



**HAL**  
open science

## Efficient models for predicting durum wheat grain Cd conformity using soil variables and cultivars

Christophe Nguyen, Agathe Roucou, Guéno   Grignon, Jean-Yves Cornu,  
Benoit Meleard

► **To cite this version:**

Christophe Nguyen, Agathe Roucou, Gu  no   Grignon, Jean-Yves Cornu, Benoit Meleard. Efficient models for predicting durum wheat grain Cd conformity using soil variables and cultivars. *Journal of Hazardous Materials*, 2021, 401, pp.1-12. 10.1016/j.jhazmat.2020.123131 . hal-02914287

**HAL Id: hal-02914287**

**<https://hal.inrae.fr/hal-02914287v1>**

Submitted on 2 Sep 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destin  e au d  p  t et    la diffusion de documents scientifiques de niveau recherche, publi  s ou non,   manant des   tablissements d'enseignement et de recherche fran  ais ou   trangers, des laboratoires publics ou priv  s.

1 Efficient models for predicting durum wheat grain Cd  
2 conformity using soil variables and cultivars

3 Christophe Nguyen<sup>a\*</sup>, Roucou Agathe<sup>b</sup>, Guénoé Grignon<sup>b</sup>,  
4 Cornu Jean-Yves<sup>a</sup>, Benoît Méléard<sup>b</sup>

5 May 27, 2020

6 <sup>a</sup> ISPA, Bordeaux Sciences Agro, INRAE. F-33140 Villenave-d'Ornon cedex, France.

7 <sup>b</sup>ARVALIS-Institut du Végétal. Station expérimentale, 91720 BOIGNEVILLE,

8 <https://www.english.arvalisinstitutduvegetal.fr>

9 \*Corresponding author phone: +33 (0)5 57 12 25 07; fax: +33 (0)5 57 12 25 15;

10 e-mail: christophe.nguyen@inrae.fr

## 11 **1 Abstract**

12 Contamination of durum wheat grain by cadmium (Cd) threatens food safety and is of in-  
13 creasing concern because regulations concerning Cd are becoming stricter due to its toxicity.  
14 This work aimed at using soil variables and cultivar types to build models to predict whether  
15 durum wheat grain Cd will conform with current and possibly lower regulatory thresholds. We  
16 combined multiple Gaussian and logistic regressions and the random forest algorithm to take  
17 advantage of their strength. Models tested using cross-validation produced excellent perfor-  
18 mances including for the lowest regulatory threshold of 0.1 mg Cd/kg, half of the current one:  
19 79-85% of the non-conformity cases were detected and the reliability of predictions was 69-82%.  
The models enabled identification of a x1.4 variability in grain Cd content between cultivars

20 that do not have the low Cd accumulation allele of the *Cdu1* gene. The models confirmed  
21 that for the grain Cd content, the between-cultivar variability had much less influence than the  
22 phytoavailability of Cd in soil, the critical contexts of which were characterized by the models.  
23 For farmers, these models are valuable tools to predict whether durum wheat production will  
24 conform with existing and future Cd regulation in foodstuffs.

25 **Keywords:** Cadmium; Durum wheat; Genetic variability; Models; Phytoavailability

## 26 **2 Introduction**

27 Cadmium (Cd) is a highly toxic and carcinogenic metal found naturally in soils. It is taken  
28 up by plant roots and transferred to edible plant parts and is therefore a major threat to food  
29 safety (Clemens et al., 2013). In 2009, the European Food Safety Authority (EFSA) published  
30 a scientific opinion recommending a tolerable weekly intake (TWI) for Cd of  $2.5 \mu\text{g kg}^{-1}$  body  
31 weight, almost three times lower than the previous one set by the World Health Organization,  
32 which was  $7 \mu\text{g kg}^{-1}$  (Alexander et al., 2009). Durum wheat concentrates more Cd in grain than  
33 bread wheat (Greger and Lofstedt, 2004). Durum wheat is a major contributor of Cd to human  
34 food intake as it is widely consumed in pasta and semolina (Clarke et al., 2010). For instance,  
35 9 Mt of durum wheat are consumed in Europe, which is also the world's main exporter at 8 Mt  
36 (FranceAgrimer, 2020). Following the downward revision of the TWI by EFSA, the European  
37 commission also revised the directive EC1881/2006, which fixes the maximum content of Cd in  
38 some foodstuffs (DGSANCO, 2011). For durum wheat, 0.1 and 0.15 mg Cd  $\text{kg}^{-1}$  were originally  
39 considered but the project was abandoned because of the strong economic negative impact the  
40 decision would have. However, the European countries were asked to conduct research and  
41 develop practices aimed at monitoring and reducing crop contamination (EC EC 488/2014,  
42 2014). Recently, new downwards revisions of the level of Cd in cereals including for durum  
43 wheat have again been the subject of discussions (ARVALIS-Institut du végétal, pers. comm.).  
44 Crop uptake of Cd depends on its phytoavailability in soils, *i.e.* on the flux of  $\text{Cd}^{2+}$  at the  
45 root surface, assuming that the free ion is the main species absorbed by root cells (Clemens,  
46 2019). In aerobic agricultural soils, Cd phytoavailability is mainly determined by sorption onto  
47 the solid organic and mineral phases (clay and oxides), by complexation with soluble ligands,

48 generally organic compounds, and by transport to plant roots by diffusion and advection (Lin  
49 et al., 2016; Antoniadis et al., 2017; Vega et al., 2010). Soil Cd, pH, organic matter and to a  
50 lesser extent clay and oxides have been found to be the major regressors of statistical models to  
51 predict the soluble soil Cd and its accumulation by plants (Adams et al., 2004; de Vries et al.,  
52 2011; Groenenberg et al., 2010; Horn et al., 2006). Except oxides, these variables are commonly  
53 measured in soil testing, and therefore, predicting plant Cd by using statistical models based  
54 on these variables is of great interest (for example, Hough et al., 2003; Tudoreanu and Phillips,  
55 2004; Viala et al., 2017). According to the literature, one major drawback of such models is  
56 that they have not been tested on new data (cross-validation). Their real predictive value for  
57 new data is consequently not known, which is an obstacle to their practical use in the field.  
58 Furthermore, because of the log-log relationship, the variance of the predictions inflates when  
59 the mean increases (Newman, 1993), making it difficult to rule on the conformity of critical grain  
60 samples with high Cd content. Binary classification models that predict conformity (yes/no)  
61 could be much more efficient because they concentrate on predicting the class, whereas log-log  
62 models are optimized to predict continuous variations in grain Cd. Among binary classification  
63 models, two have been shown to be particularly efficient: logistic regression and classification  
64 trees (Hastie et al., 2009). Logistic regression predicts the probability that an observation is  
65 positive (actually, the log of the odds) from a linear combination of predictors. If the predicted  
66 probability is greater than a cutoff threshold (generally  $p=0.5$ ), the case is classified as positive,  
67 and conversely as negative. Classification trees are a set of hierarchical decision rules that split  
68 data into two classes based on cutoff values of the most relevant predictors. In contrast to  
69 logistic regression, classification trees are a non-parametric method and may better model a  
70 complex boundary between the two conformity classes. However, they are very sensitive to  
71 the training dataset (high variance). To cope with this problem, Breiman (2001) proposed the  
72 random forest approach, which consists in aggregating the predictions of a large number of  
73 trees (*i.e.* a forest) that are uncorrelated by bootstrapping the training dataset.

74 Soil conditions and particularly soil pH, strongly influence the contamination of crops by metals  
75 by controlling the phytoavailability of the latter (Kabata-Pendias, 2004). Reducing the phy-  
76 toavailability of metals is often difficult, especially if the soil pH is already high. This is typically  
77 the case of calcareous soils, on which durum wheat is usually grown in France. Therefore, it

78 is worth taking advantage of any between-cultivar variability in Cd accumulation by crops (Li  
79 et al., 1997; Li and Zhou, 2019). Due to the strong importance of durum wheat in Canada,  
80 genetic selection of cultivars that accumulate little Cd in their grain began in the 1990s when  
81 international discussions about setting Cd limits in food products started (Clarke et al., 2010).  
82 It is known that a large part of phenotypic variability in durum wheat grain Cd is linked to  
83 the *Cdu1* locus of chromosome 5B, which is involved in Cd sequestration in roots (Knox et al.,  
84 2009; Wiebe et al., 2010). One deficient allele of the gene coding for the HMA3 transporter  
85 that transfers Cd and Zn from the root cytosol into the vacuole is thought to have been selected  
86 inadvertently during breeding because it promoted growth of durum wheat in Zn-deficient soils  
87 (Maccaferri et al., 2019). As a consequence of reduced sequestration in root vacuoles due to  
88 this deficient allele, more Cd is allocated to aboveground organs, including the grain. To our  
89 knowledge, selection of low Cd durum wheat cultivars has not yet begun in Europe but thanks  
90 to some markers of the *Cdu1* locus (AbuHammad et al., 2016; Oladzad-Abbasabadi et al., 2018;  
91 Salsman et al., 2018), many common European cultivars have been assessed, and the results  
92 show that a large proportion of all cultivars are high Cd accumulators (Zimmerl et al., 2014).  
93 Therefore, as suggested by preliminary results obtained in controlled conditions (Perrier et al.,  
94 2016), it would certainly be worth characterizing the variability of grain Cd content among  
95 high Cd cultivars to discard the highest accumulators, if this is possible with respect to other  
96 agronomic performances.

97 Based on these elements, the present work had two goals. The first was to build sensitive and  
98 reliable models to predict durum wheat grain Cd conformity using soil analysis variables. We  
99 hypothesized that soil analysis variables combined with highly efficient statistical approaches  
100 would make it possible to obtain predictive models that are sufficiently sensitive and reliable  
101 for practical use in the field. The second goal was to identify possible solutions in the case  
102 of predicted non-conformity. To this end, we investigated the variability of Cd accumulation  
103 between cultivars that do not possess the low Cd allele of the *Cdu1* gene and we used the  
104 model simulations to identify the soil conditions that could lead to an excessive contamination  
105 of durum wheat grain by Cd.

## 106 **3 Materials and Methods**

### 107 **3.1 Collection and analysis of paired soil and grain samples from** 108 **farms and of grain samples from trials comparing cultivars**

109 Between 2012 and 2018, 420 paired samples of soil and durum wheat grain were collected in  
110 farms and in ARVALIS-Institut du végétal trials across the regions that are representative of  
111 the French production. There were 192 different soil samples because several cultivars were  
112 grown on the same soil. The cultivars we studied had either been characterized as high Cd ac-  
113 cumulators (Zimmerl et al., 2014) or had not yet been characterized (Table SI1, Supplementary  
114 Information). Composite soil samples were made by mixing 12 sub-samples taken from the 0-25  
115 cm topsoil layer on a grid representative of each plot. The composite soil samples were air-dried,  
116 sieved at 2 mm before the common characteristics of agricultural soil testing were determined  
117 by a certified Inrae soil analysis laboratory (<https://www6.hautsdefrance.inrae.fr/las>). The to-  
118 tal soil Cd was quantified by ICP-MS after solubilization by fluorhydric and perchloric acids  
119 (NF X 31-147). Soil pH was measured in a 1:5 soil to water solution. Total soil organic carbon  
120 (SOC) was quantified by dry combustion and corrected for carbonate content (NF ISO 10694).  
121 The Robinson pipette method (NF X 31-107) was used for soil texture (5 size classes). Total  
122 soil CaCO<sub>3</sub> content was obtained using the acid neutralization method (NF X 31-105). In  
123 France, no other oxides are usually analyzed in soil testing even though they could play an sig-  
124 nificant role in controlling Cd phytoavailability (Sun et al., 2017). The descriptive statistics of  
125 the studied soils are listed in Table 1. Grain samples, that were selected as being representative  
126 of the plot harvest, were analyzed by Capinov certified laboratory (<https://www.capinov.fr/>)  
127 without further drying as stated by the European regulation EC466/2001 (2001). Total grain  
128 Cd was quantified by atomic absorption spectroscopy (AAS) after wet digestion in a mixture of  
129 nitric acid (10% v:v) and hydrogen peroxide (4% v:v). Both the Capinov and Inrae laboratories  
130 are certified by Cofrac for quality controls (<https://www.cofrac.fr/en/>).

131 To rank durum wheat cultivars with respect to their capacity to accumulate Cd, additional  
132 grain samples were collected in 2016, 2017 and 2018 from experimental plots located in three  
133 ARVALIS-Institut du végétal research centers. These trials are conducted annually to assess  
134 the agronomic performances of durum wheat cultivars. The trials are located in the south-west

135 (Bergerac, 43° 25' 0.12''E 0° 9' 0''), south (Montesquieu-Lauragais, 43° 25' 0.12''E 0° 9' 0'') and  
 136 center (Thizay, 47° 10' 0.12''E 0° 0' 0'') of France, the main French durum wheat production  
 137 regions. Grain Cd content was quantified using the same procedure as that described above.  
 138 No soil was collected and these samples were not used for modeling.

### 139 **3.2 Modeling durum wheat grain conformity with respect to the Cd** 140 **regulatory threshold**

141 The goal of the models was to use soil characteristics to predict if a grain sample of a particular  
 142 durum wheat cultivar will comply (0) or not (1) with the regulatory threshold (RT). Seven soil  
 143 variables were selected as predictors and combined hierarchically based on the following ranking  
 144 Cd>pH>SOC>Clay>{Fine silt, calcareous}>coarse loam. The database has a minimum of 75  
 145 positive cases for the RT of 0.2 mg Cd kg<sup>-1</sup> grain. Therefore, a maximum of 7 predictors  
 146 were considered, based on the guidelines proposed by Peduzzi et al. (1996), which state that at  
 147 least 10 positive cases per predictor are required for a correct estimation of effects in a logistic  
 148 regression. In total, nine combinations of predictors were tested, from the simplest including  
 149 only the soil Cd to the full model including the 7 soil predictors. All models also include the  
 150 cultivar as a categorical predictor.

151 Three statistical modeling approach were tested. The first approach was the mixed-effects  
 152 logistic regression (MELR), which predicted the *log*-odds of the non-conformity of the grain Cd  
 153 as a function of the *log* of the scaled soil variables (fixed effects) and of the cultivar (random  
 154 effect for the intercept).

$$\begin{aligned}
 \log\left(\frac{p}{1-p}\right) &= \log(a_0) + \sum a_i \log\left(\frac{X_i}{\text{mean}(x_i)}\right) \\
 \text{random} &= \log(a_0)|_{\text{cultivar}} \\
 p &= p(Cd_{\text{grain}} \geq RT) = p(Y = 1)
 \end{aligned}
 \tag{1}$$

155 Continuous predictors were *log* (natural) transformed to i) ensure the normality of the residuals  
 156 and to ii) model interactions between soil variables since the sum of *log* of variables is the *log*  
 157 of the product of the variables. In order to correctly estimate the model intercept, the soil  
 158 variables were first scaled by dividing them by their mean before *log* transformation. For each

159 of the nine models, the predicted *log*-odds from the calibration dataset allowed us to calculate  
160 the probabilities of non-conformity  $p(Y = 1)$ . Then, each probability was individually tested  
161 as a cutoff (*ctf*) to code all predicted probabilities in conformity classes: 0 if  $p \leq ctf$  and 1  
162 if  $p > ctf$ . The predicted conformity classes were compared to the actual classes to calculate  
163 the confusion matrix, *i.e.* the scores for the true positives (TP), true negatives (TN), false  
164 positives (FP) and false negatives (FN). The performances of the model for a given cutoff were  
165 assessed using the Youden's  $J$  statistic ( $0 \leq J \leq 1$ ) and the Matthews correlation coefficient  
166  $MCC$  ( $-1 \leq MCC \leq 1$ ):

$$J = \frac{TP}{TP + FN} - \frac{FP}{TN + FP} \quad (2)$$

$$MCC = \frac{TP.TN - FP.FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (3)$$

167 For both statistics, a value of 1 indicates a perfect model with no FN nor FP. The MCC  
168 considers the four scores of the confusion matrix (TP, TN, FP, FN) and thus makes it possible  
169 to identify models with good performances in both the detection and reliability of predictions of  
170 actual positive and negative cases. By contrast, the Youden index concentrates on the predicted  
171 positive class (TP and FP) and favors the selection of models that have the best performances  
172 in detecting positive cases. The two best cutoffs, each maximizing either the  $J$  or the  $MCC$   
173 statistics were identified and used to make future class predictions.

174 The second approach was random forest (RF) modeling (Breiman, 2001). For each of the nine  
175 models, 501 classification trees were trained with bootstrap samples that were stratified based  
176 on the (0,1) frequencies of the training dataset. The number of variables randomly sampled as  
177 candidates at each split was equal to the square root of the total number of predictors of the  
178 model. The minimum number of observations in terminal nodes was 1. For each observation  
179 of the calibration dataset, the predicted probability that the grain Cd is above the regulatory  
180 threshold was estimated by the frequency of the positive class predicted by the 501 trees for  
181 this observation. Like for logistic modeling, each predicted probability was tested as a cutoff  
182 for coding into the (0,1) classes. The Youden and MCC statistics were again used to select the  
183 best cutoffs.



184 The third approach was the mixed effects Gaussian multiple linear *log-log* regression (MEMR).  
185 The *log* of the grain Cd content was predicted from the *log* of the scaled soil predictors (fixed  
186 effects) with the cultivar as a random effect for the intercept.

$$\begin{aligned} \log(Cd_{grain}) &= \log(a_0) + \sum a_i \log\left(\frac{X_i}{\text{mean}(x_i)}\right) \\ \text{random} &= \log(a_0)|_{\text{cultivar}} \end{aligned} \tag{4}$$

187 The reasons for scaling using the median and for *log* transforming the continuous variables were  
188 the same as for the MELR. After back transformation, the predicted grain Cd ( $Cd_{grain}$ ) was  
189 coded as follows: 0 if  $Cd_{grain} \leq RT$ , otherwise 1.

### 190 **3.3 Ranking the predictive value of the models by 5-fold cross-** 191 **validation**

192 The predictive values of the 9 models x 3 modeling approaches (RF, MELR and MEMR) were  
193 evaluated by a 5-fold cross-validation with 400 repetitions. A single design was used per RT  
194 (splitting the database into 5 groups 400 times). One repetition consisted in randomly splitting  
195 the database into five groups of 84 observations, each group having the same proportions of  
196 (0,1) as the whole database. Each of the five groups was individually used to make predictions  
197 from the model trained on the four remaining groups that were pooled. For each of the 420  
198 observations of the database, the 400 repetitions were bootstrapped 400 times to generate  
199 400 combinations of the 420 predictions for the whole database. This enabled to have 400  
200 repetitions of the  $J$  and  $MCC$  statistics to further calculate their mean and standard deviation  
201 in order to compare and group the 9 models x 3 modeling approaches at  $p=0.05$  adjusted for the  
202 multiple comparisons (Tukey test). This was done for each regulatory threshold (RT): RT02,  
203 RT015, RT01. The best parsimonious models optimized by the  $J$  and  $MCC$  statistics were  
204 the models that had the lowest number of predictors while satisfying the following conditions:  
205  $J/J_{max} > 0.95$  or  $MCC/MCC_{max} > 0.95$ .

### 206 **3.4 Ranking cultivars based on their grain Cd contents**

207 The cultivars were ranked based on the random effects of the MELR best models at RT01 on  
208 the one hand, and from the field trials on the other. For field trials, grain Cd content was  
209 modeled as a function of the cultivar using a mixed-effects model with the year x location as  
210 a random effect for the intercept. The adjusted means of the grain Cd of the cultivars were  
211 grouped using the Tukey’s test at  $p=0.05$ , adjusted for multiple comparisons. For the cultivars  
212 that were the same in the modeling and in the field trials, the ranking of the two approaches  
213 was compared to examine consistency. The statistical models used in this work never accounted  
214 for an interaction between the cultivars (genotypes) and the environment *i.e* the soil variables  
215 for the modeling and the year x locations for the field trials. The interactions were tested but  
216 were not significant or led to over-parameterization of the mixed-effect models.

### 217 **3.5 Characterization of the effects of soil variables and of the cultivar** 218 **based on the conformity of grain Cd content**

219 Based on the performances of the models we tested (detailed in the results section), the random  
220 forest approach where the cutoff is optimized from the MCC statistics was used to simulate  
221 the conformity of the grain Cd content from the 7 soil variables and for the cultivar “Relief”,  
222 which was classified as the highest Cd accumulator and for the cultivar “Miradoux”, which  
223 was classified as the lowest accumulator (ranking is detailed in the results section). A factorial  
224 experimental design crossed 6 values chosen to cover realistic ranges for each of the 7 soil  
225 variables. For a given value of a given predictor  $X_i = x_{ij}$ ,  $i \in \{1, \dots, 7\}$ ,  $j \in \{1, \dots, 6\}$ , the  
226 probability of non-conformity  $p(Y = 1|X_i = x_{ij})$  was the frequency of the non-conformity class  
227 among the  $6^6 = 46656$  predictions where  $X_i = x_{ij}$  and  $X_k = x_{kj}$ ,  $k \in \{1, \dots, 7\}$ ,  $k \neq i$ . It was  
228 therefore not possible to directly derive a standard deviation for this probability. The effect of  
229 each predictor on the probability of non-conformity was characterized for the two cultivars and  
230 for the three RTs.

## 231 **3.6 Calculations and statistical software**

232 All the data processing and statistical analyses were performed with R, version 3.6.2 (R Code  
233 Team, 2019). The following specific packages were used: *randomForest* (version 4.6-14) for  
234 the random forest modeling, *lme4* (version 1.1-21) to fit the mixed effects models MELR and  
235 MEMR and the *doParallel* (version 1.0.15) for parallel computing. The residuals of MEMR were  
236 examined to check for heteroskedasticity and non-normality. When detected, heteroskedasticity  
237 was corrected by modeling the variance with the *varPower()* function of the *nlme* package  
238 (version 3.1-143)

## 239 **4 Results**

### 240 **4.1 Descriptive statistics of the modeling database**

241 The ranges of variation in soil characteristics were typical of agricultural soils (Table 1). The  
242 median soil Cd content was  $0.26 \text{ mg kg}^{-1}$ , slightly higher than the French national median of  
243  $0.20 \text{ mg Cd kg}^{-1}$  for agricultural soils (calculated from Saby et al., 2009), the contamination  
244 of which mainly originates from pedogenesis. Ninety-eight percents of soil samples had Cd  
245 contents below or equal to  $1 \text{ mg Cd kg}^{-1}$  (not shown). Most of the soils were alkaline (median  
246  $\text{pH}=8.2$ ) and only 10 had a pH of between 5.5 and 7. The soils with a pH below 7.5 did not  
247 have a significant calcareous content (Fig. SI1 in supplementary information). SOC and soil  
248 Cd were significantly correlated (Fig. SI1,  $R^2 = 0.40$ ,  $p < 0.001$ , not shown). The *log-log*  
249 linear relationship between these two variables indicates that doubling the SOC is expected  
250 to correspond to a soil Cd content increased by a factor x1.36. For the French soil survey  
251 database limited to agricultural soils, the relationship is also significant but much weaker ( $R^2 =$   
252  $0.10$ ,  $p < 0.001$ ,  $n = 5201$ ), and the magnitude of increase is x1.12 (calculated from Saby et al.,  
253 2009). The soils in the modeling database correctly reflect the fact that in France, durum  
254 wheat is mainly cultivated in alkaline and calcareous soils in the southern-half of the country  
255 (FranceAgrimer, 2020). As evidenced by a national soil survey, these soils are slightly richer in  
256 Cd, which is partly because the calcareous of these soils are naturally rich in Cd (Saby et al.,  
257 2009, 2011).

258 There were 36 distinct cultivars in the modeling database, each with between 1 to 85 observa-  
259 tions (Table 2). The cultivars with few observations were kept because they helped to estimate  
260 the fixed effects of soil variables on grain Cd conformity. By contrast, a lack of observations  
261 makes it difficult to estimate the random effect of the cultivar on the intercept of the model.  
262 As expected, the frequency of samples that did not comply with the regulation increased with  
263 decreasing RT, namely 18% for RT02, 26% for RT015 and 48% for RT01 (Table 2). These  
264 values were higher than the estimates made in a larger French national survey: 4%, 8% and  
265 22%, respectively. However, these values were in the range of some production areas where the  
266 phytoavailability of Cd is higher than average, for instance in the center of France (ARVALIS-  
267 Institut du végétal and France-AgriMer, unpublished statistics). Whatever the RT, there were  
268 marked variations in the frequency of non-conformity between the high Cd cultivars (Table 2).

## 269 4.2 Model performances resulting from the K=5-fold cross-validation

270 For RT02 and RT015, random forest (RF) was the most efficient modeling approach, followed  
271 by logistic regression (MELR) and then by Gaussian regression (MEMR) (Fig. 1). This ranking  
272 was reversed for RT01. For the three regulatory thresholds (RT02, RT015, RT01) and for both  
273 performance criteria (the Youden index  $J$  and Matthews correlation coefficient  $MCC$ ), the best  
274 models in the absolute had between 5 and 7 soil predictors, and the most parsimonious ones had  
275 between 2 and 5 (Fig.1, Table 3). Regarding the parsimonious models, the number of predictors  
276 increased with lowering of the RT, indicating that more information is required to correctly  
277 predict low RTs. The parsimonious models always produced significantly lower performances  
278 than the best models, but less than 5% by definition (see materials and methods). By combining  
279 the modeling approaches and the set of predictors, it was possible to obtain models that have  
280 very good performances. Between 79% and 85% of the cases of non-conformity were detected,  
281 the sensitivity slightly decreasing with decreasing RT (True Positive Rate, TPR, Table 2).  
282 The reliability of the model predictions (Positive Predictive Value, PPV) was barely lower,  
283 between 67% and 82 % success. Reliability increased with lowering of the RT, in contrast  
284 to sensitivity. This revealed a trade-off between the detection capacity and the reliability of  
285 predictions. Concerning cases of conformity, the model performances were a little better than  
286 for non-conformity: between 83% and 92 % of successful detection (True Negative Rate, TNR,

287 Table 3) and between 81% and 97% for reliability (Negative Predictive Value, NPV), both  
 288 decreased with lowering of the RT. Compared to the optimization of the models using the  $J$   
 289 statistics, optimization using the  $MCC$  increased the PPV (reliability) by around 2% at the  
 290 expense of a decrease in the TPR (sensitivity) also of around 2%.

291 It should be noted that the reliability performances of the models depends on the actual fre-  
 292 quencies of samples that do not comply with the RT. Hence the PPV and NPV given in Table  
 293 3 are conditioned by the percentages of samples above RT02, RT015 and RT01 in our database  
 294 (18%, 26% and 48%, respectively). If the predictions are required in a context in which the  
 295 frequency of non-conformity (prevalence:  $p$ ) differs from that in the database used to train the  
 296 models, the PPV and NPV must be corrected as follows:

$$PPV^* = \frac{pTPR}{pTPR + (1-p)(1-TNR)} \quad (5)$$

$$NPV^* = \frac{(1-p)TNR}{(1-p)TNR + p(1-TPR)} \quad (6)$$

297 Hence, in a prediction context where the prevalence could be lower than that of the modeling  
 298 database, the reliability of the models will be reduced and conversely. TPR and TNR do not  
 299 depend on prevalence.

### 300 4.3 Ranking of predictors and analysis of their effects

301 Fig. SI2 shows the estimated importance of the predictors for the random forest model with  
 302 7 soil variables and the cultivars. On average, the importance of the predictors increases with  
 303 lowering of the RT. The most influential predictors are soil Cd and pH, the least, the cultivar,  
 304 and in an intermediate position, clay, fine silt (FS) and coarse silt (CS). The relative rank of  
 305 SOC and of calcareous varied greatly depending on the RT, likely because they co-varied with  
 306 soil Cd and pH, respectively.

307 Fig. 2 shows the predicted frequencies of non-conformity as a function of the 5 most important  
 308 predictors, for the three RT and for the lowest (Miradoux) and highest (Relief) Cd accumulator  
 309 cultivar. The model predictions demonstrate the strong effect of soil Cd and pH followed by  
 310 that of clay. The risk of non-conformity increases with soil Cd with an 'S-shaped' response.

311 The effect increases with lowering of the RT. For the highest soil Cd contents, depending on  
312 the cultivar, the risk of non-conformity for RT02 is low to moderate whereas it is always certain  
313 at RT01. In the case of low soil Cd, the lag phase of the risk is severely reduced at RT01. The  
314 differences between the two cultivars is predicted to be weak at low Cd at RT02, at low and  
315 high Cd at RT01 and irrespective of soil Cd at RT015.

316 Concerning soil pH, on average, the risk was also predicted to increase with lowering of the RT.  
317 The models predicted a slight linear decrease in the risk at RT01 with little differences between  
318 cultivars. In contrast, at RT02 and RT015, the risk was predicted to be maximum when  
319 the pH is below 6.5 with marked differences between the two cultivars, and then, to strongly  
320 decrease between pH=6.5 and pH=7 with few changes at higher pH. The mean predicted risk  
321 of non-conformity was mapped for different combinations of soil Cd and pH (Fig. 3). At RT02,  
322 estimated risk was always low for Miradoux ( $p < 0.2$ ) whereas it was more than 40% ( $p > 0.4$ ) for  
323 soil Cd  $> 1.25 \text{ mg kg}^{-1}$  and pH  $< 6.5$  for Relief. With lowering of the RT, the soil Cd and pH  
324 area of high risk logically increase. At the same time, the differences between the two cultivars  
325 decrease considerably and at RT01, even for Miradoux, the lowest Cd accumulator cultivar,  
326 the safe area was very narrow. Finally, the model predicted that at RT01, at soil Cd  $> 0.2 \text{ mg}$   
327  $\text{kg}^{-1}$ , even at high soil pH, the risk of non-conformity could be high.

328 For both cultivars, SOC was predicted to have little effect at RT02 and RT015 but at RT01,  
329 increasing SOC is predicted to reduce the risk of non-conformity by around 40%, more markedly  
330 below 20 g C/kg (Fig. 2).

331 The models did not predict marked effect of clay at RT01 (Fig. 2). At RT015, the risk was  
332 predicted to first decrease at between 50-400 g clay  $\text{kg}^{-1}$  soil and to increase with higher contents  
333 with marked differences between cultivars. At RT02, an increase above 400 g clay  $\text{kg}^{-1}$  soil was  
334 also predicted. The effects of fine silt were similar but less than the effect of clay.

335 The random effects for the intercept of the model fitted to RT01 data allowed us to rank the  
336 cultivars in the modeling database (Fig. 4) showing a x1.4 factor of variation in grain Cd. On  
337 the other hand, field trials allowed us to establish groups of sensitivity to grain Cd accumulation  
338 in another set of cultivars with a x2.9 factor of variation between the two extreme groups (Fig.  
339 5). As shown in Fig. 6, for the cultivars that were used in the two approaches, ranking was  
340 consistent except for Babylone and Sculptur, which were, respectively more strongly over and

341 under classified by the models.

## 342 **5 Discussion**

### 343 **5.1 The choice of the modeling approach depends on the regulatory** 344 **threshold**

345 At the highest RTs (RT02 and RT015), the models concentrate on the highest phytoavailability  
346 of Cd in the soil, which, in these cases, is mainly controlled by the soil Cd and pH (Fig.2 and  
347 3, and see the parsimonious models, Table 3). As shown in Fig.2, at these RTs, the shape  
348 of the effects of soil Cd and pH cannot be completely modeled by the power mathematical  
349 model of the MELR and MEMR approaches. Because non-parametric classification trees are  
350 more flexible for modeling complex non-linear relationships, RF models consequently performed  
351 slightly better (Fig. 1 and Table 2) as it has also been reported in other studies (Covelo et al.,  
352 2008; Qiu et al., 2016). The use of classification trees to model complex and non monotonic  
353 responses (in this case, the boundary between non-conformity and conformity of the grain Cd  
354 content) increases the risk of obtaining over-fitted models with high variance. This pitfall is  
355 counteracted by using the RF approach, which aggregates the predictions of a large number  
356 of relatively uncorrelated trees (Breiman, 2001). In this way, the errors of prediction of some  
357 trees, in particular those due to over-fitting, are offset by the remaining good predictions. This  
358 is likely the reason why RF performed better than the logistic regression at RT02 and RT015.  
359 The cutoff optimization is one likely reason why RF and MELR performed better than MEMR  
360 at RT02 and RT015. At these two RTs, the actual frequencies of non-conformity are 18% and  
361 26%, far from 50% (Table 2). This imbalance between the conformity and non-conformity  
362 classes biases the models if a default cutoff probability of 0.5 is used and this is the reason  
363 why the bias is corrected by optimizing the cutoff (Kuhn and Johnson, 2013), as done in  
364 the RF and MELR models. By contrast, the grain Cd content predicted by the MEMR was  
365 directly transformed into conformity classes depending on whether it was above or below the  
366 RT. This is another possible explanation why MEMR performances were clearly the worst at  
367 at RT02 and RT015 (Fig. 1). At RT01, as the actual frequency of non-conformity was 48%,  
368 cutoff optimization was less necessary. Furthermore, this higher frequency provides much more

369 information to model the boundary between conformity and non-conformity with a parametric  
370 model, explaining why MEMR became the most efficient approach (Fig. 1)

## 371 **5.2 Sensitive and precise models for field prediction of durum wheat** 372 **grain Cd conformity using soil variables**

373 Our work shows that it is possible to predict the conformity of durum wheat grain Cd content  
374 with respect to the regulatory thresholds of 0.2, 0.15 and 0.1 mg Cd kg<sup>-1</sup> with high sensitivity  
375 (probability of detection) and high precision and consequently reliability (probability that a  
376 prediction is true). These models are of practical value for farmers because they only require  
377 variables already determined in soil testings plus the total soil Cd. Total soil Cd is rarely  
378 measured in France and we therefore recommend it is systematically included in all future soil  
379 analyses. On one hand, it would make it possible to use the models built in this work to predict  
380 possible risky situations and on the other hand, it would help monitor the background level of  
381 soil Cd to study in more detail should it increase as a result of agricultural practices, including  
382 fertilization with contaminated P fertilizers (Sterckeman et al., 2018; Six and Smolders, 2014).  
383 The trend in soil Cd in agricultural soil is a serious concern as shown by our results: at RT01  
384 and for sensitive cultivars, the models predict a strong risk of non-conformity for soil Cd above  
385 0.3 mg kg<sup>-1</sup> even at high soil pH (Fig. 2 and 3). The combined analysis of soil and grain  
386 Cd contents is also advisable because the models can learn and improve from new data. The  
387 very good performances of the models at RTs ranging from 0.1 to 0.2 mg Cd kg<sup>-1</sup> suggest that  
388 lowering the current regulatory threshold of 0.2 mg Cd kg<sup>-1</sup> grain would not be an obstacle  
389 to the reliable detection of the great majority of cases of non-conformity. The models will also  
390 help identify possible solutions in the case of a predicted non-conformity. Mapping the risky  
391 contexts for soil variables (Fig. 2 and 3) would help decide if it is worth taking action on the  
392 phytoavailability of Cd by shifting cultivation to another location or by increasing soil pH, for  
393 instance. Ranking cultivars is also a valuable model to reduce the risk of non-conformity of  
394 grain Cd content. The models can easily be re-calibrated for a new RT, not tested in this study.  
395 If the RT is strongly revised downwards to RT01, the global performances of the models will  
396 be reduced (see *MCC* index, Fig. 1). The models will be a little less sensitive in the detection  
397 of both conformity and non-conformity cases (Table 3). As a counterpart, their reliability will



398 increase for non-conformity but not for conformity.

399 From a practical point of view, the choice of the right model can be adjusted depending on  
400 priorities. If the cost of a prediction error is high, the model should be chosen to maximize  
401 reliability and therefore, the models optimized from the *MCC* statistics should be preferred  
402 (Table 3). On the other hand, if the priority is to maximize the detection of cases of non-  
403 conformity at the risk of overestimating the latter, the models optimized by the *J* index should  
404 be chosen. This decision rule does not concern MEMR models for RT01, for which the prediction  
405 class does not rely on an optimized cutoff.

### 406 **5.3 Phytoavailability of Cd in soil versus between-cultivar variability** 407 **to manage the conformity of grain Cd content**

408 The marked effect of soil Cd observed in this study is due to the fact that in non-polluted  
409 agricultural soils, the  $\text{Cd}^{2+}$  concentration in the soil solution is generally low, around 0.1-1 nM  
410 (Schneider et al., 2019; Sauvé et al., 2000), compared to the capacity of roots to take up the  
411 metal (Lux et al., 2011). Therefore, the factor that limits uptake is generally the supply of  $\text{Cd}^{2+}$   
412 to the root surface (Lin et al., 2016). The latter is mainly controlled by the pool of soil Cd that  
413 can be exchanged with the solution and by the speciation of soluble Cd. The strong effect of soil  
414 pH modeled in our work illustrates the competition between H and Cd for sorption sites and for  
415 association with soluble ligands. The threshold for the pH effect of 6.5-7.0 is unlikely to be due  
416 to the formation of complexes between Cd and  $\text{OH}^-$  or carbonates according to the stability  
417 constants of these compounds reported in the chemical databases (Tipping et al., 2011). The  
418 6.5-7.0 threshold could correspond to several mechanisms of Cd sorption. Regarding organic  
419 matter, based on the mean *log* of the dissociation constant of the proton ( $\log K_H$ ) for humic  
420 substances (Tipping et al., 2011; Matynia et al., 2010), the 6.5-7.0 threshold could reflect the  
421 binding of Cd to OH groups that become increasingly deprotonated at high pH. The 6.5-7.0  
422 threshold could also be due to the favored binding of Cd to variable charge sites of Fe, Mn and  
423 Al oxyhydroxides and organo-mineral complexes associated with clay (Violante et al., 2010;  
424 Rasmussen et al., 2018). The pH of most French soils are above 6.5-7.0, especially for durum  
425 wheat, which is frequently grown on calcareous soils. Hence, because the models predict little  
426 effect of pH above 7, (Fig. 2), this soil variable is not an important lever to manage the grain

427 Cd content in French durum wheat. On the other hand, because for pH above 7, Fe, Mn and Al  
428 oxides are expected to have an increasing role in binding Cd, their contents should be included  
429 in soil testing to be able to possibly improve the models.

430 SOC is generally found to strongly control Cd phytoavailability because it includes both solid  
431 and dissolved organic matter that sorbs and forms complexes with Cd. In our study, the  
432 correlation between Cd and SOC partly masked the effect of SOC but in agreement with the  
433 literature, the models predicted that SOC is expected to reduce the phytoavailability of Cd,  
434 particularly if the latter is low on average (RT01, Fig. 2). Adding organic matter to soil  
435 is encouraged for many reasons including improving fertility and storing carbon. Based on  
436 our study, this lever is also questionable because the effect is estimated to be moderate and  
437 considering the soil Cd-SOC correlation in our database, one may wonder if organic matter  
438 does not increase soil Cd, due to its own contamination or by sequestering Cd deposited in  
439 soil by atmospheric fallouts or by agricultural inputs such as P fertilizers.

440 Finally, clay and to a lesser extent fine silts were predicted to have a moderate effect on the  
441 predicted conformity of grain Cd, but only for the highest RTs. Clay and FS are involved in  
442 the reversible exchanges of Cd with the solution (buffer capacity). On one hand clay and FS  
443 are expected to reduce the background concentration of soluble Cd due to sorption but on the  
444 other hand, they facilitate the buffering of this concentration when the roots take up the Cd.  
445 However, these antagonistic effects are unlikely to explain the slight negative effect of clay < 400  
446 g kg<sup>-1</sup> soil and the positive effect above this threshold. The buffer capacity of the solid phase  
447 for Cd is all the more involved as the phytoavailability of Cd is low and therefore, the effects  
448 of clay and FS should be stronger at RT01 compared to RT02, unlike in the simulations. The  
449 effect of clay and to a lesser extent of FS are probably due to a higher abundance of soil rich  
450 in Cd for soils with high clay contents, as shown by Fig. SI1.

451 It is noteworthy that our models produced good performances although they use the total Cd  
452 content of soils and not the mobile Cd, suggesting that the relationship between the two pools  
453 was strong enough to allow correct prediction of the conformity of grain Cd. Furthermore,  
454 the error derived from approximating the mobile Cd by the total Cd is expected to have less  
455 negative impact when predicting a grain Cd binary class than when predicting grain Cd content.  
456 Modeling durum wheat grain Cd showed that soil variables governing Cd phytoavailability have

457 much more influence than the type of cultivar (Fig. SI2). This confirms previous observations  
458 (Li et al., 1997; Li and Zhou, 2019) and was probably more apparent in our work since most  
459 of the cultivars we investigated did not have the *Cdu1* low accumulation allele, thus reducing  
460 between-cultivar variability. Hence, the model estimated a x1.4 factor between the lowest and  
461 highest Cd accumulator cultivar whereas field trials, which included at least one cultivar with  
462 the *Cdu1* low Cd allele (Anvergur) showed a x3 factor of variation. The x1.4 factor of variation  
463 predicted by the models for cultivars without the low-Cd allele of the *Cdu1* gene suggests that  
464 there is also variability in some mechanisms contributing to reduce grain Cd, other than the  
465 enhanced sequestration of Cd in roots. For example, modification of the rhizosphere, including  
466 changes in pH and the release of ligands such as low molecular weight organic acids by roots can  
467 significantly differ between wheat cultivars (Cieśliński et al., 1998; Greger and Landberg, 2008).  
468 Differences in Cd sequestration in the stem and nodes of different rice cultivars has also been  
469 observed (Fujimaki et al., 2010). It was shown in sunflower and wheat, that the partitioning  
470 of plant biomass, especially aboveground biomass and plant height could partly explain the  
471 intraspecific variability in the grain Cd content (Laporte et al., 2015; Pozniak et al., 2012;  
472 Perrier et al., 2016; Álvaro et al., 2008). Hence, understanding the reasons for the variability in  
473 grain Cd contents among high Cd cultivars merits further investigations, in particular because  
474 reducing the phytoavailability of Cd in soil is not an easy task.

## 475 **6 Conclusions**

476 This work confirms the marked influence of soil Cd and pH on the transfer of Cd to durum  
477 wheat grain. For this crop, we have shown that grain Cd can exceed the current and possible  
478 lower future regulatory thresholds even in alkaline soils with moderate total Cd contents below  
479 1 mg Cd kg<sup>-1</sup> soil. Combining random forest, multiple logistic and Gaussian linear regressions  
480 enabled us to build models to predict grain Cd conformity that are both efficient and reliable.  
481 The model performances are also very good for the lowest RT that were considered by the  
482 downwards revision of the regulatory threshold. By adjusting the predicted grain Cd using  
483 the phytoavailability of Cd estimated from soil variables, the models also showed there was a  
484 x1.4 factor of variation between durum wheat cultivars that did not have the low Cd allele of  
485 the *Cdu1* gene. Because the models only require variables available from soil analyses, they

486 are valuable tools to be able to predict possible problems of non-conformity of durum wheat  
487 production with respect to the regulation concerning Cd in foodstuffs. Further, considering the  
488 current trend towards a more restrictive regulation for food products, the approach used in this  
489 work can also be used for other heavy metals such as Ni which is currently targeted in relation  
490 with the baby foods.

## 491 **7 Acknowledgments**

492 The authors are grateful to S. Thunot, T. Robert and B. Orlando for their technical assistance.  
493 They are also grateful to FranceAgriMer and to all farmers involved in this study for their com-  
494 mitment. This work was funded by ARVALIS-Institut du végétal (Cadur project CR127518)  
495 and by the Agence Nationale de la Recherche (ANR 15-CE21-0001-04, Cadon). The authors  
496 declare that they have no conflict of interest.

## 497 **References**

- 498 AbuHammad, W. A., Mamidi, S., Kumar, A., Pirseyedi, S., Manthey, F. A., Kianian, S. F.,  
499 Alamri, M. S., Mergoum, M., and Elias, E. M. (2016). Identification and validation of  
500 a major cadmium accumulation locus and closely associated snp markers in north dakota  
501 durum wheat cultivars. *Molecular Breeding*, 36(8):112.
- 502 Adams, M., Zhao, F., McGrath, S., Nicholson, F., and Chambers, B. (2004). Predicting cad-  
503 mium concentrations in wheat and barley grain using soil properties rid a-8339-2008 rid  
504 b-5127-2008. *Journal of Environmental Quality*, 33(2):532–541. WOS:000220292800011.
- 505 Alexander, J., Benford, D., Cockburn, A., Cravedi, J. P., Dogliotti, E., Di Domenico, A.,  
506 Fernández-Cruz, M. L., Fürst, P., Fink-Gremmels, J., Galli, C. L., et al. (2009). Cadmium  
507 in food. scientific opinion of the panel on contaminants in the food chain. *The EFSA Journal*,  
508 980:1–139.
- 509 Antoniadis, V., Levizou, E., Shaheen, S. M., Ok, Y. S., Sebastian, A., Baum, C., Prasad, M.  
510 N. V., Wenzel, W. W., and Rinklebe, J. (2017). Trace elements in the soil-plant interface:  
511 Phytoavailability, translocation, and phytoremediation—a review. *Earth-Science Reviews*,  
512 171:621–645.
- 513 Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.
- 514 Cieśliński, G., Van Rees, K., Szmigielska, A., Krishnamurti, G., and Huang, P. (1998). Low-  
515 molecular-weight organic acids in rhizosphere soils of durum wheat and their effect on cad-  
516 mium bioaccumulation. *Plant and Soil*, 203(1):109–117.
- 517 Clarke, J., Clarke, F., and Pozniak, C. (2010). Forty-six years of genetic improvement in  
518 canadian durum wheat cultivars. *Canadian Journal of Plant Science*, 90(6):791–801.
- 519 Clemens, S. (2019). Safer food through plant science: reducing toxic element accumulation in  
520 crops. *Journal of Experimental Botany*, 70(20):5537–5557.
- 521 Clemens, S., Aarts, M. G., Thomine, S., and Verbruggen, N. (2013). Plant science: the key to  
522 preventing slow cadmium poisoning. *Trends in Plant Science*, 18(2):92–99.

523 Covelo, E. F., Matías, J. M., Vega, F. A., Reigosa, M. J., and Andrade, M. L. (2008). A tree  
524 regression analysis of factors determining the sorption and retention of heavy metals by soil.  
525 *Geoderma*, 147(1):75–85.

526 de Vries, W., McLaughlin, M. J., and Groenenberg, J. E. (2011). Transfer functions for solid-  
527 solution partitioning of cadmium for australian soils. *Environmental Pollution*, 159(12):3583–  
528 3594. WOS:000296397300041.

529 DGSANCO (2011). Amendement of regulation (ec) n° 1881/2006 as regards cadmium.

530 EC 488/2014 (2014). Journal officiel de l’union européenne.

531 EC466/2001 (2001). Commission regulation (ec) no 466/2001 of 8 march 2001 setting maximum  
532 levels for certain contaminants in foodstuffs (text with eea relevance.).

533 FranceAgrimer (2020). Marché du blé dur. france, union européenne, monde (2019-2020).  
534 Technical Report ISSN 1779-353X.

535 Fujimaki, S., Suzui, N., Ishioka, N. S., Kawachi, N., Ito, S., Chino, M., and Nakamura, S.-I.  
536 (2010). Tracing cadmium from culture to spikelet: Noninvasive imaging and quantitative  
537 characterization of absorption, transport, and accumulation of cadmium in an intact rice  
538 plant. *Plant Physiology*, 152(4):1796–1806.

539 Greger, M. and Landberg, T. (2008). Role of rhizosphere mechanisms in cd uptake by various  
540 wheat cultivars. *Plant and Soil*, 312(1):195–205.

541 Greger, M. and Lofstedt, M. (2004). Comparison of uptake and distribution of cad-  
542 mium in different cultivars of bread and durum wheat. *Crop Science*, 44(2):501–507.  
543 WOS:000189384600019.

544 Groenenberg, J. E., Romkens, P. F. a. M., Comans, R. N. J., Luster, J., Pampura, T., Shotbolt,  
545 L., Tipping, E., and de Vries, W. (2010). Transfer functions for solid-solution partitioning of  
546 cadmium, copper, nickel, lead and zinc in soils: derivation of relationships for free metal ion  
547 activities and validation with independent data. *European Journal of Soil Science*, 61(1):58–  
548 73. WOS:000273454800006.

- 549 Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning: Data*  
550 *Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics. Springer-  
551 Verlag, New York, 2 edition.
- 552 Horn, A. L., Reiher, W., D Ring, R.-A., and G Th, S. (2006). Efficiency of pedotransfer  
553 functions describing cadmium sorption in soils. *Water, Air, and Soil Pollution*, 170(1-4):229–  
554 247.
- 555 Hough, R. L., Young, S. D., and Crout, N. M. J. (2003). Modelling of cd, cu, ni, pb and  
556 zn uptake, by winter wheat and forage maize, from a sewage disposal farm. *Soil Use and*  
557 *Management*, 19(1):19–27. WOS:000181552000003.
- 558 Kabata-Pendias, A. (2004). Soil–plant transfer of trace elements—an environmental issue.  
559 *Geoderma*, 122(2):143–149.
- 560 Knox, R. E., Pozniak, C. J., Clarke, F. R., Clarke, J. M., Houshmand, S., and Singh, A. K.  
561 (2009). Chromosomal location of the cadmium uptake gene (*cdu1*) in durum wheat. *Genome*,  
562 52(9):741–747. WOS:000271100900001.
- 563 Kuhn, M. and Johnson, K. (2013). Remedies for severe class imbalance. *Applied Predictive*  
564 *Modeling*, pages 419–443.
- 565 Laporte, M.-A., Sterckeman, T., Dauguet, S., Denaix, L., and Nguyen, C. (2015). Variability  
566 in cadmium and zinc shoot concentration in 14 cultivars of sunflower (*helianthus annuus*  
567 l.) as related to metal uptake and partitioning. *Environmental and Experimental Botany*,  
568 109(0):45–53.
- 569 Li, X. and Zhou, D. (2019). A meta-analysis on phenotypic variation in cadmium accumulation  
570 of rice and wheat: Implications for food cadmium risk control. *Pedosphere*, 29(5):545–553.
- 571 Li, Y.-M., Chaney, R., Schneiter, A., Miller, J., Elias, E., and Hammond, J. (1997). Screening  
572 for low grain cadmium phenotypes in sunflower, durum wheat and flax. *Euphytica*, 94(1):23–  
573 30.
- 574 Lin, Z., Schneider, A., Sterckeman, T., and Nguyen, C. (2016). Ranking of mechanisms govern-

575 ing the phytoavailability of cadmium in agricultural soils using a mechanistic model. *Plant*  
576 *and Soil*, 399(1):89–107.

577 Lux, A., Martinka, M., Vaculík, M., and White, P. J. (2011). Root responses to cadmium in  
578 the rhizosphere: A review. *Journal of Experimental Botany*, 62(1):21–37.

579 Maccaferri, M., Harris, N. S., Twardziok, S. O., Pasam, R. K., Gundlach, H., Spannagl, M.,  
580 Ormanbekova, D., Lux, T., Prade, V. M., Milner, S. G., Himmelbach, A., Mascher, M.,  
581 Bagnaresi, P., Faccioli, P., Cozzi, P., Lauria, M., Lazzari, B., Stella, A., Manconi, A., Gnoc-  
582 chi, M., Moscatelli, M., Avni, R., Deek, J., Biyiklioglu, S., Frascaroli, E., Corneti, S., Salvi,  
583 S., Sonnante, G., Desiderio, F., Marè, C., Crosatti, C., Mica, E., Özkan, H., Kilian, B.,  
584 Vita, P. D., Marone, D., Joukhadar, R., Mazzucotelli, E., Nigro, D., Gadaleta, A., Chao,  
585 S., Faris, J. D., Melo, A. T. O., Pumphrey, M., Pecchioni, N., Milanesi, L., Wiebe, K., Ens,  
586 J., MacLachlan, R. P., Clarke, J. M., Sharpe, A. G., Koh, C. S., Liang, K. Y. H., Taylor,  
587 G. J., Knox, R., Budak, H., Mastrangelo, A. M., Xu, S. S., Stein, N., Hale, I., Distelfeld, A.,  
588 Hayden, M. J., Tuberosa, R., Walkowiak, S., Mayer, K. F. X., Ceriotti, A., Pozniak, C. J.,  
589 and Cattivelli, L. (2019). Durum wheat genome highlights past domestication signatures and  
590 future improvement targets. *Nature Genetics*, 51(5):885–895.

591 Matynia, A., Lenoir, T., Causse, B., Spadini, L., Jacquet, T., and Manceau, A. (2010). Semi-  
592 empirical proton binding constants for natural organic matter. *Geochimica et Cosmochimica*  
593 *Acta*, 74(6):1836–1851.

594 Newman, M. C. (1993). Regression analysis of log-transformed data: Statistical bias and its  
595 correction. *Environmental Toxicology and Chemistry*, 12(6):1129–1133.

596 Oladzaad-Abbasabadi, A., Kumar, A., Pirseyedi, S., Salsman, E., Dobrydina, M., Poudel, R. S.,  
597 AbuHammad, W. A., Chao, S., Faris, J. D., and Elias, E. M. (2018). Identification and  
598 validation of a new source of low grain cadmium accumulation in durum wheat. *G3: Genes,*  
599 *Genomes, Genetics*, 8(3):923–932.

600 Peduzzi, P., Concato, J., Kemper, E., Holford, T. R., and Feinstein, A. R. (1996). A simulation  
601 study of the number of events per variable in logistic regression analysis. *Journal of Clinical*  
602 *Epidemiology*, 49(12):1373–1379.



603 Perrier, F., Yan, B., Candaudap, F., Pokrovsky, O. S., Gourdain, E., Meleard, B., Bussi re, S.,  
604 Coriou, C., Robert, T., Nguyen, C., and Cornu, J. Y. (2016). Variability in grain cadmium  
605 concentration among durum wheat cultivars: impact of aboveground biomass partitioning.  
606 *Plant and Soil*, 404(1-2):307–320.

607 Pozniak, C. J., Clarke, J. M., and Clarke, F. R. (2012). Potential for detection of marker–trait  
608 associations in durum wheat using unbalanced, historical phenotypic datasets. *Molecular*  
609 *Breeding*, 30(4):1537–1550.

610 Qiu, L., Wang, K., Long, W., Wang, K., Hu, W., and Amable, G. S. (2016). A comparative  
611 assessment of the influences of human impacts on soil cd concentrations based on stepwise  
612 linear regression, classification and regression tree, and random forest models. *PLOS ONE*,  
613 11(3):e0151131.

614 Rasmussen, C., Heckman, K., Wieder, W. R., Keiluweit, M., Lawrence, C. R., Berhe, A. A.,  
615 Blankinship, J. C., Crow, S. E., Druhan, J. L., Hicks Pries, C. E., Marin-Spiotta, E., Plante,  
616 A. F., Sch del, C., Schimel, J. P., Sierra, C. A., Thompson, A., and Wagai, R. (2018).  
617 Beyond clay: towards an improved set of variables for predicting soil organic matter content.  
618 *Biogeochemistry*, 137(3):297–306.

619 Saby, N., Marchant, B., Lark, R., Jolivet, C., and Arrouays, D. (2011). Robust geostatistical  
620 prediction of trace elements across france. *Geoderma*, 162(3–4):303–311.

621 Saby, N., Thioulouse, J., Jolivet, C., Rati , C., Boulonne, L., Bispo, A., and Arrouays, D.  
622 (2009). Multivariate analysis of the spatial patterns of 8 trace elements using the french soil  
623 monitoring network data. *Science of The Total Environment*, 407(21):5644–5652.

624 Salsman, E., Kumar, A., AbuHammad, W., Abbasabadi, A. O., Dobrydina, M., Chao, S.,  
625 Li, X., Manthey, F. A., and Elias, E. M. (2018). Development and validation of molecular  
626 markers for grain cadmium in durum wheat. *Molecular Breeding*, 38(3):28.

627 Sauv , S., Norvell, W. A., McBride, M., and Hendershot, W. (2000). Speciation and com-  
628 plexation of cadmium in extracted soil solutions. *Environmental Science & Technology*,  
629 34(2):291–296.

630 Schneider, A., Nguyen, V. X., Viala, Y., Violo, V., Cornu, J.-Y., Sterckeman, T., and Nguyen,  
631 C. (2019). A method to determine the soil-solution distribution coefficients and the con-  
632 centrations for the free ion and the complexes of trace metals: Application to cadmium.  
633 *Geoderma*, 346:91–102.

634 Six, L. and Smolders, E. (2014). Future trends in soil cadmium concentration under current  
635 cadmium fluxes to european agricultural soils. *Science of The Total Environment*, 485–  
636 486(0):319–328.

637 Sterckeman, T., Gossiaux, L., Guimont, S., Sirguey, C., and Lin, Z. (2018). Cadmium mass  
638 balance in french soils under annual crops: Scenarios for the next century. *Science of The*  
639 *Total Environment*, 639:1440–1452.

640 Sun, F., Polizzotto, M. L., Guan, D., Wu, J., Shen, Q., Ran, W., Wang, B., and Yu, G.  
641 (2017). Exploring the interactions and binding sites between cd and functional groups in soil  
642 using two-dimensional correlation spectroscopy and synchrotron radiation based spectromi-  
643 croscopies. *Journal of Hazardous Materials*, 326:18–25.

644 Tipping, E., Lofts, S., and Sonke, J. E. (2011). Humic ion-binding model vii: a revised param-  
645 eterisation of cation-binding by humic substances. *Environmental Chemistry*, 8(3):225–235.  
646 WOS:000291943800001.

647 Tudoreanu, L. and Phillips, C. J. (2004). Empirical models of cadmium accumulation in maize,  
648 rye grass and soya bean plants. *Journal of the Science of Food and Agriculture*, 84(8):845–852.

649 Vega, F., Andrade, M., and Covelo, E. (2010). Influence of soil properties on the sorption  
650 and retention of cadmium, copper and lead, separately and together, by 20 soil horizons:  
651 Comparison of linear regression and tree regression analyses. *Journal of Hazardous Materials*,  
652 174(1–3):522–533.

653 Viala, Y., Laurette, J., Denaix, L., Gourdain, E., Méléard, B., Nguyen, C., Schneider, A.,  
654 and Sappin-Didier, V. (2017). Predictive statistical modelling of cadmium content in du-  
655 rum wheat grain based on soil parameters. *Environmental Science and Pollution Research*,  
656 24(25):20641–20654.

- 657 Violante, A., Cozzolino, V., Perelomov, L., Caporale, A. G., and Pigna, M. (2010). Mobility  
658 and bioavailability of heavy metals and metalloids in soil environments. *Journal of soil  
659 science and plant nutrition*, 10(3):268–292.
- 660 Wiebe, K., Harris, N., Faris, J., Clarke, J., Knox, R., Taylor, G., and Pozniak, C. (2010).  
661 Targeted mapping of *Cdu1*, a major locus regulating grain cadmium concentration in durum  
662 wheat (*Triticum turgidum* l. var *durum*). *TAG Theoretical and Applied Genetics*, 121(6):1047–  
663 1058.
- 664 Zimmerl, S., Lafferty, J., and Buerstmayr, H. (2014). Assessing diversity in triticum durum  
665 cultivars and breeding lines for high versus low cadmium content in seeds using the caps  
666 marker usw47. *Plant Breeding*, 133(6):712–717.
- 667 Álvaro, F., Isidro, J., Villegas, D., Moral, G. d., F, L., and Royo, C. (2008). Breeding effects  
668 on grain filling, biomass partitioning, and remobilization in mediterranean durum wheat all  
669 rights reserved. no part of this periodical may be reproduced or transmitted in any form or by  
670 any means, electronic or mechanical, including photocopying, recording, or any information  
671 storage and retrieval system, without permission in writing from the publisher. *Agronomy  
672 Journal*, 100(2):361–370.

Table 1: Characteristics of the soils used for modeling grain Cd in durum wheat. SOC : soil organic carbon, Calc: soil calcareous, Cd: total soil Cd ( $\text{mg kg}^{-1}$ ). All variables but soil Cd are in  $\text{g kg}^{-1}$ . For calcareous, the values below the limit of quantification ( $<0.5 \text{ g kg}^{-1}$ ) were set to  $0.5 \text{ g kg}^{-1}$ .

	Clay	Fine silt	Coarse silt	Fine sand	Coarse sand	SOC	pH	Calc	Cd
Min	34	15	18	3	2	3.0	5.5	0.5	0.05
Q25	245	209	109	54	45	10.1	8.0	9.0	0.20
Median	300	258	139	122	104	13.3	8.2	106.8	0.26
Q75	381	310	202	181	180	18.1	8.4	276.5	0.36
Max	702	467	427	499	452	45.2	8.7	655.4	1.56

Table 2: Occurrence of observed non-conformity of grain Cd for the three regulatory thresholds of 0.2 (RT02), 0.15 (RT015) and 0.1 (RT01) mg Cd kg<sup>-1</sup> grain and for the cultivars used for modeling. Value are in % of total data (n). *Cdu1* column indicates if the cultivar has the high Cd allele of the *Cdu1* gene or if it is unknown.

Cultivar	RT02	RT015	RT01	Cdu1	n
ACTISUR	50	50	75	High Cd	4
ALEXIS	29	43	62	High Cd	21
ATOUDUR	8	8	31	High Cd	26
AURIS	33	33	67	High Cd	6
AVENTUR	40	40	60	Unknown	5
BABYLONE	8	8	31	High Cd	13
BIENSUR	0	0	0	High Cd	2
CLOVIS	0	50	50	High Cd	2
COUSSUR	0	33	50	High Cd	6
CULTUR	38	46	62	High Cd	13
DAKTER	17	25	42	High Cd	12
DAURUR	50	50	100	High Cd	2
FABULIS	10	20	60	High Cd	10
FLORIDOU	33	33	56	Unknown	9
ISILDUR	0	8	42	High Cd	12
JOYAU	0	0	33	High Cd	3
KARUR	20	36	64	High Cd	25
LIBERDUR	0	100	100	High Cd	1
LUMINUR	62	62	75	High Cd	8
MEMODUR	100	100	100	Unknown	1
MIRADOUX	7	14	25	High Cd	85
MURANO	0	0	100	High Cd	1
NEFER	0	0	50	High Cd	2
NEMESIS	100	100	100	Unknown	1
PESCADOU	25	32	57	High Cd	28
PHARAON	0	100	100	High Cd	1
PICTUR	17	50	67	High Cd	6
QUALIDOU	20	20	40	Unknown	15
RELIEF	12	31	69	Unknown	16
SANTUR	50	50	50	Unknown	2
SCULPTUR	17	26	43	High Cd	35
SY BANCO	40	40	80	High Cd	5
SY CYSCO	20	20	53	High Cd	15
SY ENZO	0	100	100	Unknown	1
TABLUR	20	32	60	High Cd	25
YELODUR	0	0	0	High Cd	1
All	18	26	48		420

Table 3: Performances of the absolute and parsimonious best models for the three regulatory thresholds of 0.2 (RT02), 0.15 (RT015) and 0.1 (RT01) mg Cd kg<sup>-1</sup> grain. FS: fine silt, CS: coarse silt, SOC : soil organic carbon, Calc: soil calcareous, RF: random forest, MELR,: mixed-effect logistic regression, MEMR: mixed-effects multi-linear regression, J: Youden statistics, MCC: Matthew’s correlation coefficient, TPR: true positive rate (% of actual positive cases that are detected), TNR: true negative rate (% of actual negative cases that are detected), PPV: positive predictive value (% of predicted positive cases that are actually positive), NPV: negative predictive value (% of predicted negative cases that are actually negative)

Regulation limit	Model	Modeling approach	Optimization Criteria	TPR	TNR	PPV	NPV				
RT02	Cd+H+SOC+Clay+FS+CS+Calc	RF	J	85.1	(2.23)	91.0	(0.92)	67.0	(2.33)	96.6	(0.49)
RT02	Cd+H	MELR	J	83.2	(1.66)	92.0	(0.62)	69.2	(1.72)	96.2	(0.36)
RT02	Cd+H+SOC+Clay+Calc	RF	MCC	83.3	(2.37)	92.0	(0.86)	69.1	(2.37)	96.2	(0.51)
RT02	Cd+H	MELR	MCC	83.2	(1.71)	91.9	(0.63)	68.9	(1.77)	96.2	(0.37)
RT015	Cd+H+SOC+Clay+FS+Calc	RF	J	81.8	(1.93)	89.4	(1.06)	72.8	(2.06)	93.4	(0.65)
RT015	Cd+H+SOC+Clay	RF	J	81.2	(1.76)	88.9	(1.11)	71.8	(2.11)	93.2	(0.6)
RT015	Cd+H+SOC+Clay+FS+CS	RF	MCC	81.0	(1.85)	90.0	(1.06)	73.6	(2.11)	93.2	(0.62)
RT015	Cd+H+SOC+Clay	RF	MCC	81.1	(1.8)	88.9	(1.07)	71.7	(1.99)	93.2	(0.61)
RT01	Cd+H+SOC+Clay+FS+CS+Calc	MEMR		78.7	(0.54)	84.9	(0.63)	82.4	(0.6)	81.6	(0.38)
RT01	Cd+H+SOC+Clay+FS	MEMR		78.4	(0.83)	83.5	(0.78)	81.0	(0.74)	81.2	(0.6)

# 1 Figure captions

Figure 1: Performances of the 9 models (y-axis) x 3 modeling approaches (different colors) for the three regulatory thresholds of 0.2 (RT02), 0.15 (RT015) and 0.1 (RT01) mg Cd kg<sup>-1</sup> grain. RF: random forest, MELR: mixed-effect logistic regression, MEMR: mixed-effect multi-linear regression. The model performances are expressed by the Youden statistics (J) or Matthew's correlation coefficient (MCC), which have both 1 as maximum value for perfect models. From right to left, the dashed lines correspond to the maximum J or MCC values (best model in absolute) and to 95% of this maximum. Letters on the right of the bars are the mean grouping by the Tukey test at p<0.05. The orange letter is the best parsimonious model, namely the model with the least predictors while having performances greater than 95% of the best model in absolute (see materials and methods).

Figure 2: Simulations of the effect of individual soil variables on the probability of non-conformity of durum wheat grain for the three regulatory thresholds of 0.2 (RT02), 0.15 (RT015) and 0.1 (RT01) mg Cd kg<sup>-1</sup> grain and for the highest (Relief) and the lowest (Miradoux) Cd accumulator cultivar of the database. The simulations were obtained from the random forest model with the following soil predictors : total soil Cd, soil pH, soil organic carbon, clay, fine and coarse silt and soil calcareous. For a given value of a soil predictor, the graphs show the frequency of non-conformity for all predictions when the other predictors vary based on a factorial design (see materials and methods).

Figure 3: Simulations of the effect of soil Cd and pH on the probability of non-conformity of durum wheat grain for the three regulatory thresholds of 0.2 (RT02), 0.15 (RT015) and 0.1 (RT01) mg Cd kg<sup>-1</sup> grain and for the highest (Relief) and the lowest (Miradoux) Cd accumulator cultivar of the database. The simulations were obtained from the random forest model with the following soil predictor : total soil Cd, soil pH, soil organic carbon, clay, fine and coarse silt and soil calcareous. For a given value of a soil predictor, the graphs show the frequency of non-conformity for all predictions when the other predictors vary based on a factorial design (see materials and methods).

Figure 4: Boxplots of the predicted grain Cd content of the different cultivars of the modeling database for the regulatory thresholds of 0.1 mg Cd kg<sup>-1</sup>(RT01) and by using the mixed-effects

30 multi-linear regression with the following soil predictors : total soil Cd, soil pH, soil organic carbon, clay, fine and coarse silt and soil calcareous. The model predicts the grain Cd content when the predictor values are set to their mean in the database and for 400 repetitions of the 5 folds cross-validation. Points are outliers extending outside 1.5 x the inter-quartile range (whiskers).

35 Figure 5: Adjusted least squared means  $\pm$  one standard deviation for the grain Cd contents of some cultivars grown in three field trials for three years. Letters and colors correspond to mean grouping by the Tukey test at  $p < 0.05$ . The trial location and year were considered as random effects whereas the cultivar was the fixed effect.

Figure 6: Ranking of some durum wheat cultivars for their grain Cd content by two approaches :  
40 data collected from field trials (y axis) and ranking from the mixed-effects multi-linear regression with the following soil predictors : total soil Cd, soil pH, soil organic carbon, clay, fine and coarse silt and soil calcareous. High rank values indicate high grain Cd as illustrated by the 1:1 arrow.



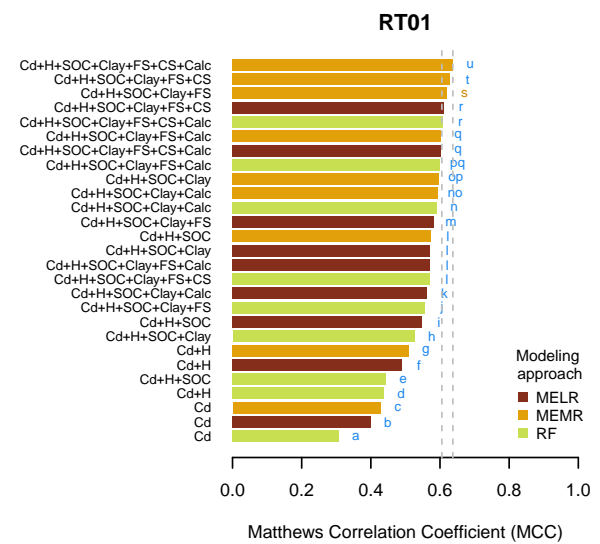
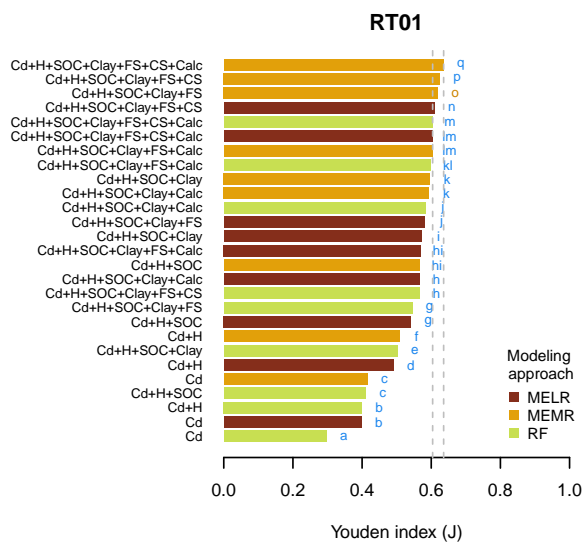
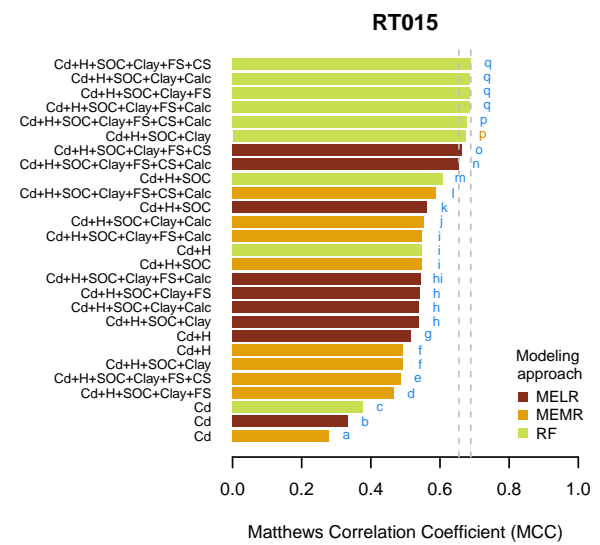
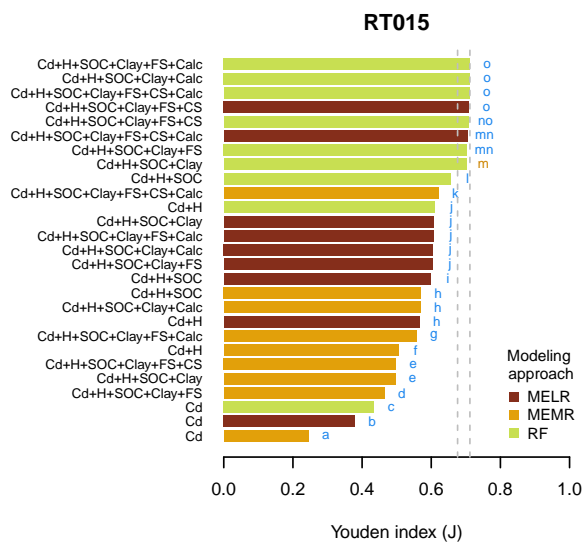
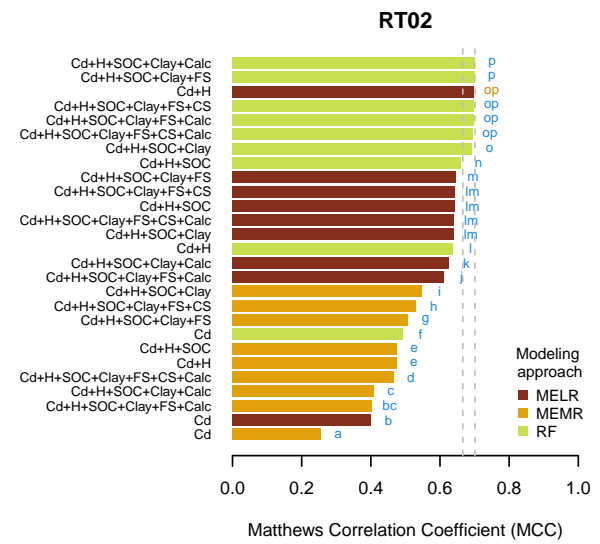
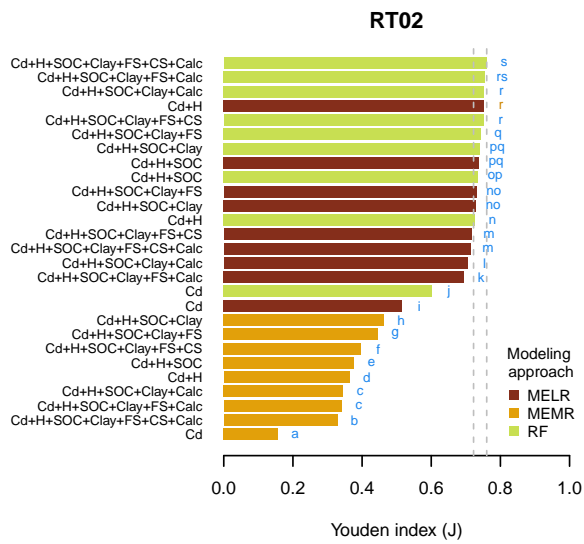


Figure 1:

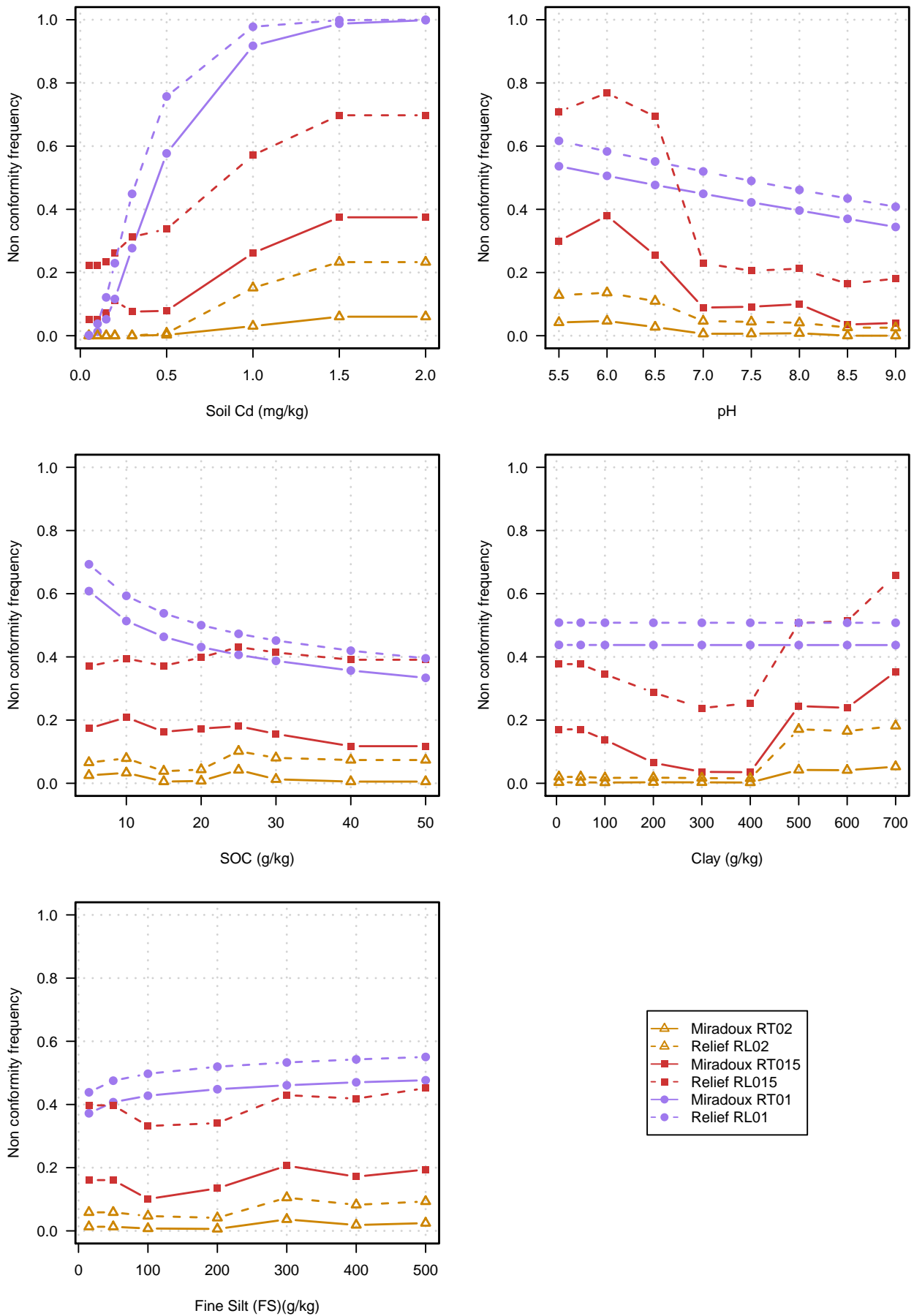


Figure 2:

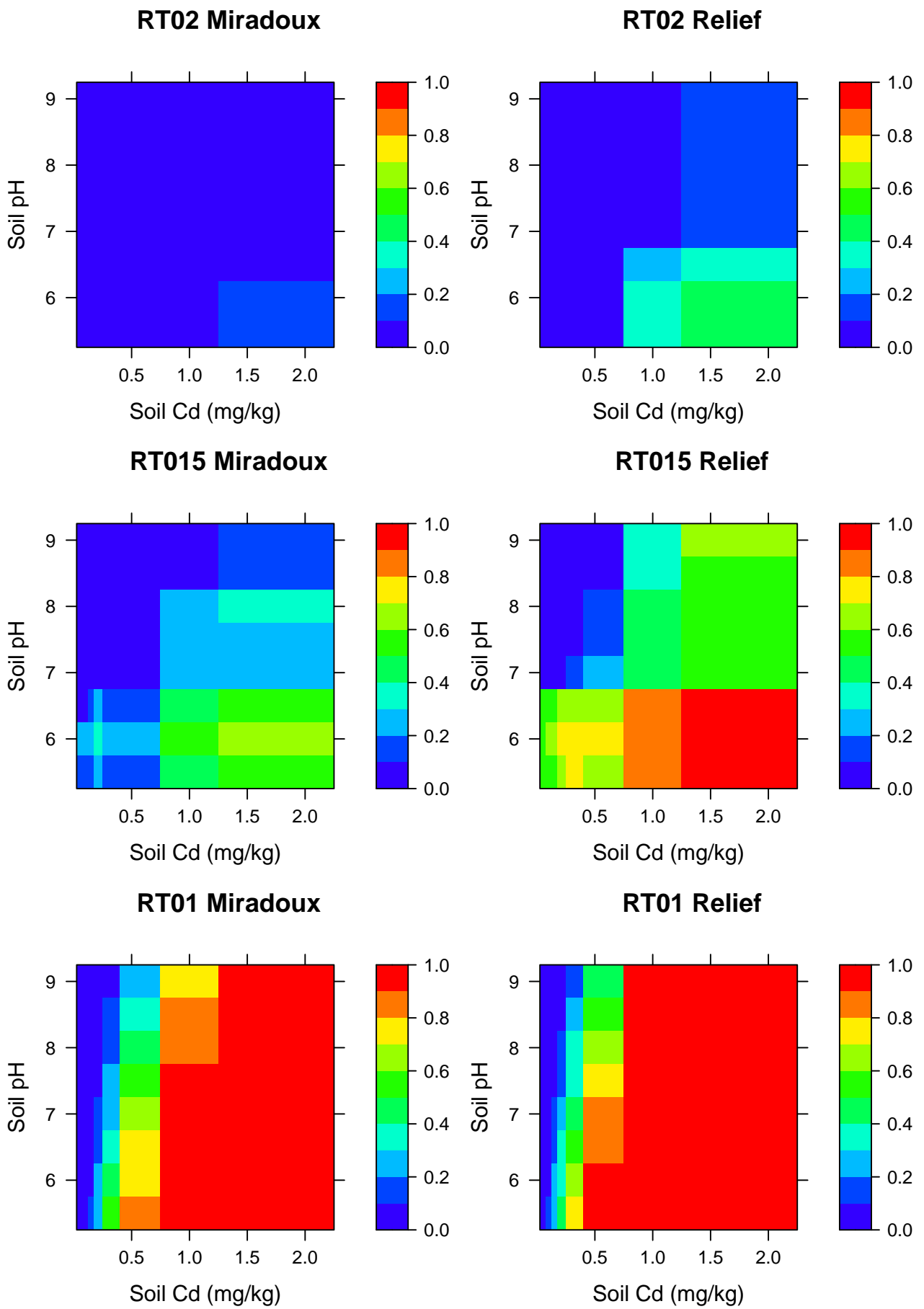


Figure 3:

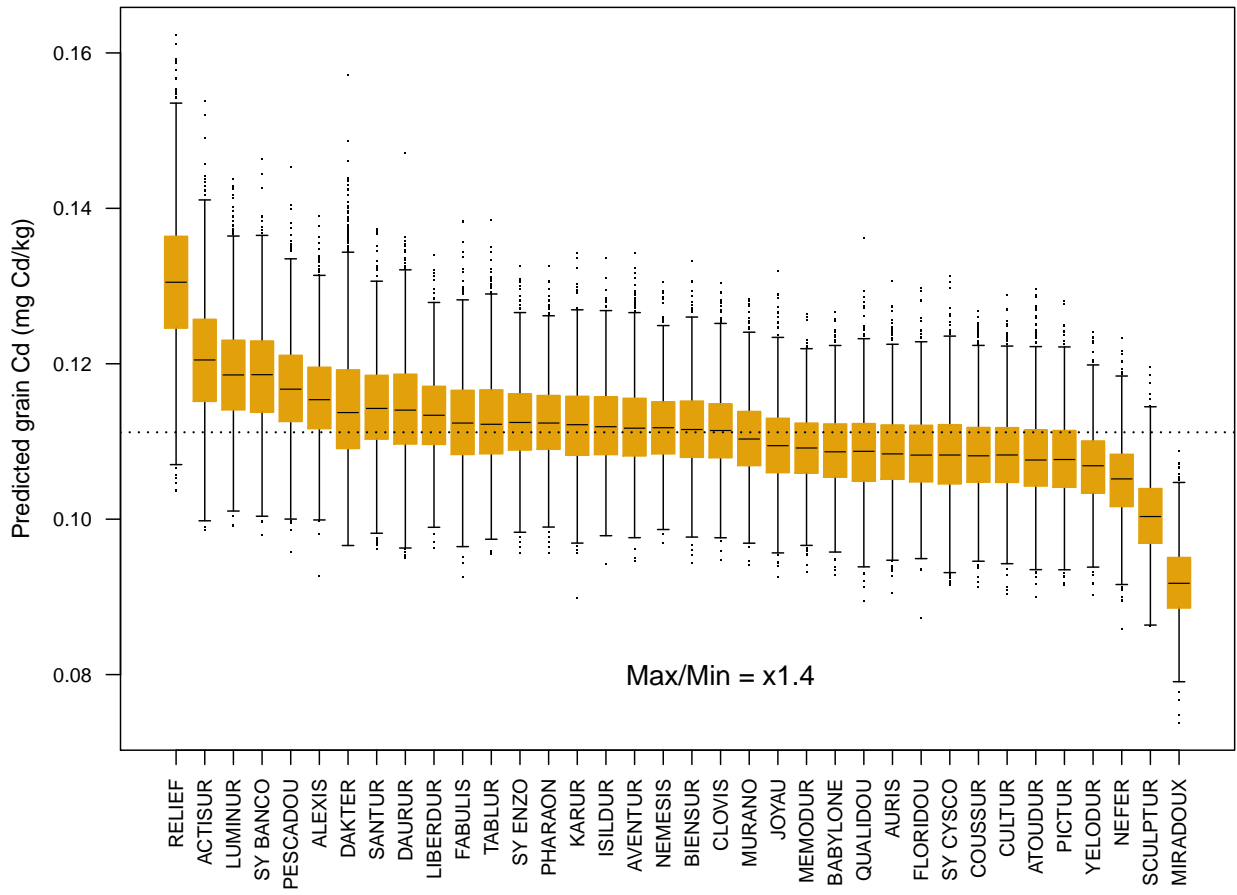


Figure 4:

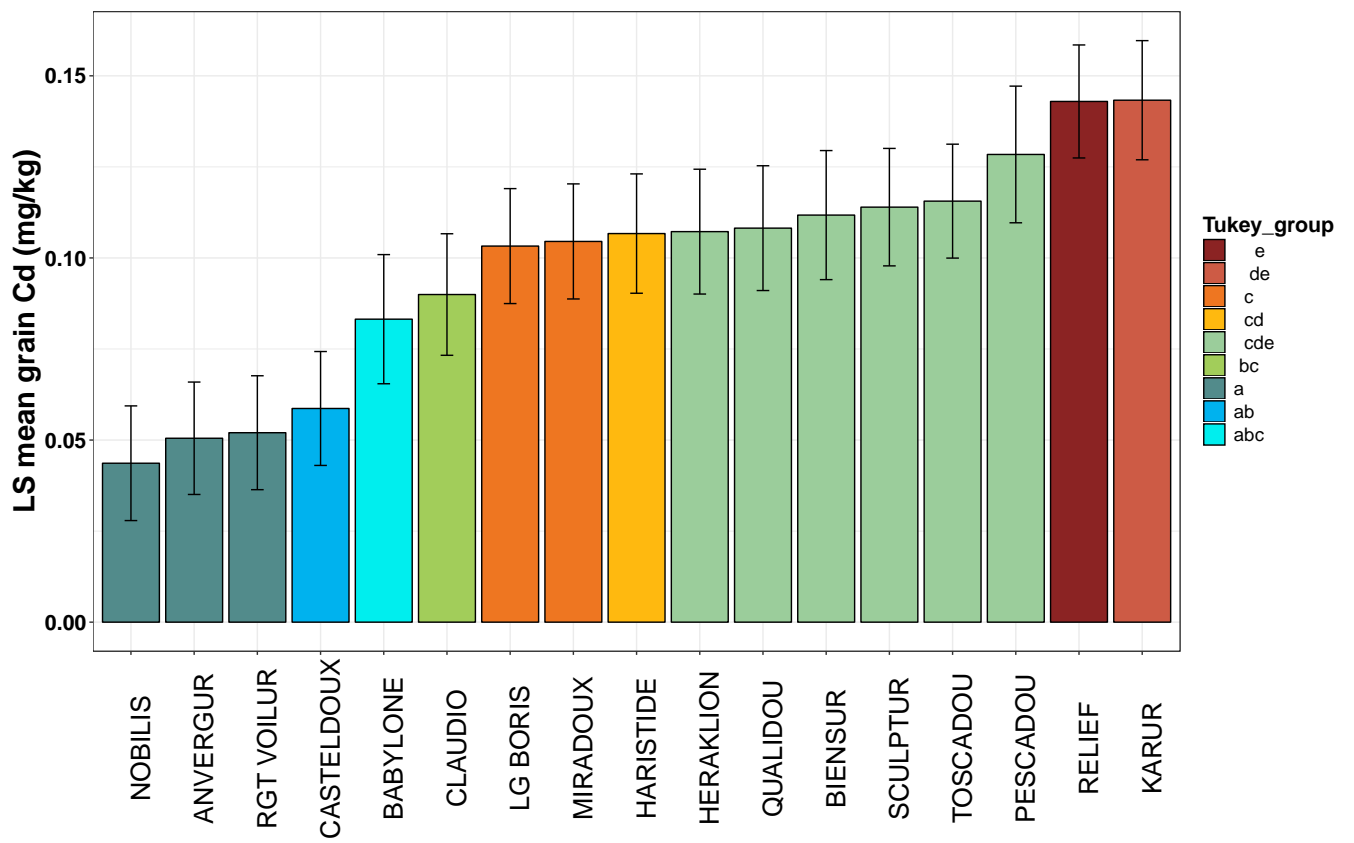


Figure 5:

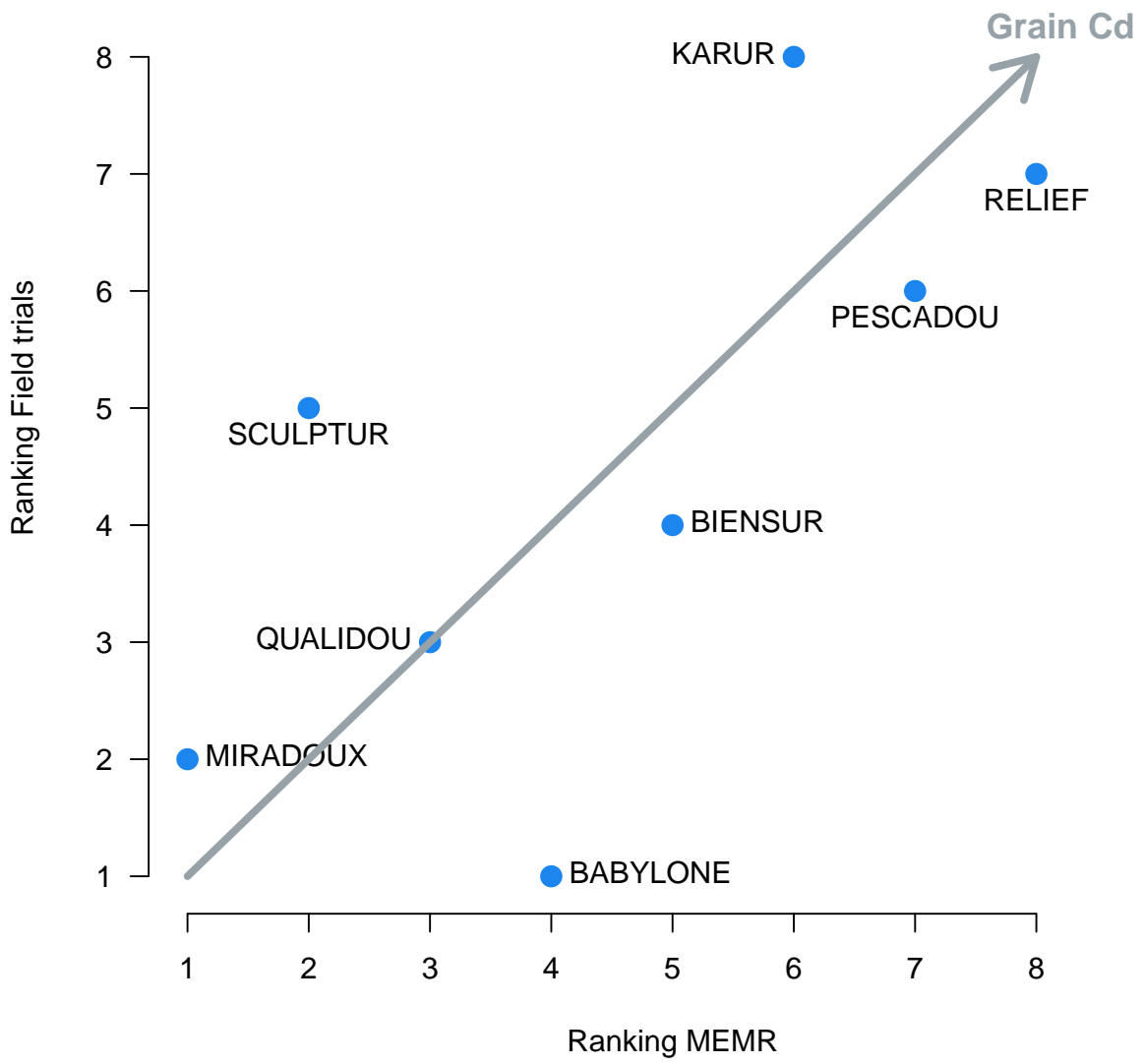
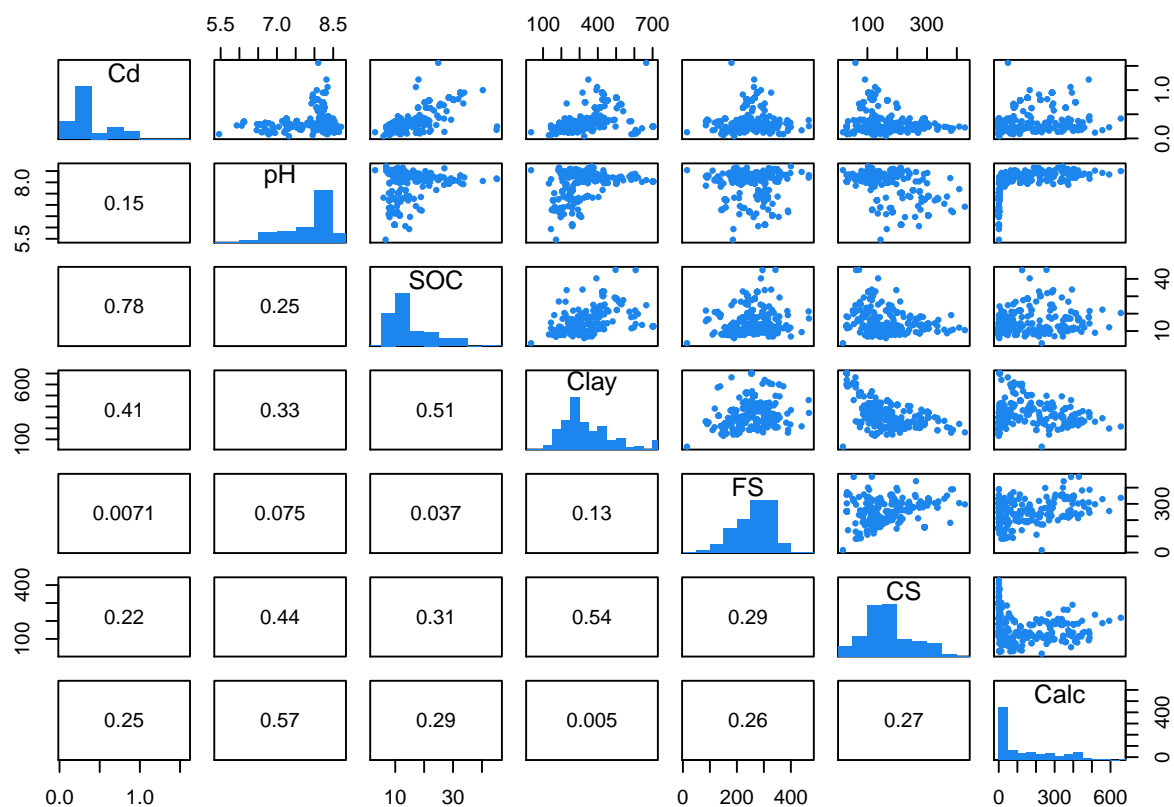


Figure 6:

## Supplementary Information



45

Figure SI1 : Matrix plot of correlations between the characteristics of the soils used for modeling the conformity of durum wheat grain Cd content. Cd: total soil Cd content, FS: fine silt, CS: coarse silt, SOC : soil organic carbon, Calc: soil calcareous. The bars on the diagonal show the distribution of the variable. Values below the diagonal are the Pearson correlation coefficients

50 with two significant digits.

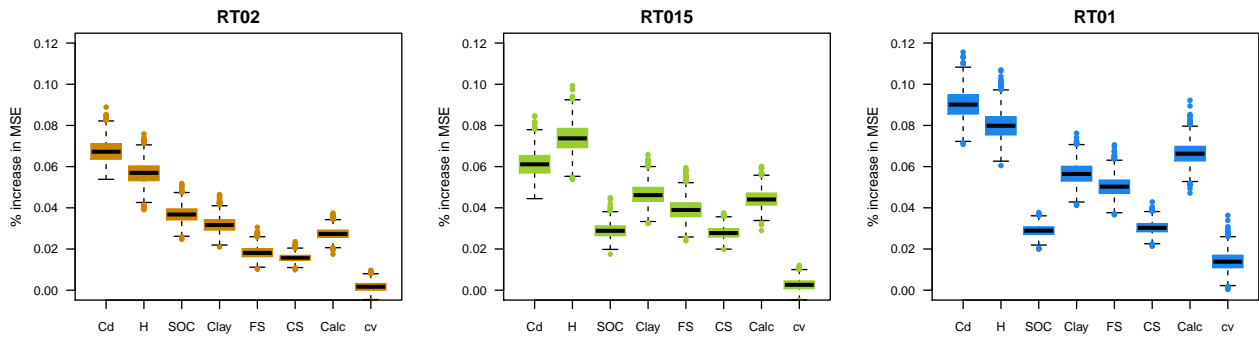


Figure SI2: Predictor importance estimated from the Random forest model with the following soil predictors : total soil Cd, soil pH, soil organic carbon, clay, fine and coarse silt and soil calcareous. The y-axis is the increase in the mean squared error of the model when the data of a given variable are shuffled while keeping other the variables unchanged.