

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

44 2.2. XGBoosted Trees

45 XGBoosted (eXtreme Gradient Boosting) trees were created using the
46 XGBoost library (Chen and Guestrin, 2016) in the Python Programming
47 language (G. van Rossum, 1995). They were chosen for this analysis as they
48 provide a both high predictive performance and ability to effectively capture
49 complex relationships. A separate XGBoosted tree was used to predict each
50 variable. As variables were both continuous, binary and multiclass, separate
51 functions were created to handle the three types of variables.

52 Following Chen and Guestrin (Chen and Guestrin, 2016), XGboosted
53 trees use a given set of data, to predict y_i from the input x_i . The method
54 of prediction is achieved through a tree ensemble model, using K additive
55 functions to predict the output.

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

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

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

64 The difference between the prediction and actual variable is a convex loss
65 function l . To optimise l , the difference is calculated for the i -th instance
66 at the t -th iteration. The function f_t is selected according to which value
67 minimises 1. Chen and Guestrin (Chen and Guestrin, 2016) illustrate, using
68 Taylor expansions and

69 how, for a fixed structure $q(x)$ the the optimal weight ω_j^* for a leaf j
70 can be derived. Furthermore they show how to successfully enumerate tree
71 structures using a the derived loss function:

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

72 The model complexity is penalised by the function Ω , this acts to smooth
73 weights in an attempt to prevent over fitting.

74 As the method uses additive tree functions to predict y_i , it can be used
75 for both classification and regression.

76 This means we greedily add the ft that most improves our model according
77 to Eq. (2). Second-order approximation can be used to quickly optimize the
78 objective in the general setting [12].

79 A common example is a linear model, where the prediction is given as
80 $y_i = \sum_j \theta_j x_{ij}$, a linear combination of weighted input features. The prediction
81 value can have different interpretations, depending on the task, i.e., regres-
82 sion or classification. For example, it can be logistic transformed to get the
83 probability of positive class in logistic regression, and it can also be used as
84 a ranking score when we want to rank the outputs.

85 The parameters are the undetermined part that we need to learn from
86 data. In linear regression problems, the parameters are the coefficients θ .
87 Usually we will use θ to denote the parameters (there are many parameters
88 in a model, our definition here is sloppy).

89 With judicious choices for y_i , we may express a variety of tasks, such
90 as regression, classification, and ranking. The task of training the model
91 amounts to finding the best parameters that best fit the training data and
92 labels

93 . In order to train the model, we need to define the objective function to
94 measure how well the model fit the training data.

95 A salient characteristic of objective functions is that they consist of two
96 parts: training loss and regularization term: $obj(\theta) = L(\theta) + \Omega(\theta)$

97 where L is the training loss function, and Ω is the regularization term.
98 The training loss measures how predictive our model is with respect to the
99 training data. A common choice of L is the mean squared error, which is
100 given by

$$101 \quad L(\theta) = \sum_i (y_i - \hat{y}_i)^2$$

102 Another commonly used loss function is logistic loss, to be used for logistic
103 regression:

104 The regularization term is what people usually forget to add. The regu-
105 larization term controls the complexity of the model, which helps us to avoid
106 overfitting. This sounds a bit abstract, so let us consider the following prob-
107 lem in the following picture. You are asked to fit visually a step function
108 given the input data points on the upper left corner of the image. Which
109 solution among the three do you think is the best fit?

110 XGBoosted Regression trees were used to predict continuous variables.
111 With data being split into 80% training data and 20% testing data.

112 XGBoosted classification trees were used to classify the binary and mul-
113 ticlass variables. Data was split into 80% training, 10% testing and 10%
114 validation data.

115 The modelled relationships are able to be scrutinised by using techniques
116 such as feature importance analysis. The use of the XGBoost library also in-
117 corporates regularisation techniques built into the software to mitigate over-
118 fitting and enhance model generalisation. The further use of cross validated
119 grid search functions allowed for the selection of better performing hyper-

120 parameters when selecting the final model.

121 *2.3. Classification Trees*

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

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

135 **3. Results**

136 *3.1. Model 1 GI Regions*

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

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

158 3.2. *Climate*

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

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

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

more prominent in cooler regions with greater rainfall due to the higher risk of disease pressure (Reynolds, 2010). This could also potentially explain the use of contractor tractor use to discern differences in grape quality, where the lack of contractor use to prevent disease could have led to lowered quality of grapes.

3.2.1. Rainfall

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

178 3.2.2. *Temperature*

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

186 3.3. *Model 3 Grape Quality*

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

198 4. Discussion

199 The difference between grape quality is most notable between warm in-
200 land regions and coastal regions such as the Riverland and Coonawarra,

201 respectively. Grape quality is only described by a singular variable within
 202 this study, however in reality it is driven by market demand and subject to
 203 complex forces such as international market pressure, fire, pests and disease
 204 (Wine Australia, 2019, 2020, 2021, 2022; Winemakers' Federation of Aus-
 205 tralia, 2015, 2016, 2017, 2018) The decision trees were able to offer some
 206 insights into the factors that influence grape quality and regional contrasts
 207 that contribute to different qualities. The most prominent being what readily
 208 available resources of each region were, particular the types of water available.
 209 Heavy water consumption is often linked to the mass production of grapes,
 210 where lower quality grapes are targeted in a quantity over quality strategy.
 211 These types of business decisions are unfortunately obfuscated by lack of in-
 212 depth data regarding vineyard business plans. Notably the literature shows
 213 that there are many complex decisions to be made on the ground depending
 214 on many compounding factors that influence both quality and yield (Abad
 215 et al., 2021; Cortez et al., 2009; Hall et al., 2011; I. Goodwin, et al., 2009;
 216 Kasimati et al., 2022; Oliver et al., 2013; Srivastava and Sadistap, 2018)
 217 . There are also further differences when comparing winegrowers to other
 218 agricultural industries as they are vertically integrated within the wine in-
 219 dustry, tying them to secondary and tertiary industries, such as wine pro-
 220 duction, packaging, transport and sales. This results in unique issues, where
 221 on-the-ground choices are influenced by other wine industry's decisions, such
 222 as the use of sustainable practices in vineyards to sell in overseas markets;
 223 notably these interactions are further complicated by some winegrowers be-
 224 ing totally integrated into wine companies, while others are not (Knight et
 225 al., 2019). It is incredibly difficult to attribute external business decisions to

226 produced grape quality but it is important to acknowledge that some growers
227 are contracted to produce grapes of a particular grade; it is difficult to know
228 whether another consumer may have graded the grape quality differently
229 paying more or less for the same grapes given the opportunity to purchase
230 them. It is difficult to untangle the contributing factors to the success of
231 winegrowers and the quality of grapes produced without further specifics of
232 choices made through out a season (Leilei He et al., 2022).

233 **5. Conclusion**

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

247 References

- 248 Abad, J., Hermoso de Mendoza, I., Marín, D., Orcaray, L., Santeste-
249 ban, L.G., 2021. Cover crops in viticulture. A systematic review (1):
250
Implications on soil characteristics and biodiversity in vineyard.
251 OENO One 55, 295–312. doi:10.20870/oeno-one.2021.55.1.3599.
- 252 Abbal, P., Sablayrolles, J.M., Matzner-Lober, É., Boursiquot, J.M., Baudrit,
253 C., Carbonneau, A., 2016. Decision Support System for Vine Growers
254 Based on a Bayesian Network. Journal of agricultural, biological, and
255 environmental statistics 21, 131–151. doi:10.1007/s13253-015-0233-2.
- 256 Agosta, E., Canziani, P., Cavagnaro, M., 2012. Regional climate variability
257 impacts on the annual grape yield in Mendoza, Argentina. Journal of
258 Applied Meteorology and Climatology 51, 993–1009.
- 259 Attorney-General’s Department, 2010. Wine Australia Corporation Act
260 1980.
- 261 Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System,
262 in: Proceedings of the 22nd ACM SIGKDD International Conference on
263 Knowledge Discovery and Data Mining, ACM, New York, NY, USA. pp.
264 785–794. doi:10.1145/2939672.2939785.
- 265 Cortez, P., Teixeira, J., Cerdeira, A., Almeida, F., Matos, T., Reis, J., 2009.
266 Using data mining for wine quality assessment, in: Discovery Science: 12th
267 International Conference, DS 2009, Porto, Portugal, October 3-5, 2009 12,
268 Springer. pp. 66–79.

- 269 Fraga, H., Costa, R., Santos, J.A., 2017. Multivariate clustering of viticul-
270 tural terroirs in the Douro winemaking region. *Ciência Téc. Vitiv.* 32,
271 142–153.
- 272 G. van Rossum, 1995. Python tutorial, Technical Report CS-R9526. Centrum
273 voor Wiskunde en Informatica (CWI),.
- 274 Hall, A., Lamb, D.W., Holzapfel, B.P., Louis, J.P., 2011. Within-season
275 temporal variation in correlations between vineyard canopy and winegrape
276 composition and yield. *Precision Agriculture* 12, 103–117.
- 277 Halliday, J.C.J.C., 2009. *Australian Wine Encyclopedia*. Hardie Grant
278 Books, VIC.
- 279 I. Goodwin,, L. McClymont,, D. Lanyon, A. Zerihun, J. Hornbuckle, M.
280 Gibberd, D. Mowat, D. Smith, M. Barnes, R. Correll, 2009. Managing soil
281 and water to target quality and reduce environmental impact.
- 282 Kasimati, A., Espejo-García, B., Darra, N., Fountas, S., 2022. Predicting
283 Grape Sugar Content under Quality Attributes Using Normalized Differ-
284 ence Vegetation Index Data and Automated Machine Learning. *Sensors*
285 22. doi:10.3390/s22093249.
- 286 Keith Jones, 2002. *Australian Wine Industry Environment Strategy*.
- 287 Knight, H., Megicks, P., Agarwal, S., Leenders, M., 2019. Firm resources and
288 the development of environmental sustainability among small and medium-
289 sized enterprises: Evidence from the Australian wine industry. *Business*
290 *Strategy and the Environment* 28, 25–39. doi:10.1002/bse.2178.

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

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

298 Oliver, D., Bramley, R., Riches, D., Porter, I., Edwards, J., 2013. Review:
 299 Soil physical and chemical properties as indicators of soil quality in Aus-
 300 tralian viticulture. Australian Journal of Grape and Wine Research 19,
 301 129–139. doi:10.1111/ajgw.12016.

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

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

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

311 Terry Therneau, Beth Atkinson, 2022. Rpart: Recursive Partitioning and
 312 Regression Trees.

- 313 Wine Australia, 2019. National Vintage Report 2019 .
- 314 Wine Australia, 2020. National Vintage Report 2020 .
- 315 Wine Australia, 2021. National Vintage Report 2021 .
- 316 Wine Australia, 2022. National Vintage Report 2022 .
- 317 Winemakers' Federation of Australia, 2015. National Vintage Report 2015 .
- 318 Winemakers' Federation of Australia, 2016. National Vintage Report 2016 .
- 319 Winemakers' Federation of Australia, 2017. National Vintage Report 2017 .
- 320 Winemakers' Federation of Australia, 2018. National Vintage Report 2018 .