

¹ Highlights

² **An analysis of interrelations between economic and environmental**
³ **variables in Australian Winegrowing.**

⁴ Author

⁵ • Highlight 1

⁶ • Highlight 2

⁷ • Highlight 3

⁸ • Highlight 4

9 An analysis of interrelations between economic and
10 environmental variables in Australian Winegrowing.

11 Author^{1,1,1}

12 1. Abstract

13 Region is a driving force in determining profits. Fuel and water are the
14 primary driving operational costs

15 The finding themselves are not that significant, it is the evidence of a
16 complex system governing these relationships that is the significant finding.

17 They were gathered by Sustainable Winegrowing Australia

18 What we add is that agricultural economics maybe complex systems that
19 are obfuscated by the expertise of growers to navigate problems such as
20 drought, pests and disease. These systems, although predictable through
21 some variables are not causal. We see this in the evidence of some years
22 not affecting outcomes even though they are known to contain the aforemen-
23 tioned issues.

24 2. Introduction

25 Historically strong demands for Australian wine have helped to create a
26 thriving industry, however recent pressures brought on by a loss of tourism
27 and labour due to the COVID-19 pandemic, the global freight crisis, war
28 in Europe, tariffs and rising inflation has negatively effected the industry's

29 outlook (Wine Australia, 2021; Australia, 2021a). The 2021-2022 financial
 30 year alone saw a decline of 19% in exports solely due to tariffs (Wine Aus-
 31 tralia, 2022). A greater understanding of the different underlying conditions
 32 leading to improved performance in agricultural productivity and sustain-
 33 ability at scale are key to introducing stronger policy and information to
 34 aid in increasing a nations agricultural sustainability (?). Specifically within
 35 the Australian Wine and vine industry there is a need to further understand
 36 the driving relationships between resource use and economic output. Where
 37 these relationships can lead to determining better and efficient methods and
 38 develop benchmarks with local growers. (?)

39 An unprecedented amount of data regarding the Australian winegrowing in-
 40 dustry has been collected through Sustainable Winegrowing Australia, offer-
 41 ing new insights into the driving economic forces of the Australian wine in-
 42 dustry. This dataset allowed insights into the economic outcome of vineyards
 43 through the incorporation of operating costs and grape revenue from grape
 44 sales within the data. We use this data to study these economic outcomes
 45 and their statistical relationships to vineyards' utilisation of the resources.
 46 Answering what the driving factors are behind vineyard economic outcomes,
 47 and linking these outcomes to predictor importance. This is done through
 48 analysing a new comprehensive nationwide data set using XGBoosted mod-
 49 els. We further compare the relationships between different resources to
 50 address the extensive collinearity found within the data (Chen and Guestrin,
 51 2016). XGBoosted models were used because they are able to overcome
 52 multicollinearity as well as highlight the level of importance that predictor
 53 variables have on response variables; with importance being able to be sta-

54 tistically defined through multiple methods.

55 **3. Methods**

56 *3.1. Data*

57 Data used in this analysis were obtained from Sustainable Winegrowing
58 Australia. Australia’s national wine industry sustainability program. The
59 program aims to facilitate grape-growers and winemakers in demonstrating
60 and improving their sustainability (SWA, 2022). Data recorded by SWA is
61 entered manually by winegrowers using a web based interface tool. A total
62 of 6091 observations were collected from 2012/2013 to 2021/2022 financial
63 years. 23 variables were used for each observation reflecting a vineyards state
64 for the given year (see Table 3.1).

65 The data originally contained only two multiclass variables: year and re-
66 gion. Variables that measured the same metric from different sources (such as
67 water collected from rivers versus water from dams) were converted into mul-
68 ticlass variables representing the source through one-hot-encoding. Changing
69 each variable class into a binary value, with one indicating the presence of
70 the class and zero indicating its absence. Occurrences of multiple sources
71 were defined as their own separate classes. Where a class variable had a
72 recorded amount the total amount used from these variables was retained
73 as a separate variable; for example water used (in Mega Litres) was also
74 included alongside water source.

75 The variable region represented one of the 65 Geographical Indicator Re-
76 gions (GI Region) used to describe different unique localised traits of vine-
77 yards across Australia (Halliday, 2009; Oliver et al., 2013; SOAR et al., 2008).

Table 1: Summary of variables used in the analysis. The recorded column indicate the number of 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
Grape Revenue	AUD	⁴ 853	
Operating Costs	AUD	853	

Each region is explicitly defined under the Wine Australia Corporation Act of 1980 (Attorney-General’s Department, 2010).

3.2. 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). XGBoosted trees are a boosted tree ensemble method that can be used to classify classes, or predict continuous response variables. They were chosen for this analysis as the data contained a mixture of class and continuous variables. And, XGBoosted trees are unaffected by multicollinearity, as well as offer high predictive performance for a wide variety of purposes (Chen and Guestrin, 2016). An XGBoosted tree was created for each variable to show how they interacted. Each tree included all but the economic variables (operating cost and revenue from grape sales), which were only included within their own trees as response variables. Separately profit (the difference between revenue and operational costs) was looked at in prior analyses (see appendix) but the results were not included due to low average loss values and model stability. This meant that every variable would have a measure of its importance to other variables (see Section 3.4), which was used to show the highly interrelated nature of variables within vineyards. The complicated interaction between variables was illustrated using Sankey and Chord diagrams; with variable importance measures being used to show the strength of connection between any two variables (see section 3.4).

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

dict the output. Each of f_k functions is a classification or regression tree, such that all functions are in the set of all decision trees, given by \mathcal{F} , is defined by $f(x) = \omega_{q(x)}(q : \mathbb{R}^m \rightarrow T, \omega \in \mathbb{R}^T)$. Where each function corresponds to an independent tree structure q of ω weights. Each tree has T leaves, which contain a continuous score, represented by ω_i for the i -th leaf. The final prediction is determined by the sum of the score of the corresponding leaves, given by:

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

The set of functions, \mathcal{F} , used by the tree is determined by minimising a regularised objective function, \mathcal{L} given by:

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

, where

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \|\omega\|^2 \quad (3)$$

As predictions are made using additive tree functions, XGboosted trees can be used for classification or regression. The difference between a prediction, $\phi(x_i)$, and actual variable, $f_k(x_i)$, is a differentiable convex loss function l . These properties of l allow the function to be versatile in which objective we choose to optimise for, which is also important in being able to process both continuous and categorical variables. To optimise l , the difference is calculated for the i -th instance at the t -th iteration.

3.3. Loss functions

The functions included as parameters in equation 2 mean that traditional optimisation methods for Euclidean space cannot be used. Chen and Guestrin

(Chen and Guestrin, 2016) illustrate, using Taylor expansions, that for a fixed structure $q(x)$ the optimal weight ω_j^* for a leaf j can be derived. Importantly a loss function can be used to fit a model iteratively to data. For this analysis several loss functions were used, as variables took the form of continuous, binary and multi-call data. The loss function for making a split within the tree structure is given by:

$$\mathcal{L}_{split} = \frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] - \gamma. \quad (4)$$

The tree structure being defined using left I_L and right I_R instance sets of nodes, with $I = I_L \cup I_R$. Instead of enumerating all possible tree structures, a greedy algorithm iteratively adds branches to the tree minimising \mathcal{L}_{split} in (4). The frequency of a variable’s occurrence within a tree is directly attributed to the minimisation of the loss function through the minimisation of \mathcal{L}_{split} .

The loss functions used for this analysis were the root-mean-square function for continuous variables. The logistic loss function for binary class variables. And, the soft max function for Multiclass variables. All objective functions are defined within the SKlearn library (Buitinck et al., 2013), which was utilised via an API to the XGBoost library (Chen and Guestrin, 2016).

3.4. Variable Importance

Due to XGBoost creating a large amount of decision trees, the interpretability of these models is obfuscated by the intricate relationships within complicated ensembles. A measure of variable importance was the technique

144 used to highlight a variables influence within the ensemble. Variable impor-
145 tance can be measured in multiple ways; we used the frequency of a variable
146 appearing as a node within the ensemble as a measure of its importance.
147 This measure was chosen as it connected a variable to the minimisation of
148 its associated objective function. The measure of a variable’s importance
149 within this study can then be interpreted as how often a variable was the
150 optimal choice in reducing the loss function of the ensemble. Importantly,
151 multiclass variables being one-hot-encoded are given an importance score
152 for each individual class; for example, each specific region will have its own
153 importance score.

154 Creating XGBoosted trees for each variable allowed the use of importance
155 to show how strongly variables were associated with each other. The impor-
156 tance of variables to one another was illustrated through the use of Sankey
157 and Chord diagrams. These diagrams were constructed using the Holoviews
158 python library (Rudiger et al., 2020). Both Chord and Sankey diagrams
159 illustrated variable importance through the size of the bands between two
160 variables. The number at the end of a connection in a Sankey diagram indi-
161 cated a variable’s importance, or the number of times it appeared within the
162 ensemble. Sankey and Chord diagrams were presented together; with Sankey
163 diagrams showing the connection of a variable to its 10 most important pre-
164 dictor variables. Chord diagrams were used alongside a Sankey diagram to
165 show the interconnectedness of the ten most prominent variables within its
166 associated Sankey diagram. Chord diagrams formed circles, with variables
167 being connected through their relative importance. The importance values
168 for the Chord diagrams were taken from the models of those individual

169 variables, with the diagram being simplified to just the ten variables in the
170 associated Sankey diagram, for readability’s sake.

171 3.5. Validation

172 As there were multiple different loss functions, multiple different forms
173 of validation were used. In each case the data was split into training data,
174 which constituted 80% of the original data. The remaining 20% was used
175 in testing and validation. Data was stratified when splitting the data into
176 these subsets to conserve the same proportion of class occurrences between
177 training, testing and validation data. For continuous variables 20% was used
178 as testing data, minimising the root-mean-square function. The final model
179 was validated using repeated k-fold cross validation for 10 folds, repeated 10
180 times. R^2 scores were used to determine the best regression models during
181 validation. For binary and multiclass variables, data was split into 80%
182 training, 10% testing and 10% validation data. For class variables, validation
183 was summarised through the accuracy, the proportion of true negatives and
184 positives.

185 The XGBoost library incorporates regularisation techniques built into
186 the software to mitigate over-fitting and enhance model generalisation. This
187 allowed us to utilise cross validated grid search functions when selecting for
188 better performing hyperparameters. The performance measure for model
189 selection was root-mean-square error for continuous variables. The receiver
190 operator characteristic’s area under the curve was used for category variables
191 (Hanley and McNeil, 1982). Multiclass variables utilised the one verse one
192 approach to minimise sensitivity to class disparity (Ferri et al., 2009; Hand
193 and Till, 2001).

194 3.6. Surrogate Models

195 The creation of more interpretable models such as linear regression in
196 parallel to AI systems has been used to explain variable's relationships within
197 black box algorithms (Molnar, 2022). As XGBoost create an ensemble of
198 decision trees, here we use classification and regression trees to gain insight
199 into intricacies of the ensembles derived through XGBoost. Decision Trees
200 were created for operating costs, revenue and region. These models describe
201 the partitions that are useful in predicting these variables; giving insight into
202 the trees that make up the ensembles created by XGBoost. These trees were
203 created using the rparts and caret packages (Kuhn, 2008; Terry Therneau
204 and Beth Atkinson, 2022) in the R statistical programming language (R
205 Core Team, 2021).

206 Decision trees were validated using K-fold cross validation. Each model
207 was validated using 10 folds, utilising a random selection of different samples
208 ten separate times to validate each of the decision trees. The same measure
209 of accuracy as the XGBoosted trees was used for comparison.

210 4. Results

211 4.1. Revenue

212 We investigated the link between revenue to other variables in the SWA
213 data by predicting it, and then linking each variable to revenue through vari-
214 able importance. The prediction of revenue performed similarly to operating
215 cost achieving an R^2 of 0.7716 (with a standard deviation of 0.1525). A
216 regression tree was used as a surrogate model to present an example of the
217 typical type of decision tree present within the XGBoost Ensemble, however

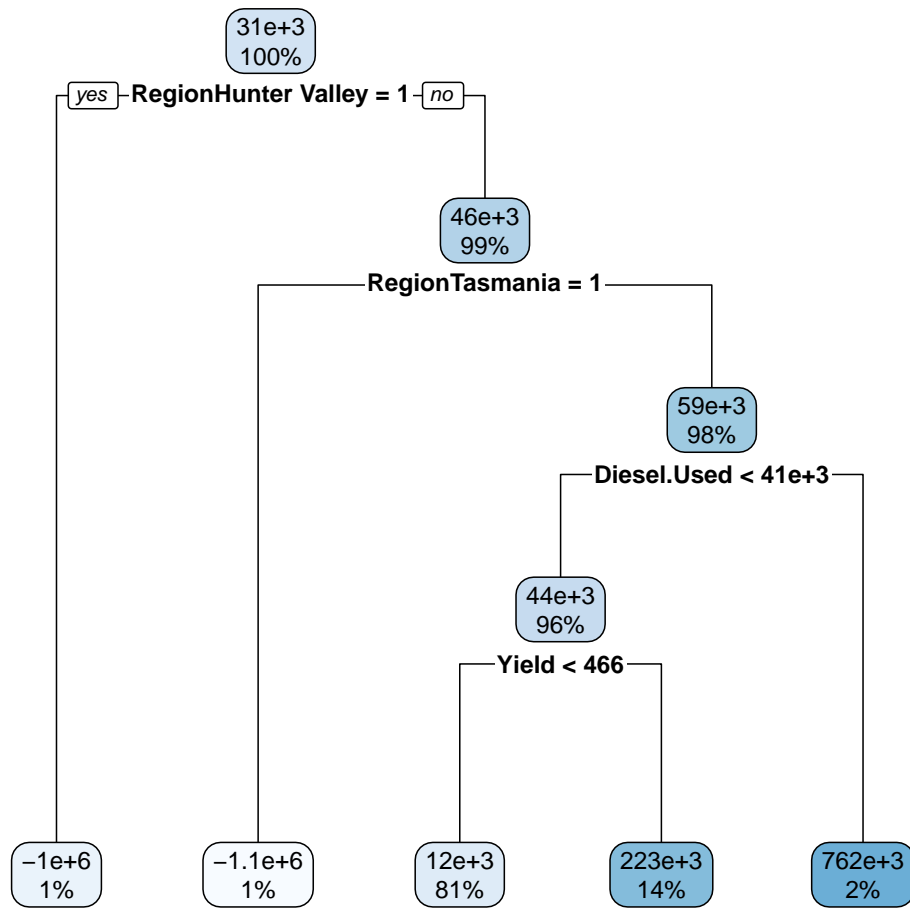


Figure 1: Decision tree predicting Profit. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

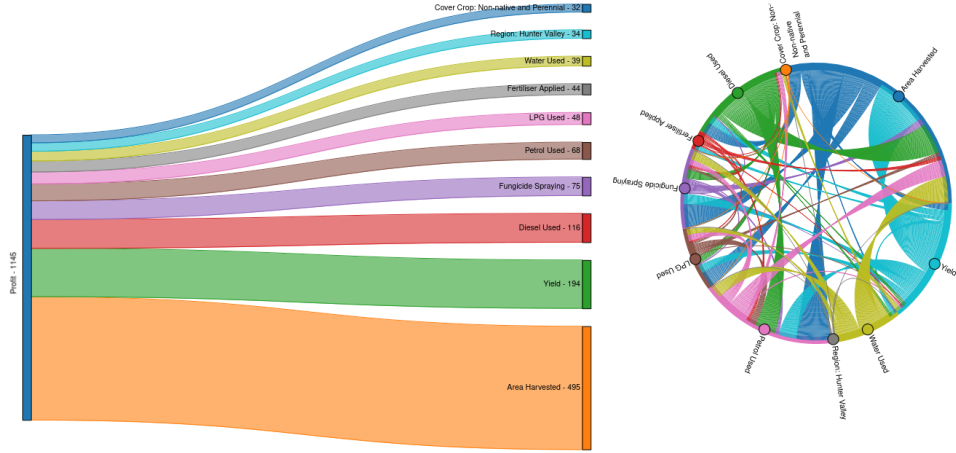


Figure 2: The left-hand side depicts the 10 most important variables in predicting Profit using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

the surrogate model only achieved an R^2 of 0.0961 (with a standard deviation of 0.0181) and the XGBoosted ensemble.

The important variables when attempting to determine revenue were size, yield, fuel and water (see 2). Due to regions being recorded separately for importance none appeared as the most important variables, overall regions contributed to 234 nodes in the ensemble making them collectively the third most important variable. Although performing poorly, the surrogate model highlights the importance of size in determining profits. Area also appearing as a variable of higher importance is show to be highly interrelated with other variables. The relation to area is likely to primarily be the effect of economies of scale, shown through its strong relations to other variables in figure 2. Area harvested is likely also an indicator of other variables such as

230 slashing passes its strongest connection presented.

231 4.2. Operating Costs

232 The relationships to operating cost through variable importance were
233 found to be similar to that of revenue. With fuel, water, area and yield
234 occurring the most (see figure 4). A surprising difference is that the most
235 important operational consideration for operating cost is the use of fungicide,
236 compared to revenue where slashing is the most important (comparing Figure
237 6 and 4). The variables that feed into these decisions are also very differing
238 with diesel being the most informative to slashing and area being the most
239 informative to the need for fungicide.

240 Again, region played a determining factor overall but not as much indi-
241 vidually with region contributing to 334 nodes within the ensemble making
242 it the most important variable when considering all regions together. It
243 was surprising that electricity, slashing and spraying passes were not more
244 prominent in operating costs due to the intrinsic nature as an agricultural
245 expense.

246 Comparatively to revenue, operating cost performed better The XG-
247 Boosted regression ensemble achieved an R^2 of 0.8025 (with a standard devi-
248 ation of 0.1033). Again the surrogate model did not perform well achieving
249 an R^2 of 0.0931 (with a standard deviation of 0.0197) but showed similarly
250 to revenue an importance placed on fungicide spraying and size (see figure
251 3).

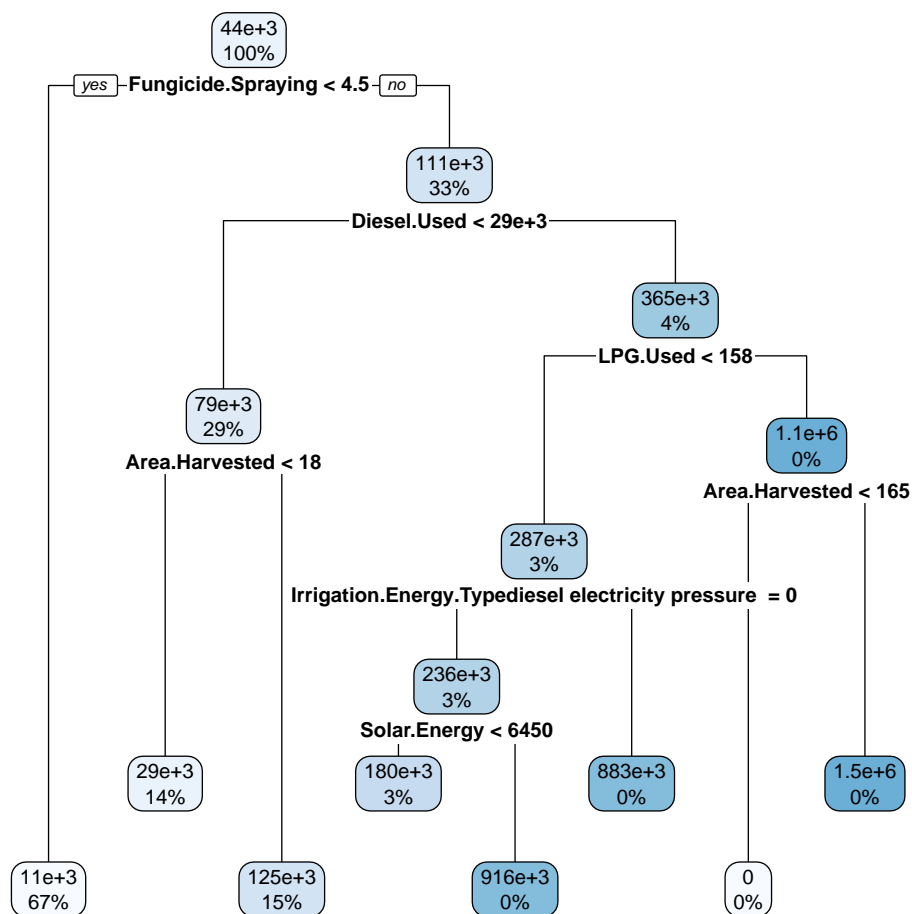


Figure 3: A surrogate model decision tree predicting operating costs. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

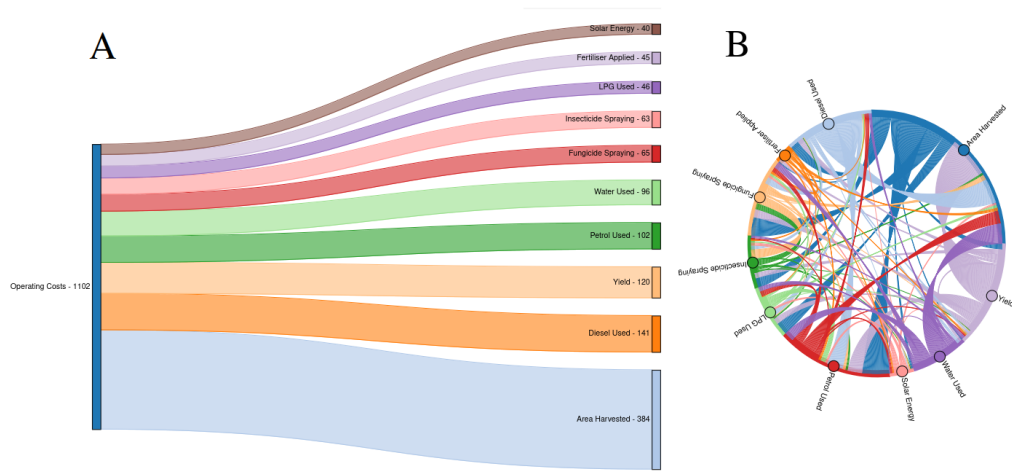


Figure 4: The left-hand side, A, depicts the 10 most important variables in predicting Operating Costs using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The number at the end of each band in the diagram is that variable's importance. The right-hand side, B, depicts the importance of the 10 variables in Sankey diagram relative to one another.

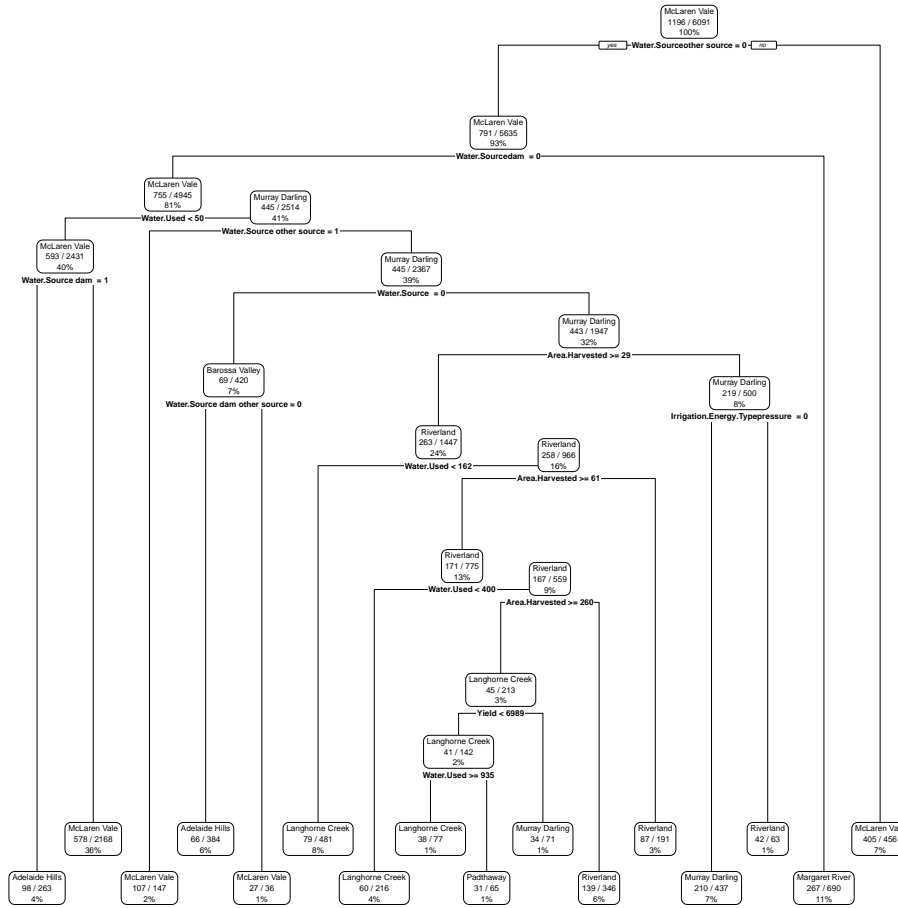


Figure 5: Decision tree predicting Region. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

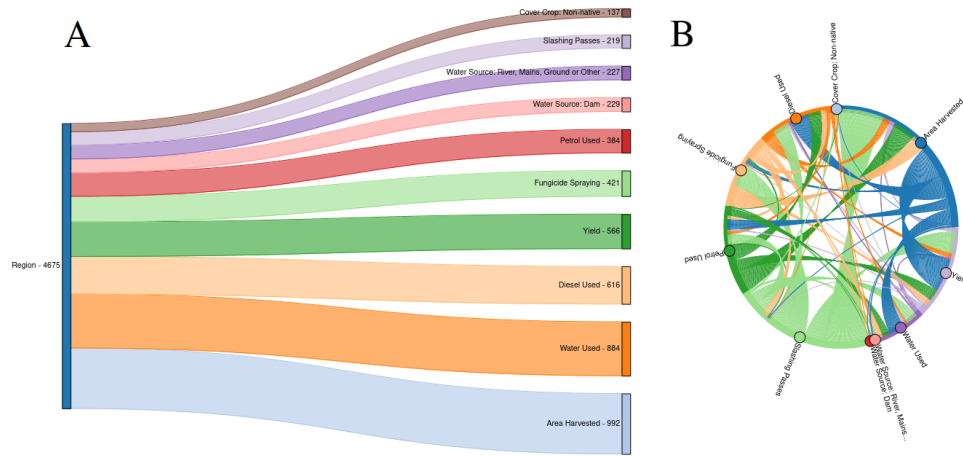


Figure 6: The left-hand side, A, depicts the 10 most important variables in predicting Region using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The number at the end of each band in the diagram is that variable's importance. The right-hand side, B, depicts the importance of the 10 variables in Sankey diagram relative to one another.

252 4.3. Region

253 When considered overall, Region was a highly informative variable through
254 measure of importance to both operating cost and revenue. Notably the
255 overall importance of region to revenue was 234 (making it the third most
256 important variable when considering all regions together). The Barossa Val-
257 ley region and Tasmania were the two most important regions in relation to
258 revenue; these two regions are considered to be some of the highest revenue
259 per hectare regions in Australia (Wine Australia, 2022). These two regions
260 are also relative opposites in winegrowing climates with the Barossa being
261 warm and dry climate focussing on Shiraz grapes and Tasmania being a cool
262 wet climate that grows Pinot.

263 When considering all regions together, it had the most node contributions
264 to determining operating costs with an importance of 334. Of all the regions,
265 again Tasmania was the most important, followed by the Adelaide Hills. In
266 contrast to revenue, both climates are considered cool and wet, and warmer
267 drier regions such as the Barossa and Hunter Valley only contributed roughly
268 half the same number of nodes to the ensemble. When looking at 6 the
269 inclusion of slashing and fungicide spraying is the likely reason; with fungal
270 and weed pressure being greater in cooler wetter regions.

271 Both diesel and petrol were of more relative importance in operating costs
272 than water was compared with region. similar to those used to classify region
273 (see Figure 6, except water used. The surrogate model relied heavily on the
274 use of water source to classify regions, which is reflective of regional access to
275 resources being a deciding factor in vineyard management (see Figure 5). A
276 major difference between region and profit was the importance given to water

277 use, with water use being a relatively more important variable in predicting
278 region than profit (considering its rank in importance to other variables).

279 The surrogate model for region performed better than other surrogate
280 models with 32.34% (3.67% standard deviation). The prominence of types
281 and use of water resources was in classifying region is reflective of difference
282 of availability of water resources is when comparing different regions (see
283 Figure 5). The XGBoost ensemble, did not perform as well as operating
284 costs or revenue with 56.82% accuracy (50.58% validation accuracy). The
285 difference in accuracy is in part due to the large number of classes (being
286 58). The ensemble did not differ greatly from the surrogate model, with a
287 continuing emphasis on Area, water, fuel and yield as determining factors
288 (see Figure (6)).

289 Many of the regions had significantly lower reporting rates, resulting in
290 much poorer classification performance. The regions with the most samples
291 performed the best. Bordering regions were routinely grouped together and
292 misclassified as the same region. Two areas that suffered the most from
293 this, specifically with the classification tree were the Limestone Coast (cool
294 coastal areas in South Australia) and the warmer inland regions along the
295 Murray Darling. The classification tree likely had more difficulty discerning
296 vineyards closer to the river using only water sources due to the greater access
297 to river water in these areas.

298 **5. Discussion**

299 The explored relationships between vineyard resource use, operations and
300 geographical properties to revenue and operating costs highlight how deci-

301 sive regional influences can be determining a vineyard’s economic outcomes.
302 Several physical parameters such as climate, geography and soil are predeter-
303 mined by a vineyard’s location; making it a widely considered key determi-
304 nant of grape yield and quality (Abbal et al., 2016; Agosta et al., 2012; Fraga
305 et al., 2017). The association between yield and region is demonstrated by
306 its rank of fourth-highest variable importance when determining region (see
307 Figure 6).

308 Warmer regions are known to be beneficial in hastening the ripening pro-
309 cess of winegrapes (WEBB et al., 2011). Warmer regions are also associated
310 with lower quality grapes, caused largely due to this hastened ripening (Bot-
311 ting et al., 1996). In general warmer regions are not associated with higher
312 yields, but if a vineyard in a warmer region is sufficiently irrigated much
313 higher yields can be achieved than in cooler regions (Camps and Ramos,
314 2012). It is likely that the combination of larger vineyards with higher water
315 use is a determining factor in classifying regions which favour larger produc-
316 tion of grapes; reflected through region using water use so prominently in the
317 XGBoost ensemble. The link to water resources in defining regions is also
318 an important consideration, as vineyards can leverage higher irrigation rates
319 given more accessible water resources. A further consideration in the link
320 between revenue and region is that grape prices are set at a regional level by
321 buyers (Wine Australia, 2022). It is also important to consider that some
322 regions carry particular fame regarding the quality of their produce such as
323 Tasmania, the Hunter Valley and Barossa Valley (Halliday, 2009). This clas-
324 sification can be contrasted with other warmer regions of higher rainfall that
325 use the warmer climate to concentrate their grapes, increasing the flavour

326 profile (and thus quality) (Goodwin I, Jerie P, 1992; MG McCarthy et al.,
327 1986).

328 In part some winegrowing strategies are restricted simply through ac-
329 cess to water resources, being reflected through the region classification tree
330 (see Figure 5). Regions are likely to have varying access to different wa-
331 ter sources, such as those along the River Murray being able to utilise river
332 water for crops, unlike most coastal regions which may be drawing from sur-
333 face or underground water sources. Similarly, the connection between region
334 and fuel use is likely an indicator of the level of infrastructure within the
335 region. Where, the need to pressurise irrigation systems from river water or
336 to generate power would require larger amounts of diesel and petrol.

337 Operational costs showed similar importance across fuel, water and trac-
338 tor use. The dominating factor of area likely played a large part in deter-
339 mining how costly a tractor pass would be, or in defining the ratio of water
340 applied to the amount of vines. The node frequency was high for area but
341 much lower in general across the other variables, which could indicate the
342 need to be more circumstantial in determining operational costs. Although
343 it was attempted to capture the complexity between how variables inter-
344 acted when determining operational costs (see Figure 4), it is likely yet more
345 complicated. An example of how interrelated operational costs can be, is
346 the optimisation of tractor passes to achieve multiple goals in a pass, being
347 shown to reduce energy use in vineyards, decreasing running costs, as well
348 as reducing soil compaction (Capello et al., 2019).

349 When determining revenue, similar variables were used to operational
350 cost; with region also being of high variable importance relative to other vari-

ables (when considering all regions together in importance). It is difficult to extrapolate the specific influence of location on a vineyard’s outcomes due to the broad and varying definition of a region. Utilising the Geographical Indicator regions defined by Wine Australia (Australia, 2021b) is a limitation in one way, as it is too broad to fully capture a vineyard’s location and how that influences variables at a more granular level. However, as buyers set prices at regional levels, it is still important to consider a vineyard’s Geographical Indicator region.

Decisions made on the ground have far-reaching effects and are difficult to completely capture. A higher number of tractor passes used as a preventative measure for occurrences such as disease, may incur higher operational costs but could be critical in preventing long term losses. With factors such as erosion and soil health being difficult to capture but also influenced by tractor use (Capello et al., 2019, 2020). Although, performing well in R^2 , the ability to predict operational costs is limited by the variables incorporated. Reductions in fuel, water and tractor use are obvious methods to reduce operational costs but not necessarily achievable decisions. Without fully capturing more granular activities such as the specifics of what fuel was used for, it is hard to determine what decisions specifically influence the operational costs. Electricity in particular is used predominantly for irrigation. Size is also a further consideration where slashing and spraying are measured in discrete tractor passes and show a surprising connection to the overall size of a vineyard, despite not being scaled to any measure of size. This would mean that, although measured as the same increment, a slashing or spraying pass in a larger vineyard would consume more fuel and wages than in a

376 smaller vineyard.

377 The reasoning for any particular decision can be widely varying. A more
378 granular definition of region may help to better discern the differences in
379 practices, and the reason for employing them. More sophisticated mod-
380 els, specifically those that utilise expert opinion, may also help to capture
381 and address the decision-making process. An example is the optimisation of
382 fungicide sprays using Bayesian models that forecast disease risk (Lu et al.,
383 2020).

384 Separately revenue and operating cost did have a greater predictability
385 than their counterpart profit. The disparity in accuracy between profit and
386 other economic outcomes is reflective of the complexity in trying to address
387 challenges such as climate change, disease and changing market demands
388 (Wine Australia, 2020, 2021, 2022). The difference between turning a profit
389 or loss is dependent on decisions made and chance. The difference between
390 vineyards that make profit and those that do not could be a multitude of fac-
391 tors including differences in farming practices not captured within this study.
392 Some decisions leading to latent effects such as large scale soil deposition in
393 extreme rain events can be caused by soil compaction due to overworking a
394 vineyard (Capello et al., 2020).

395 **6. Conclusion**

396 This study has provided valuable insights into the multifaceted dynam-
397 ics governing operational costs and revenue. The impact of different regions
398 highlighted the complex interrelatedness of variables within a vineyard. We
399 relate how factors such as water and fuel intersect to impact operational costs

400 and how different seasonal events affect these operations; as well as the signif-
401 icance of context-specific decision-making. While this investigation utilised
402 a broad regional classification, the potential benefits of adopting a more nu-
403 anced approach and incorporating expert knowledge have been highlighted.
404 Further work could pursue causal models and the creation of decision sup-
405 port systems. It is difficult to untangle the predictive and correlative nature
406 of a variable compared to the causal reasons. By delving deeper into the
407 complex interplay of variables, further advancements can be made in opti-
408 mising vineyard management strategies for lowering operational costs and
409 enhancing sustainability.

410 **References**

- 411 Abbal, P., Sablayrolles, J.M., Matzner-Lober, É., Boursiquot, J.M., Baudrit,
412 C., Carbonneau, A., 2016. Decision Support System for Vine Growers
413 Based on a Bayesian Network. *Journal of agricultural, biological, and*
414 *environmental statistics* 21, 131–151. doi:10.1007/s13253-015-0233-2.
- 415 Agosta, E., Canziani, P., Cavagnaro, M., 2012. Regional climate variability
416 impacts on the annual grape yield in Mendoza, Argentina. *Journal of*
417 *Applied Meteorology and Climatology* 51, 993–1009.
- 418 Attorney-General’s Department, 2010. *Wine Australia Corporation Act*
419 *1980*.
- 420 Australia, W., 2021a. *Australian Wine: Production, Sales and Inventory*
421 *2019–20*.
- 422 Australia, W., 2021b. *Wine Australia-Open Data*.

423 Botting, D., Dry, P., Iland, P., 1996. Canopy architecture-implications for
 424 Shiraz grown in a hot, arid climate .

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

431 Camps, J.O., Ramos, M.C., 2012. Grape harvest and yield responses to inter-
 432 annual changes in temperature and precipitation in an area of north-east
 433 Spain with a Mediterranean climate. *International Journal of Biometeo-*
 434 *rology* 56, 853–64. doi:10.1007/s00484-011-0489-3.

435 Capello, G., Biddoccu, M., Cavallo, E., 2020. Permanent cover for soil and
 436 water conservation in mechanized vineyards: A study case in Piedmont,
 437 NW Italy 15.

438 Capello, G., Biddoccu, M., Ferraris, S., Cavallo, E., 2019. Effects of Tractor
 439 Passes on Hydrological and Soil Erosion Processes in Tilled and Grassed
 440 Vineyards. *Water* 11. doi:10.3390/w11102118.

441 Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System,
 442 in: *Proceedings of the 22nd ACM SIGKDD International Conference on*
 443 *Knowledge Discovery and Data Mining*, ACM, New York, NY, USA. pp.
 444 785–794. doi:10.1145/2939672.2939785.

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

448 Fraga, H., Costa, R., Santos, J.A., 2017. Multivariate clustering of viticul-
 449 tural terroirs in the Douro winemaking region. *Ciência Téc. Vitiv.* 32,
 450 142–153.

451 G. van Rossum, 1995. Python tutorial, Technical Report CS-R9526. Centrum
 452 voor Wiskunde en Informatica (CWI),.

453 Goodwin I, Jerie P, 1992. Regulated deficit irrigation: Concept to prac-
 454 tice. *Advances in vineyard irrigation. Australian and New Zealand Wine*
 455 *Industry Journal* 7.

456 Halliday, J.C.J.C., 2009. *Australian Wine Encyclopedia.* Hardie Grant
 457 Books, VIC.

458 Hand, D.J., Till, R.J., 2001. A Simple Generalisation of the Area Under the
 459 ROC Curve for Multiple Class Classification Problems. *Machine Learning*
 460 45, 171–186. doi:10.1023/A:1010920819831.

461 Hanley, J.A., McNeil, B.J., 1982. The meaning and use of the area under a
 462 receiver operating characteristic (ROC) curve. *Radiology* 143, 29–36.

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

466 Lu, W., Newlands, N.K., Carisse, O., Atkinson, D.E., Cannon, A.J., 2020.
 467 Disease Risk Forecasting with Bayesian Learning Networks: Application
 468 to Grape Powdery Mildew (*Erysiphe necator*) in Vineyards. *Agronomy*
 469 (Basel) 10, 622. doi:10.3390/agronomy10050622.

470 MG McCarthy, RM Cirami, DG Furkaliev, 1986. The effect of crop load and
 471 vegetative growth control on wine quality. .

472 Molnar, C., 2022. Interpretable Machine Learning: A Guide for Making
 473 Black Box Models Explainable. 2 ed.

474 Oliver, D., Bramley, R., Riches, D., Porter, I., Edwards, J., 2013. Review:
 475 Soil physical and chemical properties as indicators of soil quality in Aus-
 476 tralian viticulture. *Australian Journal of Grape and Wine Research* 19,
 477 129–139. doi:10.1111/ajgw.12016.

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

480 Rudiger, P., Stevens, J.L., Bednar, J.A., Nijholt, B., Andrew, B, C., Randel-
 481 hoff, A., Mease, J., Tenner, V., maxalbert, Kaiser, M., ea42gh, Samuels, J.,
 482 stonebig, LB, F., Tolmie, A., Stephan, D., Lowe, S., Bampton, J., henri-
 483 queribeiro, Lustig, I., Signell, J., Bois, J., Talirz, L., Barth, L., Liquet, M.,
 484 Rachum, R., Langer, Y., arabidopsis, kbowen, 2020. Holoviz/holoviews:
 485 Version 1.13.3. Zenodo. doi:10.5281/zenodo.3904606.

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

- 489 SWA, S.W.A., 2022. Sustainable Wingrowing Australia.
490 <https://sustainablewinegrowing.com.au/case-studies/>.
- 491 Terry Therneau, Beth Atkinson, 2022. Rpart: Recursive Partitioning and
492 Regression Trees.
- 493 WEBB, L.B., WHETTON, P.H., BARLOW, E.W.R., 2011. Observed trends
494 in winegrape maturity in Australia. *Global change biology* 17, 2707–2719.
495 doi:10.1111/j.1365-2486.2011.02434.x.
- 496 Wine Australia, 2020. National Vintage Report 2020 .
- 497 Wine Australia, 2021. National Vintage Report 2021 .
- 498 Wine Australia, 2022. National Vintage Report 2022 .

499 *Appendix .1. Year*

500 The classification tree and XGBoosted ensemble performed similarly for
501 classifying year with 35.20% (6.28% standard deviation) and 51.81% (42.20%
502 validation accuracy) respectively. Electricity and the type of irrigation were
503 highly influential within the classification tree. Similarly, electricity was the
504 most frequently occurring node in the XGBoost ensemble. Other variables
505 such as slashing passes, and fungicide and herbicide spraying were more
506 prevalent than in the classification tree. Weed and disease outbreaks are
507 likely an influential factor when classifying different years, making the de-
508 cisions to spray and slash unique factors that differ year to year. Climatic
509 differences between years are likely tied to the influence of yield and water
510 use.

511 Over half of the interrelated importance of the predictor variables is domi-
512 nated by area harvested, yield and slashing passes. Although all the predictor
513 variables are highly connected, their relative importance is not as prominent
514 as the three major variables. It is of particular note of the relative importance
515 of slashing passes to area, fuel and yield; as these are not directly related ac-
516 tivities. The connection between the number of slashing and spraying passes
517 is that those who do a set number of spraying or slashing passes tended to
518 do that many passes for all slashing and spraying activities.

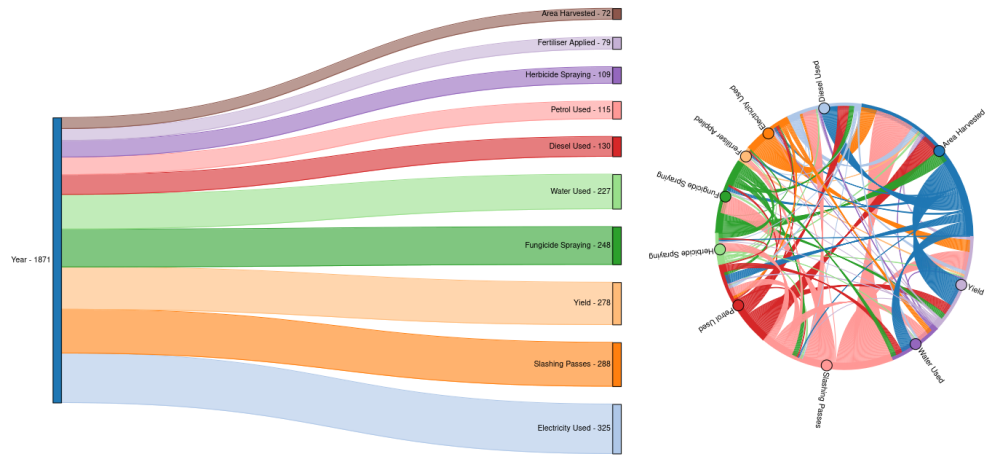


Figure .8: The left-hand side depicts the 10 most important variables in predicting Year using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.