

1 Highlights

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

4 Author

5 • ???

6 • ???

7 • ????

8 • ????



## 2. Methods

### 2.1. Data

The Australian wine industry is divided into 65 regions, known as a Geographical Indicator Regions (GI Region). Each GI Region is used to describe different unique localised traits of vineyards across Australia; with each having its own mixture of climatic and geophysical properties (Halliday, 2009; Oliver et al., 2013; SOAR et al., 2008). Each region is explicitly defined under the Wine Australia Corporation Act of 1980 (Attorney-General’s Department, 2010). The climatic properties of a GI Region are summarised by Sustainable Winegrowing Australia (2021), where regions of similar climates are amalgamated together into superset regions. The climatic regions were utilised to illustrate similar trends and explain differences between sets of regions. The data used in this analysis comes from Sustainable Winegrowing Australia and covers the period 2015 to 2022. The dataset contained 3342 samples across 52 GI Regions and 1072 individual vineyards.

### 2.2. XGBoosted Trees

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

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

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

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

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

65 The difference between the prediction and actual variable is a convex loss  
66 function  $l$ . To optimise  $l$ , the difference is calculated for the  $i$ -th instance  
67 at the  $t$ -th iteration. The function  $f_t$  is selected according to which value  
68 minimises (2). The model complexity is penalised by the function  $\Omega$ , this  
69 acts to smooth weights in an attempt to prevent over fitting.

70 As predictions are made using additive tree functions, XGboosted trees  
71 can be used for classification and regression. Due to the mixture of continu-  
72 ous, binary and multiclass variables in this analysis, both classification and  
73 regression trees were created. The difference between the trees created for

74 this analysis was the objective function used. XGBoosted regression trees  
75 were created for continuous variables, using the root-mean-square as the ob-  
76 jective function. Binary class variables utilised the logistic loss function as  
77 the objective. And, Multiclass variable used the soft max function. All objec-  
78 tive functions are defined within the SKlearn library (Buitinck et al., 2013),  
79 linked via an API to the XGBoost library (Chen and Guestrin, 2016).

80 Chen and Guestrin (Chen and Guestrin, 2016) further illustrate, using  
81 Taylor expansions, that for a fixed structure  $q(x)$  the optimal weight  $\omega_j^*$  for  
82 a leaf  $j$  can be derived. Furthermore, they show the loss reduction after the  
83 split is given by the function:

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

84 with the tree structure is defined using left  $I_L$  and right  $I_R$  instance sets of  
85 nodes, with  $I = I_L \cup I_R$ . Instead of enumerating all possible tree structures,  
86 a greedy algorithm iteratively adds branches to the tree minimising  $\mathcal{L}_{split}$   
87 in (3). The frequency of a variable’s occurrence within a tree is directly  
88 attributed to the minimisation of the objective function (or loss) through  
89 the minimisation of  $\mathcal{L}_{split}$ .

90 The frequency of a variable appearing as a node within the ensemble was  
91 used as a measure of importance. This measure was chosen as it connected  
92 a variable to the minimisation of its associated objective function, trans-  
93 lating the value into a simple count metric. Creating XGBoosted trees for  
94 each variable allowed the use of importance to show how strongly variables  
95 were associated with each other. The importance of predictor variables to

96 economic variables was illustrated through the use of Sankey diagrams con-  
97 structed using the Holoviews python library (Rudiger et al., 2020). Other  
98 variable’s interconnectedness was demonstrated through the use of a chord  
99 diagram also created using Holoviews.

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

111 The use of the XGBoost library incorporates regularisation techniques  
112 built into the software to mitigate over-fitting and enhance model generali-  
113 sation. The further use of cross validated grid search functions allowed for  
114 the selection of better performing hyperparameters when selecting the final  
115 model. The performance measure for model selection was root-mean-square  
116 error for continuous variables. The receiver operator characteristic’s area  
117 under the curve was used for category variables; with multiclass variables  
118 utilising the one verse the rest approach (Hanley and McNeil, 1982).

### 119 2.3. Classification Trees

120 Classification Trees were developed to discern the different practices within  
121 regions and climates, comparing these relationships to those linked to grape  
122 quality. This was done using the *rparts* and *caret* packages (Kuhn, 2008;  
123 Terry Therneau and Beth Atkinson, 2022) in the R statistical programming  
124 language (R Core Team, 2021).

125 Three classifications were undertaken for region, climate and grape quality.  
126 Climate was further classified into two subcategories of rainfall and tempera-  
127 ture, resulting in a total of 5 classification trees being created. Classification  
128 trees were validated using K-fold cross validation. Each model was validated  
129 using 10 folds, utilising a random selection of different samples ten separate  
130 times to validate each of the classification trees. A summary confusion ma-  
131 trix was then constructed to show the class bias and overall accuracy of each  
132 tree.

## 133 3. Results

### 134 3.1. Model 1 GI Regions

135 The first Model was used to classify GI regions and resulted in an accuracy  
136 of 36.48% across 52 classes. The most prominent features used to classify  
137 regions were the types of water resources available (see Figure 1). Two re-  
138 gions, the Riverland and Coonawarra, were the most accurate classes being  
139 92.74% and 96.97% respectively. These regions differ greatly in practice and  
140 geophysical properties, with the Riverland being a dry warm inland region  
141 and Coonawarra being a cooler, wet coastal region. However, they are both  
142 similar in operational scales, with vineyards being relatively large compared

143 with other regions. The differences in resources and practices between these  
144 regions are also significant, such as the Riverland utilising the river Murray  
145 as a water source. Many of the regions had significantly lower reporting rates,  
146 resulting much poorer classification performance. The regions with the most  
147 samples performed the best (see Table 1). Notably bordering regions were  
148 routinely grouped together and misclassified as the same region, for exam-  
149 ple the two closest regions to Coonawarra, Padthaway and Wrattenbulley,  
150 were misclassified as Coonawarra even though they had 147 and 137 samples  
151 respectively. The same case was found for the Murray Darling, with 143 sam-  
152 ples, it was misclassified as the Riverland. These misclassifications are likely  
153 due to the incredibly similar regional properties and close proximity these  
154 regions have with one another. Other misclassifications were most likely due  
155 to lower reporting rates with many regions being under represented.

### 156 3.2. *Climate*

157 Classifying the SWA climatic categorisation of the given regions had bet-  
158 ter performance than the GI Regions, with 41.66% being classified correctly.  
159 These categories were divided into 12 climatic classifications with 3 and 4  
160 separate subsets for rainfall and temperature respectively. The decision tree  
161 behaved similarly and over classified climates with higher response rates. The  
162 results posed an interesting similarity with grape quality classifier, being in-  
163 fluenced predominantly by water and area. The use of fungicide to separate  
164 regions that were 'Very dry' and 'Damp' can be considered as indicative  
165 of the different practices required due to climatic pressure; fungicides being  
166 more prominent in cooler regions with greater rainfall due to the higher risk  
167 of disease pressure (Reynolds, 2010). This could also potentially explain the



Table 1: Classification accuracy of the most prominent GI Regions.

	Accuracy	Predicted	Actual
<b>Adelaide Hills</b>	30.45%	95	312
<b>Barossa Valley</b>	51.00%	205	402
<b>Coonawarra</b>	96.97%	192	198
<b>Langhorne Creek</b>	22.84%	53	232
<b>Margaret River</b>	78.82%	201	255
<b>McLaren Vale</b>	52.89%	128	242
<b>Riverland</b>	92.74%	345	

168 use of contractor tractor use to discern differences in grape quality, where the  
 169 lack of contractor use to prevent disease could have led to lowered quality of  
 170 grapes.

### 171 3.2.1. *Rainfall*

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

### 176 3.2.2. *Temperature*

177 The classification of GI Regions by their temperatures (see Figure 5)  
 178 showed similarities to the other trees, with a heavy reliance on the types

179 of water resources used as dominant predictors. The use of contractors was  
180 again used to differentiate between warm and cool regions, likely being due  
181 to disease pressure. The temperature classification tree was only a minor  
182 improvement over the regional classification tree, with an accuracy of 49.26%  
183 as shown in the confusion matrix (see Table 4).

### 184 3.3. Model 3 Grape Quality

185 The classification of grape quality through its grade had an accuracy of  
186 55.72% across 5 separate grades. There was a notable issue with the classi-  
187 fication of B grade grapes when compared to A and C (see Table 2). The  
188 classification tree itself shows similarities to that of classifying regions in  
189 Model 1, with the type of water resource used being a prominent determiner.  
190 Although not surprising the number of contractor tractor passes is new de-  
191 ciding factor due disease and pests reducing the potential quality of a crop.  
192 The prevalence of contractor use is greater in regions such as the Barossa  
193 Valley and the McLaren Vale, this could be due to the difference in opera-  
194 tional scales, with larger sites being more likely to have ownership of their  
195 own equipment for weeding and spraying due to the cost benefit.

## 196 4. Discussion

197 The difference between grape quality is most notable between warm in-  
198 land regions and coastal regions such as the Riverland and Coonawarra,  
199 respectively. Grape quality is only described by a singular variable within  
200 this study, however in reality it is driven by market demand and subject to  
201 complex forces such as international market pressure, fire, pests and disease

(Wine Australia, 2019, 2020, 2021, 2022; Winemakers' Federation of Australia, 2015, 2016, 2017, 2018) The decision trees were able to offer some insights into the factors that influence grape quality and regional contrasts that contribute to different qualities. The most prominent being what readily available resources of each region were, particular the types of water available. Heavy water consumption is often linked to the mass production of grapes, where lower quality grapes are targeted in a quantity over quality strategy. These types of business decisions are unfortunately obfuscated by lack of in-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

227 paying more or less for the same grapes given the opportunity to purchase  
228 them. It is difficult to untangle the contributing factors to the success of  
229 winegrowers and the quality of grapes produced without further specifics of  
230 choices made through out a season (Leilei He et al., 2022).

## 231 5. Conclusion

232 The type and availability of water resources were a major contributing  
233 factor when classifying grape quality and region. This was seen in the two  
234 most accurately classified regions, Coonawarra and the Riverland, with the  
235 Riverland predominantly utilising river water. Furthermore, the study high-  
236 lighted the influence of water use, fungicide application, and contractor use in  
237 differentiating grape quality, climate and region respectively. These models  
238 provide insight into the complex dynamics between regional characteristics,  
239 sustainable practices, and grape quality in the Australian winegrowing indus-  
240 try. It is important to acknowledge that grape quality is subject to external  
241 influences such as market demands and prior established business arrange-  
242 ments. Further in-depth data and understanding are necessary to fully grasp  
243 the nuances of decision-making and the interplay of factors impacting grape  
244 quality.

## 245 References

246 Abad, J., Hermoso de Mendoza, I., Marín, D., Orcaray, L., Santeste-  
247 ban, L.G., 2021. Cover crops in viticulture. A systematic review (1):  
248 <br>Implications on soil characteristics and biodiversity in vineyard.  
249 OENO One 55, 295–312. doi:10.20870/oeno-one.2021.55.1.3599.

- 250 Abbal, P., Sablayrolles, J.M., Matzner-Lober, É., Boursiquot, J.M., Baudrit,  
251 C., Carbonneau, A., 2016. Decision Support System for Vine Growers  
252 Based on a Bayesian Network. *Journal of agricultural, biological, and*  
253 *environmental statistics* 21, 131–151. doi:10.1007/s13253-015-0233-2.
- 254 Agosta, E., Canziani, P., Cavagnaro, M., 2012. Regional climate variability  
255 impacts on the annual grape yield in Mendoza, Argentina. *Journal of*  
256 *Applied Meteorology and Climatology* 51, 993–1009.
- 257 Attorney-General’s Department, 2010. *Wine Australia Corporation Act*  
258 1980.
- 259 Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel,  
260 O., Niculae, V., Prettenhofer, P., Gramfort, A., Grobler, J., Layton, R.,  
261 VanderPlas, J., Joly, A., Holt, B., Varoquaux, G., 2013. API design for  
262 machine learning software: Experiences from the scikit-learn project, in:  
263 *ECML PKDD Workshop: Languages for Data Mining and Machine Learn-*  
264 *ing*, pp. 108–122.
- 265 Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System,  
266 in: *Proceedings of the 22nd ACM SIGKDD International Conference on*  
267 *Knowledge Discovery and Data Mining*, ACM, New York, NY, USA. pp.  
268 785–794. doi:10.1145/2939672.2939785.
- 269 Cortez, P., Teixeira, J., Cerdeira, A., Almeida, F., Matos, T., Reis, J., 2009.  
270 Using data mining for wine quality assessment, in: *Discovery Science: 12th*  
271 *International Conference, DS 2009, Porto, Portugal, October 3-5, 2009* 12,  
272 Springer. pp. 66–79.

- 273 Fraga, H., Costa, R., Santos, J.A., 2017. Multivariate clustering of viticul-  
274 tural terroirs in the Douro winemaking region. *Ciência Téc. Vitiv.* 32,  
275 142–153.
- 276 G. van Rossum, 1995. Python tutorial, Technical Report CS-R9526. Centrum  
277 voor Wiskunde en Informatica (CWI),.
- 278 Hall, A., Lamb, D.W., Holzapfel, B.P., Louis, J.P., 2011. Within-season  
279 temporal variation in correlations between vineyard canopy and winegrape  
280 composition and yield. *Precision Agriculture* 12, 103–117.
- 281 Halliday, J.C.J.C., 2009. *Australian Wine Encyclopedia*. Hardie Grant  
282 Books, VIC.
- 283 Hanley, J.A., McNeil, B.J., 1982. The meaning and use of the area under a  
284 receiver operating characteristic (ROC) curve. *Radiology* 143, 29–36.
- 285 I. Goodwin,, L. McClymont,, D. Lanyon, A. Zerihun, J. Hornbuckle, M.  
286 Gibberd, D. Mowat, D. Smith, M. Barnes, R. Correll, 2009. Managing soil  
287 and water to target quality and reduce environmental impact.
- 288 Kasimati, A., Espejo-García, B., Darra, N., Fountas, S., 2022. Predicting  
289 Grape Sugar Content under Quality Attributes Using Normalized Differ-  
290 ence Vegetation Index Data and Automated Machine Learning. *Sensors*  
291 22. doi:10.3390/s22093249.
- 292 Keith Jones, 2002. *Australian Wine Industry Environment Strategy*.
- 293 Knight, H., Megicks, P., Agarwal, S., Leenders, M., 2019. Firm resources and  
294 the development of environmental sustainability among small and medium-

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

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

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

304 Oliver, D., Bramley, R., Riches, D., Porter, I., Edwards, J., 2013. Review:  
305 Soil physical and chemical properties as indicators of soil quality in Aus-  
306 tralian viticulture. *Australian Journal of Grape and Wine Research* 19,  
307 129–139. doi:10.1111/ajgw.12016.

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

311 Rudiger, P., Stevens, J.L., Bednar, J.A., Nijholt, B., Andrew, B, C., Randel-  
312 hoff, A., Mease, J., Tenner, V., maxalbert, Kaiser, M., ea42gh, Samuels, J.,  
313 stonebig, LB, F., Tolmie, A., Stephan, D., Lowe, S., Bampton, J., henri-  
314 queribei, Lustig, I., Signell, J., Bois, J., Talirz, L., Barth, L., Lique, M.,  
315 Rachum, R., Langer, Y., arabidopsis, kbowen, 2020. Holoviz/holoviews:  
316 Version 1.13.3. Zenodo. doi:10.5281/zenodo.3904606.

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

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

323 Terry Therneau, Beth Atkinson, 2022. Rpart: Recursive Partitioning and  
324 Regression Trees.

325 Wine Australia, 2019. National Vintage Report 2019 .

326 Wine Australia, 2020. National Vintage Report 2020 .

327 Wine Australia, 2021. National Vintage Report 2021 .

328 Wine Australia, 2022. National Vintage Report 2022 .

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

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

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

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