# Satellite Dataset Visual Analysis for Remote Soil Nutrient Estimation

Andrés Isaza-Giraldo<sup>14</sup>, Manuel Pereira<sup>23</sup>, Rafael Candeias<sup>2</sup>, Lucas Pereira<sup>24</sup>

<sup>1</sup> Faculdade de Belas-Artes, U. Lisboa, 1249-058 Lisbon, Portugal giraldo@edu.ulisboa.pt

<sup>2</sup> Instituto Superior Técnico, U. Lisboa, 1049-001 Lisbon, Portugal

{manuel.afonso.pereira,

rafaelmcandeiras, lucas.pereira}@tecnico.ulisboa.pt

<sup>3</sup> Instituto de Engenharia de Sistemas e Computadores - Investigação e Desenvolvimento, 1049-001 Lisbon, Portugal

<sup>4</sup> Interactive Technologies Institute, LARSyS, 9050-100 Funchal, Portugal

Abstract. This paper proposes a methodology for visualizing satellite-based machine learning (ML) datasets to understand the visual components that will be used as inputs for developing ML models. The proposed methodology uses t-Distributed Stochastic Neighbor Embedding (T-SNE) methods to create visualizations of satellite images leveraging models that were pre-trained in ImageNet. T-SNE is a self-supervised learning tool used to transform high-dimensional spaces into two- or three-dimension embeddings, making it easier to visualize a broad dataset in a single image or space. The methodology is demonstrated using the LUCAS Topsoil Analysis dataset with satellite images from Sentinel-2. The dataset was constructed using the TerraSense Toolkit (TSTK). The T-SNE visualization tool aims to improve ML research by providing a clearer visual understanding of satellite-based datasets.

**Keywords:** Satellite imagery · t-SNE (t-Distributed Stochastic Neighbor Embedding) · Soil Sensing · Machine-Learning · Agroindustry

# 1 Introduction

Soil surveying is a process that reveals crucial information about the soil composition, such as nutrients, that could have a great impact on a more efficient and effective use of land. Traditionally this has been a time-consuming and expensive process that requires sampling on the field and intensive testing. Despite the progress made in the agricultural applications of ML using satellite images (e.g.,[1, 2]), datasets are still a challenging matter as there are still few ground-truth available and because of specific limitations of satellite photography. According to Lillesand et al. "When we look at aerial and space images, we see various objects of different sizes, shapes, and colors. (...) The images contain raw pixel data. These data, when processed by a human interpreter's brain, become usable information [3]." Following this reasoning, it becomes clear that visually analyzing the dataset and trusting human capacity for interpretation is one of the paths toward a better dataset that could ultimately improve the performance of ML algorithms.

Against this background, this paper proposes a methodology to visualize satellitebased ML datasets to understand better the challenges associated with the visual components that utterly affect the performance of ML algorithms. More precisely, the proposed methodology uses t-Distributed Stochastic Neighbor Embedding (T-SNE) [4, 5] methods to create visualizations of the satellite images leveraging models that were pretrained on ImageNet [6].

T-SNE methods have been used in many disciplines, for instance, data science, medicine, social sciences, media art. This self-supervised learning tool transforms highdimensional spaces into two- or three-dimension embeddings, making it easier to visualize a broad dataset in a single image or space. The proposed methodology is demonstrated using the LUCAS Copernicus dataset with satellite images from Sentinel-2 [7– 9]. We also analyze the dataset for two specific crops, maize, and common wheat, as these are among the most represented crops in the used dataset.

Although t-SNE has commonly been used to visualize the distance between satellite images according to their classification vectors, e.g., there are few examples of t-SNE used to plot satellite images themselves for the sake of visually embedding datasets altogether. One prominent example would be the process of creating the interactive platform Land Lines by Zach Lieberman and Google Data Arts Team for which they plotted several satellite images together in a grid before developing their interactive experience [10]. This later type of t-SNE is more commonly used in digital arts. Its application could have a significant impact on machine-learning dataset visualization, for example, for dataset cleaning purposes.

The remaining of this paper is organized as follows. Section 2 presents the datasets, tools, and methodology used to visualize the datasets. Section 3 presents and discusses the result of a simple case study that was developed to demonstrate the proposed method. The paper concludes in Section 4 with a summary of the results and future work directions.

# 2 Materials and Methods

This section first describes the datasets and tools used in this research. Then, the proposed methodology is described in detail.

#### 2.1 Dataset and Tools

The dataset was constructed using information given in the LUCAS Survey, an effort of the European Soil Data Centre (ESDAC) for which there have been collected thousands of samples of soil in the European Union described some of its properties. The ground-truth taken were the ones intersecting three different surveys, as observed in Fig. 1: the LUCAS Topsoil Analysis, which contains information about soil composition, the LUCAS Copernicus which includes information on land use; and LUCAS 2018 that contains GPS coordinates. The LUCAS Copernicus survey is not a complete dataset per se, as it does not contain aerial images of the fields. The images were obtained from Sentinel-2 satellite, a multispectral satellite orbiting Earth since 2016.



Fig. 1: Dataset construction workflow.

The dataset construction workflow was introduced by author Manuel Pereira on the TerraSense Toolkit (TSTK) [11]. For each specific sample of the Lucas Survey, several images were taken in a range of five days prior to and five days after the date on which the ground truth was taken and as long as the maximum cloud coverage was less than 80%. This is because the satellite captures every point on earth only twice or three times a week due to its orbit around Earth. In many cases, clouds might block the view of the direct soil rendering the images unusable.

The Sentinel-2 uses a multispectral camera with 12 bands ranging from the visual spectrum to short-wave infrared. The spatial resolution for the visual spectrum is 10m, while it varies from 10 to 60 m on other bands. The LUCAS Survey also contains a polygon of up to 51 meters defining an area around the ground truth where soil use does not change. The TSTK downloader downloads the polygon and an extra margin in all directions. The resolution of images may vary depending on the size of the polygon. Then images of the visible spectrum were created using red, green, and blue bands. The median resolution for the RGB images is  $162 \times 217$  px. For the first tSNE presented in this paper, all of the images were exported with a white background and then cropped to 100 x 100 px resolution. For later tSNEs, images were exported with smaller resolution were upscaled to match this resolution.



Fig. 2: RGB image with no background and image cropped to 100 x 100 px.

## 2.2 Feature Extraction and t-SNE

An overview of the proposed methodology is given in Fig. 3. To analyze the high amount of data, t-Distributed Stochastic Neighbor Embedding (tSNE) images were created. tSNE is a commonly used clustering method for visualizing high-dimensional spaces on a plane or tridimensional space. The VGG16 keras classification model [12], a convolutional neural network (CNN) pre-trained on ImageNet, was utilized for high-level feature extraction (see Fig. 4). Although the ImageNet dataset consists of images and labels of specific objects, animals, and people, the current dataset contains a different type of image often with no clear subject and smaller resolution. Whatsoever, the last layer, used for classification on the 1000 categories of ImageNet is ignored, and only the weights for the 4096 connections on the last fully connected layer, whose weights correspond to high-level features. Then a Principal Component Analysis (PCA) extraction takes place, restricting the number of features to only the most relevant.



Fig. 3: Workflow diagram of the proposed methodology.

All the images are then plotted with t-SNE to a two-dimensional plane, where neighboring relationships indicate the similarity of features. This is done through an iterative random process that ultimately optimizes the probability of distribution through a stochastic process. In the last step, the images are replotted into a grid using RasterFairy, a tool created by Mario Klingemann [13]. This tool rasters the images into a regular structure while trying to preserve the neighboring relationship.



Fig. 4: VGG16 architecture.

# **3** Results and Discussion

This section presents and discusses the obtained results. The first result is the generated tSNE for the entire dataset, it is shown on Fig. 5. This image is by its visual properties a result on itself. It is a tool to better visualize our own data, like a map that reveals the potential knowledge for the machine. The image is then described in visual terms on detail to point out the generated clusters and what it can say about our data.

In a third part, two specific crops, common wheat, and maize, are explored in more detail.

#### 3.1 Entire Dataset t-SNE

The whole dataset consists of over 50.000 captures of soil for over 15.809 different points or soil ground truth. Meaning there might be more than one image for any point. The images were created using only the red, green, and blue channels of the patches. Then they were automatically cropped to  $100 \times 100$  px due to the capture size.

The t-SNE shown in Fig. was made with a random sample of 10.000 images taken from the general set of images. Primary Component Analysis (PCA) was restricted to 300 features.



Fig. 5: t-SNE representation for the entire dataset.

Some obvious conclusions can be drawn from looking at the t-SNE. The first is the persistence of clouds on the dataset; Even though clouds were already maxed out, they still took about 20% of the dataset. In the blue-shaded section of the t-SNE at the bottom, it could be observed that heavy moisture due to atmospheric conditions blurs the view of the soil. These images add much random noise to the dataset as no specific features can be extracted from them.

On the lower right side, there are a set of images that have white lines on their margins. These images had a smaller resolution than 100px wide because the polygon provided by the LUCAS Copernicus was smaller than average. Thus, when applied the center crop, an extra white margin was left. In future experiments, images with less than 100px wide were upscaled.

It is also observed some little clusters of black images or partially black images on the bottom of the tSNE, towards the center and center-right. The absence of visual information is due to the satellite camera not covering the whole of the requested area during its orbit.

There are also some reddish images clustered on the top left. Whatsoever, some other reddish images were grouped alongside clouds on the lower center-right. Further analysis indicates that some of these images are also brighter than some of the clouds. It is possible that VGG16 assigned similar features to dim clouds and bright reddish soil, which ended up grouping them after the stochastic clustering.

The rest of the dataset seems to be clear and distinguishable. A wide variety of greens and browns is observed. But also, there is a high presence of distinct forms due to cultivation area; the presence of roads, houses, and other structures; and topographical features.

#### 3.2 Common Wheat and Maize Case Study

The TSTK was used to train individual models for each crop present on the LUCAS Survey, being that each crop has different visual properties related to soil nutrients. From the previous analysis, it was established that the model trained on Common Wheat was the most accurate one to predict soil nutrients, while the one trained on Maize was the least accurate [11]. The most significant difference for each crop is the size of the dataset. As seen in Fig. the number of samples for common wheat is 2.5 times larger than that for maize (4.133 and 1.796, respectively [7]).



Fig. 6: t-SNE representation for individual crops: common wheat (left), maize (right).

Furthermore, the t-SNE differentiates two different types of images for each crop, possibly for when the plant is fully grown, and for when it is harvested or recently planted as observed on the detail shown in Fig. . This distinction is much clearer on wheat than maize. It could be hypothesized for a study that a machine trained on these images might be able to understand these differences very easily and make separate calculus for when there is leaf presence and when there are no leaves.



Fig. 7: Details of different clusters generated for each crop showing the different stages of plant growth.

# 4 Conclusions and Future Works

To understand what is happening inside a satellite images dataset, it is essential to look at the visual information and thus to find ways to organize it to make it easier to comprehend by the human eye. Visual t-SNE help to cluster information altogether to better explore relationships within images but also to understand the possible challenges driven out of the dataset creation and further growth. It shows how task related to computer vision grow in the interdisciplinary efforts across the computational sciences, the multimedia practices and the agronomical field.

Satellite images, such as those provided by Sentinel-2, are a cost-efficient way to obtain visual information from the soil. Nonetheless, information regarding soil composition is still limited, and more robust datasets are needed to make a big picture of soil nutrients all across the European Union and the whole of the world.

# 4.1 Future Work

The current dataset is still subject to different types of visual analysis. On one hand it is needed to study the visual spectrum for each of the crops that are not cited on this article to understand the how much visual information is actually present on the dataset.

On the other hand, it is needed to be understood which of the bands of the Sentinel-2 satellite are the most helpful for the study of each crop. As suggested by Wand and Wei [14, 15], certain bands might be more valuable than others in estimating the presence of chlorophyll in leaves which can be an indication of specific nutrients on the soil, specifically nitrogen for the study. The bands depend on the specific crop but some common bands might also be determinant. Lu et al. experiments propose that for the estimation of potassium from plant canopy, certain near-infrared and shortwave infrared are more adequate for such task [16]. It would be interesting to see t-SNE experiments with the different bands provided by the satellite to actually understand how it clusters differently from the visual spectrum cluster.

According to Mandrake study to detect sulfur on the soil [17], some sulfur might be easy to differentiate through aerial imagery. Others might have similar visual properties to non-sulfur components, which might lead to false positive results. For such reason, it is necessary to study which visual characteristic might lead further algorithms to provide false positive results.

Upson visualizing the results of the t-SNE for the present study, it has been suggested several times that this kind of image clustering could be applied through an interactive interface for both dataset analysis and dataset cleaning. This kind of tool would be extremely useful for cleaning larger datasets avoiding going through thousands of files.

However, what could have the most significant impact on the success of the TSTK would be an improved dataset of soil ground truth alongside crop type and GPS coordinates. Although the effort of the LUCAS Survey has been extremely useful, the dataset is not big enough for machine-learning algorithms result to be fully trusted. Further improvement and refinement of our dataset is going to be performed, but the help of agricultural institutions might play a big role in the further development on this cost-efficient technology.

## References

- Karthikeyan N, Shashikkumar M, Ramanamurthy J (2010) A study on vegetation vigour as affected by soil properties using remote sensing approach. https://doi.org/10.1109/RSTSCC.2010.5712811
- Walshe D, McInerney D, De Kerchove RV, Goyens C, Balaji P, Byrne KA (2020) Detecting nutrient deficiency in spruce forests using multispectral satellite imagery. Int J Appl Earth Obs Geoinformation 86:101975. https://doi.org/10.1016/j.jag.2019.101975
- 3. Liu JG, Mason PJ (2016) Image Processing and GIS for Remote Sensing: Techniques and Applications, 2nd edition. Wiley-Blackwell
- 4. Cai TT, Ma R (2022) Theoretical Foundations of t-SNE for Visualizing High-Dimensional Clustered Data
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. Journal of Machine Learning Research, 9(Nov):2579–2605, 2008.
- Deng J, Dong W, Socher R, Li L-J, Kai Li, Li Fei-Fei (2009) ImageNet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, Miami, FL, pp 248–255
- d'Andrimont R, Verhegghen A, Meroni M, Lemoine G, Strobl P, Eiselt B, Yordanov M, Martinez-Sanchez L, van der Velde M (2021) LUCAS Copernicus 2018: Earth-observationrelevant in situ data on land cover and use throughout the European Union. Earth Syst Sci Data 13:1119–1133. https://doi.org/10.5194/essd-13-1119-2021
- Tóth G, Jones A, Montanarella L (2013) The LUCAS topsoil database and derived information on the regional variability of cropland topsoil properties in the European Union. Environ Monit Assess 185:7409–7425. https://doi.org/10.1007/s10661-013-3109-3
- European Commission. Joint Research Centre. (2020) Assessment of changes in topsoil properties in LUCAS samples between 2009/2012 and 2015 surveys. Publications Office, LU
- Liebermann, Zach (2016). Land Lines. https://zachlieberman.medium.com/land-linese1f88c745847
- 11. Pereira MAS (2022) TerraSenseTK: a toolkit for remote soil nutrient estimation. Master-Thesis
- 12. Simonyan K, Zisserman A (2015) Very Deep Convolutional Networks for Large-Scale Image Recognition
- 13. Klingemann M (2022) About RasterFairy-Py3. https://github.com/Quasimondo/RasterFairy
- Wang W, Yao X, Tian Y, Liu X, Ni J, Cao W, Zhu Y (2012) Common Spectral Bands and Optimum Vegetation Indices for Monitoring Leaf Nitrogen Accumulation in Rice and Wheat. J Integr Agric 11:2001–2012. https://doi.org/10.1016/S2095-3119(12)60457-2
- 15. Wang L, Wei Y (2016) Revised normalized difference nitrogen index (NDNI) for estimating canopy nitrogen concentration in wetlands. Optik 127:7676–7688. https://doi.org/10.1016/j.ijleo.2016.05.115
- Lu J, Eitel JUH, Jennewein JS, Zhu J, Zheng H, Yao X, Cheng T, Zhu Y, Cao W, Tian Y (2021) Combining Remote Sensing and Meteorological Data for Improved Rice Plant Potassium Content Estimation. Remote Sens 13:3502. https://doi.org/10.3390/rs13173502
- Mandrake L, Wagstaff KL, Gleeson D, Rebbapragada U, Tran D, Castaño R, Chien S, Pappalardo RT Hyperspectral Sulfur Detection Using An Svm With Extreme Minority Positive Examples Onboard EO-. 12