

1 An analysis of underlying relationships between factors
2 related to operating costs and revenue in Australian
3 vineyards.

4 Author^{1,1,1}

5 **1. Introduction**

6 Historically strong demands for Australian wine have helped to create a
7 thriving industry. However, recent pressures brought on by a loss of tourism
8 and labour due to the COVID-19 pandemic, the global freight crisis, war in
9 Europe, tariffs and rising inflation has negatively affected the industry's out-
10 look (Wine Australia, 2021; Australia, 2021a). The 2021-2022 financial year
11 alone saw a decline of 19% in exports solely due to tariffs (Wine Australia,
12 2022). A greater understanding of the different underlying conditions leading
13 to improved performance in agricultural productivity and sustainability at
14 scale are key to making data-informed decisions to aid in increasing a nations
15 agricultural sustainability (OECD, 2019). Specifically within the Australian
16 Wine and vine industry there is a need to further understand the driving
17 relationships between resource use and economic output, where these rela-
18 tionships can lead to determining better and efficient methods and develop
19 benchmarks with local growers (Luke Mancini, 2020).

20 An unprecedented amount of data regarding the Australian winegrowing
21 industry has been collected through Sustainable Winegrowing Australia, of-
22 fering new insights into the driving economic forces of the Australian wine

industry. This dataset allows insights into the economic outcome of vineyards through the incorporation of operating costs and grape revenue from grape sales within the data. We use this data to study these economic outcomes and their statistical relationships to vineyards’ utilisation of the resources. We further compare the relationships between different resources to address the extensive collinearity found within the data (Chen and Guestrin, 2016). We adopt XGBoosted models for this analysis because they are able to overcome multicollinearity as well as highlight the level of importance that predictor variables have on response variables.

2. Methods

2.1. Data

Data used in this analysis were obtained from Sustainable Winegrowing Australia. Australia’s national wine industry sustainability program. The program aims to support grape-growers and winemakers in demonstrating and improving their sustainability (SWA, 2022). Data recorded by SWA are entered manually by winegrowers using a web based interface tool. A total of 6049 observations were collected from 2012/2013 to 2021/2022 financial years, with each observation comprising 23 variables reflecting a vineyard’s state for the given year (see Table 2.1).

The data originally contained only two multiclass variables: year and region. Related binary variables, such as the use of river water and the use of dam water, were combined to create a single multiclass variable. This was done by first converting each combination that occurred into a unique category (such as river and dam water used, as opposed to the two separate

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	6049	10
Disease	Class	6049	2
Region	Class	6049	58
Solar	Kilowatt Hours	622	
Irrigation Type	Class	6049	20
Petrol	Litres	4309	
Slashing	Times per year	2290	
Yield	Tonnes	5935	
Irrigation Energy	Class	6049	16
Area Harvested	Hectares	6049	
Electricity	Kilowatt Hours	1014	
Insecticide Spray	Times per year	1092	
Fertiliser	KGs of Nitrogen	795	
Fungicide Spray	Times per year	2260	
Cover Crop	Class	6049	32
Water Type	Class	6049	39
Grape Revenue	AUD	853	
Operating Costs	AUD	853	

categories prior). These variables were then one-hot-encoded, changing each variable class into a binary value, with one indicating the presence of the class and zero indicating its absence. Further details about classes and their frequency is available in the appendices.

The variable region represented one of the 65 Geographical Indicator Regions (GI Region) used to describe different unique localised traits of vineyards across Australia (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).

2.2. *XGBoosted Trees*

XGBoosted (eXtreme Gradient Boosting) trees, described in more detail below (and further in the appendix), 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. Moreover, XGBoosted trees are unaffected by multicollinearity, and offer high predictive performance for a wide variety of purposes (Chen and Guestrin, 2016).

XGBoosted models were constructed with operational cost and grape revenue as the predicted variables. The analyses were aimed at uncovering what factors influenced these variables and to what extent. As the purpose of the analysis was to identify relationships between variables and to show how they interact, an XGBoosted tree was created for each of the predictor variables as well. Trees for the predictor variables did not include operational cost or

72 grape revenue as predictors. By creating an XGBoosted tree for each variable
73 it meant that every variable would have a measure of its relative importance
74 to every other variable (see Section 2.3). Together these models were used to
75 measure the interrelationships of the ten most important variables in deter-
76 mining operational cost and grape revenue using variable importance. These
77 measures of relative importance were used to illustrate the highly interrelated
78 nature of variables within vineyards. The interaction between variables was
79 depicted through the use of Sankey and Chord diagrams; with variable im-
80 portance measures being used to show the strength of connection between
81 the respective predictor variable and the response (see section 2.3).

82 Due to constraints from the XGBoost library region could only be incor-
83 porated as a one-hot-encoded variable when used as a predictor. To better
84 show what variables were related to region overall, another XGBoost tree
85 was created with Region as the predicted value. The difference for this model
86 was that relative variable importance would only be measured for each re-
87 gion specifically, as opposed to a variables overall importance in determining
88 region. Separately profit (the difference between revenue and operational
89 costs) and year was looked at in prior analyses (see appendix) but these
90 results were not included due to low average loss values and model stability.

91 XGBoosted trees are an ensemble method that combines multiple decision
92 trees together to create a more accurate predictive model. The gradient
93 boosting aspect of the ensemble is the use of a loss function used to create
94 new decision trees that add to the ensemble. Each new tree created is done so
95 using a loss function that is optimised iteratively to improve upon prior tree's
96 predictive power. The loss function can be any convex function, allowing

97 gradient descent to traverse the loss space until, no improvements can be
98 made via traversal. Because the loss function is only required to be convex,
99 both classifiers and regressors can be used. Regularisation methods can also
100 be incorporated to help prevent over fitting.

101 *2.3. Variable Importance*

102 Due to XGBoost creating a large amount of decision trees, the inter-
103 pretability of these models is obfuscated by the intricate relationships within
104 complicated ensembles. A measure of variable importance was the technique
105 used to highlight a variables influence within the ensemble. Variable impor-
106 tance can be measured in multiple ways; we used the frequency of a variable
107 appearing as a node within the ensemble as a measure of its importance.
108 This measure was chosen as it connected a variable to the minimisation of
109 its associated objective function. The measure of a variable’s importance
110 within this study can then be interpreted as how often a variable was the
111 optimal choice in reducing the loss function of the ensemble. Importantly,
112 multiclass variables being one-hot-encoded (see Section 2.1) are given an im-
113 portance score for each individual class; for example, each specific region will
114 have its own importance score.

115 The Sankey and Chord diagrams were constructed using the Holoviews
116 python library (Rudiger et al., 2020). Both Chord and Sankey diagrams
117 illustrated variable importance through the size of the bands between two
118 variables. The number at the end of a connection in a Sankey diagram indi-
119 cates a variable’s importance, or the number of times it appeared within the
120 ensemble. Sankey and Chord diagrams are presented together; with Sankey
121 diagrams showing the connection of a variable to its ten most important pre-

dictor variables. Chord diagrams were used alongside a Sankey diagram to show the interconnectedness of the ten most prominent variables within its associated Sankey diagram. Chord diagrams formed circles, with variables being connected through their relative importance. The importance values for the Chord diagrams were taken from the models of those individual variables, with the diagram being simplified to just the ten variables in the associated Sankey diagram, for readability's sake.

2.4. Validation

The predictive accuracy of each tree was assessed through a validation process. For each model the data was split into training data, which constituted 80% of the original data. The remaining 20% was used in testing and validation. Categorical data was stratified to conserve the same proportion of class occurrences between training, testing and validation data. For continuous variables 20% was used as testing data and the models were validated using 10 repetitions of the sampling process (10-fold cross validation). R^2 scores were used to determine the best regression models during validation. R^2 was used instead of RMSE to allow the comparison of models with different units to each other when considering how well each model extrapolated to further data. For binary and multiclass variables, validation was summarised through the accuracy, the proportion of true negatives and positives.

The XGBoost library incorporates regularisation techniques built into the software to mitigate over-fitting and enhance model generalisation. This allowed us to utilise cross validated grid search functions when selecting for better performing hyperparameters. The performance measure for model

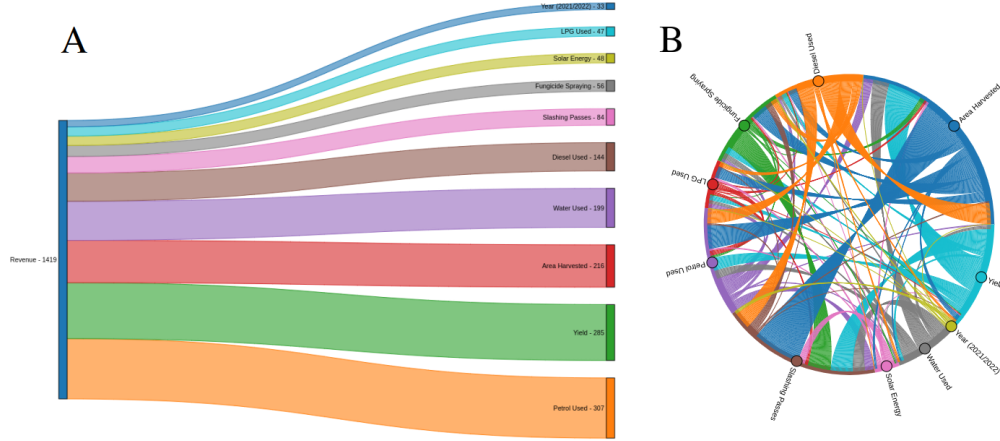


Figure 1: The left-hand side depicts the 10 most important variables in predicting revenue using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

selection was root-mean-square error for continuous variables. The receiver operator characteristic's area under the curve was used for category variables (Hanley and McNeil, 1982). Multiclass variables utilised the one verse one approach to minimise sensitivity to class disparity (Ferri et al., 2009; Hand and Till, 2001).

3. Results

3.1. Revenue

The prediction of revenue performed similarly to operating cost achieving an R^2 of 0.7716 (with a standard deviation of 0.1525). The value of predictors' relative importance was then calculated through the number of nodes used within the XGBoost. Values for relative importance were then used to

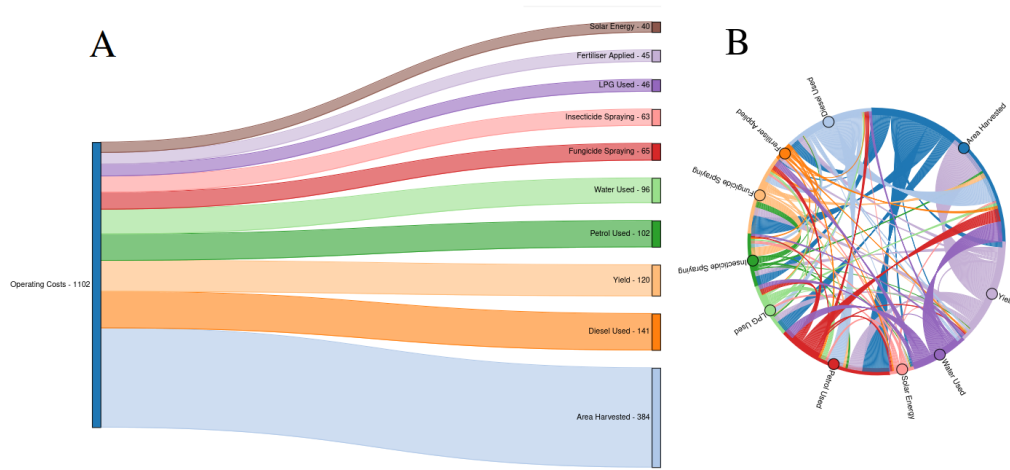


Figure 2: The left-hand side, A, depicts the 10 most important variables in predicting Operating Costs using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The number at the end of each band in the diagram is that variable’s importance. The right-hand side, B, depicts the importance of the 10 variables in Sankey diagram relative to one another.

construct Sankey and Chord diagrams to compare the contribution of each variable to predicting revenue.

In order of importance, predictors of revenue were fuel use (petrol 307 and diesel 144), yield (285), size (216) and water use (199). Here, the values in the brackets indicate the relative importance of each variable (see C.7). Overall regions contributed to 234 nodes in the ensemble making them collectively the third most important variable. The chord diagram illustrates that vineyard area is also of high relative importance to other variables especially slashing. The overall importance of area to other variables is evident by its larger circumference within the chord diagram (see B in Figure C.7).

168 3.2. Operating Costs

169 Comparatively to revenue, operating cost performed better with the XG-
170 Boosted regression ensemble achieving an R^2 of 0.8025 (with a standard
171 deviation of 0.1033). The relationships to operating cost through variable
172 importance were found to be similar to that of revenue, with fuel, water,
173 area and yield having the largest number relative importance (see figure 2).
174 A surprising difference is that the most important operational consideration
175 for operating cost is the use of fungicide, compared to revenue where slash-
176 ing is the most important (comparing Figure 3). The variables that feed into
177 these decisions are also very different with diesel having the highest relative
178 importance to slashing, and area having the greatest relative importance to
179 the need for fungicide.

180 Again, region played a determining factor overall, contributing to 334
181 nodes within the ensemble making it the most important variable when con-
182 sidering all regions together. It was surprising that electricity, slashing and
183 spraying passes were not more prominent in operating costs due to the in-
184 trinsic nature as an agricultural expense.

185 3.3. Region

186 When considered overall, Region was a highly informative variable based
187 on measures of importance for both operating cost and revenue. As noted
188 above, Region was the third most important variable for determining rev-
189 enue. The Barossa Valley region and Tasmania were the two most important
190 regions in relation to revenue; these two regions are considered to be some of
191 the highest revenue per hectare regions in Australia (Wine Australia, 2022).
192 These two regions are also relative opposites in winegrowing climates with

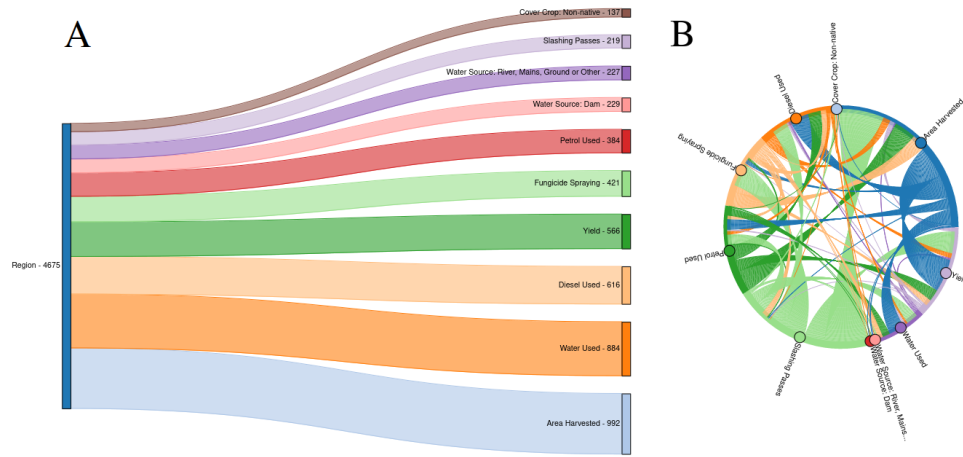


Figure 3: The left-hand side, A, depicts the 10 most important variables in predicting Region using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The number at the end of each band in the diagram is that variable's importance. The right-hand side, B, depicts the importance of the 10 variables in Sankey diagram relative to one another.

193 the Barossa having a warm and dry climate focussing on Shiraz grapes and
194 Tasmania having a cool wet climate that favours Pinot.

195 As also noted above Region was also a key determinant of operating costs.
196 Again Tasmania was the most important, followed by the Adelaide Hills. In
197 contrast to revenue, both climates are considered cool and wet, and warmer
198 drier regions such as the Barossa and Hunter Valley only contributed roughly
199 half the same number of nodes to the ensemble. Based on further analysis
200 of Regions (Figure 3) the inclusion of slashing and fungicide spraying is the
201 likely reason with fungal and weed pressure being greater in cooler wetter
202 regions.

203 The XGBoost ensemble, did not perform well when predicting operating
204 costs or revenue with 56.82% accuracy (50.58% validation accuracy). The
205 difference in accuracy is in part due to the large number of classes (58 re-
206 gions). The ensemble had a great emphasis on area, water, fuel and yield as
207 determining factors (see Figure (3)).

208 Many of the regions had significantly lower reporting rates, resulting in
209 much poorer classification performance. The regions with the most samples
210 performed the best. Bordering regions were routinely grouped together and
211 misclassified as the same region. Two areas that suffered the most from this
212 were the Limestone Coast (cool coastal areas in South Australia) and the
213 warmer inland regions along the Murray Darling.

214 4. Discussion

215 This study explored the relationships between vineyard resource use, op-
216 erations and geographical properties to revenue and operating costs. The

analysis was based on a large national study of 6049 samples collected over ten years. Three main findings were identified. First, the most important predictors of revenue and operating costs were fuel, yield and area. Secondly, area and fuel were highly interrelated to other variables (see Figure 2 and Figure C.7). Finally, the relative importance of predictor variables for region, differed from Revenue and operating costs, with water use being more prominent than yield. Region was also more prominent than illustrated in the Sankey diagrams due to the relative importance for operating cost and revenue being calculated for individual regions and not all regions together. In its entirety region contributed third most prominently in relative importance to revenue, and was of the most relative importance in determining operating costs.

Several physical parameters such as climate, geography and soil are predetermined by a vineyard's location; making it a widely considered key determinant of grape yield and quality (Abbal et al., 2016; Agosta et al., 2012; Fraga et al., 2017). The association between yield and region is demonstrated by its rank of fourth-highest variable importance when determining region (see Figure 3).

Warmer regions are known to be beneficial in hastening the ripening process of winegrapes (Webb et al., 2011). Warmer regions are also associated with lower quality grapes, caused largely due to this hastened ripening (Botting et al., 1996). In general warmer regions are not associated with higher yields, but if a vineyard in a warmer region is sufficiently irrigated much higher yields can be achieved than in cooler regions (Camps and Ramos, 2012). It is likely that the combination of larger vineyards with higher water

242 use is a determining factor in classifying regions which favour larger produc-
243 tion of grapes; reflected through region using water use so prominently in the
244 XGBoost ensemble. The link to water resources in defining regions is also
245 an important consideration, as vineyards can leverage higher irrigation rates
246 given more accessible water resources. A further consideration in the link
247 between revenue and region is that grape prices are set at a regional level by
248 buyers (Wine Australia, 2022). It is also important to consider that some
249 regions carry particular fame regarding the quality of their produce such as
250 Tasmania, the Hunter Valley and Barossa Valley (Halliday, 2009). This clas-
251 sification can be contrasted with other warmer regions of higher rainfall that
252 use the warmer climate to concentrate their grapes, increasing the flavour
253 profile (and thus quality) (Goodwin I, Jerie P, 1992; MG McCarthy et al.,
254 1986).

255 In part some winegrowing strategies are restricted simply through access
256 to water resources. Regions are likely to have varying access to different water
257 sources, such as those along the River Murray being able to utilise river water
258 for crops, unlike most coastal regions which may be drawing from surface or
259 underground water sources. Similarly, the connection between region and
260 fuel use is likely an indicator of the level of infrastructure within the region
261 because vineyards in regions without pressurised water will need to use more
262 fuel to pressurise their irrigation systems.

263 Operational costs showed similar importance across fuel, water and trac-
264 tor use. The dominating factor of area likely played a large part in deter-
265 mining how costly a tractor pass would be, or in defining the ratio of water
266 applied to the amount of vines. The node frequency was high for area but

267 much lower in general across the other variables, which could indicate the
268 need to be specific when attempting to determine the cause of a operational
269 cost. Although it was attempted to capture the complexity between how
270 variables interacted when determining operational costs (see Figure 2), it
271 is likely yet more complicated. An example of how interrelated operational
272 costs can be, is the optimisation of tractor passes to achieve multiple goals
273 in a pass, being shown to reduce energy use in vineyards, decreasing running
274 costs, as well as reducing soil compaction (Capello et al., 2019).

275 When determining revenue, similar variables were used to operational
276 cost; with region also being of high variable importance relative to other
277 variables (when considering all regions together in importance). It is difficult
278 to extrapolate the specific influence of location on a vineyard’s outcomes due
279 to the broad and varying definition of a region. Utilising the Geographical
280 Indicator regions defined by Wine Australia (Australia, 2021b) is a limitation
281 in one way, as it is too broad to fully capture a vineyards location and how
282 that influences variables at a more granular level. However, as buyers set
283 prices at regional levels, it is still important to consider this factor.

284 Decisions made on the ground have far-reaching effects and are difficult
285 to completely capture. A larger number of tractor passes used as a preven-
286 tative measure for occurrences such as disease may incur higher operational
287 costs but could be critical in preventing long term losses. Although the
288 models demonstrated a good predictive fit (via large R^2 values), the ability
289 to predict operational costs is limited by the variables incorporated in the
290 analysis. Other factors such as erosion and soil health are also influenced by
291 tractor use and would contribute to these operational costs but are difficult

292 to measure and were not available as part of the data (Capello et al., 2019,
293 2020). Reductions in fuel, water and tractor use are obvious methods to
294 reduce operational costs but not necessarily achievable decisions. Without
295 fully capturing more granular activities for example the specific reasons for
296 fuel use, it is difficult to determine what decisions specifically influence the
297 operational costs.

298 The reasoning for any particular decision can be widely varying. More
299 sophisticated models, specifically those that utilise expert opinion, may also
300 help to capture and address the decision-making process. An example is the
301 optimisation of fungicide sprays using Bayesian models that forecast disease
302 risk (Lu et al., 2020).

303 Separately revenue and operating cost did have a greater predictabil-
304 ity than their counterpart profit. The disparity in accuracy between profit
305 and other economic outcomes is reflective of the complexity in trying to
306 address challenges such as climate change, disease and changing market de-
307 mands (Wine Australia, 2020, 2021, 2022). The difference between turning a
308 profit or loss is dependent on predictable factors unforecasted factors, farm-
309 ing practice and farmers' decisions. The difference between vineyards that
310 make profit and those that do not could be a multitude of factors including
311 differences in farming practices not captured within this study. Some deci-
312 sions leading to latent effects such as large scale soil deposition in extreme
313 rain events can be caused by soil compaction due to overworking a vineyard
314 (Capello et al., 2020).

315 5. Conclusion

316 This study has provided valuable insights into the multifaceted dynam-
317 ics governing operational costs and revenue. The impact of different regions
318 highlighted the complex interrelatedness of variables within a vineyard. We
319 relate how factors such as water and fuel intersect to impact operational
320 costs and how different seasonal events affect these operations; as well as
321 the significance of context-specific decision-making. While this investigation
322 utilised a broad regional classification, the potential benefits of adopting a
323 more nuanced approach and incorporating expert knowledge have been high-
324 lighted. Further work could pursue causal models and the creation of decision
325 support systems. It is difficult to untangle the predictive and correlative na-
326 ture of a variable compared to the causal reasons. By delving deeper into
327 the complex interplay of variables, further advancements can be made in
328 optimising vineyard management strategies for lowering operational costs,
329 increasing revenue and enhancing sustainability.

330 References

- 331 Abbal, P., Sablayrolles, J.M., Matzner-Lober, É., Boursiquot, J.M., Baudrit,
332 C., Carbonneau, A., 2016. Decision Support System for Vine Growers
333 Based on a Bayesian Network. *Journal of agricultural, biological, and*
334 *environmental statistics* 21, 131–151. doi:10.1007/s13253-015-0233-2.
- 335 Agosta, E., Canziani, P., Cavagnaro, M., 2012. Regional climate variability
336 impacts on the annual grape yield in Mendoza, Argentina. *Journal of*
337 *Applied Meteorology and Climatology* 51, 993–1009.

338 Attorney-General's Department, 2010. Wine Australia Corporation Act
 339 1980.

340 Australia, W., 2021a. Australian Wine: Production, Sales and Inventory
 341 2019–20.

342 Australia, W., 2021b. Wine Australia-Open Data.

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

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

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

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

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

- 361 Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System,
362 in: Proceedings of the 22nd ACM SIGKDD International Conference on
363 Knowledge Discovery and Data Mining, ACM, New York, NY, USA. pp.
364 785–794. doi:10.1145/2939672.2939785.
- 365 Ferri, C., Hernández-Orallo, J., Modroi, R., 2009. An experimental com-
366 parison of performance measures for classification. Pattern Recognition
367 Letters 30, 27–38. doi:10.1016/j.patrec.2008.08.010.
- 368 Fraga, H., Costa, R., Santos, J.A., 2017. Multivariate clustering of viticul-
369 tural terroirs in the Douro winemaking region. Ciência Téc. Vitiv. 32,
370 142–153.
- 371 G. van Rossum, 1995. Python tutorial, Technical Report CS-R9526. Centrum
372 voor Wiskunde en Informatica (CWI),.
- 373 Goodwin I, Jerie P, 1992. Regulated deficit irrigation: Concept to prac-
374 tice. Advances in vineyard irrigation. Australian and New Zealand Wine
375 Industry Journal 7.
- 376 Halliday, J.C.J.C., 2009. Australian Wine Encyclopedia. Hardie Grant
377 Books, VIC.
- 378 Hand, D.J., Till, R.J., 2001. A Simple Generalisation of the Area Under the
379 ROC Curve for Multiple Class Classification Problems. Machine Learning
380 45, 171–186. doi:10.1023/A:1010920819831.
- 381 Hanley, J.A., McNeil, B.J., 1982. The meaning and use of the area under a
382 receiver operating characteristic (ROC) curve. Radiology 143, 29–36.

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

387 Luke Mancini, 2020. Understanding the Australian Wine Industry: A growers
388 guide to the background and participants of the wine grape industry.

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

391 OECD, 2019. Innovation, Productivity and Sustainability in Food and Agri-
392 culture.

393 Oliver, D., Bramley, R., Riches, D., Porter, I., Edwards, J., 2013. Review:
394 Soil physical and chemical properties as indicators of soil quality in Aus-
395 tralian viticulture. *Australian Journal of Grape and Wine Research* 19,
396 129–139. doi:10.1111/ajgw.12016.

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

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

- 406 SWA, S.W.A., 2022. Sustainable Wingrowing Australia.
407 <https://sustainablewinegrowing.com.au/case-studies/>.
- 408 Webb, L.B., Whetton, P.H., Barlow, E.W.R., 2011. Observed trends in
409 winegrape maturity in Australia. *Global change biology* 17, 2707–2719.
410 doi:10.1111/j.1365-2486.2011.02434.x.
- 411 Wine Australia, 2020. National Vintage Report 2020 .
- 412 Wine Australia, 2021. National Vintage Report 2021 .
- 413 Wine Australia, 2022. National Vintage Report 2022 .

414 **Appendix A. Continuous variables**

415 Table A.2 below shows the ranges of each of the continuous variables:

Table A.2: Summary statistics of continuous variables used in XGBoosted models.

	count	mean	std	min	0.25	0.5	0.75	max
Vineyard Solar	622	22916.89	104808	1	1170.75	5500	14866.25	2300000
Biodiesel	25	6635.932	11768.832104	1	200	500	10000	37216
Fungicide Spray	2260	7.724801	3.279794	1	6	7	9	68
LPG	958	327.831399	861.538804	1	40	95.835	240	11950
Petrol	4309	825.276809	1556.621119	1	135	306.66	903	38568
Insecticide Spray	1092	1.707189	1.316042	0	1	1	2	12
Water Used	5846	7301838	558206600	0.0007	13.2655	43	146.875	42680000000
Fertiliser	795	91149.89	483913.4	1	560	4759.5	45148.5	11358000
Diesel	5585	11677.070183	24380.588742	0.1267	1240	3850	12500	591000
Yield	5935	772.902449	2175.113895	0.03	68	192.3	601.8795	72305
Herbicide Spray	2026	2.646199	2.598899	0	2	2	3	103
Slashing	2290	3.311485	1.826788	1	2	3	4	26
Electricity	1014	58223.07	177626.3	0.019	2160	9637	36498.25	3000000
Area Harvested	6049	66.52604	133.4525	2.220446E-16	10.13	24.5	66.8	2436.15
Grape Revenue	875	377972	606286.8	1	76000	172964	386747	5700000
Operating Costs	853	314187.1	511522.6	1	57315	140000	327408	4482828

416 **Appendix B. Categorical Variables**

417 The tables below describe each possible class a multiclass variable could
418 have taken and the frequency that it occurred.

419 *Appendix B.1. Water Source Types*

420 Table B.3 below shows the different class types for water sources used by
421 vineyards and their frequency of occurrences.

Table B.3: Frequency and class types of water types used
by vineyards.

Water types	frequency
river water	1578
groundwater	1433
surface water dam	617
recycled water from other source	386
groundwater and surface water dam	256
not listed	235
mains water	170
river water and groundwater	147
groundwater and recycled water from	145
other source	
other water	101
river water and surface water dam	92

Continued on next page

Table B.3 – continued from previous page

Water types	frequency
groundwater and water applied for frost control	90
groundwater and mains water	76
river water and groundwater and surface water dam	70
recycled water from other source and mains water	63
groundwater and recycled water from other source and mains water	60
river water and mains water	57
surface water dam and mains water	56
groundwater and other water	33
river water and groundwater and mains water	30
groundwater and surface water dam and recycled water from other source	27
river water and water applied for frost control	27
groundwater and surface water dam and mains water	22
surface water dam and recycled water from other source	21
Continued on next page	

Table B.3 – continued from previous page

Water types	frequency
river water and recycled water from other source	19
river water and other water	19
river water and surface water dam and mains water	18
river water and groundwater and sur- face water dam and mains water	18
mains water and other water	16
groundwater and surface water dam and water applied for frost control	12
surface water dam and other water	12
groundwater and recycled water from other source and other water	11
groundwater and surface water dam and recycled water from other source and mains water	8
recycled water from other source and mains water and other water	8
river water and recycled water from other source and mains water	8
river water and surface water dam and recycled water from other source	8
Continued on next page	

Table B.3 – continued from previous page

Water types	frequency
surface water dam and mains water and other water	7
recycled water from other source and other water	7
river water and groundwater and recy- cled water from other source	6
groundwater and mains water and other water	5
groundwater and surface water dam and other water	5
groundwater and surface water dam and mains water and other water	5
river water and groundwater and re- cycled water from other source and mains water	5
river water and groundwater and wa- ter applied for frost control	5
river water and surface water dam and water applied for frost control	4
surface water dam and water applied for frost control	4

Continued on next page

Table B.3 – continued from previous page

Water types	frequency
river water and groundwater and sur- face water dam and recycled water from other source and mains water and other water	4
river water and groundwater and recy- cled water from other source and other water	3
groundwater and surface water dam and recycled water from other source and water applied for frost control	3
river water and groundwater and sur- face water dam and recycled water from other source	3
river water and recycled water from other source and other water	3
surface water dam and recycled water from other source and mains water	2
river water and recycled water from other source and mains water and wa- ter applied for frost control	2

Continued on next page

Table B.3 – continued from previous page

Water types	frequency
groundwater and surface water dam	2
and recycled water from other source	
and mains water and other water	
river water and groundwater and	2
mains water and other water	
river water and groundwater and sur-	2
face water dam and other water	
river water and surface water dam and	2
other water	
river water and mains water and water	2
applied for frost control	
river water and groundwater and sur-	2
face water dam and recycled water	
from other source and mains water	
river water and mains water and other	2
water	
river water and surface water dam and	2
mains water and other water	
river water and groundwater and	1
mains water and water applied for	
frost control	

Continued on next page

Table B.3 – continued from previous page

Water types	frequency
surface water dam and other water and water applied for frost control	1
water applied for frost control	1
groundwater and other water and wa- ter applied for frost control	1
other water and water applied for frost control	1
groundwater and surface water dam and recycled water from other source and other water and water applied for frost control	1
mains water and water applied for frost control	1
groundwater and surface water dam and recycled water from other source and other water	1
groundwater and mains water and wa- ter applied for frost control	1
river water and groundwater and sur- face water dam and mains water and other water	1

Continued on next page

Table B.3 – continued from previous page

Water types	frequency
river water and surface water dam and	1
recycled water from other source and	
mains water	

423 *Appendix B.2. Cover Crop Types*

424 Table B.4 below shows the different cover crop types used together and
425 their frequency.

Table B.4: Frequency and class types of cover crop types
used by vineyards.

Cover crop types	frequency
Cover crop types	frequency
permanent cover crop volunteer sward	1822
permanent cover crop non native	936
permanent cover crop native	490
annual cover crop	479
groundwater and surface water dam	406
annual cover crop and permanent cover crop volunteer sward	309
bare soil	225
permanent cover crop non native and permanent cover crop volunteer sward	214
annual cover crop and permanent cover crop non native	169
bare soil and permanent cover crop volunteer sward	129
Continued on next page	

Table B.4 – continued from previous page

Cover crop types	frequency
bare soil and permanent cover crop non native	115
annual cover crop and permanent cover crop non native and permanent cover crop volunteer sward	101
bare soil and annual cover crop	93
permanent cover crop native and per- manent cover crop volunteer sward	80
bare soil and permanent cover crop na- tive	78
annual cover crop and permanent cover crop native	78
permanent cover crop native and per- manent cover crop non native	68
permanent cover crop native and per- manent cover crop non native and per- manent cover crop volunteer sward	44
annual cover crop and permanent cover crop native and permanent cover crop non native and permanent cover crop volunteer sward	44

Continued on next page

Table B.4 – continued from previous page

Cover crop types	frequency
bare soil and annual cover crop and permanent cover crop volunteer sward	33
bare soil and permanent cover crop non native and permanent cover crop volunteer sward	26
annual cover crop and permanent cover crop native and permanent cover crop volunteer sward	17
bare soil and annual cover crop and permanent cover crop native	15
annual cover crop and permanent cover crop native and permanent cover crop non native	15
bare soil and annual cover crop and permanent cover crop non native	13
bare soil and annual cover crop and permanent cover crop native and per- manent cover crop non native and per- manent cover crop volunteer sward	12
bare soil and annual cover crop and permanent cover crop non native and permanent cover crop volunteer sward	11
Continued on next page	

Table B.4 – continued from previous page

Cover crop types	frequency
bare soil and annual cover crop and permanent cover crop native and per- manent cover crop non native	8
bare soil and permanent cover crop na- tive and permanent cover crop non na- tive	7
bare soil and permanent cover crop na- tive and permanent cover crop volun- teer sward	6
bare soil and permanent cover crop na- tive and permanent cover crop non na- tive and permanent cover crop volun- teer sward	4
bare soil and annual cover crop and permanent cover crop native and per- manent cover crop volunteer sward and	2

427 *Appendix B.3. Irrigation Types*

428 Below in Table B.5 are the frequency and different irrigation types.

Table B.5: Frequency and class types of irrigation types
used by vineyards.

Irrigation types	frequency
Irrigation type	frequency
dripper	4800
dripper and non irrigated	342
Not listed	319
dripper and overhead sprinkler	201
dripper and undervine sprinkler	91
non irrigated	65
undervine sprinkler	53
dripper and flood	53
overhead sprinkler	46
dripper and overhead sprinkler and undervine sprinkler	28
overhead sprinkler and undervine sprinkler	12
dripper and non irrigated and overhead sprinkler	11
flood and undervine sprinkler	10
Continued on next page	

Table B.5 – continued from previous page

Irrigation types	frequency
dripper and flood and undervine sprinkler	7
dripper and flood and non irrigated and overhead sprinkler and undervine sprinkler	3
dripper and flood and overhead sprinkler	3
non irrigated and undervine sprinkler	2
dripper and flood and non irrigated	1
dripper and non irrigated and overhead sprinkler and undervine sprinkler	1
flood and	1

430 *Appendix B.4. Irrigation Energy Type*

431 Below, Table ?? shows the different types of energy used to power vine-
 432 yards and their frequency.

Table B.6: Frequency and class types of irrigation energy types used by vineyards.

Irrigation Energy types	frequency
Irrigation energy type	frequency
electricity	2162
not listed	2053
pressure	586
electricity and pressure	396
diesel	254
diesel and electricity	227
electricity and solar	96
diesel and electricity and pressure	90
diesel and pressure	74
solar	50
electricity and pressure and solar	23
diesel and electricity and solar	14
diesel and electricity and pressure and solar	10
pressure and solar	9
Continued on next page	

Table B.6 – continued from previous page

Irrigation Energy types	frequency
diesel and solar	4
diesel and pressure and solar and	1

434 *Appendix B.5. Year*

435 Below in Table B.7 is the list of years and the number of sample collected
436 in each.

Table B.7: Frequency and class types of year

Year	frequency
Year	frequency
2021/2022	954
2020/2021	860
2019/2020	599
2012/2013	590
2013/2014	549
2015/2016	548
2014/2015	505
2017/2018	493
2016/2017	485
2018/2019	466

437

439 Below in Table B.8 are the number of collected samples for each region.

Table B.8: Frequency and class types of regions.

Regions	frequency
giregion	frequency
McLaren Vale	1195
Barossa Valley	584
Murray Darling	521
Riverland	472
Adelaide Hills	454
Langhorne Creek	347
Margaret River	344
Coonawarra	284
Padthaway	202
Wrattonbully	195
Clare Valley	149
Yarra Valley	122
Eden Valley	92
Tasmania	89
Swan Hill	83
Grampians	73
Orange	72

Continued on next page

Table B.8 – continued from previous page

Regions	frequency
Hunter Valley	70
Bendigo	53
Great Southern	51
Rutherglen	41
Robe	36
Tumbarumba	35
Mornington Peninsula	32
King Valley	32
Southern Fleurieu	30
Heathcote	29
Adelaide Plains	25
Currency Creek	24
	23
Henty	22
Canberra District	21
Southern Flinders Ranges	20
Upper Goulburn	20
Mudgee	20
Mount Benson	20
Other	19
Riverina	18
Alpine Valleys	15
Continued on next page	

Table B.8 – continued from previous page

Regions	frequency
Barossa Zone	14
Pemberton	12
Mount Gambier	11
Blackwood Valley	10
Kangaroo Island	10
Big Rivers Zone Other	9
Geographe	7
Cowra	6
Gundagai	5
Strathbogie Ranges	5
Glenrowan	4
Geelong	4
Swan District	4
Goulburn Valley	3
Beechworth	3
Southern Highlands	3
Macedon Ranges	2
Pyrenees	2
Sunbury	1

441 Appendix C. XGBoost

442 Following Chen and Guestrin (Chen and Guestrin, 2016), XGboosted
 443 trees predict a value y_i from the input x_i . The method of prediction is
 444 achieved through a tree ensemble model, using K additive functions to pre-
 445 dict the output. Each of f_k functions is a classification or regression tree, such
 446 that all functions are in the set of all decision trees, given by \mathcal{F} , is defined
 447 by $f(x) = \omega_{q(x)}(q : \mathbb{R}^m \rightarrow T, \omega \in \mathbb{R}^T)$. Where each function corresponds to
 448 an independent tree structure q of ω weights. Each tree has T leaves, which
 449 contain a continuous score, represented by ω_i for the i -th leaf. The final
 450 prediction is determined by the sum of the score of the corresponding leaves,
 451 given by:

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

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

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

454 , where

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

455 As predictions are made using additive tree functions, XGboosted trees
 456 can be used for classification or regression. The difference between a predic-
 457 tion, $\phi(x_i)$, and actual variable, $f_k(x_i)$, is a differentiable convex loss function
 458 l . These properties of l allow the function to be versatile in which objective
 459 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.

Appendix C.1. Loss functions

The functions included as parameters in equation C.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-class data. The loss function for making a split within the tree structure is given by:

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

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 (C.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).

482 *Appendix C.2. Year*

483 The classification tree and XGBoosted ensemble performed similarly for
484 classifying year with 35.20% (6.28% standard deviation) and 51.81% (42.20%
485 validation accuracy) respectively. Electricity and the type of irrigation were
486 highly influential within the classification tree. Similarly, electricity was the
487 most frequently occurring node in the XGBoost ensemble. Other variables
488 such as slashing passes, and fungicide and herbicide spraying were more
489 prevalent than in the classification tree. Weed and disease outbreaks are
490 likely an influential factor when classifying different years, making the de-
491 cisions to spray and slash unique factors that differ year to year. Climatic
492 differences between years are likely tied to the influence of yield and water
493 use.

494 Over half of the interrelated importance of the predictor variables is domi-
495 nated by area harvested, yield and slashing passes. Although all the predictor
496 variables are highly connected, their relative importance is not as prominent
497 as the three major variables. It is of particular note of the relative importance
498 of slashing passes to area, fuel and yield; as these are not directly related ac-
499 tivities. The connection between the number of slashing and spraying passes
500 is that those who do a set number of spraying or slashing passes tended to
501 do that many passes for all slashing and spraying activities.

502 *Appendix C.3. Profit*

503 Predictions of profit performed poorly compared to operating cost and
504 revenue with an average R^2 of 0.2535 and standard deviation of 0.3126. With
505 the large standard deviation being indicative of how unstable the models
506 created were.

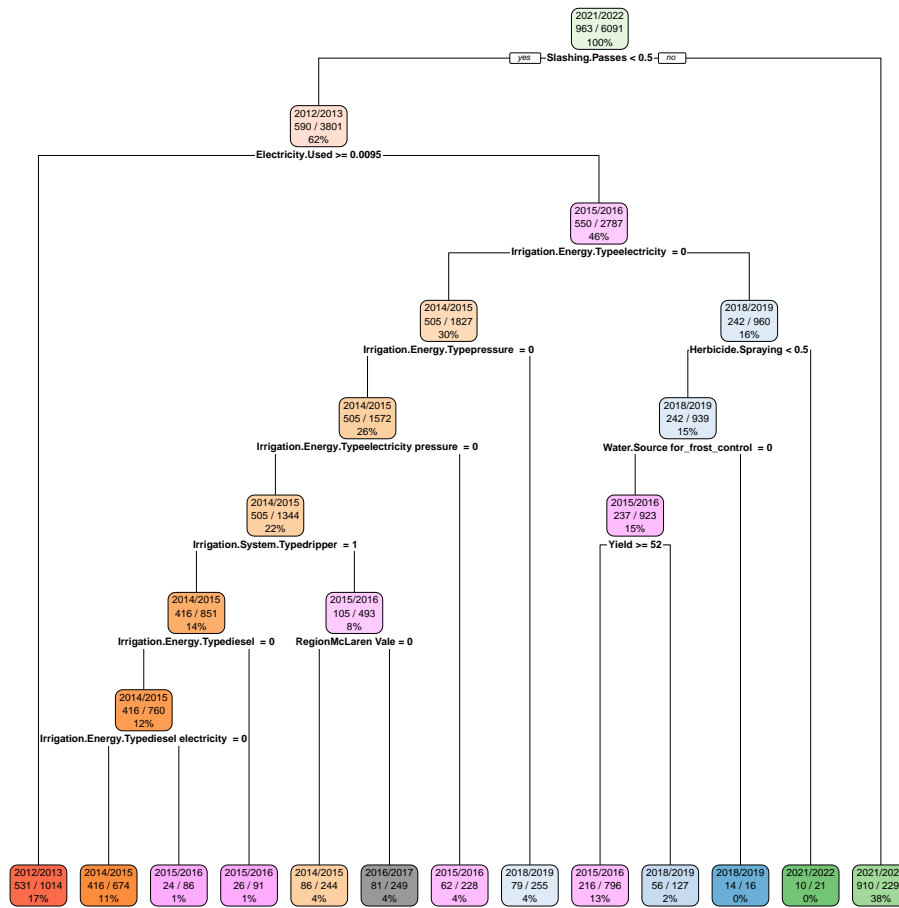


Figure C.4: Decision tree predicting Year. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

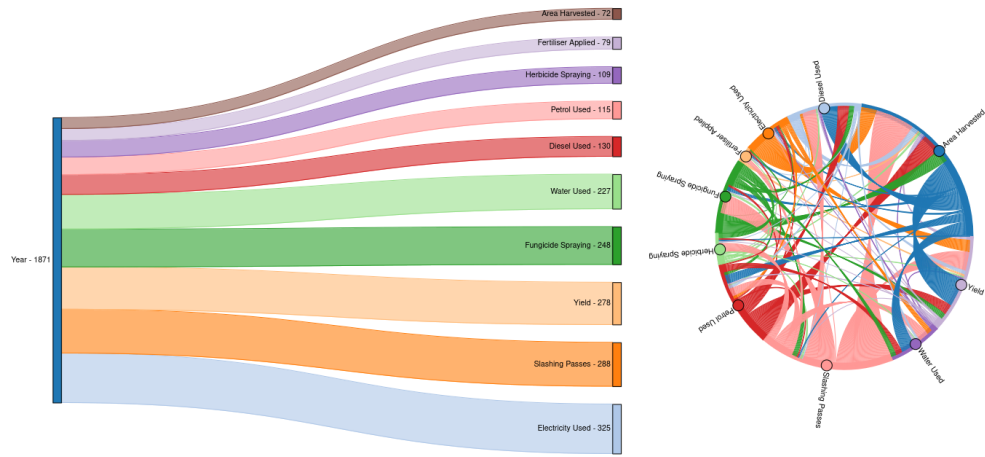


Figure C.5: The left-hand side depicts the 10 most important variables in predicting Year using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.

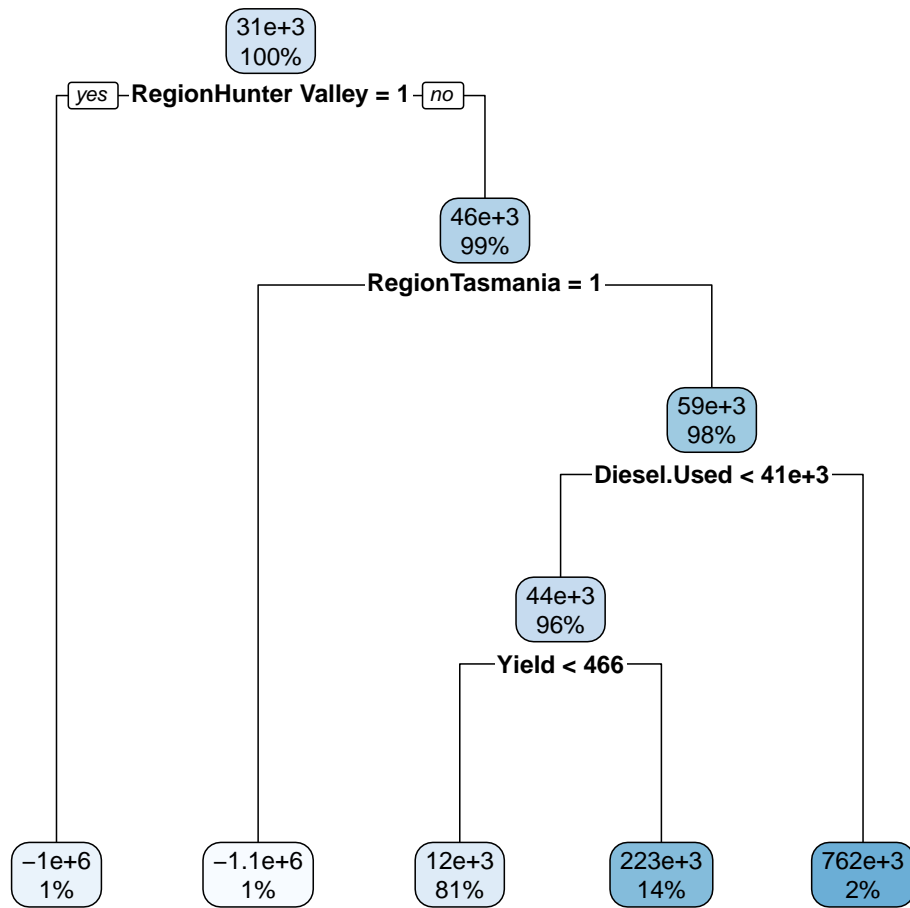


Figure C.6: Decision tree predicting revenue. Each node indicates the class predicted, and the proportion of elements agreeing with nodes partitioning, with the left direction indicating a yes to the nodes rule.

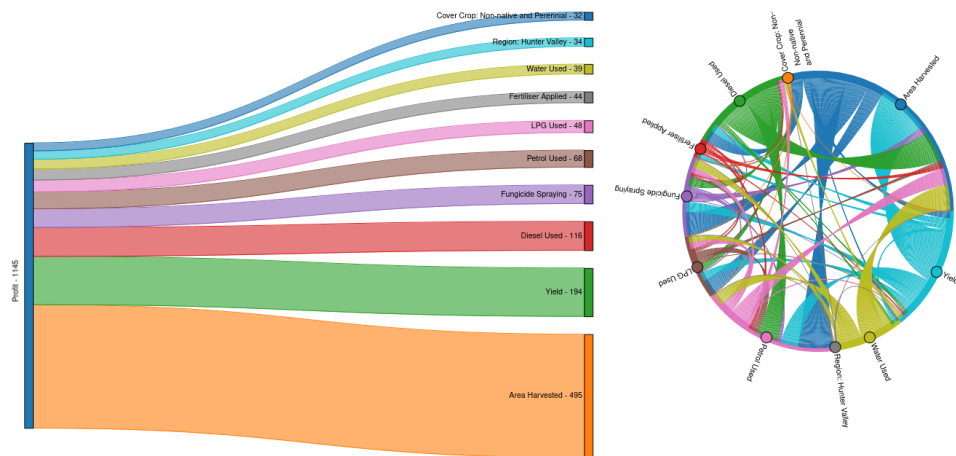


Figure C.7: The left-hand side depicts the 10 most important variables in predicting revenue using XGBoosted trees as a measure of node occurrence, using a Sankey diagram. The right-hand side depicts the interrelated importance of the ten predictor variables using a chord diagram.