive performance and ability to effectively capture
 56 complex relationships. An XGBoosted tree was created for each variable to
 57 show how they interacted. Each tree included all but the economic vari-
 58 ables (profit and operating cost), which were only included once as predicted
 59 variables.

60 Following Chen and Guestrin (Chen and Guestrin, 2016), XGboosted
 61 trees predict a value y_i from the input x_i . The method of prediction is
 62 achieved through a tree ensemble model, using K additive functions to pre-
 63 dict the output.

$$\hat{y}_i = \phi(x_i) = \sum_{k=1}^K f_K(x_i), f_K \in \mathcal{F}, \quad (1)$$

64 where each function f_K is a classification or regression tree, such that all
 65 functions are in the set of all decision trees \mathcal{F} , defined by $f(x) = \omega_{q(x)}(q :$
 66 $\mathbb{R}^m \rightarrow T, \omega \in \mathbb{R}^T)$. Where, f_K corresponds to an independent tree structure
 67 q of ω weights. Each tree has T leaves, which contain a continuous score,
 68 represented by ω_i for the i -th leaf. The final prediction is determined by the
 69 sum of the score of the corresponding leaves, given by ω . The set of func-
 70 tions used by the tree is determined by minimising the regularised objective
 71 function, given by:

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i^{t-1} + f_t(x_i)) + \sum_k \Omega(f_K). \quad (2)$$

72 The difference between the prediction and actual variable is a convex loss
 73 function l . To optimise l , the difference is calculated for the i -th instance
 74 at the t -th iteration. The function f_t is selected according to which value
 75 minimises (2). The model complexity is penalised by the function Ω , this
 76 acts to smooth weights in an attempt to prevent over fitting.

77 As predictions are made using additive tree functions, XGboosted trees
 78 can be used for classification and regression. Due to the mixture of continu-
 79 ous, binary and multiclass variables in this analysis, both classification and
 80 regression trees were created. The difference between the trees created for
 81 this analysis was the objective function used. XGBoosted regression trees
 82 were created for continuous variables, using the root-mean-square as the ob-
 83 jective function. Binary class variables utilised the logistic loss function as
 84 the objective. And, Multiclass variable used the soft max function. All objec-
 85 tive functions are defined within the SKlearn library (Buitinck et al., 2013),
 86 linked via an API to the XGBoost library (Chen and Guestrin, 2016).

87 Chen and Guestrin (Chen and Guestrin, 2016) further illustrate, using
 88 Taylor expansions, that for a fixed structure $q(x)$ the optimal weight ω_j^* for
 89 a leaf j can be derived. Furthermore, they show the loss reduction after the
 90 split is given by the function:

$$\mathcal{L}_{split} = \frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] - \gamma, \quad (3)$$

91 with the tree structure defined using left I_L and right I_R instance sets of

92 nodes, with $I = I_L \cup I_R$. Instead of enumerating all possible tree structures,
 93 a greedy algorithm iteratively adds branches to the tree minimising \mathcal{L}_{split}
 94 in (3). The frequency of a variable’s occurrence within a tree is directly
 95 attributed to the minimisation of the objective function (or loss) through
 96 the minimisation of \mathcal{L}_{split} .

97 The frequency of a variable appearing as a node within the ensemble was
 98 used as a measure of importance. This measure was chosen as it connected
 99 a variable to the minimisation of its associated objective function, trans-
 100 lating the value into a simple count metric. Creating XGBoosted trees for
 101 each variable allowed the use of importance to show how strongly variables
 102 were associated with each other. The importance of predictor variables to
 103 economic variables was illustrated through the use of Sankey diagrams con-
 104 structed using the Holoviews python library (Rudiger et al., 2020). Other
 105 variable’s interconnectedness was demonstrated through the use of a chord
 106 diagram also created using Holoviews.

107 Each variable utilised 80% of the data to train the XGBoost ensemble,
 108 with 20% reserved for testing and validation. Testing was done through the
 109 iterative minimisation of the respective objective function for the variables
 110 type. For continuous variables 20% was used as testing data, minimising the
 111 root-mean-square function. The final model was validated using repeated k-
 112 fold cross validation for 10 folds, repeated 10 times. For binary and multiclass
 113 variables data was split into 80% training, 10% testing and 10% validation
 114 data. Due to class disparity in multiclass variables (most prominently in
 115 region) data was stratified into each subset at the same ratio of class oc-
 116 currence. Validation was summarised through confusion matrices and their

117 associated accuracy

118 The use of the XGBoost library incorporates regularisation techniques
119 built into the software to mitigate over-fitting and enhance model gener-
120 alisation. The further use of cross validated grid search functions allowed
121 for the selection of better performing hyperparameters when selecting the
122 final model. The performance measure for model selection was root-mean-
123 square error for continuous variables. The receiver operator characteristic's
124 area under the curve was used for category variables (Hanley and McNeil,
125 1982). Multiclass variables utilised the one verse one approach to minimise
126 sensitivity to class disparity (Ferri et al., 2009; Hand and Till, 2001).

127 *2.3. Classification and Regression Trees*

128 Classification and Regression Trees were created for region, year, profit
129 and operating cost. These models describe the partitions that are useful
130 in predicting these variables; giving insight into the trees that make up the
131 ensembles created by XGBoost. These trees were created using the rparts
132 and caret packages (Kuhn, 2008; Terry Therneau and Beth Atkinson, 2022)
133 in the R statistical programming language (R Core Team, 2021).

134 Classification trees were validated using K-fold cross validation. Each
135 model was validated using 10 folds, utilising a random selection of different
136 samples ten separate times to validate each of the classification trees. A
137 summary confusion matrix was then constructed to show the class bias and
138 overall accuracy of each tree.

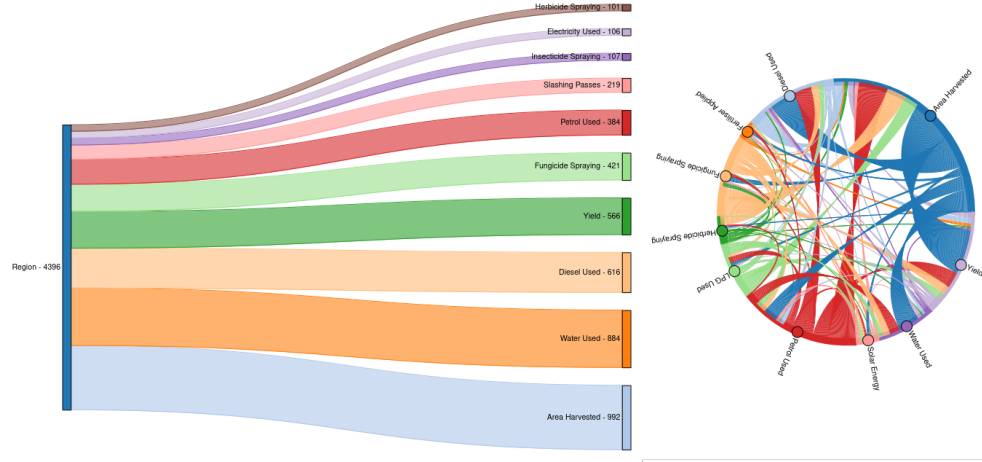


Figure 2: The left-hand side depicts the 10 most important variables in predicting Region using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

3. Results

3.1. Region

3.2. Year

3.3. Operating Costs

3.4. Profit

3.5. Validation

3.6. Model 1 GI Regions

The first Model was used to classify GI regions and resulted in an accuracy of 36.48% across 52 classes. The most prominent features used to classify regions were the types of water resources available (see Figure 1). Two regions, the Riverland and Coonawarra, were the most accurate classes being

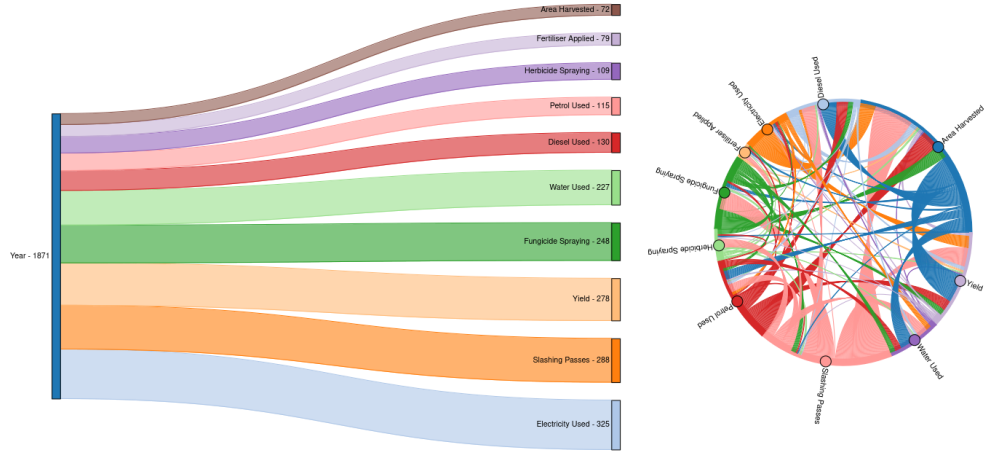


Figure 4: The left-hand side depicts the 10 most important variables in predicting Year using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

Table 2: Validation and training accuracies of each multiclass variable.

Variable	Validation	Training
cover crops	0.364086	0.396418
water type	0.742097	0.928905
profitable	0.705882	0.719737
irrigation type	0.841845	0.847554
giregion	0.505824	0.568242
irrigation energy	0.746293	0.836405
data year id	0.422003	0.518059

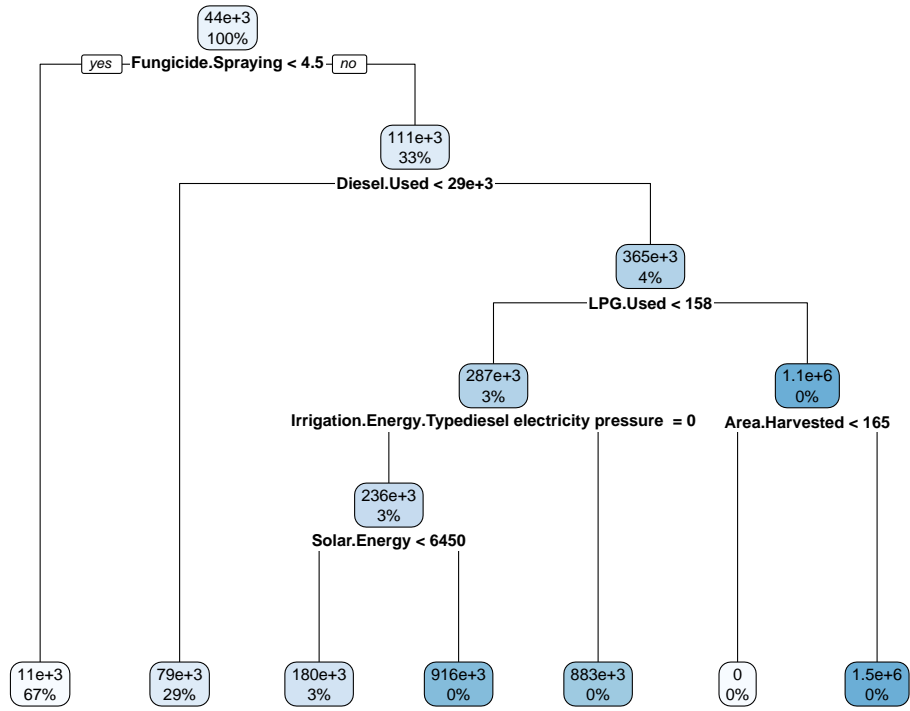


Figure 5: Decision tree predicting Operating Costs. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

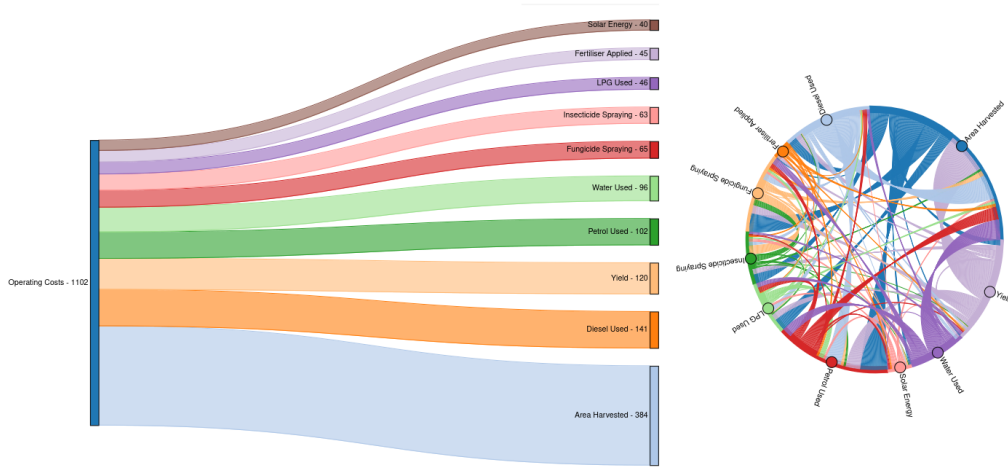


Figure 6: The left-hand side depicts the 10 most important variables in predicting Operating Costs using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

92.74% and 96.97% respectively. These regions differ greatly in practice and geophysical properties, with the Riverland being a dry warm inland region and Coonawarra being a cooler, wet coastal region. However, they are both similar in operational scales, with vineyards being relatively large compared with other regions. The differences in resources and practices between these regions are also significant, such as the Riverland utilising the river Murray as a water source. Many of the regions had significantly lower reporting rates, resulting much poorer classification performance. The regions with the most samples performed the best (see Table 1). Notably bordering regions were routinely grouped together and misclassified as the same region, for example the two closest regions to Coonawarra, Padthaway and Wratttonbulley, were misclassified as Coonawarra even though they had 147 and 137 samples

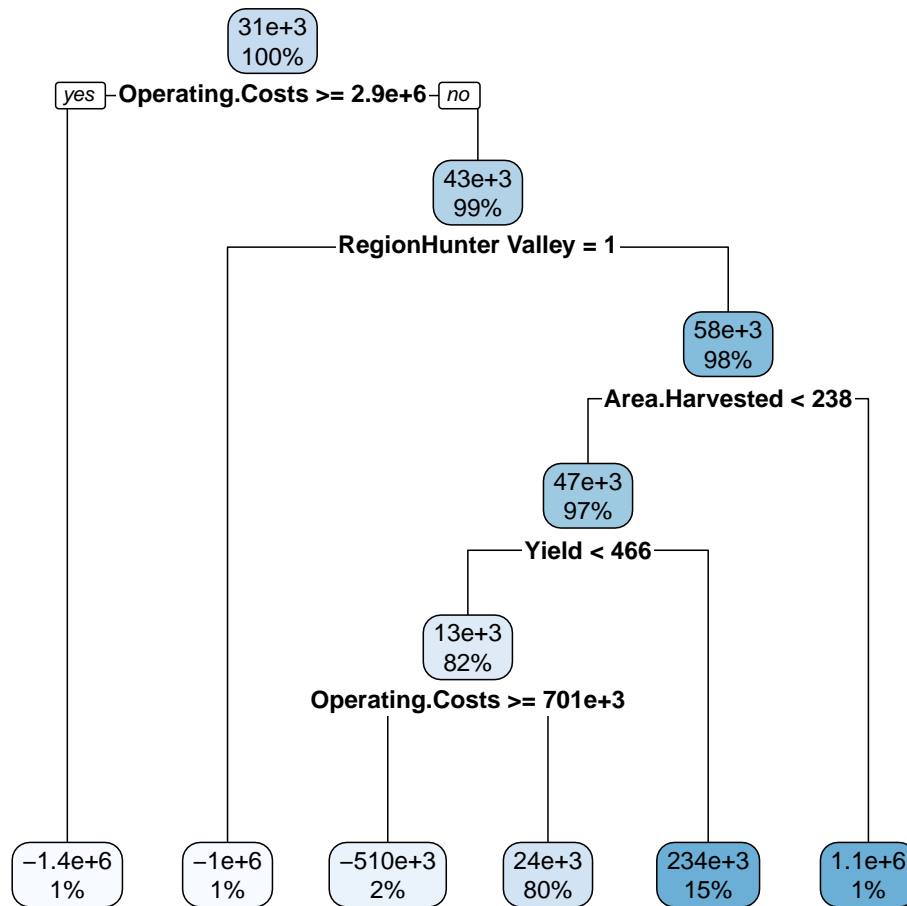


Figure 7: Decision tree predicting Profit. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

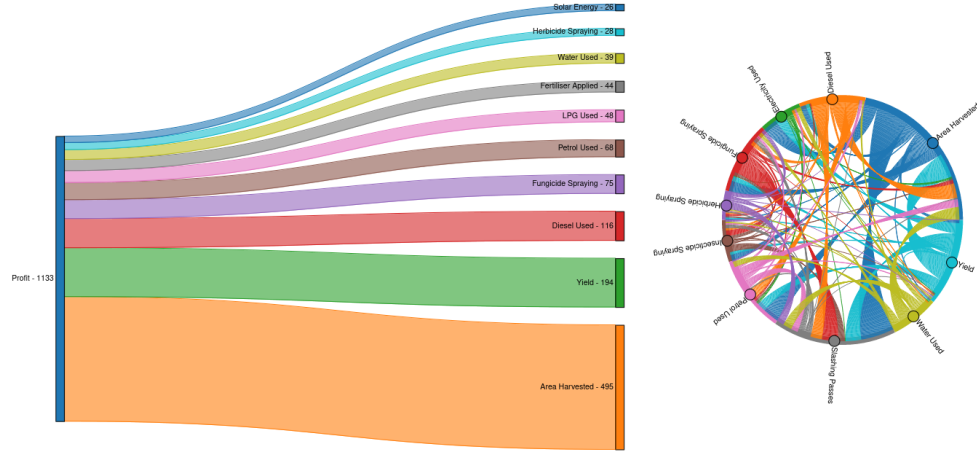


Figure 8: The left-hand side depicts the 10 most important variables in predicting Profit using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

162 respectively. The same case was found for the Murray Darling, with 143 sam-
 163 ples, it was misclassified as the Riverland. These misclassifications are likely
 164 due to the incredibly similar regional properties and close proximity these
 165 regions have with one another. Other misclassifications were most likely due
 166 to lower reporting rates with many regions being under represented.

167 3.7. Climate

168 Classifying the SWA climatic categorisation of the given regions had bet-
 169 ter performance than the GI Regions, with 41.66% being classified correctly.
 170 These categories were divided into 12 climatic classifications with 3 and 4
 171 separate subsets for rainfall and temperature respectively. The decision tree
 172 behaved similarly and over classified climates with higher response rates. The
 173 results posed an interesting similarity with grape quality classifier, being in-

174 fluenced predominantly by water and area. The use of fungicide to separate
175 regions that were 'Very dry' and 'Damp' can be considered as indicative
176 of the different practices required due to climatic pressure; fungicides being
177 more prominent in cooler regions with greater rainfall due to the higher risk
178 of disease pressure (Reynolds, 2010). This could also potentially explain the
179 use of contractor tractor use to discern differences in grape quality, where the
180 lack of contractor use to prevent disease could have led to lowered quality of
181 grapes.

182 3.7.1. *Rainfall*

183 The rainfall decision tree showed a greater use of fungicides sprays to
184 discern between damp and very Dry as shown in Figure 4; with the accuracy
185 improving to 62% but was unable to effectively discern between dry and very
186 dry regions (see Table 3).

187 3.7.2. *Temperature*

188 The classification of GI Regions by their temperatures (see Figure 5)
189 showed similarities to the other trees, with a heavy reliance on the types
190 of water resources used as dominant predictors. The use of contractors was
191 again used to differentiate between warm and cool regions, likely being due
192 to disease pressure. The temperature classification tree was only a minor
193 improvement over the regional classification tree, with an accuracy of 49.26%
194 as shown in the confusion matrix (see Table 4).

195 3.8. *Model 3 Grape Quality*

196 The classification of grape quality through its grade had an accuracy of
197 55.72% across 5 separate grades. There was a notable issue with the classi-

198 fication of B grade grapes when compared to A and C (see Table 2). The
199 classification tree itself shows similarities to that of classifying regions in
200 Model 1, with the type of water resource used being a prominent determiner.
201 Although not surprising the number of contractor tractor passes is new de-
202 ciding factor due disease and pests reducing the potential quality of a crop.
203 The prevalence of contractor use is greater in regions such as the Barossa
204 Valley and the McLaren Vale, this could be due to the difference in opera-
205 tional scales, with larger sites being more likely to have ownership of their
206 own equipment for weeding and spraying due to the cost benefit.

207 **4. Discussion**

208 The difference between grape quality is most notable between warm in-
209 land regions and coastal regions such as the Riverland and Coonawarra,
210 respectively. Grape quality is only described by a singular variable within
211 this study, however in reality it is driven by market demand and subject to
212 complex forces such as international market pressure, fire, pests and disease
213 (Wine Australia, 2019, 2020, 2021, 2022; Winemakers' Federation of Aus-
214 tralia, 2015, 2016, 2017, 2018) The decision trees were able to offer some
215 insights into the factors that influence grape quality and regional contrasts
216 that contribute to different qualities. The most prominent being what readily
217 available resources of each region were, particular the types of water available.
218 Heavy water consumption is often linked to the mass production of grapes,
219 where lower quality grapes are targeted in a quantity over quality strategy.
220 These types of business decisions are unfortunately obfuscated by lack of in-
221 depth data regarding vineyard business plans. Notably the literature shows

that there are many complex decisions to be made on the ground depending on many compounding factors that influence both quality and yield (Abad et al., 2021; Cortez et al., 2009; Hall et al., 2011; I. Goodwin, et al., 2009; Kasimati et al., 2022; Oliver et al., 2013; Srivastava and Sadistap, 2018).

. There are also further differences when comparing winegrowers to other agricultural industries as they are vertically integrated within the wine industry, tying them to secondary and tertiary industries, such as wine production, packaging, transport and sales. This results in unique issues, where on-the-ground choices are influenced by other wine industry’s decisions, such as the use of sustainable practices in vineyards to sell in overseas markets; notably these interactions are further complicated by some winegrowers being totally integrated into wine companies, while others are not (Knight et al., 2019). It is incredibly difficult to attribute external business decisions to produced grape quality but it is important to acknowledge that some growers are contracted to produce grapes of a particular grade; it is difficult to know whether another consumer may have graded the grape quality differently paying more or less for the same grapes given the opportunity to purchase them. It is difficult to untangle the contributing factors to the success of winegrowers and the quality of grapes produced without further specifics of choices made through out a season (Leilei He et al., 2022).

5. Conclusion

The type and availability of water resources were a major contributing factor when classifying grape quality and region. This was seen in the two most accurately classified regions, Coonawarra and the Riverland, with the

246 Riverland predominantly utilising river water. Furthermore, the study high-
 247 lighted the influence of water use, fungicide application, and contractor use in
 248 differentiating grape quality, climate and region respectively. These models
 249 provide insight into the complex dynamics between regional characteristics,
 250 sustainable practices, and grape quality in the Australian winegrowing indus-
 251 try. It is important to acknowledge that grape quality is subject to external
 252 influences such as market demands and prior established business arrange-
 253 ments. Further in-depth data and understanding are necessary to fully grasp
 254 the nuances of decision-making and the interplay of factors impacting grape
 255 quality.

256 References

- 257 Abad, J., Hermoso de Mendoza, I., Marín, D., Orcaray, L., Santeste-
 258 ban, L.G., 2021. Cover crops in viticulture. A systematic review (1):
 259
Implications on soil characteristics and biodiversity in vineyard.
 260 OENO One 55, 295–312. doi:10.20870/oeno-one.2021.55.1.3599.
- 261 Abbal, P., Sablayrolles, J.M., Matzner-Lober, É., Boursiquot, J.M., Baudrit,
 262 C., Carbonneau, A., 2016. Decision Support System for Vine Growers
 263 Based on a Bayesian Network. Journal of agricultural, biological, and
 264 environmental statistics 21, 131–151. doi:10.1007/s13253-015-0233-2.
- 265 Agosta, E., Canziani, P., Cavagnaro, M., 2012. Regional climate variability
 266 impacts on the annual grape yield in Mendoza, Argentina. Journal of
 267 Applied Meteorology and Climatology 51, 993–1009.

268 Attorney-General's Department, 2010. Wine Australia Corporation Act
 269 1980.

270 Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel,
 271 O., Niculae, V., Prettenhofer, P., Gramfort, A., Grobler, J., Layton, R.,
 272 VanderPlas, J., Joly, A., Holt, B., Varoquaux, G., 2013. API design for
 273 machine learning software: Experiences from the scikit-learn project, in:
 274 ECML PKDD Workshop: Languages for Data Mining and Machine Learn-
 275 ing, pp. 108–122.

276 Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System,
 277 in: Proceedings of the 22nd ACM SIGKDD International Conference on
 278 Knowledge Discovery and Data Mining, ACM, New York, NY, USA. pp.
 279 785–794. doi:10.1145/2939672.2939785.

280 Cortez, P., Teixeira, J., Cerdeira, A., Almeida, F., Matos, T., Reis, J., 2009.
 281 Using data mining for wine quality assessment, in: Discovery Science: 12th
 282 International Conference, DS 2009, Porto, Portugal, October 3-5, 2009 12,
 283 Springer. pp. 66–79.

284 Ferri, C., Hernández-Orallo, J., Modroi, R., 2009. An experimental com-
 285 parison of performance measures for classification. Pattern Recognition
 286 Letters 30, 27–38. doi:10.1016/j.patrec.2008.08.010.

287 Fraga, H., Costa, R., Santos, J.A., 2017. Multivariate clustering of viticul-
 288 tural terroirs in the Douro winemaking region. Ciência Téc. Vitiv. 32,
 289 142–153.

- 290 G. van Rossum, 1995. Python tutorial, Technical Report CS-R9526. Centrum
291 voor Wiskunde en Informatica (CWI),.
- 292 Hall, A., Lamb, D.W., Holzapfel, B.P., Louis, J.P., 2011. Within-season
293 temporal variation in correlations between vineyard canopy and winegrape
294 composition and yield. *Precision Agriculture* 12, 103–117.
- 295 Halliday, J.C.J.C., 2009. Australian Wine Encyclopedia. Hardie Grant
296 Books, VIC.
- 297 Hand, D.J., Till, R.J., 2001. A Simple Generalisation of the Area Under the
298 ROC Curve for Multiple Class Classification Problems. *Machine Learning*
299 45, 171–186. doi:10.1023/A:1010920819831.
- 300 Hanley, J.A., McNeil, B.J., 1982. The meaning and use of the area under a
301 receiver operating characteristic (ROC) curve. *Radiology* 143, 29–36.
- 302 I. Goodwin,, L. McClymont,, D. Lanyon, A. Zerihun, J. Hornbuckle, M.
303 Gibberd, D. Mowat, D. Smith, M. Barnes, R. Correll, 2009. Managing soil
304 and water to target quality and reduce environmental impact.
- 305 Kasimati, A., Espejo-García, B., Darra, N., Fountas, S., 2022. Predicting
306 Grape Sugar Content under Quality Attributes Using Normalized Differ-
307 ence Vegetation Index Data and Automated Machine Learning. *Sensors*
308 22. doi:10.3390/s22093249.
- 309 Keith Jones, 2002. Australian Wine Industry Environment Strategy.
- 310 Knight, H., Megicks, P., Agarwal, S., Leenders, M., 2019. Firm resources and
311 the development of environmental sustainability among small and medium-

312 sized enterprises: Evidence from the Australian wine industry. *Business*
313 *Strategy and the Environment* 28, 25–39. doi:10.1002/bse.2178.

314 Kuhn, M., 2008. Building Predictive Models in R Using the
315 caret Package. *Journal of Statistical Software, Articles* 28, 1–26.
316 doi:10.18637/jss.v028.i05.

317 Leilei He, Wentai Fang, Guanao Zhao, Zhenchao Wu, Longsheng Fu, Rui
318 Li, Yaqoob Majeed, Jaspreet Dhupia, 2022. Fruit yield prediction and
319 estimation in orchards: A state-of-the-art comprehensive review for both
320 direct and indirect methods 195.

321 Oliver, D., Bramley, R., Riches, D., Porter, I., Edwards, J., 2013. Review:
322 Soil physical and chemical properties as indicators of soil quality in Aus-
323 tralian viticulture. *Australian Journal of Grape and Wine Research* 19,
324 129–139. doi:10.1111/ajgw.12016.

325 R Core Team, 2021. R: A Language and Environment for Statistical Com-
326 puting. R Foundation for Statistical Computing.

327 Reynolds, A.G., 2010. Managing Wine Quality : Viticulture and Wine Qual-
328 ity. Woodhead Publishing Series in Food Science, Technology and Nutri-
329 tion ; v.1., Elsevier Science, Cambridge.

330 Rudiger, P., Stevens, J.L., Bednar, J.A., Nijholt, B., Andrew, B, C., Randel-
331 hoff, A., Mease, J., Tenner, V., maxalbert, Kaiser, M., ea42gh, Samuels, J.,
332 stonebig, LB, F., Tolmie, A., Stephan, D., Lowe, S., Bampton, J., henri-
333 queribeiro, Lustig, I., Signell, J., Bois, J., Talirz, L., Barth, L., Liquet, M.,

334 Rachum, R., Langer, Y., arabidopsis, kbowen, 2020. Holoviz/holoviews:
335 Version 1.13.3. Zenodo. doi:10.5281/zenodo.3904606.

336 SOAR, C., SADRAS, V., PETRIE, P., 2008. Climate drivers of red wine
337 quality in four contrasting Australian wine regions. Australian journal of
338 grape and wine research 14, 78–90. doi:10.1111/j.1755-0238.2008.00011.x.

339 Srivastava, S., Sadistap, S., 2018. Non-destructive sensing methods for qual-
340 ity assessment of on-tree fruits: A review. Journal of Food Measurement
341 and Characterization 12, 497–526.

342 SWA, S.W.A., 2022. Sustainable Wingrowing Australia.
343 <https://sustainablewinegrowing.com.au/case-studies/>.

344 Terry Therneau, Beth Atkinson, 2022. Rpart: Recursive Partitioning and
345 Regression Trees.

346 Wine Australia, 2019. National Vintage Report 2019 .

347 Wine Australia, 2020. National Vintage Report 2020 .

348 Wine Australia, 2021. National Vintage Report 2021 .

349 Wine Australia, 2022. National Vintage Report 2022 .

350 Winemakers’ Federation of Australia, 2015. National Vintage Report 2015 .

351 Winemakers’ Federation of Australia, 2016. National Vintage Report 2016 .

352 Winemakers’ Federation of Australia, 2017. National Vintage Report 2017 .

353 Winemakers’ Federation of Australia, 2018. National Vintage Report 2018 .