# Machine Learning for Time Series

## Introduction

Laurent Oudre
laurent.oudre@ens-paris-saclay.fr

Master MVA
2024-2025

# Contents

## 1. Organization of the course

# Organization of the course : lectures

▶ Six lectures on Thursday afternoons (except next week : course in the morning !) at ENS Paris Saclay

▶ Lectures will be in French but all material (slides, homeworks...) is in English.

▶ Attendance is mandatory

▶ Registration deadline for the mailing list : **October 15th, 2024**

https://forms.gle/KECsUsdYU7ZFtLNT9

▶ ML for time series is part of the Modelling track !

Teaching material: http://www.laurentoudre.fr/ast.html

# Organization of the course : tutorials

- For the tutorials there are two options :
  - Thursday morning : remote on Zoom
  - Thursday afternoon : onsite at ENS Paris Saclay
- Extra work for each tutorial : approximately 6 hours
- Attendance is mandatory
- Tutorial homeworks are **mandatory**

  **missing or late homeworks $\rightarrow$ fail the class**

  Teaching assistant: Charles Truong (ctruong@ens-paris-saclay.fr)

# Validation

**Validation: tutorials (25%) + mini-projects (25% report, 25% source code and 25% oral presentation)**

▶ Projects can be done in groups of two, but no more than that

▶ Students are allowed to propose additional project : ask in advance !

▶ The mini project consists in reading a research paper, implement it in Python and launch experiments on real time series

▶ Report (PDF file, $\approx$ 5 pages) + source code (Jupyter Notebook). **Deadlines : December 18th (23:59) or January 9th (23:59)**

▶ A 10 min oral presentation is scheduled on **December, 19th and 20th and January, 9th and 10th**, which will finalize the course project

▶ **Due to the large number of students, *auditeurs libres* will not be able to validate the course (no grading for tutorials, no mini-project).**

# Contents

# What is a time series?

▶ A time series is a series of data points indexed in time order
▶ In practice, array of real numbers of size $D \times N$ where $D$ is the number of dimensions and $N$ the number of samples

  ▶ Sample number $n$

  | $n$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
  |---|---|---|---|---|---|---|---|

  ▶ Time series values $x[n]$

  | $x[n]$ | 0.7 | 0.2 | 0.8 | 0.9 | 0.3 | 0.2 | 0.7 |
  |---|---|---|---|---|---|---|---|
  | | 0.4 | 0.1 | 0.6 | 0.2 | 0.5 | 0.6 | 0.3 |

  ▶ Timestamps $t[n]$

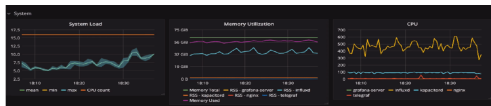  | $t[n]$ | 16:30:01 | 16:30:23 | 16:31:43 | 16:32:38 | 16:33:06 | 16:33:16 | 16:33:56 |
  |---|---|---|---|---|---|---|---|

# An un-unified field

- ▶ Different scientific communities have given different names to the same mathematical object.
    - ▶ **Time series**: mathematics, statistics, economics, finance...
    - ▶ **Signals**: signal processing, physics, engineering, simulation...
    - ▶ **Sequences**: computer sciences, bioinformatics, data mining...
- ▶ In this course, we will use indifferently one of these terms.
- ▶ Typical definition: real-valued (or at least ordered) sequential data

# Time series are everywhere



Meteorology,
Finance,
Healthcare,
Monitoring,
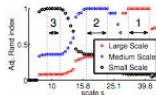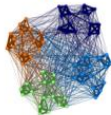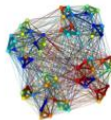Epidemiology,
Sensor networks...

# Univariate vs. multivariate

2D/3D trajectories,
Multivariate time series,
Multimodal data from
sensor networks,
Graph signals

**Sensor Data**

# Time series are complex

- ▶ Potentially massive data (e.g. sound : sampling frequency 44.1 kHz)
- ▶ Multivariate, multimodal, heterogeneous
- ▶ Noisy, missing data, trends, mixture of sources
- ▶ Often linked to an application context: data scientist is not trained to understand the data

# Annotations and ground truth

▶ Contrary to basic image processing tasks (e.g. classification of cats and dogs), annotating time series often require expertise
▶ Typical context:
  ▶ Noisy and dirty data
  ▶ A few annotated signals with blurry labels (confusing and hyper-specialized annotations that cannot be transformed into class labels)
  ▶ An expert with several years in the business, but unable to translate it into ML-compatible annotations

How to use ML in this context?

# What about time?

▶ What is the difference between regular data and time series ? Notion of sequence and chronology



▶ Each sample corresponds to the measurement of a phenomenon at a given time stamp.

▶ Time allows to study the evolution of the phenomenon and should be taken into account for processing the data

# What about time?



Same time series... but mixed up times

# World vs. Machine Learning



- ▶ Most ML algorithms do not care for time.
- ▶ How can we still use the time information to extract relevant features/patterns that can be used within a ML procedure ?

# Two visions: physics vs. statistics

▶ The notion of time have been used and modeled in physics since 18th century and before (eg. Fourier transform).
**First vision :** a time series $x[1:N]$ is the result of the digitization of a physical phenomenon $x(t)$. Physical properties of this phenomenon can be retrieved and analyzed through the study of $x[1:N]$ (and vice/versa).

▶ Randomness can also play a part to model a wider class of signals.
**Second vision :** a time series $x[1:N]$ is a realization of a stochastic process $X[1:N]$. Statistical properties of this phenomenon can be retrieved and analyzed through the study of $x[1:N]$ (and vice/versa).

In most cases, both approaches can be combined.

# Deep learning: the optimal solution?

Deep learning achieves state-of-the-art results for several tasks **BUT**...

- ▶ Good performances $\neq$ good understanding of the data (cf next slide)
- ▶ DL is a black box that may not bring satisfaction to users on the field since they cannot interpret the results
- ▶ Although some networks are able to handle time (e.g. LSTM), they still only manage at most a few hundred time samples
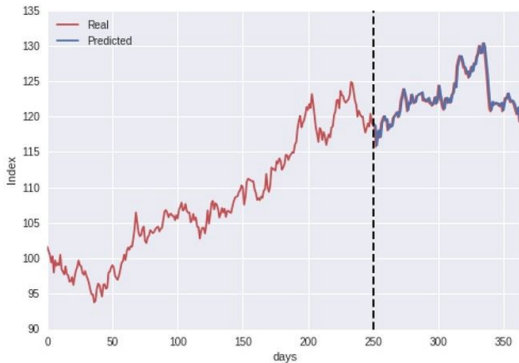- ▶ DL is bad in the context of scarce data and annotations

# Deep learning: the optimal solution?



*To validate our claim, we introduce a set of embarrass-ingly simple one-layer linear models named LTSF-Linear for comparison. Experimental results on nine real-life datasets show that LTSF-Linear surprisingly outperforms existing sophisticated Transformer-based LTSF models in all cases, and often by a large margin. Moreover, we con-*

Zeng, A., Chen, M., Zhang, L., & Xu, Q. (2022). Are Transformers Effective for Time Series Forecasting?. *arXiv preprint arXiv:2205.13504*.

**Forecasting**

Ma, Q., Zheng, J., Li, S., & Cottrell, G. W. (2019). Learning representations for time series clustering. Advances in neural information processing systems, 32.

**Clustering**

https://www.timeseriesclassification.com/results.php (June 2020)

**Classification**

In line with related work [67], we found that deep learning approaches are not (yet) competitive despite their higher processing effort on training data. We could also confirm that *"simple methods yield performance almost as good as more sophisticated methods"* [56]. Still, no single algorithm clearly performs best. We highlighted sev-

Schmidl, S., Wenig, P., & Papenbrock, T. (2022). Anomaly detection in time series: a comprehensive evaluation. *Proceedings of the VLDB Endowment*, 15(9), 1779-1797.

**Anomaly detection**

= deep learning approach

# How NOT to use DL for time series (1/4)

Blog post from V. Flovik
https://towardsdatascience.com/how-not-to-use-machine-learning-for-time-series-forecasting-avoiding-the-pitfalls-19f9d7adf424
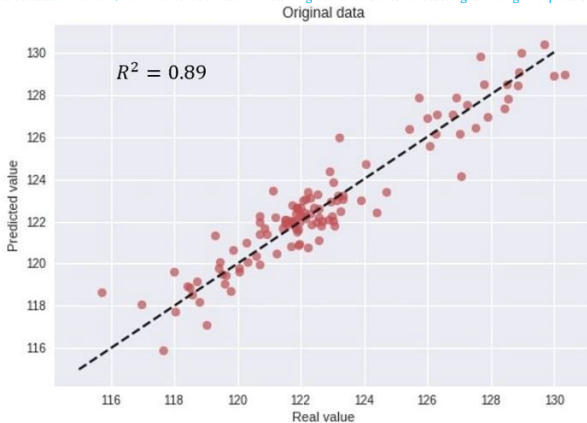


Prediction of stock index with LSTM network: use the first 250 days as training
data. Prediction seems great !!

# How NOT to use DL for time series (2/4)

Blog post from V. Flovik
https://towardsdatascience.com/how-not-to-use-machine-learning-for-time-series-forecasting-avoiding-the-pitfalls-19f9d7adf424
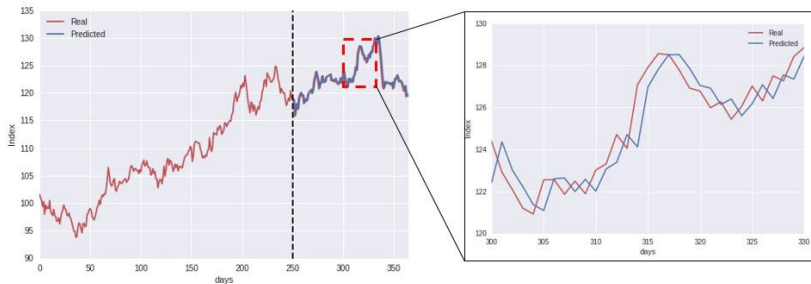


Great performances too! Accurate prediction and RMSE !

# How NOT to use DL for time series (3/4)

Blog post from V. Flovik
https://towardsdatascience.com/how-not-to-use-machine-learning-for-time-series-forecasting-avoiding-the-pitfalls-19f9d7adf424



In fact, LSTM was just repeating the previous sample...

# How NOT to use DL for time series (4/4)

Blog post from V. Flovik

https://towardsdatascience.com/how-not-to-use-machine-learning-for-time-series-forecasting-avoiding-the-pitfalls-19f9d7adf424



In fact, the data was a random walk : impossible to predict. This could have been detected by a careful pre-investigation...
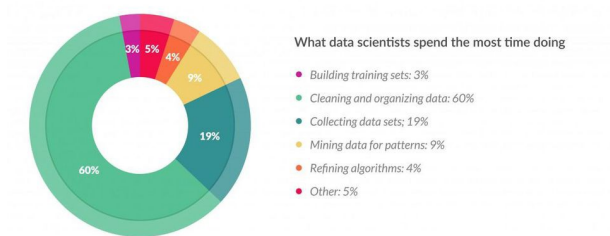
# Contents

# Data science

- ▶ Data science is not (or at least should not) attempting to obtain the best performances by launching DL packages in Python
- ▶ Data science also aims at understanding the data, interacting with experts, bring human intelligence and expertise and improve knowledge
- ▶ Artificial intelligence cannot be intelligent if the data scientist is not
- ▶ Applying complex DL methods does not prevent from a thorough preliminary phase... **and ML can also help for this!**
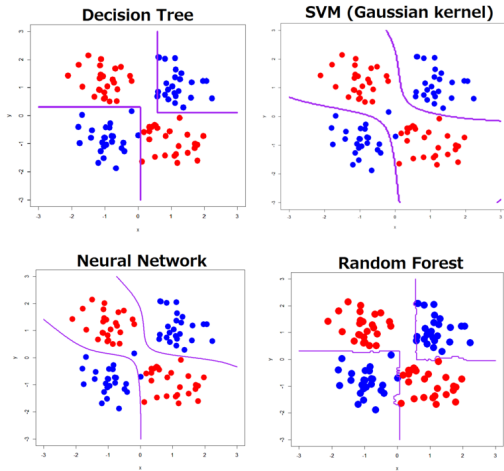
# Understanding data: complex and time-consuming task



**What data scientists spend the most time doing**

- *Building training sets: 3%*
- *Cleaning and organizing data: 60%*
- *Collecting data sets: 19%*
- *Mining data for patterns: 9%*
- *Refining algorithms: 4%*
- *Other: 5%*

Source: https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/
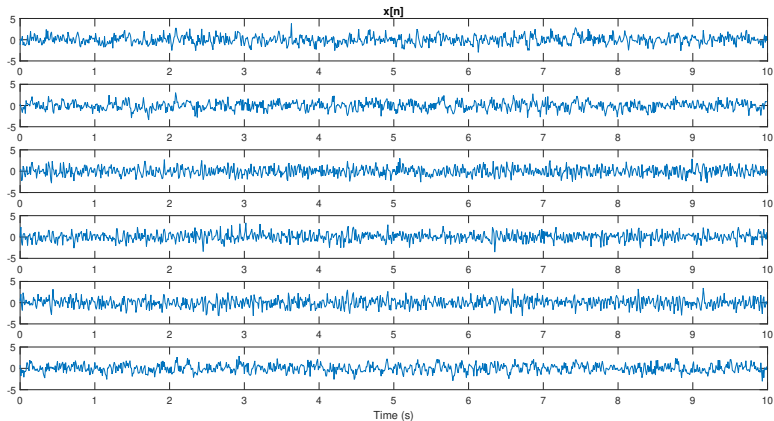
▶ Understanding the data for extracting the relevant information
▶ Understand what you do, why you do it and how you do it: interpretability
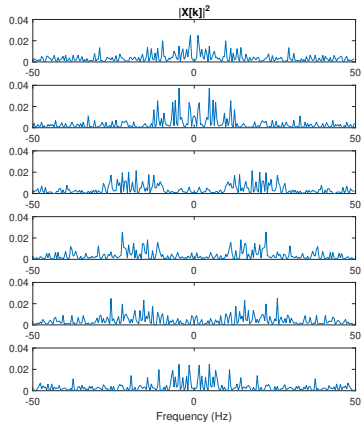
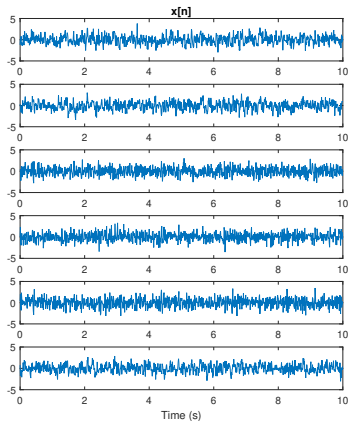# Representation vs. complexity



Data complexity often translates into algorithm and model complexity

# Importance of representation



Two classes of signals?

# Importance of representation



Trivial in the frequency domain

# Main ML tasks for time series

▶ **Prediction:** Predict the future values a time series

▶ **Completion/interpolation:** Recover missing/lost samples in a time series

▶ **Classification:** Assign a class label to a time series or to a subsequence

▶ **Clustering:** Form several groups of time series with the same properties

▶ **Query by content/indexation:** Given an input time series, retrieve the closest time series in a large database up to a given measure of fit

▶ **Segmentation/change-point detection:** Find significant abrupt changes in the time series

▶ **Anomaly detection:** Find abnormal events in a time series

▶ **Pattern extraction:** Find repetitive events in a time series

# Hidden ML tasks for time series

▶ **Understand the data:** know where they come from, how they were acquired, what are their characteristics, interact with domain-experts and understand their problems

▶ **Improve the data:** find accurate representation spaces where the events of interest can be seen, consolidate the data (denoising, detrending, detection/removal of outliers)

▶ **Model the data:** physical/statistical or expert-based models, simple, adaptive and interpretable models

▶ **Extract information from the data:** find repetitive patterns, features of interest, change-points, anomalies

# Contents

## Aim of the course

# **Machine Learning for (understanding) Time Series**

▶ Focus on the *hidden tasks*: understand, improve, model and extract information

▶ Interpretable and reproducible ML algorithms: white boxes (no Deep Learning)

▶ Unsupervised and semi-supervised ML approaches

▶ Methodology can be applied for prediction, classification, clustering etc...

# Outline of the course

- ▶ Lecture 1: Pattern Recognition and Detection
- ▶ Lecture 2: Feature Extraction and Selection
- ▶ Lecture 3: Models and Representation Learning
- ▶ Lecture 4: Data Enhancement and Preprocessings
- ▶ Lecture 5: Change-Point and Anomaly Detection
- ▶ Lecture 6: Multivariate Time Series