

Machine Learning for Earth and Environmental Sciences (5 ECTS)

Master in Environmental Sciences, FGSE, University of Lausanne

Syllabus (last updated on August 21, 2024)

Main Instructor: Tom Beucler (Assistant Professor at the FGSE, Lab Website, tom.beucler@unil.ch)

Teaching Assistants: Milton Gomez (milton.gomez@unil.ch), Kejdi Lleshi (PhD Student at IDYST, (kejdi.lleshi@unil.ch)), Haokun Liu (PhD Student at IGD, haokun.liu@unil.ch), Ayoub Fatihi (PhD Student at ISTE, ayoub.fatihi@unil.ch).

1 Summary

Our ever-improving ability to observe and model the environment produces Petabytes of data every day, which overwhelms traditional data analysis methods. Machine learning (ML) algorithms, broadly defined as novel and computational algorithms that allow computers to automatically perform a task from data without requiring explicit programming, have recently emerged as efficient tools to extract knowledge from large geoscientific datasets. Once trained, these algorithms are inexpensive to use and ideal shortcuts when time or resources prevent running a full-complexity model. In addition, the ability of ML algorithms to summarize large amounts of data makes them promising tools for environmental scientific discovery.

In this 12-week hands-on course, we will introduce common ML algorithms in the context of their application in Earth and environmental sciences. By the end of this course, you should be able to:

1. Describe common ML algorithms (listed in Sec 2) and summarize their advantages and limitations, especially in the context of environmental science,
2. Implement them in Python (using the Numpy/Scikit-Learn/Keras/Tensorflow/PyTorch libraries in Colab notebooks),
3. Know from experience which algorithms are most appropriate for environmental applications you are passionate about (e.g., your thesis research).

To achieve these three objectives, the course will combine:

- **Lectures** (≈ 2 hours/week, 15% of grade): Typical structure = 15-min answering your questions about readings, 15-min live quiz based on readings, 45-min group activities based on environmental application readings (including a 15-min break), 15-min summary of the following week's main algorithm & environmental application to help you prepare for the readings, 15-min reviewing quiz. Note that during lecture, we will favor a diversity of ML applications over mathematical foundations; if you are interested in the latter, we encourage you to take the appropriate ML courses at EPFL.
- **Readings** (≈ 3 hours/week, 15% of grade): Reading 1 is a textbook chapter covering next week's algorithms, while Reading 2 is (usually) a recently published article that successfully applied ML algorithm to tackle key environmental science issues. Both readings are posted on Moodle, and you'll get the full 15% of the grade as long as you provide thoughtful feedback on the corresponding chapter(s) of the course's e-book no more than 24 hours before the lecture.
- **Computer labs** (≈ 3 hours/week, 20% of grade): ≈ 1.5 hr applying ML covered in Reading 01 on standard ML datasets + ≈ 1.5 hr applying ML on environmental dataset covered in Reading 02. To get full credits (20%), simply push the completed Colab notebook to your fork of the course's GitHub repo no more than 24 hours before the following lab.
- **Final project** (≈ 2 hours/week, 50% of grade): The final project's goal is to answer a well-defined scientific question by applying one of the ML algorithms introduced in class on an environmental dataset of your choice (e.g., related to your thesis). You can find more specific instructions in Sec 3 and on Moodle. You may collaborate with peers on the same dataset and present your projects together in class (20% of grade), but you must choose distinct scientific questions and write separate 4-page final reports (30% of grade) by uploading this template into Overleaf (LaTeX tutorial at this link). Please submit the final report in PDF format by the deadline indicated on Moodle.
- There will be **no final exam** initiations or homework other than the readings and the final project.

For the ML components of the course, we will mainly use Géron's 2019 textbook "Hands-on Machine Learning with Scikit-Learn, Keras & Tensorflow (2nd ed)" (code/pdf on Moodle) and Chollet's 2021 textbook "Deep Learning with Python (2nd ed)" (code/pdf on Moodle), but we encourage you to use the wealth of online resources on machine learning (link to get started).

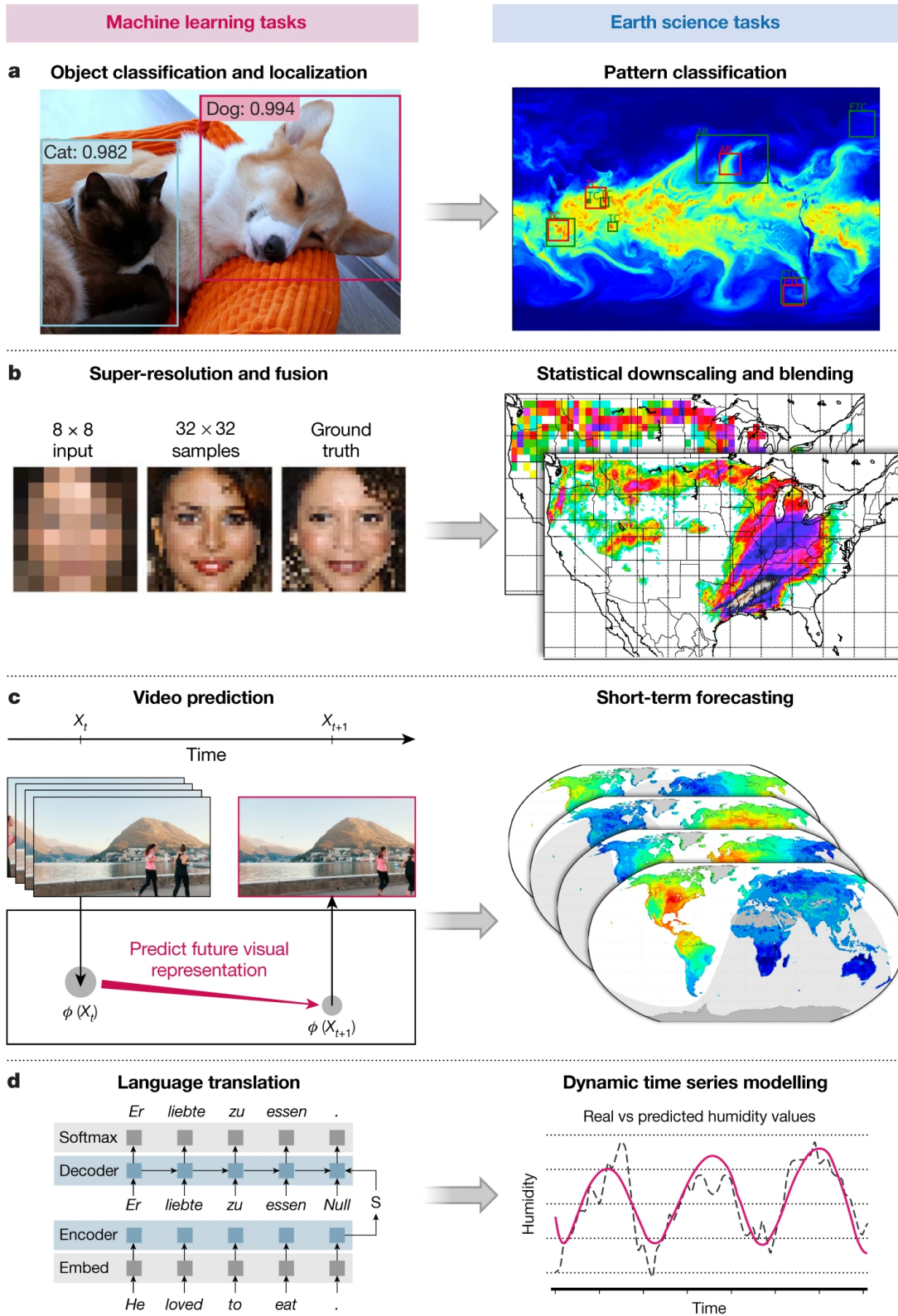


Figure 1: a) Object recognition in images links to classification of extreme weather patterns using a unified convolutional neural network on climate simulation data⁴¹. b) Super-resolution applications relate to statistical downscaling of climate model output⁷². c) Video prediction is similar to short-term forecasting of Earth system variables. d) Language translation links to modeling of dynamic time series.

Source: Figure 2 in <https://www.nature.com/articles/s41586-019-0912-1>.

2 Tentative Schedule

To stay up-to-date, consider adding the course’s e-calendar to your own calendar: Google Cal/iCal links, which contains information on dates, times, and location. We highly encourage you to complete the assigned readings and to provide feedback on the e-book (e.g., via Moodle or a pull request) no more than 24 hours before lecture to get the full 15% of the “Readings” grade.

Prerequisites: Introduction to Scientific Programming with Python

Recommended prerequisite: Completing Ch1 of e-book.

Recommended reading before the first course = Ch 1+2 of Géron

Week 1-2: Linear/Logistic Regression for Classification/Regression & Statistical Forecasting

Python/Google Colab Notebooks/Git basics, Training/Validation/Test split, Best practices for training and benchmarking.

Reading1 for W2 = Ch 3+4 of Géron; W2 Lab = Ch 2 of e-book.

Reading2 for W2 = Statistical Methods in the Atmospheric Sciences (Ch7: Statistical Forecasting 7.1-7.4+7.9) (book)

Week 2-3: Decision Trees/Random Forests/SVMs & Environmental Risk Analysis

Ensemble Learning, RVM, Gaussian Processes.

Reading1 for W3 = Ch 5+6+7 of Géron; W3 Lab = Ch 3 of e-book.

Reading2 for W3 = An ML-Based Approach for Wildfire Susceptibility Mapping. [...] (paper)

Week 3-4: Unsupervised Learning for Clustering/Dimensionality Reduction & Environmental Complexity

K-Means, DBSCAN, Hierarchical clustering, t-SNE, Gaussian Mixtures, Variational Inference.

Reading1 for W4 = Ch 8+9 of Géron; W4 Lab = Ch 4 of e-book.

Reading2 for W4 = Revealing the Impact of Global Heating on North Atlantic Circulation Using Transparent ML (article/code)

Week 4-5: Artificial Neural Networks & Surrogate Modeling

Reading1 for W5 = Ch 10+11 of Géron; W5 Lab = Ch 5 of e-book.

Reading2 for W5 = Could Machine Learning Break the Convection Parameterization Deadlock? (article/code)

Week 5-6: Generative Modeling: From Uncertainty Quantification to Stochastic Downscaling

Distributional regression, Auto-encoders, Generative adversarial networks, Diffusion models.

Reading1 for W6 = Ch 17 of Géron; W6 Lab = Ch 10 of e-book.

Reading2 for W6 = Adversarial super-resolution of climatological wind and solar data (article/code)

Week 6-7: Convolutional Neural Networks & Remote Sensing

Fully Convolutional Networks, ResNets, U-Nets.

Reading1 for W7 = Ch 14 of Géron; W7 Lab = Ch 6 of e-book.

Reading2 for W7 = Review on Convolutional Neural Networks in vegetation remote sensing (article/code)

Week 7: First In-Class Presentation

Week 7-8: Graph Neural Networks & Interconnected Systems

Graph statistics, Node embedding, Graph Convolutional Networks.

Reading1 for W8 = Ch 1+2 of Labonne, Review on Graph Neural Networks. (paper); W8 Lab = Ch 8 of e-book.

Reading2 for W8 = Regional Heatwave Prediction Using GNN and Weather Station Data (article/code)

Week 8-9: Explainable Artificial Intelligence & Understanding/Communicating Predictions

Permutation tests, Partial-dependence plots, Saliency maps, Feature visualization.

Reading1 for W9 = Extracts from “Interpretable ML” by Christoph Molnar (book); W9 Lab = Ch 9 of e-book.

Reading2 for W9 = Interpretable Deep Learning for Spatial Analysis of Severe Hailstorms (article/code)

Week 9-10: Hybrid Modeling & Knowledge-Guided Learning

Reading for W10 = Ice-flow model emulator based on physics-informed deep learning (article/code); W10 Lab = Ch 11 of e-book.

Week 10-12: Office Hours for Final Projects, In-Class Peer-Review and In-Class Presentations

1) Possible overview of student-chosen related topics not covered in class that may be relevant to final projects, such as:

Bayesian inference, Causal discovery, Data ethics, Recurrent Neural Networks, Reinforcement Learning, Symbolic Regression.

2) In-class peer review: Each class member submits the draft of their final project for review and reviews 3 drafts from peers.

3 Final Project Guidelines

The final project's goal is to answer a well-defined scientific question by applying one of the machine learning (ML) algorithms introduced in class on an environmental dataset of your choice (e.g., related to your Master's thesis or your PhD research). You may collaborate with peers on the same dataset and present your projects together in class (40% of final project's grade), but you must choose distinct scientific questions and write separate 4-page final reports (40% of final project's grade) by uploading this template into Overleaf (LaTeX tutorial at this link). To ensure a smooth peer-review process, you are encouraged to submit a draft of your written report (it does not have to be final) one day before the in-class peer review (20% of final project's grade). If you write constructive peer reviews, you can obtain as many bonus points as 30% of the final project's grade (10% per constructive peer review). Your final project will be made publicly-accessible on the course's website at the end of the semester.

3.1 Timeline

- Weeks 1-2: Pick a relatively large (if possible, more than 1000 samples) environmental dataset linked to a scientific question you are passionate about (see Section 3.4 for more guidance). **Don't hesitate to discuss possible projects with the instructor, the TA, and your classmates during labs or office hours**, especially if you are struggling to find an environmental dataset, define a scientific question, or choose an appropriate ML algorithm for the task at hand. We will schedule dedicated, one-on-one office hours to help you kickstart your final project.
- Weeks 3-7: Work on your final project at the end of labs and outside of the classroom. **Don't hesitate to discuss pitfalls and brainstorm solutions with the instructor, the TA, and your classmates during labs or office hours.**
- Week 7: First in-class presentation to get helpful feedback from the teaching team based on your *preliminary* results.
- Weeks 8-12: Work on your final project *during* labs and outside of the classroom. **Don't hesitate to discuss pitfalls and brainstorm solutions with the instructor, the TA, and your classmates during labs or office hours.**
- Week 10: Please submit a draft of your final project (even if still in-progress) one day before the in-class peer review. **This deadline is especially important** as the in-class peer-review process will not be possible without everyone's submission.
- Week 11: Final in-class presentation to share your results with your peers.
- Week 12: Please submit the final report in PDF format via Moodle (the deadline is on Moodle).

3.2 Deliverables

1. A written 4-page¹ final report addressing a well-defined scientific question by applying one of the machine learning (ML) algorithms introduced in class. Use this template on Overleaf (LaTeX tutorial at this link) and include:
 - An informative title, and an abstract with no more than 6 sentences (follow this link for tips),
 - At least 2 Figures to communicate your methodology and your results,
 - At least 2 Tables with (1) the range of your hyperparameter search and the hyperparameters you chose using your validation set; and (2) at least 2 carefully-chosen performance metrics evaluated over the training, validation, and test sets.
 - A link to the dataset you used,
 - A link to a well-documented Python notebook/script on GitHub for your project.
2. Two short² in-class presentations of your project, clearly communicating your scientific question, introducing your dataset, explaining your methodology and the reason you chose a particular ML algorithm, and summarizing your findings.

3.3 Evaluation & Peer-Review Process

The final report's evaluation and the in-class peer-review process will both use the same rubric at this link. The presentations' evaluation will assign equal importance to the scientific question, the dataset's presentation, the ML methodology's justification, the demonstrated ML knowledge, and the project's results.

3.4 Resources

We encourage you to use data from your Masters or PhD research. If you are still looking for the right dataset, consider reaching out to your Masters/PhD thesis' advisor. You may also browse Kaggle datasets, this list of benchmark datasets (maintained by Pangeo), the linked datasets from last year's final projects, the environmental datasets from the course's syllabus (refer to the course's GitHub repository), and this list of open geo datasets and final projects suggestions kindly provided by Dr. Yu.

¹excluding references

²the exact duration will be based on the total number of presentations and will include time for questions

4 Resources and Ethics

4.1 UNIL/FGSE Resources for Students

- Disability resources. If you need academic support, please email me (preferentially before the class starts) so that I can request/provide the appropriate services.
- English resources: Many of us are not native English speakers, and UNIL provides a wealth of resources to practice English, which may come in handy when writing up your final project.
- Financial support resources.
- Confidential mental health resources (first session is free) provided by the university's hospital.

4.2 Diversity and Inclusion in the Classroom

The University of Lausanne is committed to equal opportunity and stands firm against all forms of discrimination, including discrimination based on race, gender, religion, country of origin, ethnicity, socioeconomic status, sexual orientation, and disability. There are confidential resources if you feel harassed.

In the context of our classroom, this means:

- Choosing how you would like to be addressed by indicating your preferred name and pronouns in the initial course survey,
- Openly discussing and asking about concepts we struggle with to normalize difficulties in learning and applying course materials,
- Being kind and understanding towards each other: Especially in an interdisciplinary and international environment, concepts that seem obvious to you may be unknown to others or have different names depending on your sub-field,
- Emailing me or the equal opportunity office if you feel that students are not treated evenhandedly, or if the context/structure of the course is negatively impacting your learning experience and performance,
- All recognizing and working on our implicit biases by actively listening to each other.

4.3 Late Work Policy

Late work is eligible for partial credit of 50% until the official end of the semester.

4.4 Conversational Artificial Intelligence Policy

The increasing availability of high-quality, user-friendly conversational AI tools (e.g., ChatGPT, Claude, Le Chat, Bing Chat, and Google Bard) offer unprecedented opportunities to hone your writing, coding, and overall communication skills. To maximize their potential without undermining your critical thinking and subject knowledge (key reasons to attend UNIL), we offer the following guidelines: For **reading assignments**, we encourage you to read the original texts rather than rely on AI to summarize large sections, noting that AI can help clarify specific concepts. For **coding assignments (Colab)**, we recommend searching for functions within the appropriate Python library over using Gemini code assistance. For **writing assignments**, we recommend using AI as you would a third-party proofreader, e.g., ask “proofread this sentence for typos, clarity, and conciseness” rather than “write a sentence about the use of recurrent neural networks in hydrological sciences”. **Forbidden during quizzes.**

4.5 Academic Integrity

At UNIL, we all share strict rules on academic integrity, which can be found at this link (in French). In the context of this course, the following behavior can lead to an automatic failure of the class (grade of 0%):

1. Plagiarism. To avoid plagiarism, always cite your sources: at the bottom of your slides during the final presentation, including for photos/schematics, and using bibtex when writing your final report using Overleaf. Please do not take credit for someone else's work and do not have someone write in your name (this applies to all course activities).
2. Unauthorized collaboration. Even if you collaborate with some of your peers on the final project, you must answer distinct questions and write separate reports. Please transparently acknowledge any help you received from your peers (coding, research ideas, writing, proofreading, data, citations, etc.) in the acknowledgments section of your final report. During graded quizzes in class, please do not copy your peers' responses. Even if collaboration is highly encouraged, do not copy your peers' code during computer labs. Between classes, do not copy your peers' answers to the readings' guiding questions.
3. Data fabrication or falsification. Please do not fabricate the data reported in the analyses, figures, and tables of your final report. Being transparent about the shortcomings of a method or a dataset is always helpful to the community.