Multi-year multi-crop correlation analysis in Brasov area

Ioana C. Plajer¹, Alexandra Băicoianu¹, Matei Debu¹, Maria Ștefan², Mihai Ivanovici¹, Corneliu Florea^{1,3}, Adrian Ghinea², Luciana Majercsik¹

¹ Transilvania University, Brașov, Romania

² National Institute of Research and Development for Potato and Sugar Beet, Braşov, Romania ³ National University of Science and Technology Politehnica Bucharest, Romania

Keywords: Sentinel-2 data, NDVI time series, correlation analysis, crop identification.

Abstract. Artificial Intelligence (AI) models are currently deployed in smart agriculture for various applications like crop monitoring and identification or yield estimation. AI models rely on huge amounts of data; thus, the first concern is the size and quality of labeled data for training such models. The second main concern is the explainability of the results produced by AI models. In this article, based on the 5-year data set we previously produced and published, we perform a correlation analysis in the attempt of explaining the performance of an AI model in a crop identification scenario.

1. Introduction

Crop monitoring and identification are examples of applications that play a significant role in modern agriculture. With advancements in technology, monitoring systems are increasingly automated, relying on remotely-acquired data from sensors mounted on drones, airplanes, or satellites. The growing number of satellite missions has made free data more accessible to researchers and institutions, driving innovation across various fields and advancing smart agriculture (Ivanovici et al. 2024).

Multispectral satellite imagery, such the one from the EU Copernicus Earth Observation (EO) programme, provides valuable insights into crop types, growth stages, health, and soil-vegetation characteristics. One of the most widely-used and accessible satellites for vegetation analysis is Sentinel-2, which captures multispectral images across 13 spectral bands, covering both the visible and near-infrared (NIR) spectrum.

To enable automatic crop monitoring and crop identification based on satellite data, the first essential step is accurately identifying crops from this data. A common approach, accepted as the *de facto* approach (ECA, 2020), involves analyzing time series, often based on the well-known Normalized Difference Vegetation Index (NDVI) (Pettorelli, 2013; Qin et al., 2021).

In recent years, artificial intelligence (AI) and machine learning (ML) have become key tools for crop identification. However, the effectiveness of these models heavily depends on the quality and relevance of the datasets used. The accuracy of the results can be significantly impacted by data correlations and similarities, making dataset selection a critical factor in model performance. Various datasets exist (Sumbul et al., 2019), (Weikmann et al., 2021),

This paper presents a multi-year multi-crop correlation analysis across 47 parcels with 17 agricultural crop types, for the purpose of explaining the classification results based on an ML model.

1.1 Data

The current study utilizes a multispectral dataset, called μ DACIA5, constructed based on the Sentinel-2 satellite optical data (DOI: 10.5281/zenodo.14283242) collected over a time period of five years (2020–2024) in Northern Braşov area, Romania. The research focuses on parcels managed by the National Institute of Research and Development for Potato and Sugar Beet, Braşov, Romania (NIRDPSB), which provided detailed crop type data over multiple years, that allowed the labelling of the data for ML tasks. Such data enables insights into crop distribution, growth stages, and phenological events. The crop correlation analysis was performed based on the NDVI time series, as NDVI time series are often used both for monitoring and identification of agricultural crops.

1.2 Crop spectral signature analysis

As an initial step, we plotted the spectral signatures and averaged spectral signatures for each crop, generating curves for each cultivated parcel. These curves were created using Sentinel-2 data from the crop growing season, while capturing key stages such as planting, flowering, growth, and maturity. The evolution of crops over time is evident from these graphical representations. For instance, Figure 1 illustrates the average spectral signature of late potatoes, with distinct lines representing average reflectance curves for each image. This approach highlights the crop's growth period, allowing the identification of specific growth phases based on the curves. Such evolution of the spectral curves in time is characteristic for most of the types of agricultural crops.



Figure 1. Average reflectance curves for parcel no. 97 with late potato crop (code 253) in 2023.

The potato plant's phenological growth stages, as outlined by (Meier, 1997), can be summarized into four main phases: planting to sprouting, sprouting to bud flowering, bud flowering to flowering, and flowering to maturity (Munteanu et al., 2008; Nemes et al., 2008). Consider the following example: potato planting on parcel no. 97 began on May 14, 2023. Sentinel-2 data reveal key growth stages of the potato crop on parcel 97. Early reflectance curves (May 23 and June 2) in bands B5, B6, B7, B8, B8A, and B9 were low, indicating the planting phase. These bands, particularly B5-B7 for vegetation classification and B8 for chlorophyll sensitivity, show minimal vegetation activity after planting. By June 24, a slight increase in reflectance marked the second growth phase. From July 17 to August 18, the crop reached maturity, producing abundant chlorophyll and dense vegetation with high reflectance. After August 21, reflectance began to decline as chlorophyll production decreased, signaling plant senescence. A sharper decrease was observed after August 31 as the plants dried, culminating in harvest on September 10. Bands B11 and B12, used to measure soil and vegetation moisture, provided insights into crop water content. Monitoring these bands helps detect water stress, enabling timely irrigation to maintain crop health. This analysis effectively tracks crop growth, vegetation status and moisture dynamics.

1.3 NDVI monitoring

Vegetation indices play a fundamental role in crop monitoring and analysis and one of the most relevant and widely used is the NDVI (Pettorelli, 2013; Qin et al., 2021), which is calculated as:

$$NDVI = \frac{NIR - RED}{NIR + RED},$$
(1)

where, NIR represents the reflectance value of the band selected for near-infrared and RED the reflectance value of the red band. For a comprehensive list of available vegetation indices see (Ivanovici et al., 2024).

Sentinel-2 multispectral images provide the following two spectral bands: B4 (664.6 nm, RED) and B8 (832.8 nm, NIR) to calculate NDVI, an indicator of vegetation status and health ranging from -1 to 1. In summary, negative values mark bare soil and artificial, built-up areas, values in the interval (0, 0.33], indicate unhealthy or sparse vegetation, while 0.66 is the threshold between healthy and moderately healthy vegetation.

NDVI was calculated for all parcels across all acquisition dates, generating time series to track vegetation growth. Figure 2 illustrates NDVI evolution for all pixels corresponding to parcel 39, where sugar beet was cultivated in 2023. Thirty observations were made at uneven intervals due to weather conditions, thus availability of Sentinel-2 data. Sugar beet was selected as it is one of the representative crops for the region and significant for NIRDPSB research activity.



Figure 2. Temporal evolution of NDVI for each pixel in parcel 39 cultivated with sugar beet in 2023.

For each parcel, we also calculated a representative NDVI curve in time, as an average of all the NDVIs of the pixels in the parcel. A representative NDVI curve was calculated for each parcel by averaging the NDVI values of all its pixels. For parcel 39 in 2023, Figure 3 illustrates this average NDVI curve, representing the general NDVI evolution for the parcel (in red). In addition, we calculated the representative NDVI curve for each crop, by calculating the average NDVI for all the pixels, in all parcels within the dataset, with the same crop. For sugar beet in 2023 there are three parcels and the NDVI AVG curve is represented in Figure 3 with magenta. The seeding and harvesting events are indicated with blue arrows.



Figure 3. The representative average NDVI curve for parcel 39 cultivated with sugar beet in 2023 compared against the average NDVI for all parcels cultivated with sugar beet in 2023.

Representative NDVI curves illustrate crop growth, decay, and harvest periods. For sugar beet, the maximum NDVI value (\sim 0.7) slightly exceeds the 0.66 threshold for moderately healthy vegetation, highlighting possible discrepancies with in-situ measurements due to satellite image processing, sensor aggregation, and atmospheric effects. This may be also due to the characteristic canopy which does not cover the ground completely, requiring further correction of the NDVI value. NDVI time series can also be used to detect anomalous growth within a parcel by comparing deviations from the representative curve. Such curves, calculated over multiple years and correlated with meteorological data, could provide broader insights into crop health.

2. Single-year correlation analysis

In this section we perform a correlation analysis for a single year, based on the Pearson correlation coefficient computed between pairs of representative average NDVI time serios of crops. We present the analysis result as the correlation matrix, thus showing the degree of linear correlation between the chosen crop features. For the single-year we focused on 2023 as the reference year for our analysis.

In Figure 4, the correlation matrix for the per-parcel average NDVI time series of our dataset is displayed. Each label indicates the parcel index and the corresponding crop type. The correlation values between the average NDVI curves of the two parcels are represented using a colour gradient ranging from yellow to green, with yellow indicating low correlation and dark green indicating very high correlation, aligning with the interpretation of NDVI.



Figure 4. Correlation matrix of average NDVI time series per parcel.

As can be observed from Figure 4, parcels with the same crop exhibit a very large correlation, as expected. Various cereals (codes 101 and 108) or potato varieties, (codes 254 and 255), can be easily confused due to their high similarity, as indicated by the correlation coefficients. We observed that some parcels with different crops, such as sugar beet on parcel 39 (code 3017) and corn on parcel 38 (code 108), show a high correlation, despite the crops being distinct.

In order to better interpret these correlations, we calculated the correlation between all crops using the representative average curves for each crop, derived from all pixels in the respective parcels. The resulting correlation matrix is displayed in Figure 5, using the same colour scheme as in Figure 4. One can notice in Figure 5 the large correlation values for the pairs of crops with codes 253, 254, 255, and 2557, which all represent different types of potatoes, but also between sugar beet (code 3017) and corn (code 108), as observed before. Other crops, like corn (108), and autumn wheat (101), with low correlation are visibly different. In 2023 only 12 of the 17 crops present in our dataset were planted on the considered parcels and are thus represented in the correlation matrices. Such a correlation analysis shows the limitations of an analysis based on the NDVI time series, which may be very similar for different types of crops (that were planted approximately in the same period).





3. Multi-year correlation analysis

Starting from the analysis in 2023, we consider the most and the least correlated crops and verify the consistency of the observations made for 2023 for the other 4 years: 2020, 2021, 2022 and 2024. In particular, 2022 was a dryer year compared to the other 4 years. In this section we shall present the results and interpretation of our multi-year correlation analysis.

For the correlation analysis, we computed a representative average NDVI time curve for each crop in each year and estimated the correlation between time curves from two different years, such as 2023 and 2024. The 2022 and 2024 years were both dry years for the agricultural activity in the region, 2023 was normal from the point of view of precipitations. When comparing crops across years, several factors must be considered.

First, due to crop rotation and diversification policies, some crops may appear in certain years but are absent in others. To ensure a consistent correlation matrix, we included only crops present in the compared pairs of years.

Another challenge in cross-year correlation analysis is the variation in image acquisition dates and the number of observations. For instance, in 2023, images were collected on 30 different days, whereas in 2024, data was available from 40 dates. As a result, the time series differ in length, making a direct Pearson correlation calculation infeasible.

One approach to analysing the correlation between different crops across two years, as well as the autocorrelation of the same crop between two different years, is to calculate the crosscorrelation between each of the considered pair of time series. Cross-correlation measures the similarity between two sequences, even of different lengths, by shifting one sequence relative to the other.

However, a major challenge in using cross-correlation is the irregular timing of multispectral image acquisitions in our dataset. The intervals between consecutive acquisition dates vary significantly, ranging from 3–4 days to an entire month. As a result, applying cross-correlation directly to NDVI time series from different years can cause sequence misalignment, leading to inaccurate correlation results.

To ensure equal spacing between the considered dates, we interpolated the NDVI time curves, considering that NDVI values typically do not change abruptly over short time frames. After experimenting with both linear and cubic interpolation, we found that linear interpolation was the most suitable for this task. The interpolation was applied to each curve over the entire year, from January 1st to December 31st, with a fixed interval of 5 days between consecutive NDVI values.

The comparison of average NDVI time series was conducted for each crop that was present in the years 2022, 2023, and 2024 as follows. For each common crop, we computed the interpolated average NDVI time series for 2023 and evaluated different shifts of the corresponding average NDVI time series from 2022 and 2024 to determine the optimal alignment.

To perform these calculations, we utilized the cross-correlation function provided by the SciPy library in Python (SciPy. (n.d.), 2025). This function computes the correlation between two discrete sequences, x and y, representing the time series from two different years, based on the following formula:

$$Corr[k] = \sum_{l=0}^{||x||-1} x_l y_{l-k+N-1}$$
(2)

in which $k = 0, 1, ..., ||\mathbf{x}|| + ||\mathbf{y}|| - 1, ||.||$ indicates the length of the time series, $N = \max(||\mathbf{x}||, ||\mathbf{y}||)$ and $y_m = 0$ outside the range of y. The second time-series is shifted by k steps.

This calculation produces a set of values for each pair of time series, corresponding to the number of applied shifts. Larger values indicate stronger correlations; however, the results do not fall within the interval [-1, 1]. To standardize them, we applied the normalized cross-correlation, which involves first subtracting the mean from each time series and then normalizing by its standard deviation before computing the cross-correlation. Among all possible shifts between two time series, we selected the maximum correlation value obtained.

In Figure 6 are illustrated pairs average NDVI time series for corn (crop code 108), late potato (crop code 253) and alfalfa (crop code 9748) from 2023 together with the accordingly shifted time series from 2022, while in Figure 7 are presented similar pairs for 2023 and 2024.



Figure 6. Pairs of interpolated average NDVI time curves for 2022 and 2023. The curves of 2022 are shifted, as to obtain the maximum correlation with the curve of 2023.

To verify the consistency of the one-year analysis conducted on 2023, we calculated the normalized cross-correlation between the common crops for the following year pairs: (2022, 2023) and (2023, 2024). The rationale behind this choice is that 2022 and 2024 were relatively similar in terms of weather conditions, while 2023 was a normal to rainy year. We aimed to observe whether this similarity, or dissimilarity, would also be reflected in the correlation matrices. For each pair of years, we considered pairs of representative interpolated average NDVI curve in time for each of the crops present in both these years and calculated the cross-correlation as explained in the previous section. We



represented the resulted correlation matrix as a heat map like the one obtained for year 2023.

Figure 7. Pairs of interpolated average NDVI time curves for 2022 and 2023. The curves of 2022 are shifted, as to obtain the maximum correlation with the curve of 2023.

In Figure 8 we show the cross-correlation matrix for the pairs of interpolated curves in years (2022, 2023), while in Figure 9 we plotted the cross-correlation matrix for the pairs of interpolated curves in years (2023, 2024). The curves in 2022 and respectively 2024 were shifted along the curves in 2023, as to yield the best correlation. From the two figures, one can observe that the correlation values for the pair (2023, 2024) in Figure 8 are larger than those from the pair (2022, 2023) in Figure 9, suggesting that the similarity in weather conditions could influence this relationship. Furthermore, when compared with the correlation matrix of crops from 2023 presented in Figure 5, the strong and weak correlations appear between the same pairs of crops, such as the strong correlation between peas and winter wheat and the weak correlation between sugar beet and winter wheat. What is somewhat unusual, however, is the weak correlation between sugar beet in 2024 and sugar beet in 2023, indicating that the different environmental conditions in the two years considerably affected the NDVI time series-based analysis.



Figure 8. Cross-correlation matrix of average NDVI signature per crop. The pairs of interpolated signatures considered are from 2022 and 2023, where the 2022 signatures were shifted, as to obtain the best correlation.



Figure 9. Cross-correlation matrix of average NDVI signature per crop. The pairs of interpolated signatures considered are from 2024 and 2023, where the 2024 signatures were shifted, as to obtain the best correlation.

To illustrate the low correlation between NDVI curves of the same crop in different years, in Figure 10 we show, as an extreme case, the NDVI curves for alfalfa in 2021 and 2023. One can notice the very different behaviour which explains the very low correlation coefficients between identical crops considered for different years.



4. NDVI time series-based crop identification

To analyze how the correlation between NDVI time series influences the ability of a neural network (NN) to differentiate between crops, we trained a fully connected neural network (FCNN) on the NDVI time series for all the pixels in the common crops from the years 2022, 2023, and 2024. The training set was built using all the time series from 2023, and the model was tested on the time series from 2022 and 2024. All the time series were interpolated over the interval [0, 365] corresponding to a full year data at 1-day time resolution, then down sampled to a 5-day time resolution, resulting in a total of 73 values per time series. Additionally, the time series for 2022 and 2024 were shifted for each crop by a specific amount to achieve the maximum correlation with the corresponding crop in 2023. The common crops and their respective shifts are presented in Table 1. The shift value represents the number of shift steps, with each step corresponding to 5 days. A positive shift indicates a rightward shift, while a negative shift corresponds to a leftward shift.

Crop name	Crop	Shift for	Shift for	
	code	2022	2024	
Winter wheat	101	2	-2	
Spring wheat	1010	1	-3	
Corn	108	1	1	
Peas	151	0	-2	
Late potato	253	1	0	
Sugar beet	3017	1	4	
Temporal grassland	450	-3	0	
Alfalfa	9748	7	0	
Permanent grassland	606	3	0	

Table 1: Time series shifts (in number of data samples) per crop for 2022 and 2024 with respect to 2023.

Since the interpolation at the beginning and end of the year involves extrapolation, the interpolation error tends to be larger in these regions. Additionally, shifting the time series introduces some zero values at the ends of the interval. To address these issues, we removed the first 5 and the last 10 interpolated values, resulting in NDVI interpolated time series with a length of 58.

The architecture of the FCNN is presented in Figure 11 and consists of one input layer of 58 neurons, corresponding to the 58 values of the interpolated time series, three hidden layers of respectively 64, 32 and 8 neurons and of an output layer of 9 neurons corresponding to the 9 common crops of the years 2022, 2023 and 2024. The activation on the hidden layers is the Rectified Linear Unit (ReLu) and on the output layer softmax. We used as an optimizer the AdamW function, with learning rate $Ir=10^{-3}$ and a weight decay of also 10^{-3} . The loss function considered was cross entropy loss.



Figure 11: FCNN trained on the interpolated NDV time series of year 2023. On the right side the class labels for each crop code are indicated.

The FCNN was trained over 400 epochs in several folds and reached an average accuracy of 85% on the train set. But we observed that after about 300 epochs the accuracy on the 2022 and 2023 sets stabilize or even get worse, indicating an overfit on the training set.

Figure 12 presents the confusion matrix for the training set (year 2023) after 280 epochs, with a total accuracy of 83%. From the confusion matrix, one can see that 5 crops are correctly recognized (codes 101, 108, 151, 606 and 1010), while late potato (code 253) and sugar beet (code 3017) are consistently misclassified as corn. Additionally, temporal grassland is misclassified as permanent grassland. The latter confusion is not entirely unexpected, as both crops share significant similarities. As we can observe in Figures 6 and 7 the average NDVI curves of corn and late potato are very similar, this can be an explanation of the misclassification done by the deployed ML model.

When using the trained network to determine the crops in 2022, the detection accuracy is, at best, around 60%. The confusion matrix in Figure 13 shows that while some crops are classified very accurately, others exhibit significant confusion with different classes. As seen in Figure 13, winter wheat, corn, peas, and spring wheat are classified with an accuracy exceeding 87%. However, temporal grassland (code 450) and permanent

grassland (code 606) are consistently misclassified as winter wheat. Looking at the correlation matrix in Figure 8, we notice a strong correlation between these three crops, which could explain this misclassification. FCNN misclassifies late potato and sugar beet as corn, a pattern that also occurs in 2023 for these crops.











Figure 14: The confusion matrix obtained when classifying the interpolated and shifted NDVI curves of 2024 by the FCNN trained on 2023 data.

The results of the FCNN-based classification on the interpolated and shifted NDVI time series of 2024 are illustrated by the confusion matrix presented in Figure 14. The average best accuracy in the different folds for 2024 was approximately 35%. One can notice that in the case of 2024, only winter wheat and corn are correctly classified. The classification pattern observed in Figure 13 corresponds, to some extent, to the correlation matrix pattern which can be seen in Figure 9.

5. Conclusions

In this article we presented a single- and multi-year correlation analysis for agricultural crops, for the purpose of explaining the ML-based crop classification results. The µDACIA5 dataset was used in our experiments, which is based on the Sentinel-2 multispectral data collected over 5 years from 47 parcels with 17 crop types, in the Brasov North area, Romania. Such multi-spectral datasets are very useful for the development and validation of crop analysis tools, such as AI- or ML-based crop identification, in smart agriculture. The µDACIA5 dataset provides reliable, accurately labelled data, verified in-situ by the owner of the land, NIRDPSB Brasov. The correlation analysis was performed on the NDVI time series of the under-study crops, considering the representative average curves for each parcel and each culture. Given the availability of the Sentinel-2 data in different years, interpolation and resampling of the NDVI time series was performed. Given the difference in seeding times in different years, due to the weather conditions, required shifting the time series in time for increasing the efficacy of the analysis. Some consideration on the correlation of the different cultures, as represented by the specific NDVI was provided. The computed correlation matrices showed high correlation for both identical and different crops, while for certain identical crops, the correlation was low. The results also pointed out the impact of the weather conditions on the correlation. The correlation analysis was followed by a ML-based classification of the crops, using an FCNN model, in a crop identification scenario. The experimental results using a FCNN showed good similarity with the results of the correlation analysis, pointing out the limitations of the NDVI time series-based, ML-based identification of agricultural crops.

6. Acknowledgment

Funded by the European Union. The AI4AGRI project entitled "Romanian Excellence Center in Artificial Intelligence on Earth Observation Data for Agriculture" received funding from the European Union's Horizon Europe research and innovation program under grant agreement no. 101079136.



References

ECA – European Court of Auditors (2020). Using new imaging technologies to monitor the Common Agricultural Policy: steady progress overall, but slower for climate and environment monitoring.

Ivanovici, M., Olteanu, G., Florea, C., Coliban, RM., Ștefan, M., Marandskiy, K. (2024). Digital Transformation in Agriculture. In: Ivascu, L., Cioca, LI., Doina, B., Filip, F.G. (eds) Digital Transformation. Intelligent Systems Reference Library, vol 257. Springer, Cham. https://doi.org/10.1007/978-3-031-63337-9 9 Meier, U. (1997). Growth stages of mono-and dicotyledonous plants (BBCH Monograph). Blackwell. https://doi.org/10.5073/20180906-074619

Munteanu, L. S., Cernea, S., Morar, G., Duda, M. M., Vârban, D. I., & Muntean, S. (2008). Fitotehnie. AcademicPres.

Nemes, Z., Baciu, A., Popa, D., Mike, L., Petrus-Vancea, A., & Danci, O. (2008). The study of the potato's life-cycle phases important to the increase of the individual variability. Analele Unversitatii Oradea Fasc. Biol., 15, 60-63.

Pettorelli, N. (2013). The normalized difference vegetation index. Oxford University Press. https://doi.org/10.1093/acprof:osobl/9780199693160.001.0001

Qin, Q., Xu, D., Hou, L., Shen, B., & Xin, X. (2021). Comparing vegetation indices from Sentinel-2 and Landsat 8 under different vegetation gradients based on a controlled grazing experiment. Ecological Indicators, 33, 108363. https://doi.org/10.1016/j.ecolind.2021.108363

SciPy. (n.d.). (2025) *scipy.signal.correlate* — *SciPy* v1.10.1 *Manual*. Retrieved March 18, 2025, from https://docs.scipy.org/doc/scipy/reference/generated/scipy.signa l.correlate.html

Sumbul, G., Charfuelan, M., Demir, B., & Markl, V. (2019). BigEarthNet: A Large-Scale Benchmark Archive for Remote Sensing Image Understanding. In IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium (pp. 5901-5904).

Weikmann, G., Paris, C., & Bruzzone, L. (2021) TimeSen2Crop: A Million Labeled Samples Dataset of Sentinel 2 Image Time Series for Crop-Type Classification. In IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 14, 4699-4708.

Appendix

In the Table 2 we show the name of the crops and their corresponding codes used by APIA (Agency for Payments and Interventions in Agriculture), Romania.

Fable 2.	Crop	names	and	corres	ponding	APIA	codes.
----------	------	-------	-----	--------	---------	------	--------

Crop name	APIA code		
Common winter wheat	101		
Common spring wheat	1010		
Corn	108		
Peas	151		
Late potatoes	253		
Other potato crop	254		
Potatoes for seed	255		
Potatoes for seed	2557		
Sugar beets	3017		
Temporal grassland	450		
Alfalfa	9748		
Permanent grassland	606		
Corn silage	131		
Soybean	2037		
Alfalfa	9747		
Winter rapeseed	202		
Sunflower	123		