

1 Comparing Machine-Learning models of different  
2 levels of complexity for crop protection : A look  
3 into the complexity-accuracy tradeoff

4 Olivier Gauriau<sup>1</sup>, Luis Galárraga<sup>2</sup>, François Brun<sup>3</sup>, Alexandre  
5 Termier<sup>4</sup>, Loïc Davadan<sup>5</sup>, and François Joudelat<sup>6</sup>

6 <sup>1</sup>Univ.Rennes, Inria, Irisa, France

7 <sup>2</sup>Univ.Rennes, Inria, Irisa, France

8 <sup>3</sup>ACTA, INRAE, UMR AGIR, Castanet Tolosan, France

9 <sup>4</sup>Univ.Rennes, Inria, Irisa, France

10 <sup>5</sup>IFV, UMT SEVEN, Cognac, France

11 <sup>6</sup>ITB, Paris, France

12 **1 Abstract**

13 Crop diseases and pests constitute significant causes of yield losses for crops.  
14 To limit the harm incurred by those events, farmers resort to plant protec-  
15 tion products. Such products are known to have adverse effects both on the

16 environment and on human health. Agronomists make continuous efforts  
17 to limit the usage of plant protection products to situations where those  
18 products are strictly necessary. To determine such situations, agronomists  
19 and policy-makers often rely on decision support tools to model and predict  
20 the dynamics of plant diseases. Decision support tools are based either on  
21 mechanistic models or on statistical approaches learned from large datasets of  
22 biotic (e.g., disease incidence, plant phenological stage) and abiotic (meteo-  
23 rological, soil characteristics) observations in cultures. The surge of powerful  
24 machine learning (ML) methods in the last decade makes such approaches a  
25 natural pathway to model the dynamics of plant diseases.

26 Machine learning models can reveal the factors that contribute the most to  
27 disease and pests outbreaks, provided that those models are simple enough  
28 for human inspection. Simplicity, however, may come at the price of lower  
29 prediction performances when compared to more complex models.

30 In this paper, we offer a deep look at the performance of ML models of differ-  
31 ent complexity when used on two use cases of crop disease prediction: downy  
32 mildew in the grapevine, and *Cercospora* leaf spot in the sugar beet.

33 We compare model accuracy and complexity using a year-based cross-validation  
34 approach. Our results suggest that interannual meteorological variations are  
35 a very important factor in plant disease prediction. Moreover, in line with the  
36 observations of the research community in interpretable ML, model complex-

ity stands in clear trade-off with accuracy. This makes models of intermediate complexity appealing for predicting the dynamics of crop diseases as they can provide explicit insights about the rationale of their predictions.

## 2 Introduction

Crop protection against plant diseases is crucial to secure crop yields. To this end, farmers and agronomists make use of plant protection products, i.e., pesticides, to combat plant diseases and pests in cultures. It is well-known, however, that the usage of such products has multiple downsides. Besides their impact on farmers' health, and their polluting effect on the environment, such products incur an economic cost on both farmers and consumers, not to mention their role in the development of pesticide-resistant breeds [Heap, 2014] and the indirect contamination in other stages of the food supply chain [Parsons et al., 2021]. It follows that minimizing the usage of pesticides in cultures incurs countless benefits. One way to reduce our dependence on such products is to adapt their usage to local factors [Chen, 2019] such as the climate/weather, the soil type, or the farming practices. This can be achieved through the deployment of models that can predict disease incidence or risk of outbreak. Such tools help farmers and agronomists avoid the usage of pesticides when they are not necessary.

There have been multiple efforts to model and predict the risk of outbreak

57 and the incidence of plant diseases in cultures [Edwards-Jones, 1993]. Exist-  
58 ing methods can be categorized into two families. On the one hand, mecha-  
59 nistic models are constructed based on prior observations and knowledge of  
60 the diseases or pests’ life cycles. These models require extensive agronomical  
61 studies and experts intervention, and were the preferred approach to model  
62 plant diseases for long time. An inflexion point arose with the emergence  
63 of large amounts of data including past observations of diseases in cultures  
64 – as human annotations or as images –, but also information about abiotic  
65 factors such as the characteristics of the soil and meteorological data. This  
66 data abundance has made statistical models, in particular machine learning  
67 models, more appealing in the last decade, and has nurtured their steady  
68 increase in accuracy and sophistication.

69 ML models used in crop protection are usually trained for a single type of  
70 crop and disease. This is due to the fact that different crops develop in  
71 different ways, and so do diseases and pests. It is also known that models  
72 are typically trained for a given region, and are less accurate when used on  
73 data from other regions [Velasquez-Camacho et al., 2023]. Some approaches  
74 rely on image classification with deep learning [van Klompenburg et al.,  
75 2020, Ip et al., 2018, Liakos et al., 2018, van Evert et al., 2017] for disease  
76 diagnosis. Other models are designed to predict or forecast the incidence  
77 of a disease at a particular period of the year, e.g., before harvest, based  
78 on human annotations. This forecast can take the form of an incidence

79 prediction (regression) or a risk of outbreak (classification) [Chen, 2019]. In  
80 those cases the models' outcomes help agronomists decide whether to apply  
81 or not plant protection products in their cultures.

82 But besides forecasting the incidence of crop diseases, ML models can also  
83 help agronomists understand which factors contribute to the development  
84 of those diseases. This is possible, however, if the model is simple and  
85 interpretable enough to be understood by humans. Examples of interpretable  
86 ML models are linear functions and shallow decision trees.

87 A simple, yet effective proxy to model interpretability is model complex-  
88 ity [Galárraga et al., 2021]. Complexity is usually measured as the number  
89 of relevant parameters that play a role in the model's answers, and it is  
90 known to be correlated with interpretability. To see why, it suffices to com-  
91 pare the effort of interpreting a linear model with 5 variables versus a linear  
92 model with 300 variables.

93 While complex models such as neural networks or gradient boosting tend  
94 to be less interpretable than transparent simple methods such as linear re-  
95 gression or shallow trees, in some cases this complexity pays off in terms  
96 of prediction performance [Mori and Uchihira, 2019, Johansson et al., 2011,  
97 Galárraga et al., 2021]<sup>1</sup>. This trade-off between complexity and prediction  
98 accuracy can happen because more parameters or weaker assumptions en-

---

<sup>1</sup>As shown by Rudin [2019], Bell et al. [2022], the accuracy-interpretability trade-off is not necessarily observed in every application domain and depends on multiple factors such as the quality of the data.

99    dow models with more expressiveness and flexibility to capture subtle inter-  
100    actions. Simpler models make assumptions that may **not encapsulate** the  
101    complexity of real data. For example linear models assume there exists a  
102    linear relationship between the input features and the target variable i.e.,  
103    the variable we want to model or predict. This, for instance, excludes any  
104    potential interactions between the input features as predictors for the target  
105    variable. Between complex approaches and simple models lie pattern-based  
106    models [Galárraga et al., 2021, Dong and Taslimitehrani, 2015] that strike an  
107    interesting trade-off because they remain relatively simple white boxes that  
108    exhibit higher predictive power than linear regression or decision trees.

109    Existing works that use ML methods for crop protection have paid little  
110    attention to the potential needs for interpretability and the complexity-  
111    interpretability trade-off [Fenu and Mallocci, 2021, Ip et al., 2018, van Evert  
112    et al., 2017]. We therefore contribute to the state of the art by studying  
113    this trade-off in the context of crop protection. We train different popular  
114    machine learning models of varied complexity for two typical crop protec-  
115    tion tasks: (i) disease incidence prediction, and prediction of the symptoms  
116    appearance date. We predict these target variables for the downy mildew  
117    in grapevine cultures, and for the Cercospora leaf spot in sugar beet crops,  
118    both in France. In both cases we resort to biotic (e.g., past disease inci-  
119    dences) and abiotic (e.g., meteorological data) predictors. Our tasks are  
120    classical regression problems, therefore the studied models include (i) black-

121 box ensemble methods such as random forests and gradient boosting trees;  
122 (ii) white boxes such as linear regression; and (iii) HiPaR [Galárraga et al.,  
123 2021], a pattern-based regression method of intermediate complexity. Our  
124 experiments confirm a clear complexity-accuracy trade-off in our use cases,  
125 and also show different techniques to distill agronomical insights from both  
126 white- and black-box ML models. Our results suggest that despite the dif-  
127 ference in prediction accuracy and model architecture, the models agree on  
128 some common insights. Moreover, interannual effects play a very impor-  
129 tant role, which makes very difficult to have a single model that can predict  
130 disease incidence for any arbitrary year.

131 Section 3 describes the datasets used for our study, the methods trained on  
132 those datasets as well as their performance. This is followed by a discussion of  
133 the different agronomical findings we extracted from the trained ML models  
134 in Section 4. Section 5 concludes the paper with avenues for future research  
135 in the prediction of disease incidence in cultures.

### 136 **3 Material and Methods**

137 We now describe the agronomical datasets used in our study as well as the  
138 machine learning models trained on those datasets.

## 139 3.1 Data

140 Our study case builds upon four datasets covering two major plant diseases  
141 observed in French cultures: Grape downy mildew and Sugar beet Cercospo-  
142 ria.

### 143 3.1.1 Sugar beet *Cercospora* epidemiologic data

144 Sugar beet *Cercospora* (SBC) incidences were observed in several vineyards  
145 located in France by different extension services, including the ITB (Insti-  
146 tut Technique de la Betterave). The experimental observations have been  
147 collected from 2009 to 2020 in different regions in France.

148 For each monitored site, a specific part of the area, further referred to as the  
149 “plot”, was observed throughout a specific year. Weekly visual inspections  
150 were performed on leaves covering one hundred plants in order to assess dis-  
151 ease incidence. The incidence was calculated as the proportion of sugar beet  
152 leaves displaying symptoms of *Cercospora* leaf spot (*Cercospora beticola*).  
153 Weekly inspections were conducted in each plot from leaf emergence (which  
154 happens in mid-May) until harvest (after mid-September). The collected  
155 dataset adds up to 1235 individual plots. We highlight that no plot was  
156 observed every year, and that conversely, not all plots can be monitored in  
157 a single year.

158 For each plot, we define the date of SBC onset (yearly symptoms apparitions



159 date) as the first day in which the proportion of infected leaf exceeded 10%.  
160 The end of season incidence for SBC was defined as the maximum incidence  
161 for the period going from the 25th of August to the 15th of September.

### 162 3.1.2 Grape downy mildew epidemiologic data

163 Grape downy mildew (GDM) incidence were observed in several vineyards  
164 located in France by different wine extension services including the IFV  
165 (Institut Français de la Vigne et du Vin). The data have been collected from  
166 2010 to 2017.

167 For each considered plot, an untreated row of vines was observed. Each  
168 untreated row was surrounded by two other untreated rows to ensure that  
169 they were not unintentionally sprayed with fungicides. In the monitored  
170 central row, weekly visual inspections were performed on leaves in order to  
171 measure disease incidence. The incidence was calculated as the proportion  
172 of vine leaves displaying downy mildew symptoms caused by *Plasmopara*  
173 *viticola*. Weekly inspections were conducted in each vineyard from budburst  
174 (early March) until at least bunch closing (mid-late July) or stopped when  
175 the incidence was close to 100%. The observations consist of around 9407  
176 weekly datapoints corresponding to 713 plots.

177 For each plot, date of GDM onset (yearly symptoms apparitions date) was de-  
178 fined as the first week in which the proportion of infected vines leaf exceeded  
179 1%. The end of season incidence for GDM was defined as the maximum

180 incidence for each plot.

### 181 3.1.3 Meteorological data

182 Meteorological variables were provided by the SAFRAN weather database  
183 constructed and maintained by the French national meteorological service  
184 (Météo-France). SAFRAN organizes the French territory into a grid of size  
185  $8 \times 8$  Km and stores meteorological data for each cell in the grid [Quintana-  
186 Seguí et al., 2008]. Daily observations on humidity, mean temperature, wind,  
187 amount of rainfall, and solar radiation were used to compute different mete-  
188 orological variables for both diseases.

189 For SBC, each meteorological variable covers a period of half a month (15  
190 days) from January to June. Features in the dataset follow a given conven-  
191 tion. The first part describes the temporal characteristics of the feature with  
192 the first three letters of the corresponding month, followed by an ‘A’ for the  
193 first half of a month or a ‘B’ for the second half. The second part describes  
194 the climatic nature of the feature and how this information was calculated.  
195 The feature suffixes are described in Table 1. For example, the variable  
196 named *JanA-ndRHm60* corresponds to the number of days (**nd**) such that  
197 the relative humidity was higher than 60 percent (**RHm60**) during the first  
198 half (**A**) of January (**Jan**).

Name	Feature
RHmX	Mean Relative Humidity lower than X ( $X = \{60, 65, 80, 90\}$ )
H87	Humidity index equals to 87
H87Y	Humidity index equals to 87 for at least ( $Y = \{6, 10\}$ ) hours
TmX	Mean Temperature higher than ( $X = \{15, 20\}$ )
TmXTinfYZ	Mean Temperature higher than ( $X = \{15\}$ ) but lower than ( $Y = \{10\}$ ) for at least ( $Z = \{3\}$ ) hours
TbloX	Number of days where temperatures were defined as <i>inhibiting</i> to SBC growth for more than ( $X = \{3, 6\}$ ) hours.

Table 1: Description of the meteorological variables used to model the dynamics of the Sugar beet Cercosporia (SBC). Temperatures are considered as *inhibiting* below 10°C or above 38°C

For GDM, features either describe meteorological conditions at the date of recording or its sum for the four previous weeks before recording. For example, the predictive variable ETP gives us the evapotranspiration at the time of recording. ETP-4w is the sum of evapotranspiration for the four previous weeks. Two exceptions are the number of rainy and dry days, which are counted from the beginning of January. This length of four weeks was chosen based on expert insights about the growth speed of downy mildew.

#### 3.1.4 Four prediction targets

From both diseases data and associated climatic variables, we finally obtained 4 data sets corresponding to our 4 prediction targets.

- Sugar beet Cercosporia (SBC) end of season incidence (% of leaves with diseases) with 1235 plots and 367 variables including one categorical variable and 366 numeric ones. The categorical feature is the

212        *risk-exposure*, an indicator defined by agronomists based on their own  
213        knowledge of each plot’s sensitivity to SBC. The numerical variables  
214        correspond to the one described in Subsection 3.1.3.

215        • Sugar beet Cercosporia (SBC) symptoms appearance date (day number  
216        of year) with 1235 plots and 367 variables.

217        • Grape downy mildew (GDM) end of season incidence (% of sick leaves)  
218        with 359 plots and 22 variables including two categorical and 20 nu-  
219        meric.

220        • Grape downy mildew (GDM) symptoms appearance date (week num-  
221        ber of year) with the same 359 plots and 22 variables.

222        Thus, the target variables are numerical. We are thus confronted to a re-  
223        gression problem in all cases.

## 224    3.2    Regression Methods

225        We assume that the goal is to predict the values of a real variable, that we  
226        call the *target variable*, using observations from another set of variables that  
227        we call the *predictive variables*. Examples of target variables are given in  
228        Subsection 3.1.4. Conversely, the predictive variables constitute the set of  
229        meteorological indicators (see Table 1). This scenario constitutes a classical  
230        regression problem. We first introduce some notation and then survey the  
231        most popular regression methods used in crop protection on tabular data.

232 We extend the discussion with the description of a pattern-aided regression  
233 method that deals with the complexity-accuracy trade-off introduced in pre-  
234 vious sections.

### 235 3.2.1 Problem Formulation and Notation

236 Let us assume that we count on a set of  $n$  target observations represented as  
237 a column vector  $\mathbf{y} \in \mathbb{R}^n$ . Those target observations are associated to a set  
238 of observations on the predictive variables, organized in a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ .  
239 Each row  $\mathbf{x}_i^\top \in \mathbb{R}^d$  in the matrix stores the observed values of the  $d$  predictive  
240 variables associated to a target observation  $y_i$ . From now on, we denote  
241 vectors and matrices with names in bold to distinguish them from scalars  
242 and functions. Moreover, matrices are denoted with capital letters. If a  
243 predictive variable is categorical, e.g., plant variety, we assume its values  
244 have been encoded as real numbers, for instance, by resorting to strategies  
245 such as one-hot encoding or dimensionality reduction.

246 The goal of regression analysis is to learn a function  $f$  such that  $\mathbf{y} = f(\mathbf{X}) + \boldsymbol{\epsilon}$   
247 and  $\boldsymbol{\epsilon}$  is minimal. The function  $f$  is a model of the data designed to predict  
248 the target variable for unseen instances  $\mathbf{x}^\top \in \mathbb{R}^d$  of the predictive variables.  
249 The term  $\boldsymbol{\epsilon}$  is the error of the regression model and accounts for potentially  
250 unobserved predictors of  $\mathbf{y}$ . The model  $f$  is learned on a set of training and  
251 validation observations.

### 252 3.2.2 Classical Regression Methods

253 **Linear Regression.** This method assumes that the relation between the  
254 target variable  $\mathbf{y}$  and the predictive variables  $\mathbf{X}$  is linear, that is,

$$\mathbf{y} = \boldsymbol{\beta} \mathbf{X}' + \boldsymbol{\epsilon} \text{ with } \boldsymbol{\beta} = \underset{\hat{\boldsymbol{\beta}}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}' \hat{\boldsymbol{\beta}}\|_2^2 \quad (1)$$

$$\boldsymbol{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (2)$$

255 where  $\mathbf{X}' = \mathbf{1} \oplus \mathbf{X}$ , i.e.,  $\mathbf{X}' \in \mathbb{R}^{n \times (d+1)}$  and  $\boldsymbol{\beta} \in \mathbb{R}^{d+1}$  are the parameters  
256 of the model (the operator  $\oplus$  denotes column concatenation), namely the  
257 linear coefficients associated to the each of the  $d$  predictive variables plus  
258 the intercept coefficient  $\beta_0$ . The parameters of the model can be computed  
259 by minimizing the loss function  $\mathcal{L}_l(\hat{\boldsymbol{\beta}}) = \|\mathbf{y} - \mathbf{X}' \hat{\boldsymbol{\beta}}\|_2^2$  with the method of  
260 ordinary least squares (OLS) as illustrated in Equation 2. Linear models  
261 are among the most popular regression methods due to their simplicity and  
262 interpretability. This is because the magnitude of the coefficients tells us ex-  
263 plicitly how much a predictive variable contributes to the model's prediction.  
264 On the downside, the linearity assumption may come at the expense of low  
265 prediction accuracy, which is why linear models are often used as baseline  
266 methods.

267 **Lasso.** To reduce the risk of over-fitting in linear regression, Lasso [Tibshi-  
 268 rani, 1994] proposes an L1-regularization of the loss function, which favors  
 269 models with few non-zero coefficients. This is achieved by minimizing the  
 270 following objective:

$$\beta = \operatorname{argmin}_{\hat{\beta}} \mathcal{L}_l(\hat{\beta}) + \theta \|\hat{\beta}\|_1. \quad (3)$$

271 By minimizing the L1-norm of  $\beta$  we can obtain sparse models that can not  
 272 only prevent or mitigate over-fitting, but that are less complex and therefore  
 273 easier to inspect by humans. The penalization term  $\theta$  is a hyper-parameter  
 274 that controls the importance of the sparsity constraint in the optimization  
 275 process. The Lasso method selects the set of parameters  $\hat{\beta}$  that achieves the  
 276 highest performance in cross-validation.

277 **Decision/Regression Trees.** A decision tree is a binary tree where each  
 278 internal node evaluates a Boolean condition on a predictive variable. The  
 279 children of a node are decision trees associated to an evaluation outcome,  
 280 i.e., true or false. Leaves (also called final nodes) are linked to a prediction of  
 281 the model for the target variable. When the target variable is numerical, we  
 282 talk about *regression trees* [Kramer, 1996]. Regression trees are white-box  
 283 models because the model’s prediction on a particular instance  $\mathbf{x}^\top \in \mathbb{R}^d$  can  
 284 be explained by following the path from the root to the leaf node that pre-  
 285 dicts the outcome for  $\mathbf{x}^\top$ . This makes regression trees interpretable models,

286 provided that the tree is not too deep for human inspection. Despite their  
287 interpretability, decision trees are prone to over-fitting if not properly param-  
288 eterized, and are usually outclassed in terms of predictive performance by  
289 ensemble methods such as random forests and gradient boosting trees.

290 **Random Forests.** Random Forests are ensembles of weak decision tree  
291 estimators [Breiman, 2001]. Predictions are computed by averaging the pre-  
292 dictions of each tree in the ensemble. The weak estimators are learned by  
293 applying bagging and random feature selection. In bagging, each tree is  
294 learned by sampling from  $\mathbf{X}$  and  $\mathbf{y}$  uniformly and with replacement. More-  
295 over, the trees are trained on different subsets of the features, which gives  
296 each tree a “partial” but “unique” view of the data. These techniques make  
297 random forests very robust to over-fitting, and a very popular choice for crop  
298 protection [Elavarasan et al., 2018]. On the downside, random forests are  
299 not interpretable because the aggregation step makes it very difficult to trace  
300 the outcome of the model back to the input features – without resorting to  
301 post-hoc inspection approaches as we will show later.

302 **Gradient Boosting.** Another popular ensemble method is gradient boost-  
303 ing [Mason et al., 1999]. Like random forests, the basic principle is to  
304 compute a robust prediction from the predictions of a set of weak learn-  
305 ers. Different from random forests, learning is based on an additive model  
306 where each learner  $h_m$  is fit on the error of the previous learner  $h_{m-1}$  –



307 technically on the negative gradient of the minimized loss function. Put dif-  
 308 ferently, each new learner is trained to correct the errors of the previous one:

$$309 \quad f_m(\mathbf{X}) = f_{m-1}(\mathbf{X}) + \gamma_m h_m(\mathbf{X}) \quad (4) \quad \gamma_i = \mathcal{L}(\mathbf{y}, f_m(\mathbf{X})) \quad (5)$$

310

311 The individual learners can be of any type, however decision trees are a  
 312 common choice [Mason et al., 1999]. Gradient boosting models are very  
 313 robust to over-fitting, and like random forests, behave pretty much like black  
 314 boxes.

### 315 3.2.3 Hierarchical Pattern-aided Regression (HiPaR)

316 **Pattern-aided Regression.** Pattern-based regression models consist of a  
 317 set of local models trained on regions of the data. Those regions are charac-  
 318 terized by interpretable patterns, namely logical conditions on the predictive  
 319 variables, e.g., *wind-speed* > 50. The local models are usually interpretable  
 320 functions, e.g., linear functions, that capture local relationships between the  
 321 target and the predictive variables that cannot be observed at the “global  
 322 level”. As shown in the literature [Galárraga et al., 2021, Dong and Taslim-  
 323 itehrani, 2015], these methods exhibit higher predictive performance than lin-  
 324 ear regression at the price of a manageable increase in complexity. Examples  
 325 of pattern-aided regression methods include piecewise regression [McGee  
 326 and Carleton, 1970], regression trees [Breiman, 2001], model trees [Wang

and Witten, 1997]<sup>2</sup>, Contrast pattern-aided regression (CPXR) [Dong and Taslimitehrani, 2015], and HiPaR [Galárraga et al., 2021]. We elaborate on the latter method in the following.

**HiPaR.** Hierarchical Pattern-aided Regression Galárraga et al. [2021] estimates the values of the target variable via a compact set of local hybrid rules on the predictive variables. These rules have the form:

$$p = C_1 \wedge \dots \wedge C_m \Rightarrow \mathbf{y} = f_p(\mathbf{X}_p). \quad (6)$$

In this expression, the pattern  $p$  is a conjunction of conditions on the predictive variables such as  $wind-speed > 50 \wedge humidity > 30$ . Those conditions define subsets or regions of the data  $\mathbf{X}_p \subset \mathbf{X}$ . A hybrid rule is associated to a local linear model  $f_p$  that has been trained on  $\mathbf{X}_p$ , and that refines the predictions of a global linear model  $f$  trained on  $\mathbf{X}$ . The model  $f$ , called the *default* model, is used to make predictions whenever none of the local hybrid rules applies. After having learned the default model, HiPaR mines a compact set of hybrid rules by means of two phases:

1. During the enumeration phase, the learning algorithm explores the space of patterns  $p$  in a depth-first hierarchical fashion. When a pattern  $p$  is visited, HiPaR learns a hybrid rule of the form  $p \Rightarrow \mathbf{y} = f_p(\mathbf{X}_p)$

---

<sup>2</sup>These are regression trees such that some nodes, usually the leaves, are linear models on the target variable

344 on  $\mathbf{X}_p$  – the set of observations that satisfy  $p$  –, and then explores  
 345 the sub-regions of  $\mathbf{X}_p$ . Since the search space is exponential in the  
 346 number of features, a set of pruning strategies reduces it by avoiding  
 347 the exploration of unpromising sub-regions; for example a minimum  
 348 support threshold is enforced to avoid sub-regions with very few points.

349 2. Despite the pruning strategies carried out during the enumeration  
 350 stage, the set of resulting hybrid rules can still be very large. For  
 351 this reason, HiPaR carries out a selection phase that retains a small  
 352 set of hybrid rules with good performance and minimal overlap. This  
 353 phase is governed by two hyper-parameters: the support and the over-  
 354 lap bias. They determine, respectively, to which extent very specific  
 355 rules are preferred over general rules, and how much overlap between  
 356 the selected rules is allowed.

357 Contrary to tree-based models, HiPaR’s hybrid rules are extracted from a  
 358 hierarchy with potentially overlapping regions as depicted in Figure 1. When  
 359 a new observation  $\mathbf{x}^\top$  satisfies more than one hybrid rule, the final prediction  
 360 is the weighted average of the predictions of the individual rules. The weight  
 361 is inversely proportional to the rule’s error on a validation subset. This  
 362 makes HiPaR models more robust than linear functions and regression trees,  
 363 but significantly more complex. That said, HiPaR hybrid rules remain white-  
 364 box models that allow for simple inspection of the most important predictive  
 365 variables in the prediction for an observation  $\mathbf{x}^\top \in \mathbb{R}^d$ .

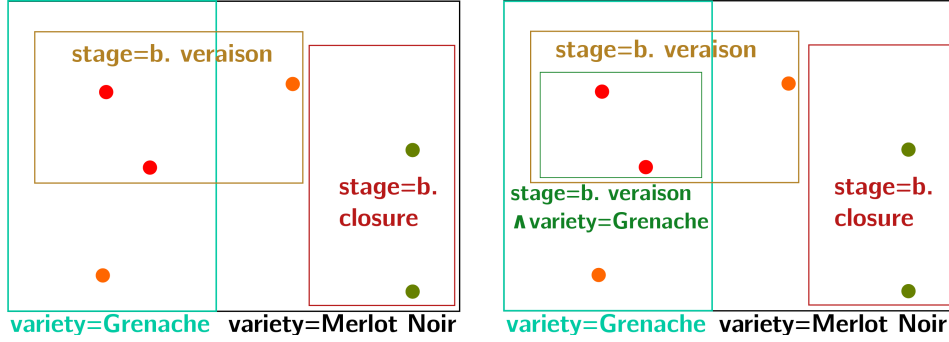


Figure 1: A depiction of the regions explored by HiPaR for two steps of the enumeration phase. Each rectangle defines a region described by a pattern, on which HiPaR learns a local regression model. Regions can overlap; an example are the regions *stage* = “*b. veraison*” and *variety* = “*Grenache*”. Once a region is explored, e.g., *stage* = “*b. veraison*”, HiPaR will look at its sub-regions in a depth-first-search manner (figure on the right).

366 Table 2 summarizes the strengths and weakness of the methods discussed in  
367 this section.

Models	Characteristics	Advantages	Disadvantages
Lasso	Sparse linear regression	Simple/interpretable	Baseline method
HiPaR	Pattern-based	Medium-complexity	High computation time
Random Forests	Ensemble-, tree-based	High accuracy, Built-in feature importance values	Black-box model
Gradient Boosting	Ensemble-based	High accuracy	Black-box model

Table 2: Overview of the machine learning methods used in this study.

### 368 3.3 Training and testing procedures

#### 369 3.3.1 Optimization and performance evaluation

370 One of the challenges of evaluating different machine learning models is to  
371 select the best configuration so that comparisons are fair and meaningful. In

an agronomical scenario the important interannual differences make standard cross-validation unadapted. Therefore, we use cross-validation by year, that is, each year is used as a fold in the process. The data from a given year is separated from the rest of the dataset for testing, whereas the observations from remaining years are used to train the algorithm. That way we are able to estimate the actual capacity of the algorithms to predict for unseen scenarios, e.g., for a new year.

Inside each fold, we select the best model by optimizing the hyper-parameters of each method. HiPaR’s enumeration phase can take long for very low support thresholds. Therefore we run the enumeration phase with a support threshold of 30% the size of the dataset once – i.e., regions covering fewer points are not explored –, and we then optimize the hyper-parameters of the selection phase to pick the most performing set of rules.

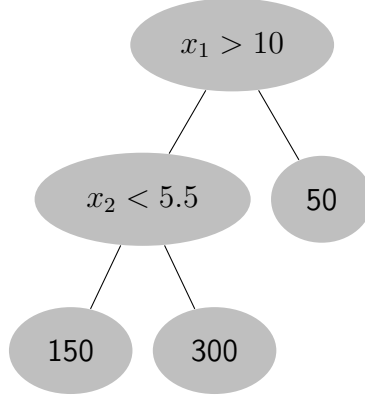
We use the coefficient of determination ( $R^2$ ) as prediction performance metric. The  $R^2$  score is defined as the proportion of the variance in the predicted target variable explained by the independent variables. Contrary to the root mean square error (RMSE),  $R^2$  values can be compared among different prediction tasks (e.g., disease incidence and symptoms appearance). Indeed, the closer to 1 the  $R^2$  is, the better the model fits the data. Values close to zero denote a performance comparable to predicting the mean of the target variable, whereas negative scores mean the model is worse than a mean-based simple predictor.

### 394 3.3.2 Complexity measure

395 To measure the complexity of the studied machine learning models, we resort  
396 to the complexity measure for pattern-based models proposed by Galárraga  
397 et al. [2021] that counts the number of elements in the model. An element  
398 is either a non-zero coefficient or a condition on a predicting variable. We  
399 remark that this measure is also applicable to tree-based methods such as  
400 random forests or gradient boosting trees because each node of each tree of  
401 the ensemble defines either a condition on one attribute or a linear model  
402 – for simple regression trees this linear model is a single constant. The  
403 number of elements can be very large when the ensemble consists of many  
404 trees, which points out the complexity of such models.

405 Under this principle, a Lasso model is generally less complex than a HiPaR  
406 model with several rules. This is the case because for Lasso we only need to  
407 count the non-zero coefficients in the linear function, whereas for HiPaR we  
408 must consider both the number of conditions and the coefficients of each of  
409 the local models.

410 If we consider the following regression tree  $T$  :



411

412 Then its complexity  $c(T)$  is 5. Likewise, if we consider the rule  $R$  :

$$C_1 \wedge C_2 \wedge C_3 \wedge C_4 \Rightarrow \mathbf{y} = 3x_1 + 4x_2 - 4x_3 + 8. \quad (7)$$

413 then its complexity  $c(R)$  is 8 because the rule consists of 4 conditions and 4  
 414 linear coefficients.

## 415 3.4 Results

### 416 3.4.1 Performance-Complexity Trade-off

417 Figure 2 depicts the trade-off between the complexity and accuracy of the  
 418 studied machine learning methods. On the x-axis we show the complexity of  
 419 the models (in log scale). The y-axis corresponds to the median  $R^2$  coefficient  
 420 of each model in cross-validation. Models located in the top-left part of  
 421 the space strike a better accuracy-complexity trade-off as they predict the  
 422 data more accurately with fewer elements. As suggested by Galárraga et al.  
 423 [2021], more complex models such as random forests or gradient boosting

424 trees achieve the best performance at the price of high complexity. Lasso  
425 regression, our baseline, is often the least accurate model. HiPaR positions  
426 itself in between linear regression and ensemble methods striking a very  
427 interesting trade-off for 3 of the 4 prediction tasks.

428 We highlight that accuracy varies drastically across tasks: All models strug-  
429 gle when it comes to predicting the date of apparition of downy mildew in  
430 vine cultures, as the median  $R^2$  for all methods is negative (bottom-right fig-  
431 ure). We observe  $R^2$  scores between 0.12 and 0.26 for the final downy-mildew  
432 incidence (on the bottom-left) with gradient boosting as the winner. HiPaR  
433 lies close to Lasso, which means that it did not find many regression rules  
434 improving performances marginally over the baseline. The performance dif-  
435 ferent between the two target variables in the downy-mildew dataset could  
436 be explained by the relatively low number of observations for the date of  
437 symptoms apparition – 359 versus 700 observations for the end-of-season  
438 incidence.

439 The reach of the aggregated variables is relatively limited too. By this we  
440 mean most of these variables over a range of 4 weeks before data collec-  
441 tion. While this confirms the trade-off, the low  $R^2$  makes this dataset less  
442 interesting to study further.

443 The results for the sugar beet *Cercospora* are more encouraging. The  $R^2$   
444 median scores for the apparition date vary between 0.13 and 0.18 with gra-



445 dient boosting leading the rank and followed by HiPaR (top-right figure).  
446 For the prediction of the end-of-season incidence performance ranges from  
447 0.05 to 0.35. In this use case HiPaR outperforms all methods and finds a  
448 large number of rules that improve performance significantly when compared  
449 to a single linear model, and without incurring as much complexity as the  
450 ensemble methods.

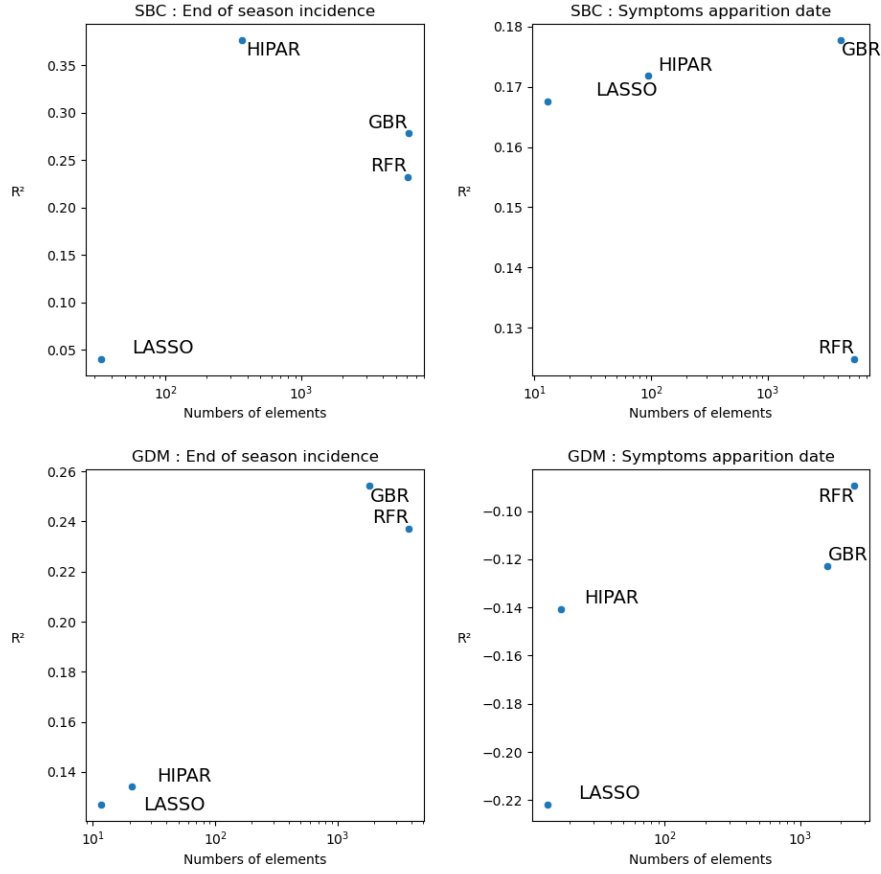
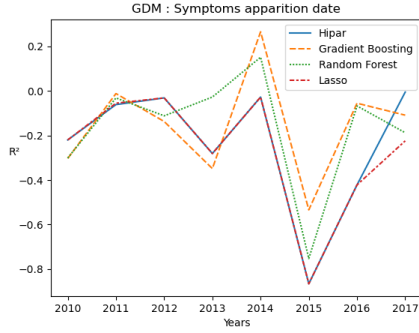
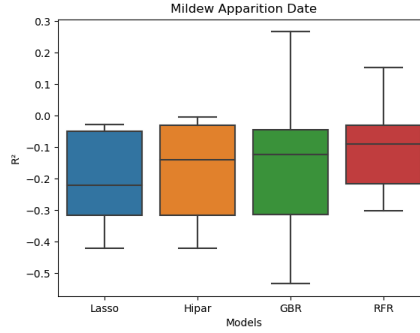


Figure 2:  $R^2$  of different machine learning models compared against their complexity. The x-axis correspond to the number of elements that compose each model (log scale). The y-axis are the median  $R^2$  values in cross-validation. GBR stands for gradient boosting regression, and RFR for random forests regression.

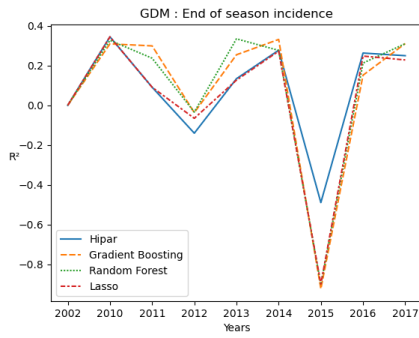


(a)  $R^2$  year-to-year cross-validation

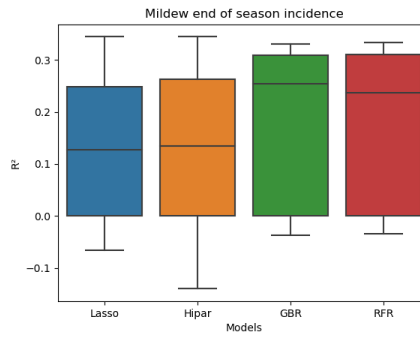


(b)  $R^2$  year-to-year distribution

Figure 3: Mildew symptoms apparition date

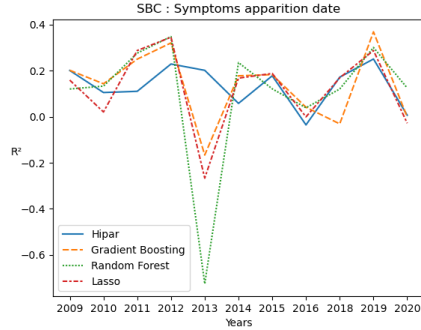


(a)  $R^2$  year-to-year cross-validation

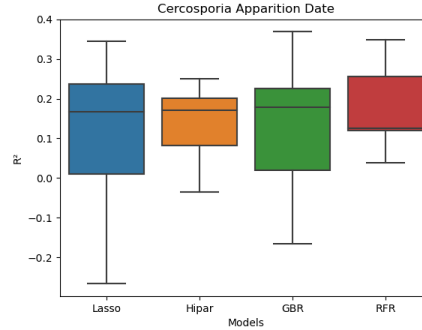


(b)  $R^2$  year-to-year distribution

Figure 4: Mildew end of season incidence

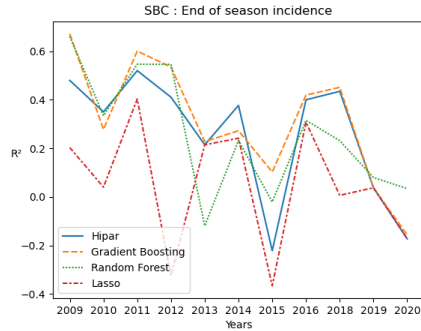


(a)  $R^2$  year-to-year cross-validation

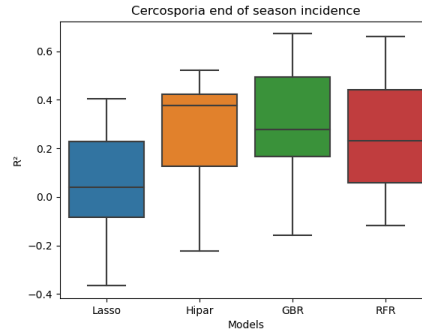


(b)  $R^2$  year-to-year cross-validation distribution

Figure 5: Cercosporia symptoms date of apparition



(a)  $R^2$  year-to-year cross-validation



(b)  $R^2$  year-to-year cross-validation distribution

Figure 6: Cercosporia end of season incidence

451 When we look at the performance of the methods per year (Figures fig. 2),  
 452 we notice that performance can vary drastically from one year to another,  
 453 and that both end-of-season incidence and date of apparition are very hard  
 454 to model for some years. This is true for all methods. As a general trend,  
 455 we can observe that Cercospora end of season incidence predictions seem to  
 456 follow a downward trend in performance. The performance variability across

457 folds (Figures fig. 2) for the different methods is comparable and does not  
458 seem to follow a noticeable pattern.

459 Now that we have illustrated the accuracy-complexity trade-off present in our  
460 use cases, we delve into the knowledge captured by the different methods.  
461 To do so we analyze the models trained to predict year 2009 for the end-  
462 of-season incidence of the sugar beet *Cercospora*, as these models exhibit  
463 the highest explained variance across all years ( $R^2$  scores of 0.67 and 0.66  
464 for gradient boosting trees and random forests, 0.3 for Lasso, and 0.47 for  
465 HiPaR). For white-box models such as Lasso and HiPaR, we conduct direct  
466 inspection of the models' elements. For the complex black-box approaches,  
467 we resort to classical model inspection techniques and assess whether our  
468 models agree on the relationships between the predictive variables and the  
469 target variables.

### 470 **3.4.2 Use case: Incidence of the Sugar Beet *Cercospora***

471 In this section we carry out an inspection phase aimed to distill agronomical  
472 insights from the experimental machine learning models trained to predict  
473 the incidence of sugar beet *Cercospora*. These models were trained on all  
474 years except 2009 and correspond to the most performing cross-validation  
475 round of our experiments. We resort to classical interpretation techniques  
476 including feature importance rankings, partial dependence plots, and simple  
477 rule inspection. The first technique tells us which are the most important

variables that play a role in the prediction. PDPs and rule inspection allow us to identify *threshold effects* on the predicting variables, that is, cases when the behavior of the target variable varies in a piece-wise manner, i.e., according to thresholds on the predicting variables. Pattern-based regression methods such as HiPaR are good at detecting such kind of effects. Moreover, such methods allow us to study more fine-grained interactions among the predicting variables present in the rules. Our observations set the ground for the discussion in Section 4.

**Feature importance.** A simple way to interpret the knowledge captured by a machine learning model is to construct a feature-importance ranking that tell us how much the model’s input variables affect the model’s output. This ranking can be based on the actual contributions of a variable to the answers of a model, e.g., the coefficients of a linear regression, or on model-aware scores computed a posteriori for black-box models. In this spirit we contrast the feature-importance rankings of Lasso, RFR and GBR and depicted them in Figure 7. Lasso’s linear coefficients encode the actual contributions of the input features to the answers of the model. They are therefore signed. To turn the linear coefficients into importance scores, we take their absolute value. Conversely, RFR and GBR are based on tree ensembles for which different importance scores have been developed. We choose the permutation feature importance method as implemented in the

scikit-learn library. This approach estimates the importance of a feature by shuffling its values across rows in  $\mathbf{X}$ . The resulting decrease in accuracy is then used to determine how much the model relies on a feature to make predictions – the higher the decrease, the informative the feature is for predicting the target variable.

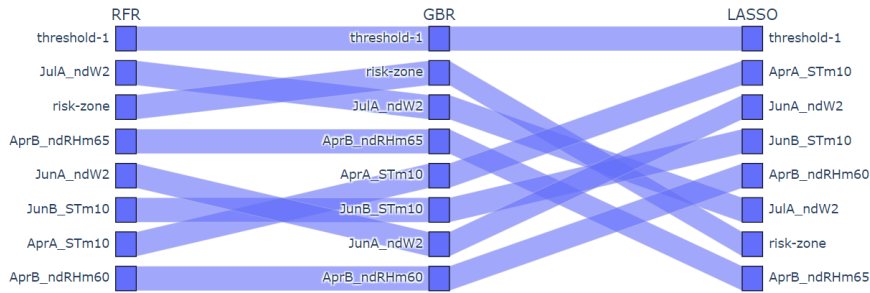


Figure 7: Parallel coordinates chart comparing feature-importance rankings for Lasso, random forests, and gradient boosting trees when predicting end-of-season incidence for the sugar beet *Cercospora*. For each model we chose the top-4 most important features. For each model, features below the 4 can be further down the importance order than what is displayed

As we can see, RFR and GBR yield very similar rankings – their top-4 variables are the same even though the order is not identical. The variable *threshold-1* is the most important feature for all three models. This variable represents the day in which the first symptoms of *Cercospora* were detected in the culture. Conversely RFR and GBR’s accuracy rely on the *risk-zone* expert-based indicator, which is less important for linear regression. While importance scores tells us which information the model is looking at, it does not tell whether those features tend to increase or decrease the model’s incidence prediction. We can, however, obtain this information by looking

513 at the linear coefficients learned by Lasso.

514 Table 3 shows the top-5 most important linear coefficients:

Variable	Coefficient
Threshold-1	-41
AprA-STm10	27.64
JunA-ndW2	-23.39
JunB-STm10	22.24
AprB-ndRHm60	14.73

Table 3: Top-5 important linear coefficients learned by Lasso

515 We remind the reader the meaning of these variables:

516 **Threshold-1** : The symptoms apparition date

517 **AprA-STm10** : The sum of the daily average temperatures of the days  
518 above 10°C during the first half of April.

519 **JunB-ndW2** : The number of days in the first half of June such that the  
520 average wind speed was higher than 2 m/s.

521 **JulB-STm10** : The sum of the daily average temperatures above 10°C  
522 during the second half of June.

523 **AprB-ndRHm60** : The number of days in the second half of April such  
524 that the relative humidity is higher than 60%.

525 Table 3 tells us that the later symptoms appear, the lower the final incidence  
526 tends to be. The predicted incidence tends to increase when temperature and  
527 humidity in June and April increase, whereas faster winds seem to hinder



528 the development of Cercospora. This results must be taken with a grain  
529 of salt given the fact that our baseline Lasso model can explain only 30%  
530 of the target variable’s variance. That said, these variables are used by  
531 more accurate models such as RFR and GBR, which means that we are not  
532 uncorrelated with the target variable.

533 **Threshold Effects.** As stated before, pattern-based regression methods  
534 are constructed to detect predicting variable threshold effects on the target  
535 variable. In HiPaR such effects are explicitly stated in the rule conditions.  
536 To observe whether our models captured such effects we have a deep look  
537 at the hybrid rules learned by HiPaR on our studied use case, and con-  
538 trast those thresholds to those learned by the more complex models, namely  
539 RFR and GBR. Since those models are actually based very large ensembles  
540 of threshold-based estimators, we observe those threshold effects by means  
541 of partial dependence plots (PDP). This widely-used inspection technique  
542 allows us to visualize the behavior of a model’s prediction (y-axis) for the  
543 different values of a predicting variable (x-axis).

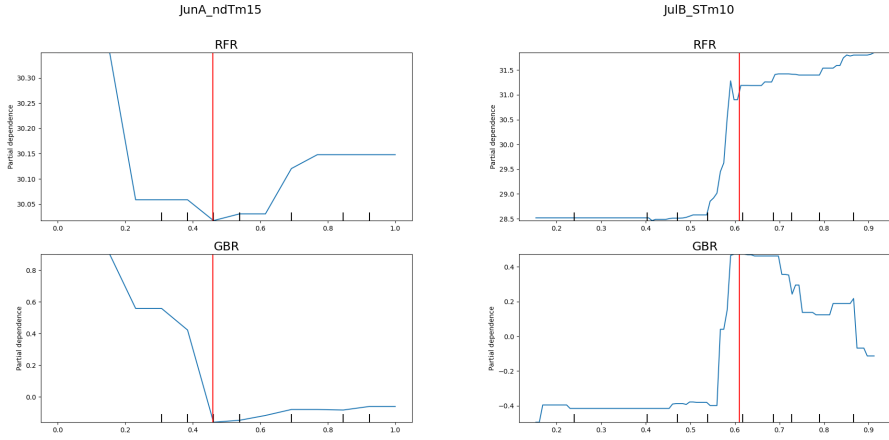
544 In our use case, HiPaR learned 3 hybrid rules whose conditions are listed  
545 in Table 4. As displayed before, thresholds (in red) used in HiPaR’s rules  
546 roughly fits with changes in the PDPs behaviour. While they not the most  
547 important features as seen before, it seems to indicate that these thresholds  
548 are not insignificant (according to RFR and GBR models). We suppose that

549 these features might act as proxies for other features, or simply have an indi-  
550 rect influence on the final result that is not detectable by using PDPs.

Rule 1	$JunA-ndTm15 < 8, risk-zone=false$
Rule 2	$JulB-STm10 < 324$
Rule 3	$MarB-ndW2 \geq 4, JunB-ndTm15 < 13$

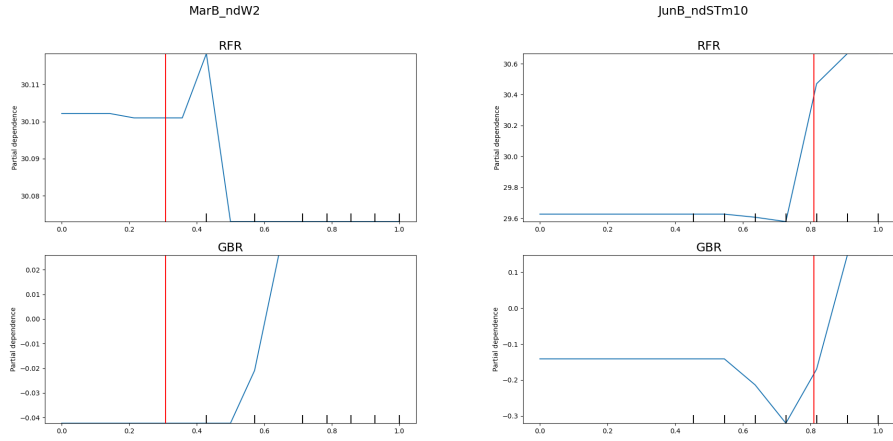
Table 4: Conditions of the hybrid rules learned by HiPaR when predicting the end-of-season incidence of the sugar beet *Cercospora*.

551 In other words, HiPaR detected different linear behaviors based on whether  
552 a plot lies or not within a region deemed risky by experts (*risk-zone*), or  
553 whether the number of hot days in June and July are below certain thresholds  
554 ( $JunA-ndTm15, JulB-STm10, JunB-ndTm15$ ), or whether the second half of  
555 March was windy ( $MarB-ndW2$ ). We now construct PDPs for the numerical  
556 variables on RFR and GBR, which we depict in Figure 8.



(a) Number of days where average temperatures were higher than  $15^\circ$  during the first half of July

(b) Number of hours where average temperatures were higher than  $10^\circ$  during the second half of June



(c) Number of days where average wind speed were higher than 2km/h during the 2nd half of March

(d) Number of days where average temperatures were higher than  $10^\circ$  during the 2nd half of June

Figure 8: Partial Dependence Plots for the predicting variables  $MarB\_ndW2$ ,  $JulB\_STm10$ ,  $JunA\_ndTm15$  and  $JunB\_STm10$  on random forests and gradient boosting trees. The red line — represents a threshold learned by HiPaR.

557 **Feature Interactions.** Each of the conditions listed in Table 4 is associ-  
 558 ated to a local linear model (learned using Lasso). Those models reveal local

interactions between the variables in the conditions and the linear coefficients, and are designed to refine the baseline linear (called also the default) model learned on the entire dataset. Out of 368 features used as input in the models, Lasso selects between 25 and 55. This represents between 6.7% and 15% of the available features. Moreover, local models are systematically less complex than the default model as Table 5 shows.

	Rule 1	Rule 2	Rule 3	Default Model
Rule 1	25	8	3	6
Rule 2		28	6	16
Rule 3			26	12
Default Model				55

Table 5: Number of common non-zero coefficients of the linear models learned by HiPaR for the prediction of the end-of-season sugar beet Cercospora.

We can also observe that coefficients overlap between the different hybrid rules is low. This means that each local model is relying on different signals to make predictions on the end-of-season incidence. Figure 9 depicts the intensity and polarity of 16 of those coefficients for both the default and local models.

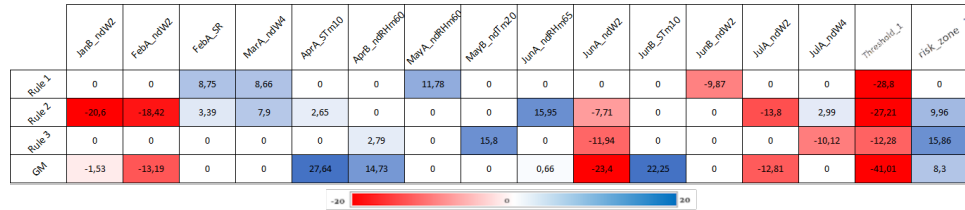


Figure 9: A color encoding for the linear coefficients of the three hybrid rules learned by HiPaR. Cells in white  $\square$  denote features with a linear coefficient strictly equal to 0, which means those features aren't used by the model

Our first observation is that the apparition date (*threshold-1*) is consistently important across all models – and always correlated negatively with the predict incidence. The features *risk-zone* and *JunA-ndW2* (wind speed in the first half of June) are used in all models except the first rule because these variables appear in the conditional part of this rule (Table 4).

This rule can be interpreted as follows: Plots with lower disease exposure (*risk-zone=false*) and lower temperatures in the first half of June (*JunA-ndTm15* < 8), experience an aggravated development of Cercosporia as humidity in May (*MayA-ndRHm60*), wind speed in March (*MarA-ndW4*), and rainfall in February (*FebA-SR*) increase. Wind during June (*JunB-ndW2*) is associated to a slow down of the disease.

The second rule suggests that lower temperatures in July (*JulB-STm10* < 324) make Cercosporia sensitive to wind in January, February, June, and July (*JanB-ndW2*, *FebA-ndW2*, *JunA-ndW2*, *JulA-ndW2*). Conversely, a wet June (*JunA-ndRHm65*) or a windy March (*MarA-ndW4*) appear as aggravating factors. A windy July (*JulA-ndW4*), a rainy February (*FebA-SR*) and a hotter April (*AprA-STm10*) have a mitigated effect on the development of Cercosporia.

The third rule triggers when the month of March is windy (*MarB-ndW4* ≥ 4) June is not very hot (*JunB-ndTm15* < 13). In that case, higher temperatures in May (*MayB-ndTm20*) and a wet April (*AprB-ndRHm60*) are

591 correlated with the growth of *Cercosporia* growth. Conversely, wind in June  
592 (*JunA-ndW2*) and July (*JulA-ndW4*) exhibit a negative correlation with  
593 growth.

594 Finally, we observe that the default model combines signals from all the lo-  
595 cal rules, even though it does not always relies on the same variables. This  
596 happens because the learning objective of this model must fit the observa-  
597 tions from all the sub-regions. This translates into selecting variables (such  
598 as *JunB-STm10*) that explain incidence for all the observations, i.e., at the  
599 global level, but that have little to no explanation power when limited to  
600 subsets of the data such as the observations on regions not deemed risky by  
601 the experts (*risk-zone=false*).

## 602 4 Discussions

603 We structure our discussion along three axes: (a) the complexity-accuracy  
604 trade-off discussed in Subsection 3.4.1, (b) the implications of complexity  
605 in interpretability, and (c) the agronomical insights offered by the models  
606 trained.

607 **Complexity and Accuracy.** Our results go in line with what has been ob-  
608 served in other works on model complexity [Dong and Taslimitehrani, 2015,  
609 Galárraga et al., 2021], that is, the tendency of complex models to outper-  
610 form simple models in terms of prediction accuracy. It is crucial to highlight

though, that such a trend holds under the assumption that the models have been properly parameterized and trained. For instance, a complex model trained on very little data will surely over-fit that data specially if there are as many or more parameters than data points. Conversely if the data adheres to the learning hypothesis of a simple model, e.g., linearity, such model will surely shine in terms of performance regardless of its complexity. Finally, even if a model was trained under a reasonable learning hypothesis, testing it on data that diverges from the training distribution will result in unsatisfactory prediction performance. We can observe such a phenomenon for the models tested on years 2013 and 2015 for the prediction of both the incidence and the date of apparition in both cultures. The observations collected those years are atypical because some of the predicting variables exhibited measures outside the amplitudes observed other years. This translated into a clear under-fitting with the lowest  $R^2$  scores registered in our experiments.

**Interpretability.** It is widely-assumed that interpretability and model complexity are positively correlated. An illustration of such phenomenon can be observed from our use case. Both linear and pattern-based model allowed us to distill insights easily and directly from the structures of the models themselves. For more complex models such as random forests and gradient boosting trees we had to resort to external inspection tools such as

the permutation-based accuracy decrease and the partial dependence plots (PDPs). Albeit effective, those techniques have limitations. Importance scores do not tell us if a variable is positively or negatively correlated with the prediction of the model. PDPs can be applied to up to two variables at the same time, and make independence assumptions that often do not hold on the data. This happens because each point in the curve is the result of averaging the model answers over all possible values of the remaining predicting variables. Since some combinations of values may be unlikely, PDPs must be taken with a grain of salt, specially when the predicting variables exhibit some correlation. That said, the PDPs for RFR and GBR in our experiments were in concordance with the threshold effects observed when using HiPaR. It should be noted that while RFR, GBR, and HiPaR resort to thresholds on the predicting variables, the fact they all outperform Lasso significantly suggest that threshold effects are a reasonable hypothesis for the prediction of plant diseases based on meteorological data.

**Agronomical insights.** Based on our use case study on the sugar beet *Cercospora*, we observe that aggregating the meteorological indicators according to the seasons, i.e., winter, spring, and summer can effectively explain some of the variation in disease incidences.

Winter defines the initial conditions: This is the period in which the primary inoculum of *Cercospora* lies in the soil in the form of spores. Spring defines



653 the development period for both crops and the Cercospora. Finally, summer  
654 encompasses the end of the season, and the moment in which the disease's  
655 symptoms, as well as its effects, are obvious.

656 As a general rule, dry summers seem to hinder the growth of Cercospora.  
657 This follows from the importance assigned by the models to the wind and  
658 temperature factors during June and July. Dry winters also seem to mitigate  
659 the disease's spread. Conversely, a hot and humid spring stands as the  
660 main aggravating factor in Cercospora's incidence. Thanks to the hybrid  
661 rules provided by HiPaR, we can obtain more nuanced relationships between  
662 the incidence and the predicting variables. Rule 2 in Table 4 tells us that  
663 a mild month of July should make us focus the attention on the initial  
664 conditions (winter), in particular the wind and the sun exposure and the  
665 temperature – the two latter factors contributing positively to the presence of  
666 the primary inoculum. Moreover, a windy spring with mild temperatures in  
667 June should target our attention towards the development phase (spring) in  
668 particular towards temperature and humidity, which are positively correlated  
669 with incidence. In all cases, the date of apparition is the best predictor of the  
670 final incidence, which means that early detection is the best weapon against  
671 Cercospora.

672 We could not draw insights from the prediction of the downy-mildew on the  
673 vine because the transparent models explain no more than 14% of the ob-  
674 served variance for the incidence – the results for the date of apparition are

worse. We think this performance gap is due to the fact that the dataset relies only on meteorological measures for the four weeks that precede the end of the season. In other words, this dataset lack the richness of the meteorological signals available for the sugar beet Cercospora dataset. This observation confirms the importance of accurate and complete meteorological measurements when modeling the dynamics of plant diseases. We also believe that the studying the impact of the granularity of the meteorological indicators in such tasks remains an interesting research avenue.

## 5 Conclusions

In this paper, we have shown the interest of exploring the complexity trade-off for machine learning models when applied to predicting the incidence of plant diseases. It is accepted that in some applications, complex models such as neural networks or gradient boosting generally perform better than simpler ones such as linear regression. This comes, however, at the cost of interpretability, which (a) is vital when we need to draw insights from the prediction model, and (b) fosters transparency, which can in turn favor acceptability by users. Post-hoc explanation methods can help us extract insights from accurate black-box models, but they are not the only solution as we have shown in this work: medium-complexity models based on pattern-aided regression can achieve competitive prediction performance while remaining simple and interpretable. Moreover, our experiments with

696 post-hoc explainability techniques such as partial dependence plots suggest  
697 that pattern-aided regression can reveal threshold effects that are also ex-  
698 ploited by the more accurate black-box ensemble methods. Using those mod-  
699 els, we have also shown that medium-complexity methods are well suited to  
700 extract more pertinent information compared to simpler models. Likewise  
701 medium-complexity models are easier to interpret compared to more com-  
702 plex methods. This shows the utility of pattern-aided regression and makes  
703 it appealing for crop prediction. Since there is a direct correlation between  
704 interpretability and acceptability, evaluating the complexity of a model is  
705 not trivial and should be taken into account. This aspect has been already  
706 addressed from the angle of learning complexity [Kearns, 1990] or from the  
707 perspective of data complexity [Dwivedi et al., 2020], but rarely in terms  
708 of the complexity of the resulting model. Finally, our study suggests that  
709 the meteorological inter-annual variations make disease incidence prediction  
710 very challenging, and that predicting disease incidence for any year requires  
711 more research as well as more historical (quality) data.

712 In the future we envision to study whether increasing the temporal and  
713 spatial granularity of the meteorological attributes can help us improve the  
714 quality of our predictions. An interesting research avenue could be to apply  
715 representation learning techniques in order to learn novel and useful mete-  
716 orological indicators that predict disease incidence more accurately. Given  
717 the inter-annual variations of weather patterns, future approaches should be

718 able to categorize prediction models based on the meteorological profile of the  
719 data used to train them. We believe that unsupervised learning techniques  
720 could be adapted in that regard. Such approaches may be even necessary in  
721 the light of a climate that will keep changing in the upcoming years.

## 722 6 Acknowledgements

723 We thank the French applied agricultural research organization for sugar  
724 beet (ITB - Institut Technique de la Betterave) (ITB) and French wine and  
725 vine Institut (IFV – institut francais de la Vigne et du Vin) for providing  
726 the meteorological and agronomical data used in this study.

727 We also thank all the experts from each institutes who helped us in inter-  
728 preting our results with their insight : Fabienne Maupas, Ghislain Malatesta,  
729 Gouwie Céline (ITB) and Marc Raynal, Christian Debord, Xavier Burgun,  
730 Marc Vergnes (IFV). We thank Lucile Vallet (Acta) for her work and the  
731 preparation of the sugar beet dataset.

732 This work was funded by the DigitAg institute (ANR-16-CONV-0004) and  
733 the RegEpi project (ECOPHYTO R&D program, French Biodiversity Agency  
734 – OFB). This project data was also part of the network data science and mod-  
735 eling methods for agriculture and agri-food sector ([www.modelia.org](http://www.modelia.org), funded  
736 by CASDAR grants of the French ministry of agriculture).

## 737 References

- 738 Ian Heap. Global perspective of herbicide-resistant weeds. *Pest Management*  
739 *Science*, 70(9):1306–1315, 2014. ISSN 15264998. doi: 10.1002/ps.3696.
- 740 Paul Parsons, Elaine Freeman, Ryan Weidling, Gary L. Williams,  
741 Philip Gill, and Neil Byron. Using existing knowledge for the  
742 risk evaluation of crop protection products in order to guide ex-  
743 posure driven data generation strategies and minimise unneces-  
744 sary animal testing. *Critical Reviews in Toxicology*, 51(7):600–621,  
745 2021. ISSN 15476898. doi: 10.1080/10408444.2021.1987384. URL  
746 <https://doi.org/10.1080/10408444.2021.1987384>.
- 747 Mathilde Chen. *Analyse du risque de mildiou de la vigne dans le Bordelais à*  
748 *partir de données régionales et d’informations locales collectées en cours*  
749 *de saison*. PhD thesis, Université Paris-Saclay (ComUE), 2019.
- 750 Gareth Edwards-Jones. Knowledge-based systems for crop protection: the-  
751 ory and practice. *Crop Protection*, 12(8):565–578, 1993. ISSN 02612194.  
752 doi: 10.1016/0261-2194(93)90119-4.
- 753 Luisa Velasquez-Camacho, Marta Otero, Boris Basile, Josep Pijuan, and  
754 Giandomenico Corrado. Current Trends and Perspectives on Predictive  
755 Models for Mildew Diseases in Vineyards. *Microorganisms*, 11(1):1–19,  
756 2023. ISSN 20762607. doi: 10.3390/microorganisms11010073.
- 757 Thomas van Klompenburg, Ayalew Kassahun, and Cagatay Catal. Crop

758 yield prediction using machine learning: A systematic literature re-  
 759 view. *Computers and Electronics in Agriculture*, 177(January):105709,  
 760 2020. ISSN 01681699. doi: 10.1016/j.compag.2020.105709. URL  
 761 <https://doi.org/10.1016/j.compag.2020.105709>.

762 Ryan H.L. Ip, Li Minn Ang, Kah Phooi Seng, J. C. Broster, and  
 763 J. E. Pratley. Big data and machine learning for crop protection.  
 764 *Computers and Electronics in Agriculture*, 151(November 2017):376–  
 765 383, 2018. ISSN 01681699. doi: 10.1016/j.compag.2018.06.008. URL  
 766 <https://doi.org/10.1016/j.compag.2018.06.008>.

767 Konstantinos G. Liakos, Patrizia Busato, Dimitrios Moshou, Simon Pearson,  
 768 and Dionysis Bochtis. Machine learning in agriculture: A review. *Sensors*  
 769 *(Switzerland)*, 18(8):1–29, 2018. ISSN 14248220. doi: 10.3390/s18082674.

770 F. K. van Evert, S. Fountas, D. Jakovetic, V. Crnojevic, I. Travlos, and  
 771 C. Kempenaar. Big Data for weed control and crop protection. *Weed*  
 772 *Research*, 57(4):218–233, 2017. ISSN 13653180. doi: 10.1111/wre.12255.

773 Luis Galárraga, Olivier Pelgrin, and Alexandre Termier. HiPaR: Hierarchical  
 774 Pattern-Aided Regression. In *Advances in Knowledge Discovery and Data*  
 775 *Mining*, pages 320–332, Cham, 2021. Springer International Publishing.  
 776 ISBN 978-3-030-75762-5.

777 Toshiki Mori and Naoshi Uchihira. Balancing the Trade-off between Accu-  
 778 racy and Interpretability in Software Defect Prediction. *Empirical Softw.*

779 *Engg.*, 24(2):779–825, apr 2019. ISSN 1382-3256. doi: 10.1007/s10664-  
780 018-9638-1. URL <https://doi.org/10.1007/s10664-018-9638-1>.

781 Ulf Johansson, Cecilia Sönströd, Ulf Norinder, and Henrik Boström. Trade-  
782 off between accuracy and interpretability for predictive in silico modeling.  
783 *Future Medicinal Chemistry*, 3(6):647–663, 2011. doi: 10.4155/fmc.11.23.  
784 URL <https://doi.org/10.4155/fmc.11.23>. PMID: 21554073.

785 Cynthia Rudin. Stop Explaining Black Box Machine Learning Models for  
786 High Stakes Decisions and Use Interpretable Models Instead. *Nature Ma-*  
787 *chine Intelligence*, 1:206–215, 05 2019. doi: 10.1038/s42256-019-0048-x.

788 Andrew Bell, Ian Solano-Kamaiko, Oded Nov, and Julia Stoyanovich. It’s  
789 just not that simple: An empirical study of the accuracy-explainability  
790 trade-off in machine learning for public policy. In *Proceedings of the 2022*  
791 *ACM Conference on Fairness, Accountability, and Transparency*, FAccT  
792 ’22, page 248–266, New York, NY, USA, 2022. Association for Computing  
793 Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533090. URL  
794 <https://doi.org/10.1145/3531146.3533090>.

795 G. Dong and V. Taslimitehrani. Pattern-Aided Regression Modeling and  
796 Prediction Model Analysis. *IEEE Transactions on Knowledge and Data*  
797 *Engineering*, 27(9):2452–2465, 2015.

798 Gianni Fenu and Francesca Maridina Mallocci. Review forecasting plant  
799 and crop disease: An explorative study on current algorithms. *Big*

800 *Data and Cognitive Computing*, 5(1):1–24, 2021. ISSN 25042289. doi:  
801 10.3390/bdcc5010002.

802 P. Quintana-Seguí, P. Le Moigne, Y. Durand, E. Martin, F. Habets,  
803 M. Baillon, C. Canellas, L. Franchisteguy, and S. Morel. Analysis of  
804 near-surface atmospheric variables: Validation of the safran analysis  
805 over france. *Journal of Applied Meteorology and Climatology*, 47(1):  
806 92 – 107, 2008. doi: <https://doi.org/10.1175/2007JAMC1636.1>. URL  
807 <https://journals.ametsoc.org/view/journals/apme/47/1/2007jamc1636.1.xml>.

808 Robert Tibshirani. Regression Shrinkage and Selection Via the Lasso. *Jour-*  
809 *nal of the Royal Statistical Society, Series B*, 58:267–288, 1994.

810 Stefan Kramer. Structural regression trees. In *AAAI/IAAI, Vol. 1*, pages  
811 812–819. Citeseer, 1996.

812 Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.

813 Dhivya Elavarasan, Durai Raj Vincent, Vishal Sharma, Albert Y.  
814 Zomaya, and Kathiravan Srinivasan. Forecasting yield by inte-  
815 grating agrarian factors and machine learning models: A sur-  
816 vey. *Computers and Electronics in Agriculture*, 155(October):257–282,  
817 2018. ISSN 01681699. doi: 10.1016/j.compag.2018.10.024. URL  
818 <https://doi.org/10.1016/j.compag.2018.10.024>.

819 Llew Mason, Jonathan Baxter, Peter Bartlett, and Marcus Frean. Boosting



- 820 algorithms as gradient descent. *Advances in neural information processing*  
821 *systems*, 12, 1999.
- 822 Victor E. McGee and Willard T. Carleton. Piecewise Regression. *Journal of*  
823 *the American Statistical Association*, 65(331):1109–1124, 1970.
- 824 Yong Wang and Ian H. Witten. Inducing Model Trees for Continuous Classes.  
825 In *ECML Poster Papers*, 1997.
- 826 Michael J Kearns. *The computational complexity of machine learning*. MIT  
827 press, 1990.
- 828 Raaz Dwivedi, Chandan Singh, Bin Yu, and Martin J Wainwright. Re-  
829 visiting complexity and the bias-variance tradeoff. *arXiv preprint*  
830 *arXiv:2006.10189*, 2020.