# Using Advanced Regression Models for Determining Optimal Soil Heterogeneity Indicators

1
2
3

**Georg Ruß, Rudolf Kruse, Martin Schneider, and Peter Wagner**

4

**Abstract** Nowadays in agriculture, with the advent of GPS-based vehicles and sensor-aided fertilization, large amounts of data are collected. With the importance of carrying out effective and sustainable agriculture getting more and more obvious, those data have to be turned into information – clearly a data analysis task.

Furthermore, there are novel soil sensors which might indicate a field's heterogeneity. Those sensors have to be evaluated and their potential usefulness should be assessed. Our approach consists of two stages, of which the first stage is presented in this article.

The data attributes will be comparable to the ones described in Ruß (2008). In the first stage, we will build and evaluate models for the given data sets. We will present a comparison between results using neural networks, regression trees and SVM regression. Results for an MLP neural network have been published in Ruß et al. (2008). In a future second stage, we will use the model information to evaluate and classify new sensor data. We will then assess their usefulness for the purpose of (yield) optimization.

5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

## 1 Introduction

20

Due to the modernization and better affordability of state-of-the-art GPS technology and a multitude of available sensors, a farmer nowadays harvests not only crops but also growing amounts of data. These data are small-scale and precise – which is essentially why the combination of GPS, agriculture and data has been termed *precision agriculture*.

However, collecting large amounts of data often is both a blessing and a curse. There is a lot of data available containing information about a certain asset – here: soil and yield properties – which should be used to the farmer's advantage. This is a

21
22
23
24
25
26
27
28

G. Ruß (✉)
Otto-von-Guericke-Universität Magdeburg, Germany
e-mail: russ@iws.cs.uni-magdeburg.de

common problem for which the term *data mining* or *data analysis* has been coined. 29
Data analysis techniques aim at finding those patterns or information in the data that 30
are both valuable and interesting to the farmer. 31

A common specific problem that occurs is yield prediction. As early into the 32
growing season as possible, a farmer is interested in knowing how much yield he is 33
about to expect. In the past, this yield prediction has usually relied on farmers' 34
long-term experience for specific fields, crops and climate conditions. What if 35
a computational model could be generated that allows to predict current year's 36
yield based on past data and current year data? This problem of yield predic- 37
tion encountered here is one which intelligent data analysis should be applied to. 38
More specifically, multi-dimensional regression techniques could be used for yield 39
prediction. 40

Nowadays, we can collect small-scale, precise data in-season using a multitude 41
of sensors. These sensors essentially aim to measure a field's heterogeneity. In future 42
work, these sensors will be assessed as to how useful they are for the purpose of yield 43
prediction. For this work, this article should serve as an overview on the capabilities 44
of different regression techniques used on agricultural yield data. 45

## 1.1 Research Target

46

The overall research target is to find those indicators of a field's heterogeneity which 47
are suited best to be used for a yield prediction task. Since this should be done 48
in-season, the sub-task here is one of multi-dimensional regression – predicting 49
yield from past and in-season attributes. At a later stage, when multi-year data are 50
available, models from past years could be used to predict present year's yield. 51

Therefore, this work aims at finding suitable data models that achieve a high 52
accuracy and a high generality in terms of yield prediction capabilities. Since multi- 53
year data are not yet available, the prediction can only be done in space using cross- 54
validation, instead of in time. As soon as multi-year data are available, the models 55
can be trained using these data for prediction in time. We will evaluate different 56
types of regression techniques on different data sets. Since these models usually are 57
strongly parameterized, an additional question is whether the model parameters can 58
be carried over from one field to other fields which are comparable in (data set) size. 59
This issue will be addressed in this work. This is especially useful when new 60
data have to be evaluated using one of the presented models. 61

## 1.2 Article Structure

62

Section 2 lays out the data sets that this work builds upon. The attributes and 63
their properties will be presented shortly. Section 3 briefly presents four selected 64
regression techniques from the data mining area which will be used for yield 65
prediction. 66

Section 4 shows the results of the modeling/regression stage and provides ans- 67
wers to the aforementioned research questions. 68

At the end of this article, future work is pointed out and implementation details 69
are provided. 70

## 2   Data Description 71

The data available in this work have been obtained in the years 2003–2006 on 72
three fields near Köthen, north of Halle, Germany (GPS coordinates: Latitude N 73
51 40.430, Longitude E 11 58.110). All information available for these 65-, 72- 74
and 32-hectare fields was interpolated using kriging (Stein 1999) to a grid with 10 75
by 10 meters grid cell sizes. Each grid cell represents a record with all available 76
information. During the growing season of 2006, the latter field was subdivided into 77
different strips, where various fertilization strategies were carried out. For an exam- 78
ple of various managing strategies, see e.g. Schneider and Wagner (2006), which 79
also shows the economic potential of PA technologies quite clearly. The fields grew 80
winter wheat, where nitrogen fertilizer was distributed over three application times 81
during the growing season. 82

Overall, for each field there are seven attributes – accompanied by the respective 83
current year's yield (2004 or 2006) as the target attribute. Those attributes have been 84
described in detail in Ruß et al. (2008), an overview is given below. In total, for the 85
F04 field there are 5241 records, for F131 there are 2278 records, for F330 there are 86
4578 records, thereof none with missing values and none with outliers. In addition, 87
a subset for F131 was available: in this subset, a special fertilization strategy was 88
carried out which used a neural network for prediction and optimization – this data 89
set is called F131net and has 1144 records. 90

In this work, data sets from three different fields are evaluated. A brief summary 91
of two of the available data sets is given in Tables 1a and 1b. On each field, dif- 92
ferent fertilization strategies have been used. One of those strategies is based on a 93
technique that uses a multi-layer perceptron (MLP) for prediction and optimization. 94
This technique has been presented and evaluated in, e.g., Ruß et al. (2008); Ruß 95
(2008) or Weigert (2006). For each field, one data set will contain all records, thus 96
containing all the different fertilization strategies. In addition, a subset of F131 has 97
been chosen to serve as a fourth data set to be evaluated. 98

## 3   Advanced Regression Techniques 99

As mentioned in the introduction, the task of yield prediction is essentially a task 100
of multi-dimensional regression. Therefore, this section will serve as an overview 101
about different regression techniques that are applicable to the yield data sets. We 102
aim to evaluate these techniques on the data sets presented in the preceding section. 103

t1.1 **Table 1** Overview of the *F04* and *F131* data sets. The additional data sets *F330* and *F131net*, which is a subset of *F131*, are not shown as their statistical properties are very similar to those of *F04* and *F131*

|       | (a) Data overview, F04 | | | |       | (b) Data overview, F131 | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| **F04** | *min* | *max* | *mean* | *std* | **F131** | *min* | *max* | *mean* | *std* |
| YIELD03 | 1.19 | 12.38 | 6.27 | 1.48 | YIELD05 | 1.69 | 10.68 | 5.69 | 0.93 |
| EM38 | 17.97 | 86.45 | 33.82 | 5.27 | EM38 | 51.58 | 84.08 | 62.21 | 8.60 |
| N1 | 0 | 100 | 57.7 | 13.5 | N1 | 47.70 | 70 | 64.32 | 6.02 |
| N2 | 0 | 100 | 39.9 | 16.4 | N2 | 14.80 | 100 | 51.71 | 15.67 |
| N3 | 0 | 100 | 38.5 | 15.3 | N3 | 0 | 70 | 39.65 | 13.73 |
| REIP32 | 721.1 | 727.2 | 725.7 | 0.64 | REIP32 | 719.6 | 724.4 | 722.6 | 0.69 |
| REIP49 | 722.4 | 729.6 | 728.1 | 0.65 | REIP49 | 722.3 | 727.9 | 725.8 | 0.95 |
| YIELD04 | 6.42 | 11.37 | 9.14 | 0.73 | YIELD06 | 1.54 | 8.83 | 5.21 | 0.88 |

The regression task can be formalized as follows: the training set          104

$$T = \{\{x_1, \ldots, x_n\}, y_i\}_{i=1}^{N} \tag{1}$$

is considered for the training process, where $x_i, i = 1, \ldots, n$ are continuous input   105
values and $y_i, i = 1 \ldots, N$ are continuous output values. Given this training set,   106
the task of the regression techniques is to approximate the underlying function   107
sufficiently well.                                                              108

### 3.1  Introduction to Regression Techniques                                  109

Since one particular technique, namely MLPs, has been used successfully in previ-   110
ous work (Ruß et al. 2008; Ruß 2008), it is used as a reference model here. Three   111
additional modeling techniques, namely RBF networks, regression trees, and sup-   112
port vector regression, will be presented, which are suitable for the task of yield   113
prediction. The aforementioned techniques have, to the authors' knowledge, not   114
been compared to each other when used with different data sets in the agriculture   115
context. This section presents some of the background for each of the techniques   116
before they will be evaluated in Sect. 4.                                        117

### 3.2  Neural Networks                                                        118

In previous work multi-layer perceptrons (MLPs), a type of neural networks, have   119
been used for a modeling task (Ruß et al. 2008; Ruß 2008) similar to the one   120
encountered here. Furthermore, neural networks have shown to be quite effective   121
in modeling yield of different crops (Drummond et al. 1998; Serele et al. 2000).   122
The MLP model was established as a reference model against which further regres-   123

sion techniques would have to compete. For a more detailed and formal description 124
of MLP neural networks, it is referred to Hagan (1995) or Haykin (1998). The net- 125
work layout and the parameters will be given in Sect. 4. In this work, the matlab 126
implementation for the MLP network was used: `newff`. 127

Furthermore, a different type of neural network, a radial basis function (RBF) 128
network, will be evaluated, since it is well-suited to the regression task. For this 129
network, matlab's `newrb` function has been utilized. 130

### 3.3 Regression Tree 131

Regression as well as decision trees are usually constructed in a top-down, greedy 132
search approach through the space of possible trees (Mitchell 1997). The basic algo- 133
rithms for constructing such trees are CART (Breiman et al. 1984), ID3 (Quinlan 134
1986) and its successor C4.5 (Quinlan 1993). The idea here is to ask the question 135
"which attribute should be tested at the top of the tree?" To answer this question, 136
each attribute is evaluated to determine how well it is suited to split the data. The 137
best attribute is selected and used as the test node. This procedure is repeated for the 138
subtrees. For further information on the construction details and possible problems 139
(such as overlearning) the reader is referred to Mitchell (1997). For this work the 140
standard matlab implementation of `classregtree` has been utilized. 141

### 3.4 Support Vector Regression 142

Support Vector Machines (SVMs) are a supervised learning method discovered 143
by Boser et al. (1992). However, the task here is regression, so the focus is on sup- 144
port vector regression (SVR). A more in-depth discussion can be found in Gunn 145
(1998). Given the training set, the goal of SVR is to approximate a linear function 146
$f(x) = \langle w, x \rangle + b$ with $w \in \mathbb{R}^N$ and $b \in \mathbb{R}$. This function minimizes an empirical 147
risk function defined as 148

$$R_{emp} = \frac{1}{N} \sum_{i=1}^{N} L_{\varepsilon}(\hat{y} - f(x)), \tag{2}$$

where $L_{\varepsilon}(\hat{y} - f(x)) = \max((|\xi| - \varepsilon), 0)$. $|\xi|$ is the so-called slack variable, which 149
has mainly been introduced to deal with otherwise infeasible constraints of the opti- 150
mization problem, as has been mentioned in Smola and Schölkopf (1998). By using 151
this variable, errors are basically ignored as long as they are smaller than a prop- 152
erly selected $\varepsilon$. $L_{\varepsilon}$ is called $\varepsilon$-insensitive loss function. Other kinds of functions 153
can be used, some of which are presented in Chap. 5 of Gunn (1998). To estimate 154
$f(x)$, a quadratic problem must be solved. See Mejía-Guevara and Kuri-Morales 155

(2007) for the dual form of this problem. In this work, the SVMtorch implementa- 156
tion from Collobert et al. (2001) has been utilized. Its documentation also points out 157
further details of the SVR process.                                                    158

### 3.5  Linear Regression and Naive Estimator                                         159

For comparison reasons, two further prediction methods are employed to compare 160
the advanced regression techniques against. The first of these is a simple multi- 161
linear regression estimator. The second is a naive estimator which simply reports 162
the previous year's yield as the output yield of the current year.                     163

### 3.6  Model Parameter Estimation                                                    164

Each of the aforementioned four different models will be evaluated on the same 165
data sets. One of the research goals here is to establish whether a model which has 166
been used on one data set can be used on a different data set without changing 167
its parameters. This would lead us to believe that comparable fields could use the 168
same prediction model. Hence, the *F04* data set is used to determine the model 169
parameters experimentally. Afterwards, the models are re-trained on the remaining 170
data sets using the settings determined for *F04*. The parameter settings are given in 171
Sect. 4.                                                                               172

For training the models, a cross-validation approach is taken. As mentioned in 173
e.g. Hecht-Nielsen (1990), the data will be split randomly into a training set, a vali- 174
dation set and a test set. The model is trained using the training data and after each 175
training iteration, the error on the validation data is computed. During training, this 176
error usually declines towards a minimum. Beyond this minimum, the error rises – 177
overlearning (or overfitting) occurs: the model fits the training data perfectly but 178
does not generalize well. Hence, the model training is stopped when the error on 179
the validation set starts rising. A size ratio of 8:1:1 for training, validation and test 180
sets is used. The data sets are partitioned randomly 20 times and the models are 181
trained. The models' performance will be determined using the root mean squared 182
error (RMSE) and the mean absolute error (MAE) on the test set. It is assumed that 183
the reader is familiar with these measures.                                            184

## 4  Regression Results                                                               185

The models are run with the parameter settings given below. Those were determined 186
experimentally on *F04* using a grid search, and carried over to the remaining data 187
sets.                                                                                  188

MLP    A relatively small number of 10 hidden neurons is used and the network is 189
       trained until a minimum gradient of 0.001 is reached, using a learning rate of 190
       0.25 and the *tangens hyperbolicus* sigmoid activation function.               191

t2.1 **Table 2** Results of running different models on different data sets. The best predictive model for each data set is marked in **bold** font

| Model/Dataset | MAE | | | | RMSE | | | |
|---|---|---|---|---|---|---|---|---|
| | F04 | F131 | F131net | F330 | F04 | F131 | F131net | F330 |
| MLP | 0.3706 | 0.2468 | 0.2300 | 0.3576 | 0.4784 | 0.3278 | 0.3073 | 0.5020 |
| RBF | 0.3838 | 0.2466 | 0.2404 | 0.3356 | 0.5031 | 0.3318 | 0.3205 | **0.4657** |
| REGTREE | 0.4380 | 0.2823 | 0.2530 | 0.4151 | 0.5724 | 0.3886 | 0.3530 | 0.6014 |
| SVR | **0.3446** | **0.2237** | **0.2082** | **0.3260** | **0.4508** | **0.3009** | **0.2743** | 0.4746 |
| LINREG | 0.4285 | 0.3257 | 0.2766 | 0.3820 | 0.5578 | 0.4392 | 0.3871 | 0.5330 |
| NAIVE | 2.9061 | 0.6135 | 0.6492 | 4.7157 | 3.1253 | 0.7613 | 0.7847 | 4.8308 |

RBF     For the radial basis function network, a radius of 1 is used for the radial basis 192
neurons in the hidden layer. The algorithm, which incrementally adds neurons 193
until the error goal of 0.001 is met, uses a maximum number of 70 neurons. 194

RegTree     For the regression tree, the default settings of classregtree perform opti- 195
mal; the full tree is pruned automatically and the minimum number of training 196
examples below which no split should be done is 10. 197

SVR     For the support vector regression model, the radial basis function kernel 198
yields the best results, using the parameters $C = 60, \sigma = 4.0$ and $\xi = 0.2$. 199

Considering the results in Table 2, support vector regression obviously performs 200
best on all but one of the data sets, regarding both error measures. Furthermore, 201
SVR also is the model taking the least amount of computation time. Hence, the 202
slight difference between the RMSE of SVR and RBF on the *F330* data set may 203
be considered insignificant in practice when computational cost is also taken into 204
account when deciding for a model. Regarding the understandability of the gener- 205
ated models, it would certainly be desirable to have the regression tree as the best 206
model since simple decision rules can easily be generated from the tree. However, 207
the regression tree performs worst in all of the cases. On the other hand, when com- 208
paring the hitherto reference model MLP with the current best model SVR, there is 209
not much difference in the understandability of both models. 210

## 5   Conclusion 211

The results clearly show that support vector regression can serve as a better reference 212
model for yield prediction than MLP. Even if the improvement should be statistically 213
insignificant, the advantages of SVR over MLP remain. It is computationally less 214
demanding, at least as understandable as the MLP and, most importantly, mostly 215
produces better yield predictions. Furthermore, the comparison against a linear 216
regression baseline and a naive estimator shows that the additional effort for using 217
SVR is worth it. 218

Furthermore, the results also show that model parameters which have been estab- 219
lished on one data set can be carried over to different (but similar with respect 220
to the attributes) data sets. A model for identifying the most useful heterogeneity 221
indicators is currently being evaluated. 222

## 5.1 Future Work

223

Due to the relatively high spatial resolution of the input data, the possible issue 224
of spatial autocorrelation arises. This influences the modeling during the cross- 225
validation stage. This will be investigated in future work. 226

# References

229

Boser, B. E., Guyon, I. M., & Vapnik, V. N. (1992). A training algorithm for optimal margin clas- 230
sifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory* 231
(pp. 144–152). New York: ACM Press. 232
Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). *Classification and Regression Trees.* 233
Monterey, CA: Wadsworth and Brooks. 234
Collobert, R., Bengio, S., & Williamson, C. (2001). Svmtorch: Support vector machines for large- 235
scale regression problems. *Journal of Machine Learning Research, 1,* 143–160. 236
Drummond, S., Joshi, A., & Sudduth, K. A. (1998). Application of neural networks: precision 237
farming. In *International Joint Conference on Neural Networks, IEEE World Congress on* 238
*Computational Intelligence* (Vol. 1, pp. 211–215). 239
Gunn, S. R. (1998). *Support vector machines for classification and regression.* Technical Report, 240
School of Electronics and Computer Science, University of Southampton, Southampton, U.K. 241
Hagan, M. T. (1995). *Neural network design (electrical engineering).* Thomson Learning. 242
Haykin, S. (1998). *Neural networks: A Comprehensive Foundation* (2nd ed.). Englewood, Cliffs, 243
NJ: Prentice Hall. 244
Hecht-Nielsen, R. (1990). *Neurocomputing.* Reading, MA, USA: Addison-Wesley. 245
Mejía-Guevara, I., & Kuri-Morales, A. (2007). Evolutionary feature and parameter selection in 246
support vector regression. In *Lecture Notes in Computer Science* (Vol. 4827, pp. 399–408). 247
Berlin, Heidelberg: Springer. 248
Mitchell, T. M. (1997). *Machine learning.* NY, USA: McGraw-Hill Science/Engineering/Math. 249
Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning, 1*(1), 81–106. 250
Quinlan, R. J. (1993). *C4.5: Programs for Machine Learning (Morgan Kaufmann Series in* 251
*Machine Learning).* Los Altos, CA: Morgan Kaufmann. 252
Ruß, G., Kruse, R., Schneider, M., & Wagner, P. (2008). Estimation of neural network parameters 253
for wheat yield prediction. In M. Bramer (Ed.), *Artificial Intelligence in Theory and Practice II* 254
of *IFIP International Federation for Information Processing* (Vol. 276, pp. 109–118). Berlin: 255
Springer. 256
Ruß, G., Kruse, R., Schneider, M., & Wagner, P. (2008). Optimizing wheat yield prediction using 257
different topologies of neural networks. In J. L. Verdegay, M. Ojeda-Aciego, & Magdalena, L. 258
(Eds.), *Proceedings of IPMU-08* (pp. 576–582). University of Málaga. 259
Ruß, G., Kruse, R., Wagner, P., & Schneider, M. (2008). Data mining with neural networks for 260
wheat yield prediction. In P. Perner (Ed.), *Advances in Data Mining (Proc. ICDM 2008)* (pp. 261
47–56). Berlin, Heidelberg: Springer Verlag. 262
Schneider, M., & Wagner, P. (2006). Prerequisites for the adoption of new technologies – the 263
example of precision agriculture. In *Agricultural Engineering for a Better World*, Düsseldorf: 264
VDI Verlag GmbH. 265
Serele, C. Z., Gwyn, Q. H. J., Boisvert, J. B., Pattey, E., Mclaughlin, N., & Daoust, G. (2000). 266
Corn yield prediction with artificial neural network trained using airborne remote sensing and 267

topographic data. In *2000 IEEE International Geoscience and Remote Sensing Symposium, 1*, 384–386.

Smola, A. J., & Schölkopf, B. (1998). *A tutorial on support vector regression.* Technical report, Statistics and Computing.

Stein, M. L. (1999). *Interpolation of Spatial Data : Some Theory for Kriging (Springer Series in Statistics).* Berlin: Springer.

Weigert, G. (2006). *Data Mining und Wissensentdeckung im Precision Farming - Entwicklung von ökonomisch optimierten Entscheidungsregeln zur kleinräumigen Stickstoff-Ausbringung.* PhD thesis, TU München.