



The International Society of Precision Agriculture presents the

15th International Conference on Precision Agriculture

26–29 JUNE 2022

Minneapolis Marriott City Center | Minneapolis, Minnesota USA

A Bayesian Network for wheat yield prediction using topographic, soil and historical data

M. Karampoiki¹, L.C. Todman², S. Mahmood², A.J. Murdoch², J. Hammond², E. Ranieri³, D.S. Paraforos¹

¹ University of Hohenheim, Institute of Agricultural Engineering, Technology in Crop Production, Garbenstr. 9, 70599 Stuttgart, Germany

² University of Reading, School of Agriculture, Policy and Development, Earley Gate, PO Box 237, Reading RG6 6EU, United Kingdom

³ Agricolus SRL, Via Settevalli 320, 06129 Perugia, Italy

A paper from the Proceedings of the
15th International Conference on Precision Agriculture
June 26-29, 2022
Minneapolis, Minnesota, United States

Abstract.

Bayesian Network (BN) is the most popular approach for modelling in the agricultural domain. Many successful applications have been reported for crop yield prediction, weed infestation, and crop diseases. BN uses probabilistic relationships between variables of interest and in combination with statistical techniques the data modelling has many advantages, although it is crucial to reduce data overfitting to improve the model accuracy.

*In this study, electrical conductivity (EC) from Veris iScan sensor and Dualem scanner, yield data from a combine harvester (John Deere), and Sentinel 2 (S2) imageries were collected for ten winter wheat (*Triticum aestivum* L.) fields in Germany and ten fields in the UK for the 2020 and 2021 seasons. The combine harvester data were analysed using ArcGIS software. For the German fields, the topographic wetness index (TWI), a good indicator for soil moisture, was calculated based on the digital elevation model (DEM) using ArcGIS software. Additionally, drone imageries were collected. An unmanned aerial vehicle (UAV) (DJI Mavic 2 Zoom) was equipped with a compact multispectral sensor (Parrot Sequoia+) and flew 60 m above ground level. The UAV data were analysed in Pix4D software. For the UK fields, samples of soil organic matter were collected. S2 imageries with 10 m spatial resolution and 5-day temporal resolution were downloaded from Europe's Copernicus website for the German and the UK fields respectively. The obtained imageries were analysed using SNAP toolbox.*

This paper focuses on developing a Machine Learning Approach (MLA) based on BN model to predict wheat yield using two novel parameters which are Prior Inherent Potential (PIP) and Inherent Potential (IP). The model has been developed using Netica (Norsys software),

The authors are solely responsible for the content of this paper, which is not a refereed publication. Citation of this work should state that it is from the Proceedings of the 15th International Conference on Precision Agriculture. EXAMPLE: Last Name, A. B. & Coauthor, C. D. (2018). Title of paper. In Proceedings of the 15th International Conference on Precision Agriculture (unpaginated, online). Monticello, IL: International Society of Precision Agriculture.

categorizing each node within each field from high to low PIP based on the data available for a given fields. A high PIP leads to a higher IP and a higher IP leads in turn to a high yield. Yield predictions are based on the probabilities of 50%. The actual and predicted yields of 50% probability maps had similar patterns of spatial variation and the correlations for the testing fields in Germany and UK were 0.42 and 0.35 respectively.

Keywords.

Yield predictions, Bayesian Networks, Prior Inherent Potential, Inherent Potential.

Introduction

Crop yield prediction is one of the most challenging problems in precision agriculture so far, and thus numerous models have been developed and tested. To solve this problem, multi-source data may be required considering that crop yield depends on different factors such as weather, soil, fertilization and seed varieties (Xu et al., 2019). However, crop yield prediction models can predict the actual yield, but still, better performance is needed (Filippi et al., 2019a).

Machine Learning (ML) has been studied by several authors (Elavarasan et al., 2018, Liakos et al., 2018, Somvanshi and Mishra, 2015) and can be used to provide promising yield predictions (Willcock et al., 2018). ML algorithms have the advantage of modeling non-linear relationships between different sources of data (Chlingaryan et al., 2018) and the model performance is improved when more data for training are available (Goodfellow et al., 2016). The most popular machine learning approach for modelling in the agricultural domain is the Bayesian Network (BN). A BN is a directed acyclic graph with nodes that represent variables and the links represent the relationships between nodes. This relation is specified by conditional probability tables. Handling missing data, and learning the relationships between variables are the two main advantages of BNs (Uusitalo, 2016). Moreover, BNs are a useful tool for combining expert knowledge with multi-sources of data (Walters and Martell, 2004). On the other hand, BNs cannot deal with continuous data (Jensen, 2001), and thus they need to be discretized.

Several studies, such as the yield of malting barley in absence of pesticides (Kristensen, 2002), the yield of energy crops in Western Canada (Newlands et al., 2010), and the influence of climate on sweet potato yield (Villordon et al., 2011), rice crop yield in India (Gandhi et al., 2016), country-level corn yield in Iowa (Chawla et al., 2016) have applied BNs to the problem of yield prediction. Most of the studies used available climatic data and satellite data. For instance, Fu et al., 2020, have created a prediction model using six machine learning methods to improve the accuracy of the model. Normalized difference vegetation index (NDVI) was used to construct the model. The correlation between actual and predicted yields was 0.78 in the random forest regression (RFR). Wang et al. 2020, created a two-branch deep learning model to predict wheat yield on a country level. The model performance reached an overall R^2 of 0.75. However, few previous studies have applied BNs for wheat yield prediction using a combination of multi-source of data i.e, topographic, soil, historical and weather data.

This study aims to use collected multi-source data to develop an algorithm that predicts the future wheat yield using the minimum number of variables in combination with a probable weather effect. The novelty of this algorithm is the introduction of the concept of 'Prior Inherent Potential' (PIP) and 'Inherent Potential' (IP) to reduce the number of parameters of the developed model to avoid overfitting and increase the accuracy.

Materials and methods

Study area

This study was conducted in winter wheat fields located in two countries, i.e., Germany and the UK. In 2020 and 2021, a total of twenty fields, ten fields in each of the two countries involved.

In Germany, the acquired data were obtained from ten fields near Gauersheim (49°40'38.39"N, 8°03'22.93"E). All ten fields are from Füge and Landfried Farm (Fig. 1). The criteria for the field selection were based on the known spatial variability in soil type and topography (mainly farmer knowledge), the availability of historical data (yield maps of previous crops, soil types, soil electrical conductivity), and the farmer's access to precision agriculture technology especially yield mapping and VRA (variable-rate application) techniques.



Fig. 1. Google Earth image showing the study sites and the locations of the German winter wheat fields in Füge and Landfried Farm based in Gauersheim, Rhineland-Palatinate, Germany.

In the UK, the ten fields are located near Reading (51°29'58.02"N, 0°56'40.63"E). All ten fields are from Coppid Farming Enterprises Iip (Fig. 2) and were selected using the same criteria as in Germany.



Fig. 2. Google Earth image showing the locations of the UK winter wheat fields (Coppid Farming Enterprises Iip, Reading, UK).

Data sources

Multi-source data with a various spatial and temporal resolutions, including data on current and historical yields, soils, satellite and drone imagery and weather, were collected for the twenty fields.

Yield data for the current and previous crops and elevation above sea level were provided by John Deere combine harvester at 2 m intervals with a 9 m swath. Thousands of observations per field were collected, providing a clear picture of the spatial variability. The number of observations varied from one field to another based on their areas. Sentinel-2 (S2) images were downloaded from the European Space Agency (ESA) website (<https://scihub.copernicus.eu/dhus/#/home>) and the data were used to calculate NDVI at 10 m spatial resolution for different times during the growing season. Additionally, UAV-based imagery was collected 2 weeks before and after nitrogen application using a Mavic 2 Zoom UAV (DJI, Nanshan, Shenzhen, China Mavic), equipped with a compact Parrot Sequoia+ multispectral camera (Sensefly, Lausanne, Switzerland). The UAV with the camera flew 60 m above ground level. The camera provides a set of four bands: green (550 ± 40 nm); red (660 ± 40 nm); red edge (735 ± 10 nm); and near-infrared (790 ± 40 nm). The NDVI was calculated based on the following formula:

$$\text{NDVI} = \frac{\text{NIR} - \text{Red}}{\text{NIR} + \text{Red}} \quad (1)$$

Soil Electrical Conductivity (EC) was provided by Veris iScan-sensor. Weather data was provided by Weierhof, a weather station located close to Füge and Landfield farm. An overview of the data is shown in Table 1.

Table 1. Overview of available and collected data from the German fields.

Fields	Yield data	Year
Hinteres Tal (12 ha)	Wheat	2020
	Barley	2017
Morgen (9.8 ha)	Wheat	2020
	Wheat	2017
Rosengarten (5.5 ha)	Wheat	2020
	Wheat	2016
Schanzgewanne (2.7 ha)	Wheat	2020
	Wheat	2016
Wiederschein (4.8 ha)	Wheat	2020
	Wheat	2017
Alzeyer (5.76 ha)	Wheat	2021
Birnbäum (4.52 ha)	Wheat	2021
	Wheat	2018
	Wheat	2016
Brunnenwiese (3.13ha)	Wheat	2021
	Barley	2020
	Barley	2018
	Wheat	2017
Horn (3.12 ha)	Wheat	2021
	Wheat	2018
	Wheat	2016
Morgen Unten (3.61 ha)	Wheat	2021

For the UK fields, yield data, grain and elevation were collected by a combine harvester from John Deere. S2 images were downloaded from the ESA website and the data were used to calculate NDVI at 10 m spatial resolution for different dates during the growing season. Soil Electrical Conductivity (EC) was provided by SOYL and was measured using a Dualem scanner at 6 m intervals along rows and 24 m between rows. Data of soil available nutrients were also provided by SOYL based on one sample per hectare and the maps of soil types were created by SOYL based on EC data and one soil sample/ha data. An irregular grid sampling scheme with some nested samples was designed by the University of Reading to collect soil samples for organic

matter analysis by SOYL. Weather data was provided by the University of Reading on a daily basis from a weather station located approximately 5 km from Coppid Farm. An overview of the aforementioned data is presented in Table 2.

Table 2. Overview of available and collected data from the UK fields.

Fields	Yield data	Year
Camp A (9.6 ha)	Wheat	2020
	Barley	2017
Chalkhouse (12 ha)	Wheat	2020
	Wheat	2017
Top Lane Right (7.6 ha)	Wheat	2020
	Wheat	2016
Top Lane Left (4.4 ha)	Wheat	2020
	Wheat	2016
Homewood (4.8 ha)	Wheat	2020
	Wheat	2017
Audreys (13 ha)	Wheat	2021
	Barley	2019
	Wheat	2018
Crowsley (14.8 ha)	Wheat	2021
	Barley	2019
	Wheat	2018
Ladyshaw (14.7 ha)	Wheat	2021
	Barley	2020
	Barley	2019
	Wheat	2018
Cliff 1 (12 ha)	Wheat	2021
	Wheat	2018
	Oilseed rape	2017
	Barley	2017
Homestead (13 ha)	Wheat	2021
	Barley	2019
	Wheat	2018
	Beans	2017

Data processing and integration

The data came from a different number of sources, including sensors on the combine harvester, soil electrical conductivity data, and satellite and drone data. The points at which measurements were taken varied in their spatial distribution and densities within the field. Consequently, it was necessary to create a regular grid of data points in which all data sources were interpolated or extrapolated to a common grid in ArcGIS software. A 6×6 m grid was chosen because the section control of agricultural machinery used in precision agriculture applications is 6 m. Interpolation maps were created for each variable (current and previous yields, elevation, and electrical conductivity) based on the Kriging method of interpolation. The interpolation maps were converted to a raster map on which the 6×6m grid was overlaid and the value at each grid point extracted. For the topographic wetness index (TWI), a raster map based on the digital elevation model (DEM) and the slope was created and the value at each grid point was extracted. For vegetative indices (NDVI), the 6×6 m grid was overlaid onto the raster map of each satellite band, and the value of any pixel was assigned to all points located within that pixel. However, for variable rate N fertilizer application, it was not possible to create a map and extract the value at points, as the value at the neighbouring points for a 30 m swath will be the same, therefore, a “Spatial Join with K-nearest neighbour” method as described in the following paragraph, was used to extract the values of these points.

Yield data provided by combine harvesters were filtered based on an automated data cleaning protocol proposed by Natale *et al.* (2020). This method filters the data in three steps. First, null points, where yield is equal to 0, are eliminated from the dataset. Second, removing overlapping points with the same yield value. Third, a data cleaning method, which is based on the concept of a “moving window” is applied such that each point is classified as an “outlier” or “non-outlier” based on values occurring in neighbouring points included in a circle of radius=R. The

validation of points in each circle includes the following steps: 1. Definition of maximum acceptable coefficient of variation (CV_{max}- threshold of acceptability of the coefficient of variation: 20% was used). 2. Definition of R of the neighbourhood. - defined here as 1.5 times the working swath width of the machine meaning that circles with a diameter of three working widths were used to define a window. 2.1 Calculation of the number of points included in the neighbourhood (N), 2.2 Calculation of coefficient of variation (CV) of the points included in the neighbourhood, 2.3 Calculation of the number of points with CV > CV_{max} within the neighbourhood (NCV_{max}), 2.4 Definition of the outlier: points where N=NCV_{max} are identified as outliers. Then, removing the points not included in the numeric interval between ($\mu-3\sigma$) and ($\mu+3\sigma$) where μ is the population mean and σ is the standard deviation. The protocol was applied for the yield data collected in both countries.

Learning from the soil, topographic and yield data

The MLA uses a Bayesian Network, in which the ‘nodes’ of the network represent each variable and links between nodes represent the interactions between variables. The resulting structure can be used for probabilistic inference, that is the probability that a yield at a specified grid point will equal or exceed a given value. It also allows the conditional dependence of variables to be represented so that variables are not considered to contribute additive effects, but rather a Bayesian Network can combine the information from correlated variables. For example, to infer that in situations when a particular location in a field repeatedly provides high relative yields over two or more years there is likely to be higher Inherent Potential than a location that has a more variable yield in different years.

To “learn” what we have called the ‘Prior Inherent Potential’ (PIP) of each grid point in each field and also its ‘Inherent Potential’, a machine learning algorithm (MLA) works on the logic that locations with a high Prior Inherent Potential typically led to higher Inherent Potential and a high Inherent Potential lead in turn to high yield. The PIP and the IP are what are known as ‘latent variables’ which means they are not properties that can be measured or observed in the traditional sense. They summarize several factors that influence the potential of the crop to develop. Many of these factors are likely to be due to soil properties. For instance, the soil structure affects how water is stored and regulated and soil nutrient cycling affects nutrient availability to the crop. Observations of soil properties such as texture and carbon content may, therefore, improve the inference of the IP. IP also captures other factors influencing crop growth such as weeds.

In this study, the Bayesian Network used with the MLA aims to categorise each node within each field from high to low PIP based on the data available for a given field. Thus, in Germany (Fig. 3), electrical conductivity and topographic wetness index were used to characterise the PIP whereas electrical conductivity and soil organic matter was used in the UK (Fig. 4). Learning was performed using the Expectation Maximisation approach in Netica, a Bayesian Belief Network software (<http://www.norsys.com>).

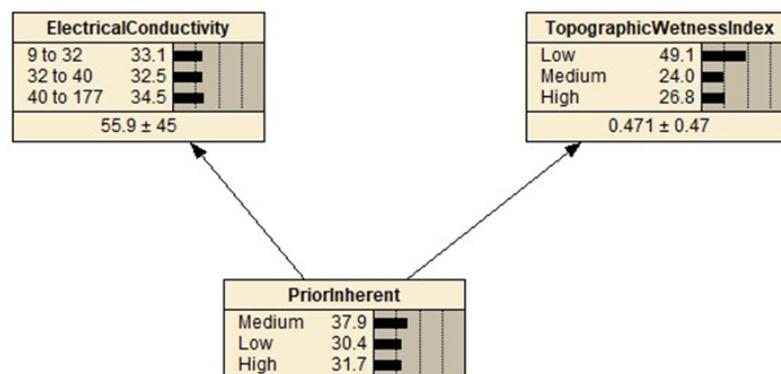


Fig. 3. Diagram for the MLA of 'Prior Inherent Potential' in the ten German fields. Topographic Wetness Index (TWI) and electrical conductivity (EC) have been separated into three discrete states where each state shows the range of the measured values. The probability distribution across states is shown as a %-probability and visualized with black horizontal bars. At the bottom of the EC and TWI nodes, the numbers show the mean and the SD of the values for each node. The model was applied to grid point locations in all wheat fields being studied for the 2020 and 2021 crop seasons on a given farm. There were up to ten fields per farm.

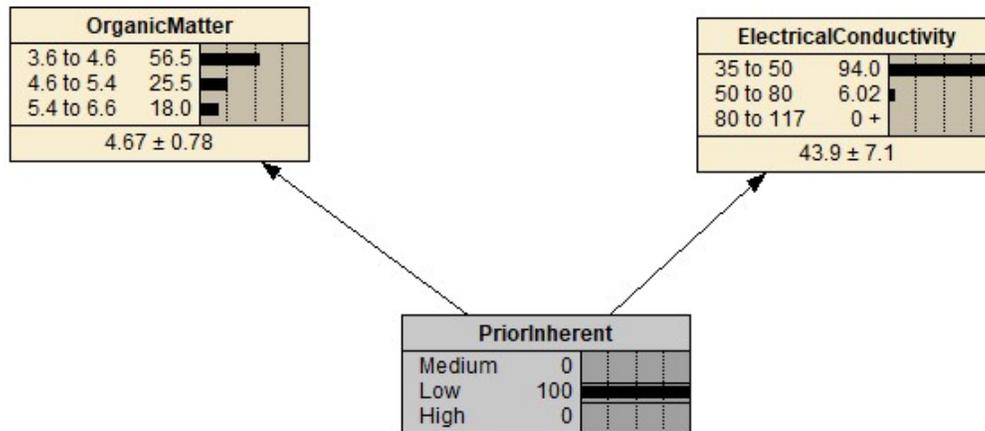


Fig. 4. Diagram for the MLA of 'Prior Inherent Potential' in the ten UK fields. Organic matter (OM) and electrical conductivity (EC) have been separated into three discrete states where each state shows the range of the measured values. The probability distribution across states is shown as a %-probability and visualized with black horizontal bars. At the bottom of the EC and OM nodes, the numbers show the mean and the SD of the values for each node. The model was applied to grid point locations in all wheat fields being studied for the 2020 and 2021 crop seasons on a given farm. There were up to ten fields per farm.

For the German fields, data were available for two crops. In two fields these data corresponded to the observed yields in 2016 and 2018, in one other field yield data were observed in three years (2017, 2018 and 2020). To facilitate comparisons between yields of different crops in different fields in different years, the yield data for a given field was normalized as percentages of the mean yield for that field. The Inherent Potential for all ten fields was learned using the data from all fields simultaneously. As an example, the IP learned from the previous yields for the German fields is shown in Fig. 5.

For the UK fields, data were available for five crops. In four fields this data corresponded to observed yields in 2018 and 2019, in one other field yield data were observed in 2017, 2019, and 2020. IP model is presented in Fig. 6.

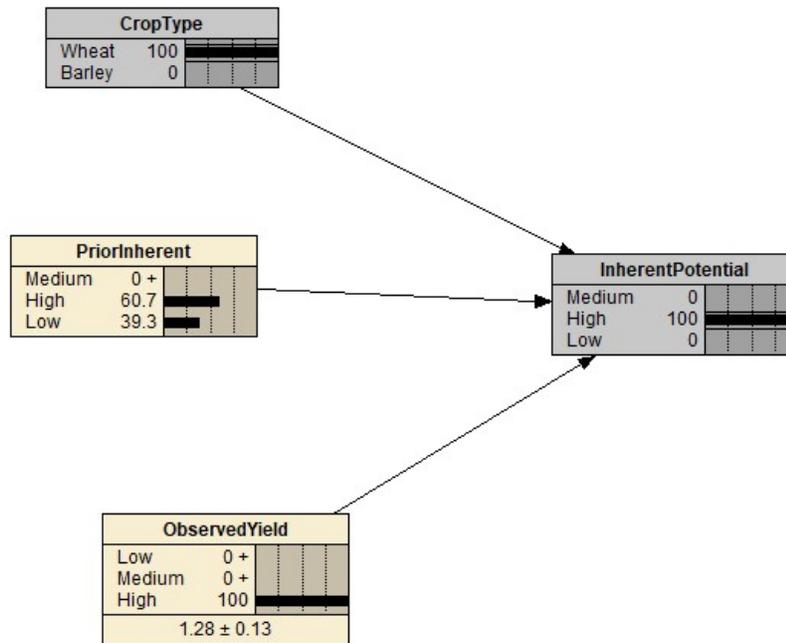


Fig. 5. Diagram for the MLA of 'Inherent Potential' from yields of previous crops (2016-2020) in ten German fields. The observed Yield node has been discretized into three categories, where each category represents a range of measured data. The probability distribution across states is shown as a %-probability and visualized with black horizontal bars. At the bottom of the Observed yield node, the numbers show the mean and the standard deviation (SD) of the values. The model was applied to grid point locations in all wheat fields being studied in 2020 and 2021 on a given farm. There were up to ten fields per farm. Crops varied and so yields were normalized for each year (taking the percentage of mean).

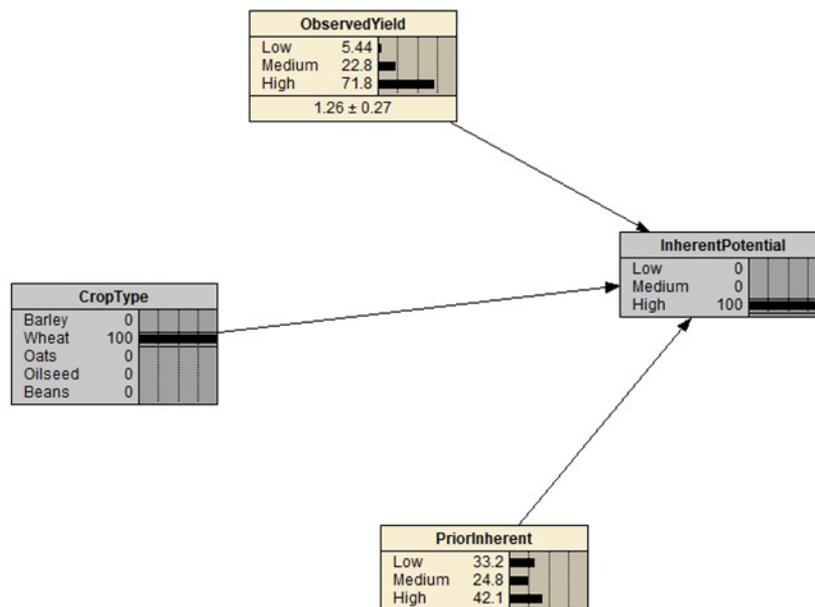


Fig. 6. Diagram for the MLA of 'Inherent Potential' from yields of previous crops (2016-2020) in ten UK fields. The observed Yield node has been discretized into three categories, where each category represents a range of measured data. The probability distribution across states is shown as a %-probability and visualized with black horizontal bars. At the bottom of the Observed yield node, the numbers show the mean and the standard deviation (SD) of the values. The model was applied to grid point locations in all wheat fields being studied in 2020 and 2021 on a given farm. There were up to ten fields per farm. Crops varied and so yields were normalized for each year (taking the percentage of mean).

Representing weather variability

A significant variation in yields of winter wheat from season to season is caused by the variability in weather conditions. The weather variables that were expected to have more influence on wheat yield are rainfall and mean temperature (Addy et al., 2020).

To represent this is the model we used a Reference Yield node that described the expected average yield due to the weather conditions throughout the season.

Model development

The output of IP and the entire weather algorithm were used with in-season crop biomass (base on the NDVI in May) to predict grain yield. Additionally, the variable-rate nitrogen fertilizer application (VRA-N) which was available for German and UK fields was also used only as an input variable to predict grain yield. The model structure (Fig. 7) showed how variables were linked to predict yield for the UK fields. The model was, however, applied on a farm basis which included ten winter wheat fields for each country, five fields in the 2020 season and five fields in the 2021 season. The conditional interdependencies of these variables were learned using 75% of the data from eight fields (four fields from each season) using expectation maximization within Netica software and then applied to the 25% of the data and to the other two fields (one from each season) to test the model. The eight fields were split randomly using Matlab R2020a software (www.matlab.mathworks.com). The yields were normalized based on farm level. In Germany, Morgen and Birnbaum were used as a testing set. In the UK the fields that were selected as a testing set were Camp A and Crowsley. Probabilities for each grid point were extracted from Netica and then the output file was processed in Matlab to calculate the predicted values with a 50% probability. The predicted grain yield was mapped for each field and compared to the observed yield.

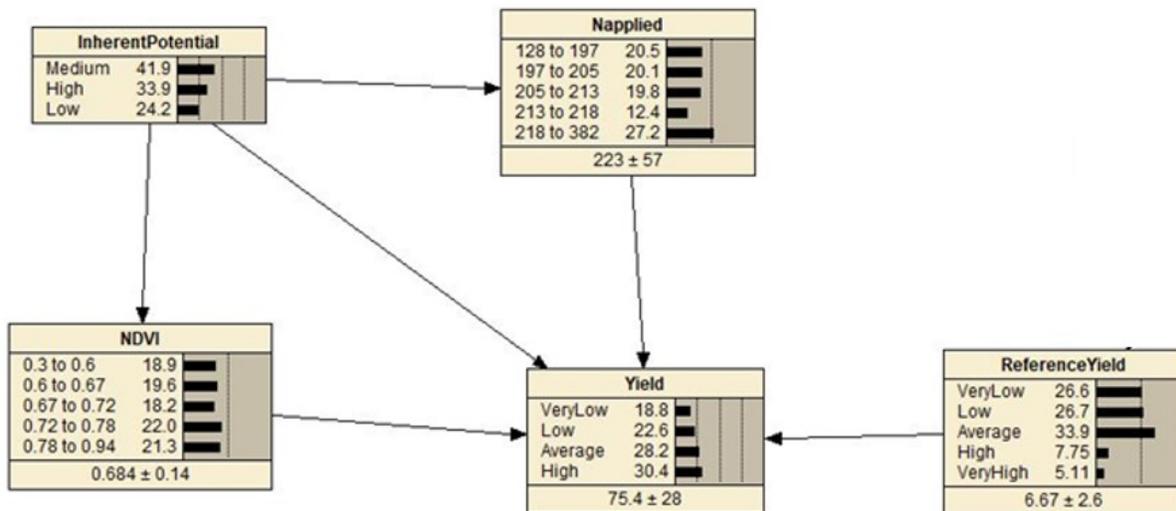


Fig. 7. Diagram for MLA learning model for predicting wheat yield based on Inherent Potential, NDVI in May, weather algorithm and total nitrogen application. The model was applied on a farm basis. Eight fields were used to learn the model and two to test it.

Results and discussion

Validation of predictions against ground-truthing field observations

For each country, some fields were used to train the model and some to test the model. Testing the model, a confusion matrix (Table 3) helps to evaluate the goodness of predictions.

Table 3. Confusion matrix from the German and UK fields.

Germany	Predicted	Very Low	Low	Medium	High	Very High	Actual
		243	95	117	51	14	Very Low
		78	172	216	68	101	Low
		45	96	169	70	129	Medium
		24	134	189	136	183	High
		18	49	174	300	1104	Very High

UK	Predicted	Very Low	Low	Medium	High	Very High	Actual
		243	95	117	51	14	Very Low
		78	172	216	68	101	Low
		45	96	169	70	129	Medium
		24	134	189	136	183	High
		18	49	174	300	1104	Very High

To evaluate the confusion matrix the following metrics were considered:

$$\text{Accuracy} = \frac{TP+TN}{TN+TN+FP+FN} \quad (2)$$

$$\text{Precision} = \frac{TP}{TP+FP} \quad (3)$$

where TP = true positive, TN = true negative, FP = false positive and FN = false negative.

In Table 4, the outcomes of the confusion matrix to calculate the precision and the accuracy of the model are presented.

Table 4. Outcomes from the confusion matrix for Germany and UK.

Germany	Very Low	Low	Medium	High	Very High
TP	243	172	169	1104	1824
TN	3290	2966	2770	2820	1903
FP	277	463	340	530	427
FN	165	374	696	489	541

UK	Very Low	Low	Medium	High	Very High
TP	255	147	145	108	271
TN	1278	1202	1225	1293	1187
FP	106	238	168	206	109
FN	114	166	215	146	186

The accuracy and the precision of the German and UK model were calculated using the equations mentioned above. The accuracy and precision for the model are 77%, 67% and 81%, 51% for Germany and UK respectively.

However, the predictions in the fields used to test the model were not as good as those in the fields used to train the model. For example, using two of the German fields, the Kriged maps of predicted and observed showed similar patterns of spatial variation (Fig. 8 and Fig. 9). Alzeyer field was used as a training field and Birnbaum field as a testing field.

Alzeyer field (5.76 ha)

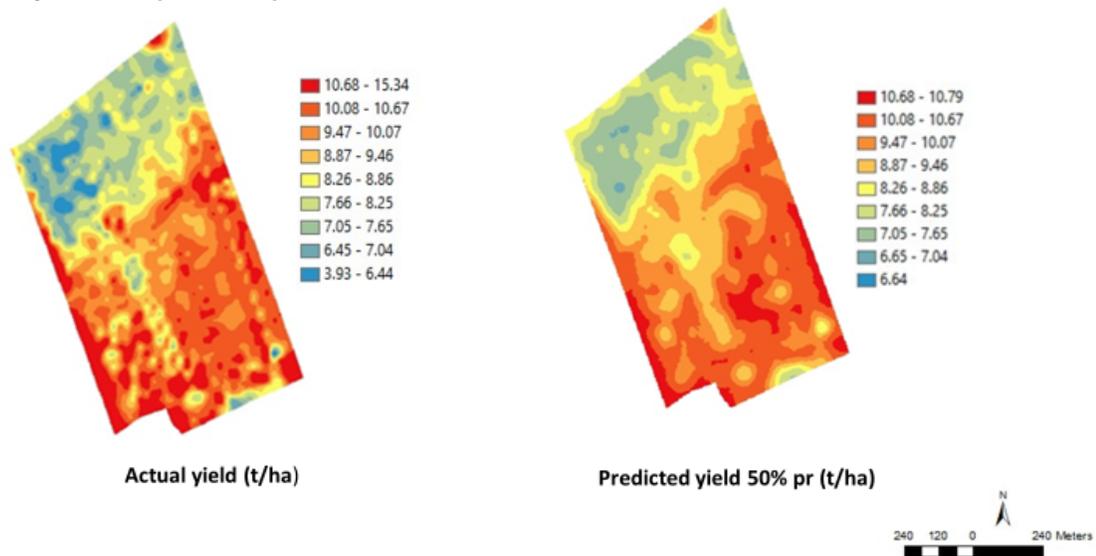


Fig. 8. Maps of actual and predicted wheat yield (t/ha) in Alzeyer (5.76 ha) field in Germany which was used to train the model. The predicted map on the right side is based on 50% probability. The values on the colour ramp cover slightly different ranges.

Birnbaum field (4.52 ha)

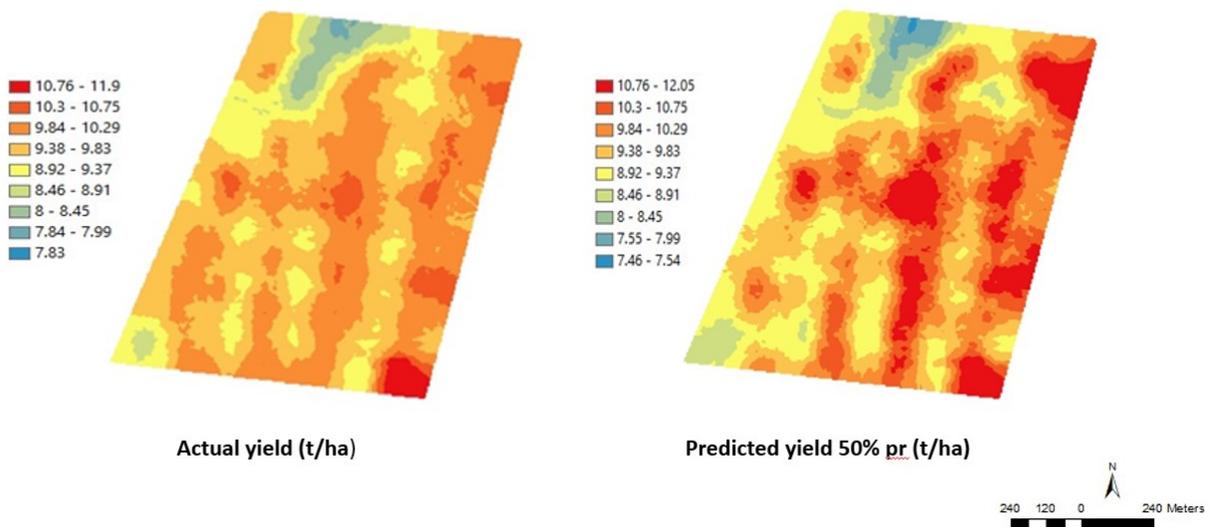


Fig. 9. Maps of actual and predicted wheat yield (t/ha) in Birnbaum (4.52 ha) field in Germany which was used to test the model. The predicted map on the right side is based on 50% probability. The values on the colour ramp cover slightly different ranges.

In the Birnbaum field which is one of the fields used for testing the model the correlation between actual and predicted values is 0.42 compared to the Alzeyer field that was used for training where the correlation between actual and predicted yield is 0.78. The fitted 1:1 line shows that the model

overestimated the yield in low areas of the field and underestimated the yield in high areas of the fields (Fig. 10). Sources of inaccuracy in predictions can be explained by the weed infestation, diseases or the model being overfitted.

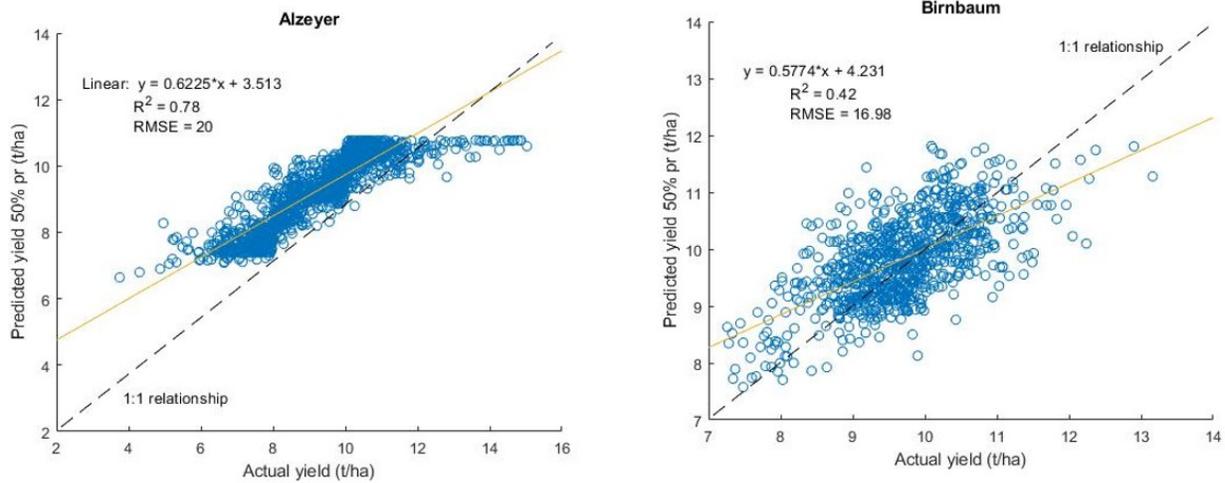


Fig. 10. Linear regression (yellow line) between actual and predicted grain yield in Alzeyer (left graph) and Birnbaum (right graph) fields. Data points are for each 6x6 m grid square in the fields. Predicted yields are at 50% probability. The dashed line is a 1:1 relationship.

In the UK, the predicted yield is again better for the fields that were used to train the model than the fields that were used to test the model. Fig. 11 shows similar patterns for the observed and predicted yield in the testing field. The correlation between actual and predicted yield for the testing field is 0.35 (Fig. 12).

Crowsley field (14.8 ha)

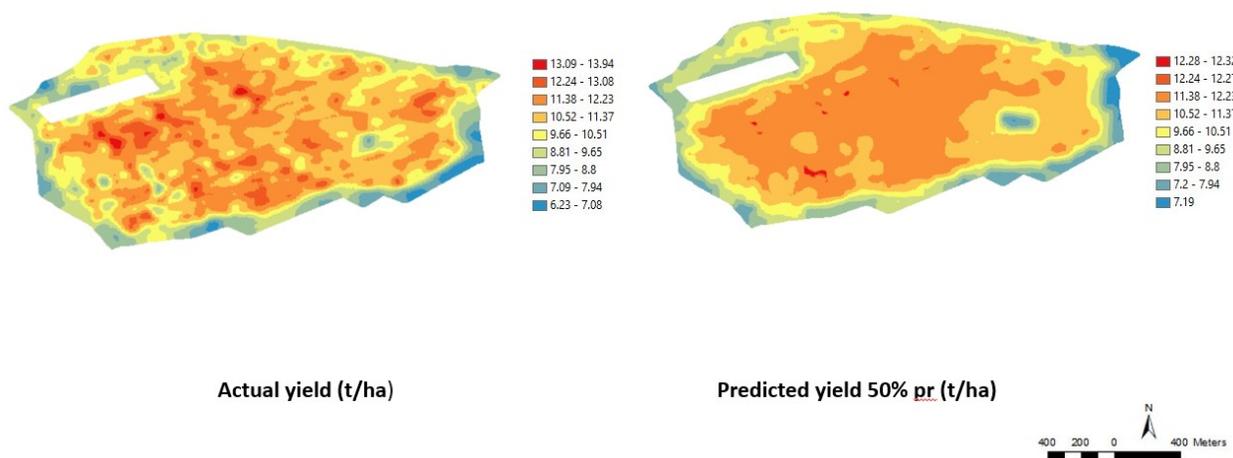


Fig. 11. Maps of actual and predicted wheat yield (t/ha) in Crowsley (14.8 ha) field in the UK which was used to test the model. The predicted map on the right side is based on 50% probability. The values on the colour ramp cover slightly different ranges.

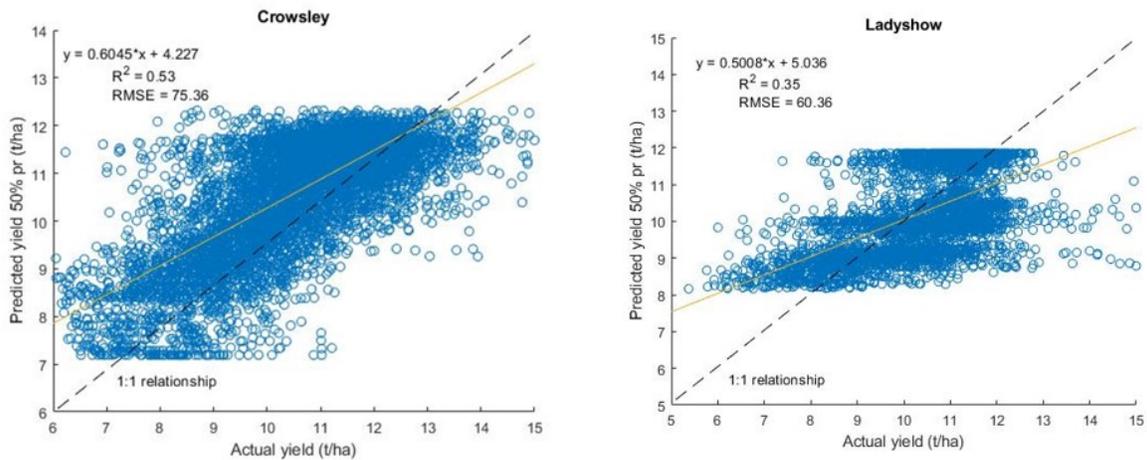


Fig. 12. Linear regression (yellow line) between actual and predicted grain yield in Crowsley (left graph) and Ladyslaw (right graph) fields. Data points are for each 6x6 m grid square in the fields. Predicted yields are at 50% probability. The dashed line is a 1:1 relationship.

Conclusion

The proposed model which is based on a Machine Learning Approach and more specifically on Bayesian Networks showed a good performance for both study areas, i.e Germany and the UK. The driving variables which were used for the yield predictions were Inherent Potential (IP), the normalized difference vegetation index (NDVI), the variable rate of nitrogen application and the weather data. The introduction of IP was an indicator of the representation of different spatial locations to support wheat growth based on Prior Inherent Potential (PIP) and historical yield data. The PIP was calculated based on topographic and soil data. The accuracy and precision of the model reached 77%, 67% and 81%, 51% for Germany and UK respectively. The correlation between actual and predicted yield of 50% probability was better for the fields that were used for training. For example, in Germany, the correlation between actual and predicted yield of 50% probability was 0.78 for the training field and 0.42 for the testing field. In the UK, the training and the testing field showed a correlation of 0.53 and 0.35 respectively.

Future work is needed to improve the model accuracy as the final algorithm will be integrated into the Agricolus platform. Thus, more fields will be added in 2022. In this platform the users could upload individual field data and depend on the output, they could decide their next strategies.

Acknowledgements

The authors are most grateful to the whole LINKDAPA team in Agricolus, John Deere, the University of Hohenheim and the University of Reading. We also thank our eight “key” farmers for 2020, and 2021 and technical staff, without whom, the research could never take place. Finally, we thank EIT Food, which is co-funded by the European Union, for supporting this work both financially and also for Kerstin Burseg’s encouragement and interest in the project.

References

- Addy, J.W., Ellis, R.H., Macdonald, A.J., Semenov, M.A. and Mead, A. (2020). Investigating the effects of inter-annual weather variation (1968–2016) on the functional response of cereal grain yield to applied nitrogen, using data from the Rothamsted Long-Term experiments. *Agricultural and Forest Meteorology* 284: article 107898.
- Chawla, V., Naik, H. S., Akintayo, A., Hayes, D., Schnable, P., Ganapathysubramanian, B., Sarkar, S. (2016). A Bayesian Network approach to County-Level Corn Yield Prediction using historical data and expert knowledge.

- Chlingaryan, A., Sukkarieh, S., Whelan, B. (2018). Machine learning approaches for crop yield prediction and nitrogen status estimation in precision agriculture: a review. *Comput. Electron. Agric.* 151, 61–69.
- Dhivya, E., Durai Raj, V., Vishal, S., Albert Y., Z., Kathiravan, S. (2018). Forecasting yield by integrating agrarian factors and machine learning models: A survey, *Computers and Electronics in Agriculture*, 155, 257-282, ISSN 0168-1699, <https://doi.org/10.1016/j.compag.2018.10.024>.
- Filippi, P., Jones, E.J., Wimalathunge, N.S. (2019). An approach to forecast grain crop yield using multi-layered, multi-farm data sets and machine learning. *Precision Agric* 20, 1015–1029.
- Fu, Z., Jiang, J., Gao, Y., Krienke, B., Wang, M., Zhong, K., Cao, Q., Tian, Y., Zhu, Y., Cao, W. (2020). Wheat growth monitoring and yield estimation based on multi-rotor unmanned aerial vehicle. *Remote Sens.*, 12, 508.
- Gandhi, N., Armstrong, L: J, Petkar, O. (2016). Predicting Rice crop yield using Bayesian networks, International Conference on Advances in Computing, Communications and Informatics (ICACCI), 795-799, DOI: 10.1109/ICACCI.2016.7732143.
- Goodfellow, I., Bengio, Y., Courville, A. (2016). *Deep Learning*. MIT Press. <http://www.deeplearningbook.org> (Last accessed: April 21, 2022).
- Jensen, F.V. (2001). *Bayesian networks and decision graphs*. Springer-Verlag, New York. ISBN 0-387-95259-4.
- Kristensen, K., Rasmussen, Ilse A. (2002). The use of a Bayesian network in the design of a decision support system for growing malting barley without use of pesticides, *Computers and Electronics in Agriculture*, 33, 197-217.
- Kuschner, K., Malyarenko, D., Cooke, W., & Cazares, L., Semmes, O., Tracy, E. (2010). A Bayesian Network Approach to Feature Selection in Mass Spectrometry Data. *BMC bioinformatics.* 11, 177.
- Liakos, K. G., Busato, P., Moshou, D., Pearson, S., Bochtis, D. (2018). Machine learning in agriculture: a review *Sensors (Switzerland)*, 18 (8).
- Natale, A., Antognelli, S., Ranieri, E., Cruciani, A., Boggia, A. (2020). A novel cleaning method for yield data collected by sensors: A case study on winter cereals. 20th International Conference, Cagliari, Italy, July 1–4, 2020, Proceedings, Part V.
- Newlands, N., Townley-Smith, L. (2010). Predicting energy crop yield using bayesian networks. *Computational Intelligence*.
- Somvanshi P, Mishra B. N., (2015). *Machine learning techniques in plant biology PlantOmics: The Omics of Plant Science*, Springer India, New Delhi, 731-754.
- Uusitalo, L. (2007). Advantages and challenges of Bayesian networks in environmental modeling. *Ecological Modelling*. 203. 312-318.
- Villordon, A., Sheffield, R., Rojas J., and Chiu, Y. (2011). Development of simple Bayesian belief and decision networks as interactive visualization tools for determining optimal in-row spacing for 'Beauregard' sweet potato.
- Walters, C.J., and Martell, S.J.D. (2004). *Fisheries Ecology and Management*. Princeton University Press. ISBN 0-691-11545-1.
- Wang, X., Jianxi, H., Quanlong F., Dongqin, Y. (2020). Winter Wheat Yield Prediction at County Level and Uncertainty Analysis in Main Wheat-Producing Regions of China with Deep Learning Approaches. *Remote Sensing* 12(11), 1744.
- Willcock, S., Hooftman, D.A., Bagstad, K.J., Balbi, S., Marzo, A., Prato, C., Sciandrello, S., Signorello, G., Voigt, B. (2018). Machine learning for ecosystem services. *Ecosystem services* 33, 165–174
- Xiangying, X., Ping, G., Xinkai, Z., Wenshan, G., Jinfeng, D., Chunyan, L., Min, Z., Xuanwei, W. (2019). Design of an integrated climatic assessment indicator (ICAI) for wheat production: A case study in Jiangsu Province, China, *Ecological Indicators*, 101, 943-953.