

# Machine Learning Segmentation Case Study: Land Usage and Land Coverage

Leo Richard Hermann Giesen\*  
University of Munster  
Munster, North Rhine-Westphalia, Germany  
leo.giesen10@uni-muenster.de

Johannes Kauffmann  
University of Munster  
Munster, North Rhine-Westphalia, Germany  
johannes.kauffmann@uni-muenster.de



## ABSTRACT

In this case study, a supervised deep learning algorithm is presented with the objective of semantically segmenting a landscape into ten predefined land usage and coverage classes. This is achieved by a sliding window approach with a three-dimensional convolutional neural network resulting in an accuracy of 73%. At the hand of seven steps of machine learning projects, the functionality and characteristics of the algorithm are outlined and how it may be improved upon.

## CCS CONCEPTS

• Computing methodologies → Neural networks.

## KEYWORDS

machine learning, deep learning, semantic segmentation, convolutional neural networks

## ACM Reference Format:

Leo Richard Hermann Giesen and Johannes Kauffmann. 2021. Machine Learning Segmentation Case Study: Land Usage and Land Coverage. In *Proceedings of Deep Learning with Python (Specialization Module'21)*. Deep Learning with Python, Munster, Germany, 5 pages.

\*Both authors contributed equally to this research.

**Unpublished working draft. Not for distribution.**

Specialization Module'21, August 2021, Munster, Germany  
© 2021 Association for Computing Machinery.

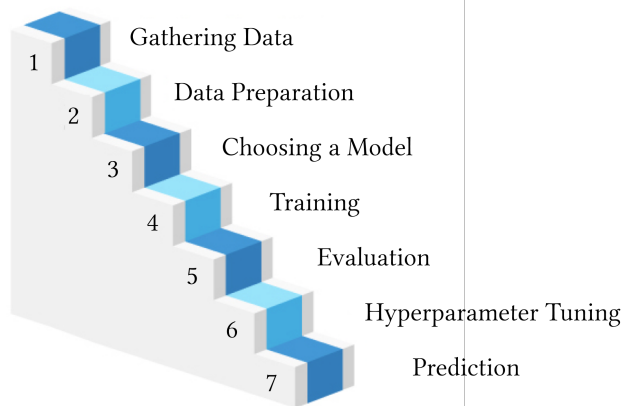
## 1 CONTEXTUALIZATION

The Information Systems specialization module 'Deep Learning with Python' held by Prof. Fabian Gieseke and Moritz Seiler, is supplemented by a machine learning case study. Its project work offered a realistic glimpse into machine learning tasks through the practical application of various theoretical concepts presented in the specialization module. However, advanced concepts such as U-Net were applied, because of the limited scope and time of the project.

The machine learning algorithm utilizes a sliding window approach in a three-dimensional convolutional neural network. It is trained on the training set, tested on the public test set, and finally evaluated on the hidden test set. The latter two include larger images for a segmentation of a landscape. The task falls under supervised learning because the algorithm draws on labeled data. More specifically, it belongs to the category of semantic segmentation classes are assigned to pixels of an image.

## 2 METHODOLOGY

The Jupyter Notebook of this project is written in Python and executed in Google Colab, which accessed the training and testing data via a mounted version of Google Drive. Frameworks such as TensorFlow and Keras were used in the machine learning implementation, i.e., they facilitate one-hot encoding, data augmentation for the sequential model and its layers and callbacks. Further libraries include for instance, NumPy, Matplotlib and scikit-learn for supporting high-level mathematical operations, plotting, and defining optimal class weights.



**Figure 2: Adjusted from Seven Steps of Machine Learning.** Web Illustration by Vaishali Advani, via Great Learning [3]. (<https://bit.ly/3C0fBqp>).

A machine learning project can be divided into seven steps as visualized in Figure 2, which are explained first and then examined in the context of the project. First, data needs to be gathered so that the algorithm can learn from it. Second, the data must be prepared to eliminate any unintended biases, for example the data is shuffled as the data instance order may affect the model's decisions. Third, a various models are tested and the best one is chosen. The validation set is used for this purpose. Fourth, the chosen model is trained on the training data to be able to perform an accurate prediction. Fifth, the model is evaluated regarding its generalization, which expresses its proficiency in application to new data. Sixth, the hyperparameters, which control and steer the learning process and define the amount of regularization, are tuned to further increase the model's accuracy. Seventh, the model is fully developed and applied in practice by performing a prediction [18].

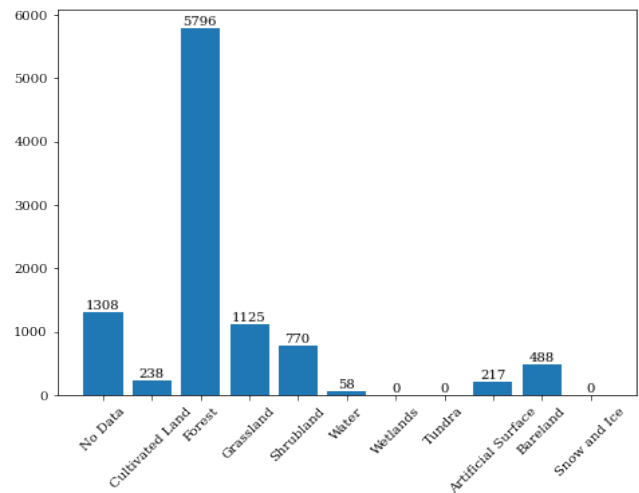
### 3 APPLICATION AND RESULTS

#### 3.1 Gathering Data and Data Exploration

These seven steps guide this project and their application in the case study is examined in-depth. The first step of gathering data has already been performed before the project. The data preparation requires and exploration of the data, hence the shape of the training dataset is inspected. The 'bands' object with a shape of (10000, 12, 33, 33, 6) comprises 10.000 patches, where the 33x33 pixel images are portrayed in red, green, blue and three near-infrared spectroscopic channels resulting in a total of six color channels. Each satellite image is covered over the course of twelve months, which may increase the overall accuracy. The central pixel of each of these images is labeled, which is stored in the 'lulc' data with a shape of (10000,).

A profound understanding of the gathered data was achieved by visualizing the landscapes and taking the various color channels and the eleven different class labels into account. The class no. 0 called 'no data' represents a label, which could not be assigned to a specific kind of biome. This could be the case when a cloud

covers the landscape. It turns out that class two 'cultivated land' is largely overrepresented with almost 6.000 instances, and classes six 'wetlands', seven 'tundra' and ten 'snow and ice' do not occur at all. This means that there is no training data available for them and consequently they cannot be predicted. The other classes occurred between 58 and 1308 times, which is observed in Figure 3.



**Figure 3: Class Label Frequency of Occurrence in Data Set.**

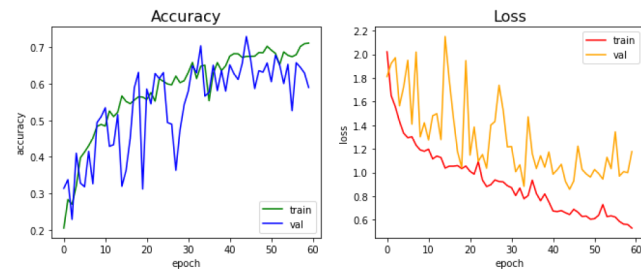
#### 3.2 Data Preparation

The imbalance of class occurrence is fatal because the accuracy is not representative of the actual performance. I.e., the algorithm would guess the overrepresented class two 'cultivated land' most of the time and be correct about  $(5796/10.000 \Rightarrow 57,96\%)$  of the time. However, this would not be a realistic prediction as other classes are neglected, and the model is not penalized accordingly. To punish the model correctly, class weights are automatically calculated and initialized in a dictionary with a predefined function. The not representative classes can be completely ignored in the prediction or hold a class weight of zero.

Before working with the data set, it is imperative to split the data set into a training (64%), validation (16%) and test set (20%), because if the same instances of the training set are also used in testing, the model predict correctly without demonstrating its learning. After that, the data is converted to a TensorFlow dataset object to make the runtime TPU (Tensor Processing Unit) of Google Colab work more efficiently with the data.

The datasets are subsequently prepared for the model fitting: First, the label datasets are one-hot encoded in order to use the categorical cross entropy as a loss function. After that, the TensorFlow prepare function was used to shuffle, augment, batch and prefetch all data [23]. As mentioned above, shuffling the data is necessary for the model to avoid the recognition of unwanted biases. The data augmentation TensorFlow preprocessing layers randomly flip the image horizontally and vertically and rotate it to provide a higher quantity of data for the model to learn from.

### 3.3 Choosing a Model, Training and Evaluation



**Figure 4: Accuracy and Loss of Machine Learning Algorithm.**

Various architectures were considered and inspired the final model architecture, which was designed from scratch to provide high flexibility. A three-dimensional convolutional neural network is required because of the additional dimension of months in the data set. As displayed in Figure 4, the lower layers consist of three main stacks of two 3D convolutional layers followed by a 3D max pooling and batch normalization layer. The higher layers comprise a flattening and three dense layers in combination with dropouts. The model is compiled with the categorical cross entropy loss function and the Adam optimizer. Three callbacks are defined for model training and fitting. The first one is a learning rate scheduler managing and reducing the learning rate by multiplying it with  $e^{-0.1}$  after the 10<sup>th</sup> epoch [24]. The second is an early stopping callback, which prevents overfitting by stopping the training process when the validation accuracy is not improving after a patience of 15 epochs. The third callback is a checkpoint, which automatically saves the best model in the h5-format. This is a hierarchical data format, which stores and organizes large amounts of data in a compact form [25]. The model is trained for a maximum of 100 epochs with a batch size of 32. The training stopped at the 68<sup>th</sup> epoch due to early stopping. After training, the fitted model is evaluated using the test set. The accuracy results in 72,44% and the loss is 54,75%.

### 3.4 Hyperparameter Tuning and Prediction



**Figure 5: Comparison of Prediction and Solution in the Public Test Set.**

The step of hyperparameter tuning is skipped as the focus of the case study lied on the general approach of the given semantic segmentation task and not on achieving the highest accuracy

possible. However, it is advised to perform hyperparameter tuning if the machine learning algorithm is planned to be deployed. The final prediction is performed with a sliding window approach on the public and hidden test set, whose images measure 500x500 and 1500x1500 pixels in size. This approach predicts the central pixel of a 33x33 pixel block around it, exactly as in the training set. The sliding aspect is added, so that this block for the prediction slides across the image to predict every pixel of every row and column. The problem arises that the edge pixels are not surrounded by pixels on every side, which is required for the prediction box of 33x33 pixels. This is resolved by zero padding, which places a  $((33 - 1)/2 =)$  16-pixel margin of zero values ('no data') around the whole image. As a result, the sliding window approach produces an output prediction with the same size as the input. The output 2D prediction array is visualized using an RGB-color-map, where each label has a corresponding color.

## 4 DISCUSSION

### 4.1 Data Preparation

There are various options and alternative approaches, which may be considered in each of the previously mentioned seven steps of the machine learning process. For instance, three classes cannot be predicted because of their non-existence of data. This could be solved by gathering data from further satellite images. The class instance imbalance was solved with initializing corresponding class weights, but it could have also been achieved by duplicating data instances to account for the difference in occurrence. However, the former is simpler to implement and faster to train because of fewer training instances.

The inspection of data also revealed that clouds pose a problem as they complicate the prediction. The color channel values are scaled between zero and one to foster the model's learning from the data. Though, there are satellite images with color channel values greater than one, which might the case because of the cloud's high reflection of the sun light. In these cases, normal scaling would scale down other color channel values so much that their information is lost, which would prevent the model from learning. These data instances could be manually removed using data cleaning, which could potentially increase the model's accuracy. However, this is not tested because this step is very time consuming, and it would decrease the quantity of data.

The data augmentation is performed with four TensorFlow pre-processing layers and are part of the model, which results in a longer training duration, but the model also benefits from GPU acceleration because the data augmentation layers run synchronously with the rest of the model's layers.

### 4.2 Choosing a Model

Alternatively, a pretrained model could have been utilized, which would improve the model's performance. This is because of transfer learning, which enables the model to start off with a much higher accuracy because the improved model's architecture and structure including the neuron's weights, which are initialized closer to their optimum. Though, the unusual input shape of the segmentation task hinders transfer learning as no comparable 3D convolutional neural networks could be found online during our research. Hence,



one could evade this situation by using very accurate and high performant model architectures such as ResNet or U-Net with no pretrained weights in a future version of this machine learning project.

Initially, the model only comprised few layers and relatively few parameters, hence further model testing revealed that more convolutional 3D layers make the model more powerful and increase accuracy, because the previous model version with just under 700.000 parameters underfitted and did not have the capacity to capture the landscape's complexity.

As mentioned above, the categorical cross entropy is used as a loss function. Instead, the model could have also used sparse categorical cross entropy as a loss function, which produces a category index of the most likely matching *class*. In comparison, categorical cross entropy displays each *class probability* in a one-hot encoding vector. This leads to two major advantages of sparse categorical cross entropy as it saves memory space and may reduce the training time. The probability vector is beneficial in situations, where the classes are not mutually exclusive, and one instance may have multiple classes. However, this is not the case in this project, thus sparse categorical cross entropy is advised to be used.

### 4.3 Evaluation

The evaluation discloses a significant fluctuation of the accuracy and loss. This may originate from consecutive cloudy images, which may negatively impact the model's weights and lower accuracy.

In alternative to accuracy, other types of performance measurement could have been selected, which are not affected by the imbalance of class instances. For instance, a confusion matrix or a receiver operating characteristic curve (ROC) are better in evaluating the quality of the model in this situation without applying class weights.

### 4.4 Prediction

The naive sliding window approach takes significantly longer than two hours to predict the public test image because each pixel is predicted separately. The prediction process could be speeded up, for instance a whole row could be predicted by passing a row of pixels and returning an array of predictions. In that case, NumPy arrays slow down the prediction because the whole array needs to be loaded, which is not the case for a python list. Therefore, the prediction of a row could be appended to the final prediction.

Apart from the naive sliding window approach, a convolutional sliding window approach could be utilized to lower the prediction duration. This is the case, because it is a one-shot approach, which does not iterate over the rows and columns in two for-loops. Moreover, one could make use of innovations such as the framework Mask R-CNN, which are significantly more performant than the current approach in this project.

## 5 CONCLUSION

In conclusion, the objective of semantic segmentation of landscapes from satellite images is achieved with a satisfactory accuracy. However, there are still options to be tested out to improve the performance. Especially, alternatives to the model's architecture and possible frameworks should be shed a light on.

## 6 CITATION EXAMPLES

online document / world wide web resource [1, 19, 26]  
journal articles: - Paginated [2] - enumerated [8]  
reference to entire issue [7]  
monograph (whole book) [17] whole book in a series (see 2a in spec. document) [13] chapter in a divisible book ... - [22] - in a series [9]  
multi-volume work as book [16]  
couple of articles in a proceedings (of a conference, symposium, workshop for example) (paginated proceedings article) [4, 11],  
proceedings article with all possible elements [21] enumerated  
proceedings article [10] informally published work [12] doctoral  
dissertation [6] master's thesis: [5]  
work accepted for publication [20]  
multi-volume works as books [15] and [14]

## ACKNOWLEDGMENTS

Regular exchange of ideas and possible solutions with Marten Jostmann, Gerrit Kruse and Tim Heibel.

## REFERENCES

- [1] Rafal Ablamowicz and Bertfried Fauser. 2007. *CLIFFORD: a Maple 11 Package for Clifford Algebra Computations*, version 11. Retrieved February 28, 2008 from <http://math.tntech.edu/rafal/cliff11/index.html>
- [2] Patricia S. Abril and Robert Plant. 2007. The patent holder's dilemma: Buy, sell, or troll? *Commun. ACM* 50, 1 (Jan. 2007), 36–44. <https://doi.org/10.1145/1188913>.
- [3] Vaishali Advani. 2021. What is Machine Learning? How Machine Learning Works and Future of It? Retrieved 2021-07-27 from <https://www.mygreatlearning.com/blog/what-is-machine-learning/>
- [4] Sten Andler. 1979. Predicate Path expressions. In *Proceedings of the 6th. ACM SIGACT-SIGPLAN symposium on Principles of Programming Languages (POPL '79)*. ACM Press, New York, NY, 226–236. <https://doi.org/10.1145/567752.567774>
- [5] David A. Anisi. 2003. *Optimal Motion Control of a Ground Vehicle*. Master's thesis. Royal Institute of Technology (KTH), Stockholm, Sweden.
- [6] Kenneth L. Clarkson. 1985. *Algorithms for Closest-Point Problems (Computational Geometry)*. Ph.D. Dissertation. Stanford University, Palo Alto, CA. UMI Order Number: AAT 8506171.
- [7] Jacques Cohen (Ed.). 1996. Special issue: Digital Libraries. *Commun. ACM* 39, 11 (Nov. 1996).
- [8] Sarah Cohen, Werner Nutt, and Yehoshua Sagie. 2007. Deciding equivalences among conjunctive aggregate queries. *J. ACM* 54, 2, Article 5 (April 2007), 50 pages. <https://doi.org/10.1145/1219092.1219093>
- [9] Bruce P. Douglass, David Harel, and Mark B. Trakhtenbrot. 1998. Statecards in use: structured analysis and object-orientation. In *Lectures on Embedded Systems*, Grzegorz Rozenberg and Frits W. Vaandrager (Eds.). Lecture Notes in Computer Science, Vol. 1494. Springer-Verlag, London, 368–394. [https://doi.org/10.1007/3-540-65193-4\\_29](https://doi.org/10.1007/3-540-65193-4_29)
- [10] Matthew Van Gundy, Davide Balzarotti, and Giovanni Vigna. 2007. Catch me, if you can: Evading network signatures with web-based polymorphic worms. In *Proceedings of the first USENIX workshop on Offensive Technologies (WOOT '07)*. USENIX Association, Berkley, CA, Article 7, 9 pages.
- [11] Torben Hagerup, Kurt Mehlhorn, and J. Ian Munro. 1993. Maintaining Discrete Probability Distributions Optimally. In *Proceedings of the 20th International Colloquium on Automata, Languages and Programming (Lecture Notes in Computer Science, Vol. 700)*. Springer-Verlag, Berlin, 253–264.
- [12] David Harel. 1978. *LOGICS of Programs: AXIOMATICS and DESCRIPTIVE POWER*. MIT Research Lab Technical Report TR-200. Massachusetts Institute of Technology, Cambridge, MA.
- [13] David Harel. 1979. *First-Order Dynamic Logic*. Lecture Notes in Computer Science, Vol. 68. Springer-Verlag, New York, NY. <https://doi.org/10.1007/3-540-09237-4>
- [14] Lars Hörmander. 1985. *The analysis of linear partial differential operators. III. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*, Vol. 275. Springer-Verlag, Berlin, Germany. viii+525 pages. Pseudodifferential operators.
- [15] Lars Hörmander. 1985. *The analysis of linear partial differential operators. IV. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*, Vol. 275. Springer-Verlag, Berlin, Germany. vii+352 pages. Fourier integral operators.

- [16] Donald E. Knuth. 1997. *The Art of Computer Programming, Vol. 1: Fundamental Algorithms (3rd. ed.)*. Addison Wesley Longman Publishing Co., Inc.
- [17] David Kosiur. 2001. *Understanding Policy-Based Networking* (2nd. ed.). Wiley, New York, NY.
- [18] Juan Cruz Martinez. 2020. Seven Steps of Machine Learning. Retrieved 2021-07-28 from <https://livecodestream.dev/post/7-steps-of-machine-learning/>
- [19] Poker-Edge.Com. 2006. Stats and Analysis. Retrieved June 7, 2006 from <http://www.poker-edge.com/stats.php>
- [20] Bernard Rous. 2008. The Enabling of Digital Libraries. *Digital Libraries* 12, 3, Article 5 (July 2008). To appear.
- [21] Stan W. Smith. 2010. An experiment in bibliographic mark-up: Parsing metadata for XML export. In *Proceedings of the 3rd. annual workshop on Librarians and Computers (LAC '10, Vol. 3)*, Reginald N. Smythe and Alexander Noble (Eds.). Paparazzi Press, Milan Italy, 422–431. <https://doi.org/99.9999/woot07-S422>
- [22] Asad Z. Spector. 1990. Achieving application requirements. In *Distributed Systems* (2nd. ed.), Sape Mullender (Ed.). ACM Press, New York, NY, 19–33. <https://doi.org/10.1145/90417.90738>
- [23] TensorFlow. 2021. Data Augmentation. Retrieved 2021-07-31 from [https://www.tensorflow.org/tutorials/images/data\\_augmentation](https://www.tensorflow.org/tutorials/images/data_augmentation)
- [24] TensorFlow. 2021. Early Stopping. Retrieved 2021-07-31 from [https://www.tensorflow.org/api\\_docs/python/tf/keras/callbacks/EarlyStopping](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping)
- [25] The HDF Group. 2021. The HDF5 Library and File Format. Retrieved 2021-07-31 from <https://www.hdfgroup.org/solutions/hdf5>
- [26] Harry Thornburg. 2001. *Introduction to Bayesian Statistics*. Retrieved March 2, 2005 from <http://ccrma.stanford.edu/~jos/bayes/bayes.html>