

1 Quantifying individual and collective influences of soil properties on 2 crop yield

3 Rebecca Whetton^a, Yifan Zhao^{b,*}, Abdul M Mouazen^{a,c}

4 ^a *Cranfield Soil and AgriFood Institute, Cranfield University, Bedfordshire MK43 0AL, UK.*

5 ^b *Through-life Engineering Services Centre, Cranfield University, Bedfordshire MK43 0AL, UK*

6 ^c *Department of Soil Management, Ghent University, Coupure 653, 9000 Gent, Belgium*

7 **E-mail of corresponding author: Abdul.Mouazen@UGent.be*

8 **Abstract**

9 Quantifying the agronomic influences of soil properties, collected at high sampling
10 resolution, on crop yield is essential for site specific soil management. This study
11 implements a novel Volterra Non-linear Regressive with eXogenous inputs (VNRX-
12 LN) model, to quantify causal factors to explain yield using high resolution data on key
13 soil properties affecting wheat yield in a 22 ha field with waterlogging problem in
14 Bedfordshire, UK. A total of eight soil properties including total nitrogen (TN), organic
15 carbon (OC), pH, available phosphorous (P), magnesium (Mg), calcium (Ca), moisture
16 content (MC), and cation exchange capacity (CEC) were collected with an on-line
17 (tractor mounted) visible and near infrared spectroscopy (vis-NIR) sensor and used as
18 multiple-input to the VNRX-LN model, while crop yield represented the single-output
19 in the system.

20 Results showed that the largest contributors to wheat yield were CEC, Mg and TN, with
21 error reduction ratio contribution (ERRC) values of 14.6%, 4.69% and 1%, respectively.
22 The overall contribution (SEER) of the soil properties considered in this study totals a
23 value of 23.21%. This was attributed to a large area of the studied field having been

24 waterlogged, which masked the actual effect of soil properties on crop yield. It is
 25 recommended to validate the introduced concept on a larger number of fields, where
 26 other crop yield affecting parameters e.g., crop disease, pests, drainage, topography and
 27 microclimate conditions should be taken into account.

28 **Keywords:** Yield limiting factors; proximal soil sensing; VNRX; nonlinear parametric
 29 modelling.

30 **Table of abbreviation**

AFOLS	Adaptive-forward-orthogonal least squares
Ca	Calcium
CAA	Circle-based average approximation
CEC	Cation exchange capacity
DGPS	Differential global positioning system
ERR	Error reduction ratio
ERRC	Error reduction ratio contribution
LAI	Leaf area index
MC	Moisture content
Mg	Magnesium
NDVI	Normalised difference vegetation index
NFIR	nonlinear finite impulse response
OC	Organic carbon
OLS	Orthogonal least squares
P	Phosphorus
SDA	Shortest distance approximation
SERR	Sum of error reduction ratio
TN	Total nitrogen
vis-NIRS	Visible and near infrared
UTM	Universal transverse Mercator
VNRX	Volterra non-linear regressive with eXogenous
VNRX-LN	Volterra non-linear regressive with eXogenous, accounting for both linear and non-linear variability

31

32

33 **1. Introduction**

34 The world's population is expected to rise to 9 billion by 2050 and, based on the current
35 land available, an increase in crop yield of 60% will be required. Precision management
36 of farm resources (e.g., fertilisers, seeds, water, etc.) is one potential way to increase
37 crop yield. The spatial variability in agricultural fields exists at different spatial scales
38 (Raun 1998; Dhillon *et al.* 1994), which requires precise management with the aim to
39 increase yield at reduced input cost and related environmental impacts. This is hardly
40 achievable by conventional agriculture that relies on homogeneous applications of
41 external inputs. For example, current fertiliser applications are made based on a bulked
42 composite soil sample collected per field or 1-3 ha in the best scenario, which ignore
43 within field variability. This may result in over-application in rich zones, and under-
44 applications in poor zones in the field. In this context, recent years have seen a surge of
45 variable-rate application technologies where external farm inputs are applied in
46 response to input data from normalised difference vegetation index (NDVI), leaf area
47 index (LAI), high resolution soil properties or a combination of these (Lowenberg-
48 DeBoer and Aghib 1999; Maleki *et al.* 2008; Mouazen *et al.* 2009; Halcro *et al.* 2013;
49 Mouazen and Kuang 2016). Although variable rate fertilisation is a strategy to increase
50 crop yield, understanding and quantifying the yield limiting factors is still a crucial
51 research question to be answered, before variable rate applications can be optimised.

52 Since spatial variability in the majority of agricultural fields exist, proximal sensor
53 technologies are invaluable to measure this variability accurately. This will require
54 robust and reliable sensing platforms of crop and soil. Proximal (e.g., Crop Circle ACS
55 470, Holland Scientific, Lincoln, NE USA) and remote sensing (e.g., satellite imagery,

56 unmanned aerial vehicles or aircrafts) both can provide high resolution data on crop
57 canopy characteristics indicated e.g., as NDVI or LAI (Mulla 2013; Kipp *et al.* 2014)
58 and they are commercially available. However, remote sensing methods based on
59 spectral reflectance provide data on the top millimetres of soil and require a bare soil
60 surface. Furthermore, due to the complex nature and vast variability of agricultural
61 soils, the majority of proximal soil sensors are still premature to fulfil this requirement.
62 Kuang *et al.* (2012) concluded in an extensive review that the most promising proximal
63 sensing technologies for quantifying soil properties are electrochemical technique and
64 optical visible and near infrared (vis-NIR) spectroscopy. Although they are limited to
65 particular research groups worldwide, on-line (tractor mounted) vis-NIR sensors
66 (Shibusawa *et al.* 2001; Mouazen *et al.* 2006; Christy 2008) enable the collection of
67 high sampling resolution (e.g., >500 samples per ha) of key soil properties (Kuang *et al.*
68 2012; Kuang and Mouazen 2013; Marin-González *et al.* 2013; Kodaira and Shibusawa
69 2013; Kweon *et al.* 2013), which are valuable sources of information to manage the
70 within field spatial variability.

71 Nonlinear parametric modelling approaches offer novel tools for the quantification and
72 better understanding of the influences of soil related yield limiting factors, collected at
73 high sampling resolution with on-line soil sensors, which cannot be obtained with the
74 traditional soil sampling and laboratory analytical methods. One of these parametric
75 methods is Volterra Non-linear Regressive with eXogenous inputs (VNRX-LN) model,
76 which was broadly used in the engineering sector, but not common in agriculture.

77 The aim of this work was to use the VNRX-LN model to quantify causal factors to
78 explain yield using high resolution data on key soil properties affecting wheat yield in a
79 22 ha field with waterlogging problem in Bedfordshire, UK.

80 **2. Materials and Methods**

81 *2.1 Study site*

82 The study site was one field designated as Horns End, and located at a commercial
83 farm, called Duck end farm, in Wilstead, Bedfordshire UK (52°5'52.087''W latitude and
84 0°27'19.76''N longitude). The field is about 22 ha area, with an average annual rainfall
85 of 598 mm. According to the UK meteorology Office
86 (<http://www.metoffice.gov.uk/climate/uk/summaries>), May was particularly wet in
87 2013, and spring was cooler than average, whilst summer was the driest for the UK
88 since 2003. Nevertheless, there were some notably wet days, particularly in July and
89 August. The farm has a crop rotation of barley (*Hordeum vulgare*), wheat (*Triticum*
90 *aestivum*) and oil seed rape (*Brassica napus*). The soil texture over the field down to
91 0.20 m is non-homogeneous, including three textures of sandy loam, loam, and sandy
92 clay loam according to the United State Department of Agriculture (USDA) texture
93 classification system. Wheat was cultivated during the experiment in 2013.

94 *2.2 On-line collected data*

95 The on-line vis-NIR sensor (Mouazen 2006) was used (Figure 1) to carry out the field
96 measurement. It consists of a subsoiler that penetrates the soil to the required depth,
97 making a trench, whose bottom is smoothed due to the downwards forces acting on
98 the subsoiler (Mouazen *et al.* 2005).

99

[Figure 1]

100 The optical probe, housed in a steel lens holder, was attached to the rear of the subsoiler
101 chisel to acquire soil spectra in reflectance mode from the smooth bottom of the trench.
102 The subsoiler, retrofitted with the optical unit, was attached to a frame that was
103 mounted onto the three point hitch of the tractor. An AgroSpec mobile, fibre type, vis-
104 NIR spectrophotometer (tec5 Technology for Spectroscopy, Germany) with a
105 measurement range of 305-2200 nm was used to measure soil spectra in diffuse
106 reflectance mode. A differential global positioning system (DGPS) (EZ-Guide 250,
107 Trimble, USA) was used to record the position of the on-line measured spectra with
108 sub-metre accuracy. On-line soil measurement occurred in summer 2012 after the
109 harvest of the previous crop, at parallel transects of 15 m space, with an average
110 forward speed of the tractor of 2 km h⁻¹ and the measurement depth set at 150 mm. A
111 few on-line collected vis-NIRS spectra are shown in Figure 2, as an example.

112

[Figure 2]

113 During on-line measurement, two or three soil samples per line were collected from the
114 bottom of a trench and the sampling positions were carefully recorded with the DGPS.
115 These samples were analysed for calcium (Ca), magnesium (Mg), cation exchange
116 capacity (CEC), phosphorous (P), pH, moisture content (MC), organic carbon (OC) and
117 total nitrogen (TN), using the following laboratory analytical methods:

- 118 • Exchangeable Ca and Mg were determined by Agilent 240 FS AA atomic
119 absorption spectrophotometry (Agilent Technologies, Inc. USA).
- 120 • CEC was determined using a Flame Photometer (Chapman 1965).

- 121 • Available P concentration was determined by an ascorbic acid method (Olsen *et*
122 *al.* 1954).
- 123 • pH was measured potentiometrically on a suspension of soil to water ratio
124 (1:2.5) (DEFRA 1986).
- 125 • MC was determined by oven drying of samples at 105° for 24 h.
- 126 • OC was determined using a combustion method (British Standard BS 7755
127 Section 3.8 1995).
- 128 • TN was determined by the Dumas method, where soil samples are heated to 900
129 °C in the presence of oxygen gas (British Standard BS EN 13654-2:2001).

130 The selection of these eight soil properties was attributed to the fact that these
131 properties are considered important in explaining crop yield response and can be
132 measured with the on-line vis-NIRS sensor with appreciable accuracy (Kuang and
133 Mouazen 2013; Marin-González *et al.* 2013).

134 Partial least squares regression (PLSR) based calibration models, developed with
135 Unscrambler V9.8 software (Camo Software, Norway) were used to predict all eight
136 soil properties using the on-line collected soil spectra (>500 samples per ha). The on-
137 line prediction accuracy of properties with direct spectral responses (i.e., MC, OC and
138 TN) indicated as residual prediction deviation (the ration of standard deviation divided
139 by root mean square error of prediction (RMSEP) ranged between 1.96 and 3.06 (good
140 to excellent predictions). For the soil properties with indirect spectral responses (i.e.,
141 Ca, Mg, CEC, P and pH), RPD ranged between 1.30 and 2.14 (moderate to good
142 predictions). More details about the on-line vis-NIR sensor and accuracy of

143 measurement can be found in Kuang and Mouazen (2013) and Marin-González *et al.*
144 (2013).

145 Wheat yield data was collected in August, 2013 by the on-board yield sensor and GPS
146 system of the farmer's combine harvester (New Holland, CX8070 model), with a header
147 width of 7.25 m commonly used for barley and wheat harvest. In addition, the harvest
148 was optimised to: I) record wheat yield when the machine header was full for the full
149 length of the study area, and II) avoid the bare soil in the tramlines. Total yield was
150 calculated from the mean yield (tonnes per hectare) of an area, multiplied by the size of
151 the area (m²), which was derived using ArcGIS (Esri, USA).

152 *2.3 Data processing*

153 Features in the environment, are the product of many interacting processes, including
154 physical, chemical and biological. They are determined with exceedingly complex
155 interactions, which along with incomplete understanding can make the occurrence seem
156 random. Due to this, a way of overcoming the prediction of distribution is to treat the
157 variation as if it is random (Matheron 1963). The measurement points from the on-line
158 soil sensor and yield sensor required a method of interpolation, to provide a continuous
159 data set across the locations. Kriging was selected as a non-biased approach to predict
160 the values between the sample points, where semi-variograms were first produced and
161 then applied in Kriging predictions. The interpolated data were then converted into a
162 common 5 m raster grid in ArcGIS (Esri, USA) in order to assist data fusion (Frogbrook
163 and Oliver 2007). The raster squares of the layers were converted into this common grid
164 of points by extracting the value at the midpoint of each raster square. A smaller
165 resolution has no practical implementation, due to the limitations of the size and

166 response time of the precision farming equipment. The 5 m grid size provided a balance
167 between adequately characterising the spatial variation and practical farm management.
168 These steps ensured that all layers consisted of a common set of 5 m grid point-values,
169 to allow the application of parametric modelling to be carried out. This method allowed
170 data from a diverse range of soil and crop property surveys, measured at different
171 resolutions, to be merged (Khosla *et al.* 2008). The different soil and crop layers of a 5
172 by 5 m grid were subjected to the VNRX-LN detailed in the following section.

173 *2.4 Volterra Non-linear Regressive with exogenous Model*

174 In this study, the simplified VNRX-LN model, also known as NFIR model, was
175 implemented, which represent a multi-inputs and single-output system:

176

177

$$178 \quad y = f(u_1, u_2, \dots, u_R) + \varepsilon \quad (1)$$

179

180

181 where R is the number of the system inputs, f is some unknown linear or non-linear
182 mapping, which links the system output y to the system inputs u_1, u_2, \dots, u_R ; ε denotes
183 the model residual.

184 The on-line measured soil properties (i.e., TN, OC, pH, P, Mg, Ca, MC, and CEC) were
185 normalised and used as inputs ($R = 8$) to the VNRX-LN model, whereas the model
186 output was wheat yield. The analysis also included the interaction between pairs of soil

187 properties and their contribution to crop yield. The aim was to investigate the
188 contribution of each soil property and their pairwise interaction on crop yield.

189 Parameters are estimated based on the observations, and these are determined by the
190 structure, using the orthogonal least squares (OLS) estimation procedures. Adaptive-
191 forward-orthogonal least squares (AFOLS) was employed not only to determine the
192 model structure but also to estimate the unknown parameters. More detailed description
193 of this method can be found in Zhao *et al.* (2012).

194 Performance of VNRX-LN model output was evaluated by considering the value of
195 error reduction ratio (ERR) for each parameter to the prediction of yield (system
196 outputs). Values of ERR always range from 0% to 100%. The larger the ERR is, the
197 higher the dependence is between this term and the output. It is, therefore, a useful
198 index to indicate the contribution of each term to the output. To calculate the
199 contribution of each input variable to the output, the sum of ERR values (SERR) of all
200 selected terms is used to describe the percentage explained by the identified model to
201 the system output. If the considered inputs can fully explain the variation of system
202 output, the value of SERR is equal to 100%. It is an indicator of model performance and
203 uncertainty. The contribution of the i^{th} input variable to the variation of the system
204 output, denoted as $ERRC_i$, is defined as the sum of ERR values of the terms that include
205 this input variable. The value of $ERRC_i$ should be always between 0% and 100%.

206 2.5 Significance Test

207 To determine the statistical significance of the contribution from each input to the
208 system output, a threshold τ_i , representing the level of contribution, above which value

209 had less than a 5% probability of occurring by chance, requires being determined. The
210 conventional 95% confidence interval is not suitable for this study because the
211 distribution of ERRC value is unknown. For this purpose, the following surrogate data
212 technique was used.

213 Assuming the signal Y is a function of the signal X , this sort of dependence is destroyed
214 when Y is ordered randomly in some way while X keeps the same order. For this
215 purpose, the order of the data in Y was randomised by a shuffle procedure that saves the
216 distribution properties of the Y signal, but destroys the spatial relationship between X
217 and Y . This procedure was repeated 100 times and then the 95% quantile was
218 determined as the threshold. A significance threshold for each term is firstly calculated,
219 and then the significance threshold for each input can then be derived by the same way
220 to calculate $ERRC_i$.

221 *2.6 Optimal spatial resolution of soil properties versus yield*

222 Since the spatial sampling resolutions of soil properties and crop yield are different,
223 before applying the proposed VNRX-LN modelling method, the data must be re-
224 sampled to establish the correspondence between the inputs and the output. Two re-
225 sampling techniques have been used in this study. In the first technique, for each crop
226 yield data $y(e_i, w_i)$ on a location (e_i, w_i) , the corresponding soil properties were
227 approximated by the properties on the location that has the shortest distance to (e_i, w_i) ,
228 which must be smaller than a radius r . It is possible that some crop yield data cannot
229 find corresponding soil properties if r is too small, for which scenario this yield data
230 will be discarded. In the second technique, for each crop yield data $y(e_i, w_i)$, each
231 corresponding soil property was approximated by the averaging value of all values of

232 this soil property inside a circle with a radius r . A small value of r refers to more
233 accurate correspondence between yield and soil properties, but a lower number of
234 samples included in the analysis. The former method of re-sampling is designated here
235 as ‘shortest distance approximation (SDA)’, whereas the latter method is designated as
236 ‘circle-based average approximation (CAA)’.

237 **3. Results and discussion**

238 *3.1 Pearson correlations*

239 Pearson coefficient (r) values between pairs of soil properties suggest collective
240 (positive) linear relationships to exist between Ca and CEC, MC, Mg, OC, pH and TN
241 ($r = 0.519 - 0.747$) and between CEC and Ca, Mg, MC and pH ($r = 0.590 - 0.748$). This
242 may indicate that although Ca has no direct spectral response in the NIR range, it is
243 measured with vis-NIR spectroscopy through covariation with MC and OC, both having
244 direct spectral response (Stenberg *et al.* 2010; Kuang *et al.* 2012). However, CEC is
245 measured through covariation with MC only. As expected, TN correlated with OC,
246 which is a similar result to that reported elsewhere (Carlyle 1993; Kuang and Mouazen
247 2011).

248 Examining r values between the eight on-line measured soil properties and yield,
249 reveals negligible (negative) relationships (Table 1) between laboratory measured soil
250 properties and yield. The highest linear correlation is calculated between CEC and yield
251 ($r = -0.349$). This again proves the complexity of the system and necessitates the need
252 for more advanced modelling techniques that account for both linear and nonlinear
253 interactions.

254 **[Table 1]**

255 *3.2 Model output*

256 The detailed correspondence between inputs variables and soil properties are described
257 in Table 2. The initial full model, based on quadratic terms, was chosen in this paper,
258 which can be written as follows:

259

260

261
$$y = \theta_0 + \sum_{i=1}^8 \theta_i u_i + \sum_{i=1}^8 \sum_{j=i}^8 \theta_{ij} u_i u_j + \varepsilon \quad (2)$$

262

263

264 This model has 45 terms. All inputs and output were normalised by removing the mean.
265 The proposed method was then applied to calculate the ERRC of each term. Table 3
266 lists the first 10 terms selected using the SDA re-sampling technique with a 3 m radius.
267 From this calculation it was observed that the contribution of CEC to the wheat yield
268 variability was the largest (e.g. ERRC = 15.68%) among the 45 terms, including all soil
269 properties and their interactions. This was followed successively by Mg (ERRC =
270 3.57%) and Ca * CEC (ERRC = 1.13%) terms. This is explained by the fact that
271 although CEC is not a nutrient, it is a widely accepted measure to assess the fertility of
272 the soil. In fact, CEC represents the soil ability to hold positively charged ions e.g.,
273 exchangeable cations, which is directly linked to nutrients, hence, it is an important
274 indicator of soil fertility (Hazelton and Murphy 2007). Its significant contribution to

275 crop yield could be due to the quantity of nutrients in the field being variable through
276 the field. Furthermore, CEC is an important indicator influencing soil structure stability,
277 nutrient availability, soil pH and the soil's reaction to fertilisers and other ameliorants
278 (Hazelton and Murphy, 2007), which as a result will have a positive influence of crop
279 growth and yield. Furthermore, CEC is also related to potassium content and clay
280 particles, which affect available water content (Bergaya and Vayer 1997), hence,
281 influencing crop growth and development.

282 **[Table 2]**

283 **[Table 3]**

284 By comparing the contribution of each soil property to the wheat yield with the
285 corresponding significance threshold, the soil properties having significant contribution
286 to the crop yield can then be highlighted as shown in Table 4. Amongst the eight studied
287 soil properties, CEC, Mg, TN, Ca, OC and MC all have significant influence on the
288 crop yield, with declining order. However, the largest influence is attributed to CEC,
289 followed successively by Mg and TN. It is worth noting that pH is normally associated
290 with soil fertility and CEC (Hazelton and Murphy 2007) has the lowest influence on
291 yield. But, pH level directly affects nutrient availability and crop nutrient uptake
292 (HGCA 2014). With acidic soils (soil pH is smaller than 5), the pH would have negative
293 influence on nutrient uptake. It is commonly stated in farmer's guides that the optimum
294 pH for soils under continuous arable cropping of cereal crops is between 6 and 7 with
295 6.5 being the ideal. However, in the Horns End experimental field, the pH value of the
296 majority of the field area ranged between 5.6 and 8, which may explain the low
297 contribution of pH to yield prediction (Bruulsema 2015). Similar observation can be

298 made for P. Although P is a key nutrient for crop growth and development, no
299 significant contribution to wheat yield was observed. One explanation could be that P is
300 not a limiting property in Horns End field, as manure is being frequently applied
301 (Mouazen and Kuang 2016). Another reason might be the fact that a part of the field
302 i.e., the north-west part experienced a waterlogging problem associated with a poor
303 drainage system for many years. This is also reflected on the poor yield harvested in
304 2013, as shown in Figure 3, where low harvest can be observed particularly on the
305 northern and south western parts of the field, coinciding well with areas with the
306 waterlogging problem.

307 **[Figure 3]**

308 **[Table 4]**

309 A multiple linear regression analyses with least square estimation conducted by
310 Kravchenko and Bullock (2000) found OC as the main and most consistent, positively
311 correlated parameter with corn and soybean yield. Interestingly, they found that the
312 contribution from K, CEC and P was mostly negligible, and this was attributed to K and
313 P being ample in abundance in the soils. This finding is in line with those of the current
314 work regarding P only. However, Kravchenko and Bullock (2000) stated that the
315 performance of crop prediction models varies from field to field across different
316 cropping seasons.

317 After CEC and Mg, TN ranked as the third largest contributor to wheat yield, a result
318 which is supported by previous research suggesting that nitrogen supply is a large
319 limiting factor of crop yield (Agegnehu *et al.* 2016) and is strongly linked with soil TN

320 content before planting and uptake rate by plants during the growing season.
321 Surprisingly, Mg has the second largest contribution to the wheat yield variance. Mg is
322 an essential plant nutrient for plant growth, as it has well-known roles in photosynthesis
323 process and chlorophyll building (Mengel and Kirkby 1987). Deficiency in Mg by
324 leaching may take place in highly acidic sandy soils. However, this is not the case of
325 the current experimental field, where pH varies between 5.6 and 8 in a mixture of
326 medium soil texture classes of sandy loam, loam, and sandy clay loam. This could
327 explain the high contribution of magnesium distribution to crop yield variation.

328 Due to the waterlogging problem associated with the poor drainage system in the north-
329 west part of the field (Figure 4), MC had only a minor influence on crop yield as it is
330 ranked sixth among the eight soil properties included in the analysis. There is an
331 optimum for soil moisture (varying with crop growth stage) being beneficial to crop
332 yield. As MC increases it may become a hindrance to crop yield after reaching a
333 threshold. The waterlogged areas are of high MC and nutrient concentrations but low in
334 yield due to the water stress, which affects crop establishment, growth and yield.
335 Waterlogging causes the crop roots to be unable to respire and when there is too little
336 oxygen in the soil pores, the demand for oxygen varies with crop and crop growth stage
337 (Boyer 1982). Waterlogging at grain filling stages can cause a significant loss in grain
338 yield (Condon and Giunta 2003).

339 **[Figure 4]**

340 *3.3 Model sensitivity to sampling technique*

341 All results discussed above are based on the SDA re-sampling technique with a 3 m
342 radius. To evaluate the sensitivity of the results to the selection of re-sampling
343 technique and the size of radius, more tests were performed, whose results are shown in
344 Table 5, in which only the top 3 significant soil properties are presented. Inspection of
345 Table 5 reveals that the top two soil properties (e.g., CEC and Mg) showed exactly
346 same response for all tests, appearing at first and second factors affecting yield,
347 respectively, whereas TN appears three times and Ca appears once in the third position.
348 Additionally, the CAA re-sampling technique consistently had a larger total
349 contribution (SERR = 22.97% for 3 m radius) to wheat yield than that of the SDA re-
350 sampling technique (SERR = 20.29% for 3 m radius), which indicates the CAA
351 technique may be more suitable for the high resolution soil and yield data, because the
352 identified model explains more of the system output. Also, the total contribution
353 decreases following the increase of the sample number, which is expected because more
354 samples indicate more spatial variations of the underlying rule (Billings 2013). This is
355 also true for the radius, because with a larger radius, larger samples are included in the
356 analysis.

357 **[Table 5]**

358 Results showed that the overall contribution of the eight soil properties to wheat yield is
359 23.21%. One would expect that the contribution of soil properties to yield should be
360 larger than the overall calculated contribution in this study. However, the results
361 obtained confirmed this to be a significant contribution, but also shows that there is
362 variability still at play, influencing the crop yield (e.g., crop disease, pests, topography,

363 micro-climatic conditions etc.). For example, whilst TN and OC should have significant
364 effects, and both are required by the crop for healthy growth and grain production, they
365 can also increase and prolong the leaf area index of the crop, which in turn increases
366 humidity, making the plant more susceptible to disease, hence, crop yield is negatively
367 affected (Bryson *et al.* 1997). Therefore, there is a need for a future work to expand on
368 the current data mining approach to quantify yield limiting factors, under larger number
369 of fields with different crops and different agricultural systems. The study should also
370 account for the other affecting factors of crop yield including crop disease, pests,
371 topography, micro-climatic conditions etc.

372 **4. Conclusions**

373 A volterra non-linear regressive with eXogenous inputs (VNRX) model accounting for
374 the linear and non-linear variability (VNRX-LN) was used to quantify yield limiting
375 factors of wheat in one field in Bedfordshire, the UK. The input data were eight soil
376 properties (e.g. OC, TN, CEC, Mg, MC, Ca, pH and P), collected at a high sampling
377 resolution rate (>500 sample per ha), with an on-line visible and near infrared
378 spectroscopy (vis-NIRS) sensor, whereas crop yield represented the single-output in the
379 system. Based on the results obtained the following conclusions can be drawn:

- 380 1. The VNRX-LN model can be successfully used to quantify the influence of
381 multi-soil properties, collected at high sampling resolution with an on-line soil
382 sensor, on crop yield.
- 383 2. The effect of soil properties on crop yield varied with soil property, with the
384 largest contribution observed for CEC, Mg and TN, with error reduction ratio
385 contribution (ERRC) values of 14.6%, 4.69% and 1%, respectively.

386 3. The overall contribution of the eight soil properties sums up to an ERRC value
387 of 23.21%. This value was found to be surprisingly low, but was explained by
388 the fact that a large part of the studied field suffers of a drainage problem, which
389 masked the actual effect of soil properties on crop yield.

390 It was recommended to validate the concept introduced in this study on a larger number
391 of fields, where other affecting parameters (e.g. crop diseases, pests, topography,
392 microclimate conditions) of crop growth and yield should be taken into account.

393 **Acknowledgements**

394 We acknowledge the funding received for FarmFUSE project from the ICT-AGRI
395 under the European Commission's ERA-NET scheme under the 7th Framework
396 Programme, and the UK Department of Environment, Food and Rural Affairs (contract
397 no: IF0208).

398 **References**

- 399 Agegnehu G, Nelson P, Bird M (2016) Crop yield, plant nutrient uptake and soil
400 physicochemical properties under organic soil amendments and nitrogen fertilization on
401 Nitisols. *Soil & Tillage Research* **160**, 1–13.
- 402 Bergaya F, Vayer M (1997) CEC of clays: measurement by absorption of a copper
403 ethylenediamine complex. *Applied clay science* **12(3)**, 275-280.
- 404 Billings SA (2013) 'Nonlinear system identification: NARMAX methods in the time,
405 frequency, and spatio-temporal domains.' (London: John Wiley & Sons)
- 406 Boyer J (1982) Plant productivity and environment. *Science* **218**, 443–448.
- 407 Bruulsema T (2015) 'Plant Nutrition TODAY.' (Georgia USA: International Plant
408 Nutrition Institute (IPNI))

409 Bryson RJ, Paveley ND, Clark WS, Sylvester-Bradley R, Scott RK (1997) Use of in-
410 field measurements of green leaf area and incident radiation to estimate the effects of
411 yellow rust epidemics on the yield of winter wheat. *European Journal of Agronomy* **7**,
412 53-62.

413 Carlyle J (1993) Carbon in forested sandy soils: properties, processes, and the impact of
414 forest management. *New Zealand Journal of Forestry Science* **23(3)**, 390-402.

415 Christy C (2008) Real-time measurement of soil attributes using on-the-go near infrared
416 reflectance spectroscopy. *Computers and Electronics in Agriculture* **61**, 10-19.

417 Condon A, Giunta F (2003) Yield response of restricted-tillering wheat to transient
418 waterlogging on duplex soils. *Australian Journal of Agricultural Research* **54(10)**, 957-
419 967.

420 Corwin DL, Lesch S, Shouse PJ, Soppe R, Ayars JE (2003) Identifying soil properties
421 that influence cotton yield using soil sampling directed by apparent soil electrical
422 conductivity. *Agronomy Journal* **95(2)**, 352-364.

423 Dhillon N, Samra J, Sadana U, Nielson D (1994) Spatial variability of soil test values in
424 a typic Ustochrept. *Soil Technology* **7**, 163–171.

425 Frogbrook ZL, Oliver MA (2007) Identifying management zones in agricultural fields
426 using spatially constrained classification of soil and ancillary data. *Soil Use and*
427 *Management* **23(1)**, 40-51.

428 Halcro G, Corstanje R, Mouazen A (2013) Site-specific land management of cereal
429 crops based on management zone delineation by proximal soil sensing. Precision
430 agriculture '13. (Wageningen, Neitherlands, Wageningen Academic Publishers) pp.
431 475-481.

432 Hazelton PA and Murphy BW (2007) 'Interpreting soil test results: what do all the
433 numbers mean Australia.' (Melbourne, Australia.: CSIRO publisher)

434 HGCA (2014). Oilseed rape guide, s.l.: Agriculture and Horticulture Development
435 Board.

436 Khosla R, Inman D, Westfall, DG, Reich, RM, Frasier M, Mzuku M, Koch B, Hornung
437 A (2008) A synthesis of multi-disciplinary research in precision agriculture: site-
438 specific management zones in the semi-arid Western Great Plains of the USA. *Precision*
439 *Agriculture* **9**, 85-100.

440 Kipp S, Mistele B, Schmidhalter U (2014) The performance of active spectral
441 reflectance sensors as influenced by measuring distance, device temperature and light
442 intensity. *Computers and Electronics in Agriculture* **100**, 24–33.

443 Kodaira M, Shibusawa S (2013) Using a mobile real-time soil visible-near infrared
444 sensor for high resolution soil property mapping. *Geoderma* **199**, 64-79.

445 Kravchenko A, Bullock D (2000) Correlation of corn and soybean grain yield with
446 topography and soil properties. *Agronomy Journal* **92(1)**, 75-83.

447 Kuang B, Mahmood HS, Quraishi Z, Hoogmoed WB, Mouazen AM, Henten E (2012)
448 Sensing soil properties in the laboratory, in situ, and on-line: a review. *Advances in*
449 *Agronomy* **114**, 155-223.

450 Kuang B, Mouazen A (2011) Calibration of a visible and near infrared spectroscopy for
451 soil analysis at field scales across three European farms. *European Journal of Soil*
452 *Science* **62(4)**, 629- 636.

453 Kuang B, Mouazen AM (2013). Effect of spiking strategy and ratio on calibration of
454 on-line visible and near infrared soil sensor for measurement in European farms. *Soil &*
455 *Tillage Research* **128**, 125-136.

456 Kweon G, Lund E, Maxton C (2013). Soil organic matter and cation-exchange capacity
457 sensing with on-the-go electrical conductivity and optical sensors. *Geoderma* **199**, 80-
458 89.

459 Lowenberg-DeBoer J, Aghib A (1999). Average return and risk characteristics of site
460 specific P and K management: eastern Corn Belt on-farm trial results. *Journal of*
461 *Production Agriculture* **12**, 276–282.

462 Maleki MR, Mouazena AM, Ketelaerea BD, Ramona H, Baerdemaekera, JD (2008).
463 On-the-go variable-rate phosphorus fertilisation based on a visible and near infrared soil
464 sensor. *Biosystems Engineering* **99(1)**, 35-46.

465 Marin-González O, Kuang B, Quraishi MZ, Munóz-García MA, Mouazen AM (2013).
466 On-line measurement of soil properties without direct spectral response in near infrared
467 spectral range. *Soil & Tillage Research* **132**, 21-29.

468 Matheron G (1963). Principles of geostatistics. *Economic Geology* **58**, 1246–1266.

469 Mouazen A (2006). Soil survey device. BE, Patent No. WO/2006/015463.

470 Mouazen AM, De Baerdemaeker J, Ramon H (2005) Towards development of on-line
471 soil moisture content sensor using a fibre-type NIR spectrophotometer. *Soil & Tillage*
472 *Research* **80(1-2)**, 171-183.

473 Mouazen AM, De Baerdemaeker J, Ramon H (2006) Effect of wavelength range on the
474 measurement accuracy of some selected soil constituents using visual-near infrared
475 spectroscopy. *Journal of Near Infrared Spectroscopy* **14**, 189–199.

476 Mouazen AM, Kuang B (2016) On-line visible and near infrared spectroscopy for in-
477 field phosphorous management. *Soil & Tillage Research* **156**, 471-477.

478 Mouazen AM, Maleki MR, Cockx L, Van Meirvenned M, Van Holm LHJ, Merckx R,
479 De Baerdemaeker J, Ramon H (2009) Optimum three-point linkage set up for
480 improving the quality of soil spectra and the accuracy of soil phosphorous measured
481 using an on-line visible and near infrared sensor. *Soil & Tillage Research* **103(1)**, 144-
482 152.

483 Mulla D (2013) Twenty Five Years of Remote Sensing in Precision Agriculture: Key
484 Advances and Remaining Knowledge Gaps. *Biosystems Engineering* **114**, 358-371.

485 Raun WR, Johnson GV, Lees HL, Sembiring H, Phillips SB, Solie JB, Stone ML,
486 Whitney RW (1998). Microvariability in soil test, plant nutrient and yield parameters in
487 Bermudagrass. *Soil Science Society of America Journal* **62(3)**, 683–690.

- 488 Shibusawa S, Imade Anom SW, Sato S, Sasao A, Hirako S (2001). Soil mapping using
489 the real-time soil spectrophotometer. In ‘Third European Conference on Precision
490 Agriculture’. pp. 497–508. (Montpellier, France: ECPA).
- 491 Zhao Y, Billings SA, Wei H, Sarrigiannis P (2012) Tracking time-varying causality and
492 directionality of information flow using an error reduction ratio test with applications to
493 electroencephalography data. *Physical Review E* **86**, 051919.



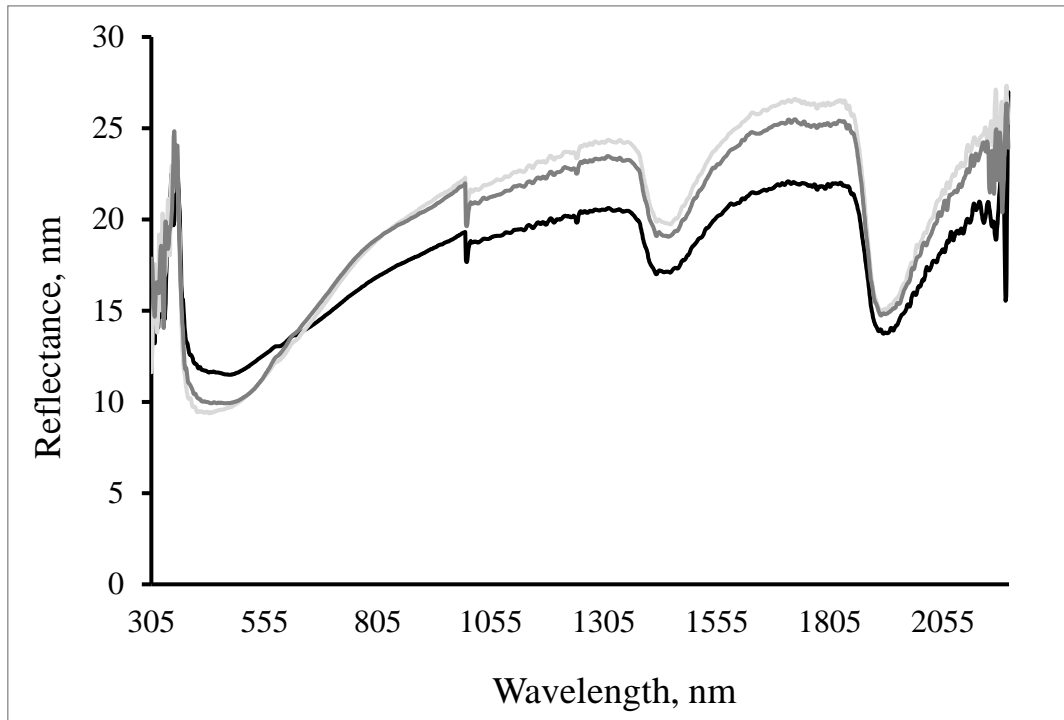
494

495

496 Figure 1. Illustrated image of the tractor mounted on-line visible and near infrared

497

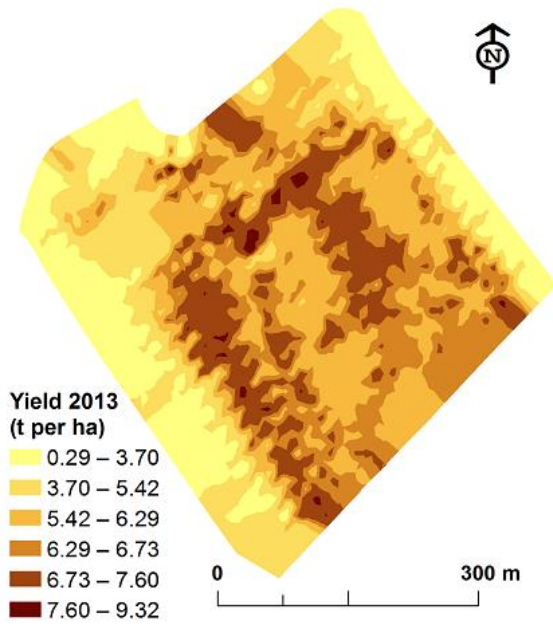
spectroscopy (vis-NIRS) sensor (Mouazen 2006).



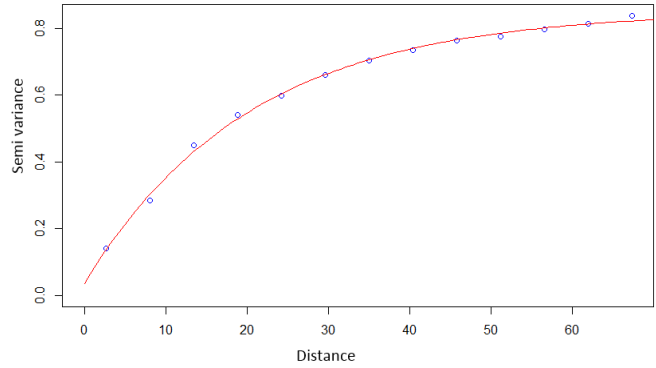
498

499 Figure 2. Examples of the raw on-line soil visible and near infrared (vis-NIR) spectra,
500 collected with the on-line sensor. Showing slight deviations in relative absorbance,
501 across the wavelengths, which is dependent on the soil properties.

502
503
504
505
506
507
508
509
510



(a)



(b)

511
512 Figure 3. Interpolated yield map (a) and exponential semi-variogram of 0.036, 0.817
513 and 20.358, representing, nugget, sill and range, respectively (b) based on the 2013
514 harvest of wheat grain in tons per hectare. Lighter areas representing lower yield.

515

516

517

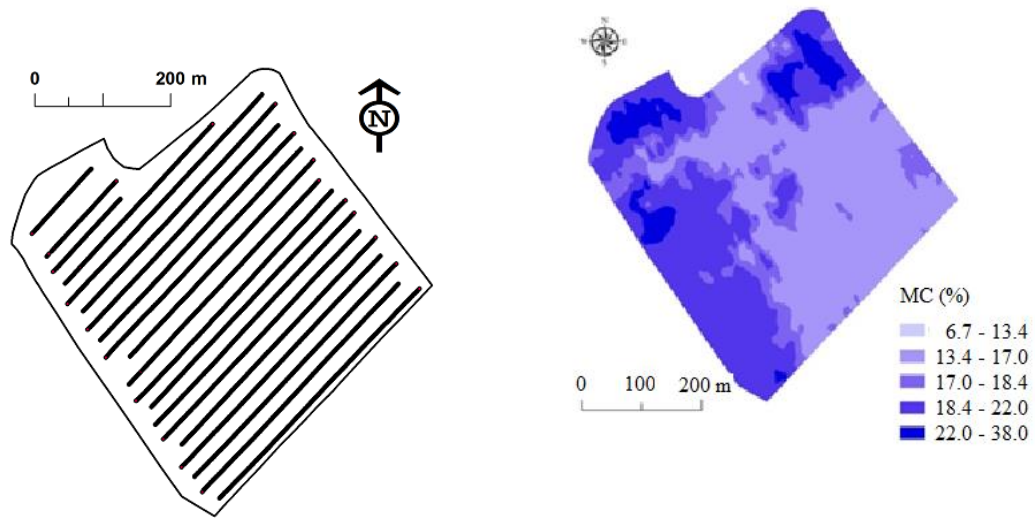
518

519

520

521

522



(a)

(b)

523

524

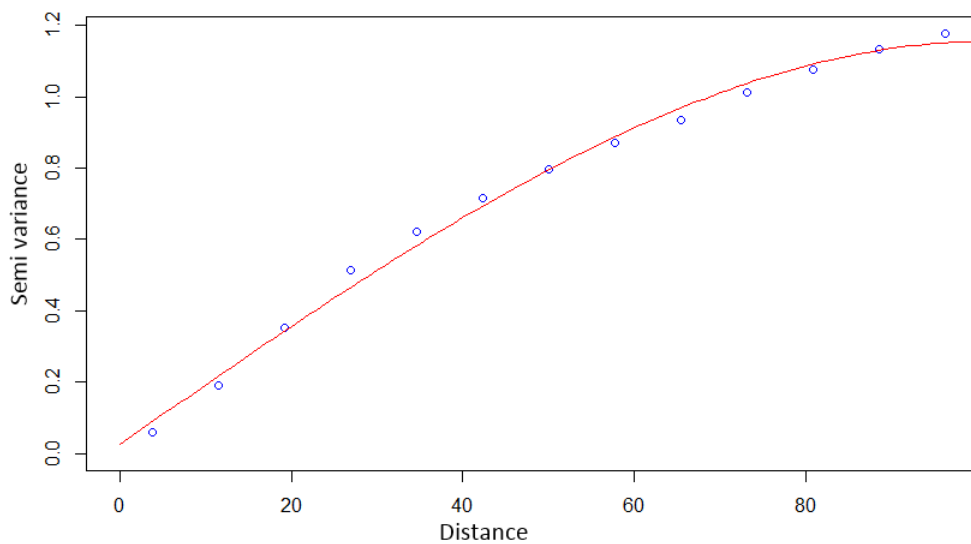
525

526

527

528

529



(c)

530

531 Figure 4. Measured transects (a), map of the soil moisture content (MC) measured with
532 the on-line visible and near infrared spectroscopy (vis-NIRS) sensor after crop harvest
533 in August, 2012 (b), and the spherical semi-variogram used for krigging of MC map
534 with nugget, sill and range values of 0.036, 0.817 and 20.358, respectively.

535

536 Table 1: Pearson correlation (r) between on-line measured soil properties in 2012 and
 537 wheat yield harvested in 2013.

	Ca	CEC	MC	Mg	OC	P	pH	TN	Yield
Ca	1.000								
CEC	0.733	1.000							
MC	0.519	0.748	1.000						
Mg	0.628	0.586	0.476	1.000					
OC	0.650	0.441	0.436	0.176	1.000				
P	0.163	0.216	0.019	0.042	0.027	1.000			
pH	0.747	0.590	0.492	0.348	0.432	-0.013	1.000		
TN	0.596	0.411	0.269	0.167	0.543	0.556	0.307	1.000	
Yield	-0.321	-0.349	-0.209	-0.320	-0.199	-0.000	-0.152	-0.057	1.000

538 OC is organic carbon in %; P is extractable phosphorous in mg/l; MC is moisture
 539 content in %; TN is total nitrogen in %, CEC is cation exchange capacity in meq/100g;
 540 Ca is calcium in mg/l; Mg is magnesium in mg/l; and pH the log measurement of
 541 acidity.

542 Table 2: The correspondence between inputs variables in Volterra Non-linear
 543 Regressive with eXogenous inputs (VNRX) model and soil properties

Input	Property	Input	Property	Input	Property	Input	Property
u_1	Ca	u_2	CEC	u_3	MC	u_4	Mg
u_5	OC	u_6	P	u_7	pH	u_8	TN

544 OC is organic carbon in %; P is extractable phosphorous in mg/l; MC is moisture
 545 content in %; TN is total nitrogen in %, CEC is cation exchange capacity in meq/100g;
 546 Ca is calcium in mg/l; Mg is magneium in mg/l; and pH the log measurement of acidity.
 547

548 Table 3: The first ten terms with corresponding error reduction ratio contribution
 549 (ERRC) values and coefficients based on the shortest distance approximation (SDA) re-
 550 sampling technique with a three m radius

Rank	Term	ERRC	Coefficient θ_i
1	CEC	15.68%	-0.0948
2	Mg	3.57%	-0.4840
3	Ca*CEC	1.13%	-0.0025
4	MC*Mg	0.72%	-0.0558
5	OC	0.78%	-0.2056
6	Mg*P	0.34%	-0.9615
7	Mg*TN	0.78%	5.0750
8	pH*pH	0.39%	-0.0670
9	constant	0.82%	0.1917
10	TN*TN	0.37%	-8.5096

551 OC is organic carbon in %; P is extractable phosphorous in mg/l; MC is moisture
 552 content in %; TN is total nitrogen in %, CEC is cation exchange capacity in meq/100g;
 553 Ca is calcium in mg/l; Mg is magnesium in mg/l; and pH the log measurement of
 554 acidity.

555 Table 4: Error reduction ratio contribution (ERRC) contribution of each soil property
 556 (input) to the crop yield (system output) with corresponding significance threshold
 557 based on the shortest distance approximation (SDA) re-sampling technique with a three
 558 m radius

Rank	Input	ERRC (%)	Significance threshold (%)	Significant
1	CEC	14.60	0.60	Yes
2	Mg	4.69	0.52	Yes
3	TN	1.00	0.50	Yes
4	Ca	0.98	0.43	Yes
5	OC	0.68	0.49	Yes
6	MC	0.62	0.47	Yes
7	pH	0.34	0.46	No
8	P	0.30	0.56	No
Total		23.21	4.03	

559 OC is organic carbon in %; P is extractable phosphorous in mg/l; MC is moisture
 560 content in %; TN is total nitrogen in %, CEC is cation exchange capacity in meq/100g;
 561 Ca is calcium in mg/l; Mg is magnesium in mg/l; and pH the log measurement of
 562 acidity.

563 Table 5: Contribution of the top three significant soil properties in terms of the sum of
 564 error reduction ratio (SERR) on the crop yield, based on shortest distance
 565 approximation (SDA) and (CAA) sampling techniques calculated for different radius
 566 values.

Re-sampling technique	Re-sampling radius	Sampled number	Top three inputs		Total Contribution (SERR)
			Inputs	Contribution	
SDA	3	1377	CEC	14.60%	20.29%
			Mg	4.69%	
			TN	1.00%	
CAA	3	1377	CEC	16.54%	22.97%
			Mg	4.00%	
			TN	2.43%	
SDA	5	3605	CEC	9.20%	13.61%
			Mg	2.45%	
			TN	1.96%	
CAA	5	3605	CEC	12.90%	15.87%
			Mg	3.02%	
			Ca	2.65%	

567 TN is total nitrogen in %, CEC is cation exchange capacity in meq/100g; Mg is
 568 magnesium in mg/l.