# Multivariate Time Series:
## Challenges, missing data, and forecasting

Ana Filipa Almeida
anaa@ua.pt

# Introduction

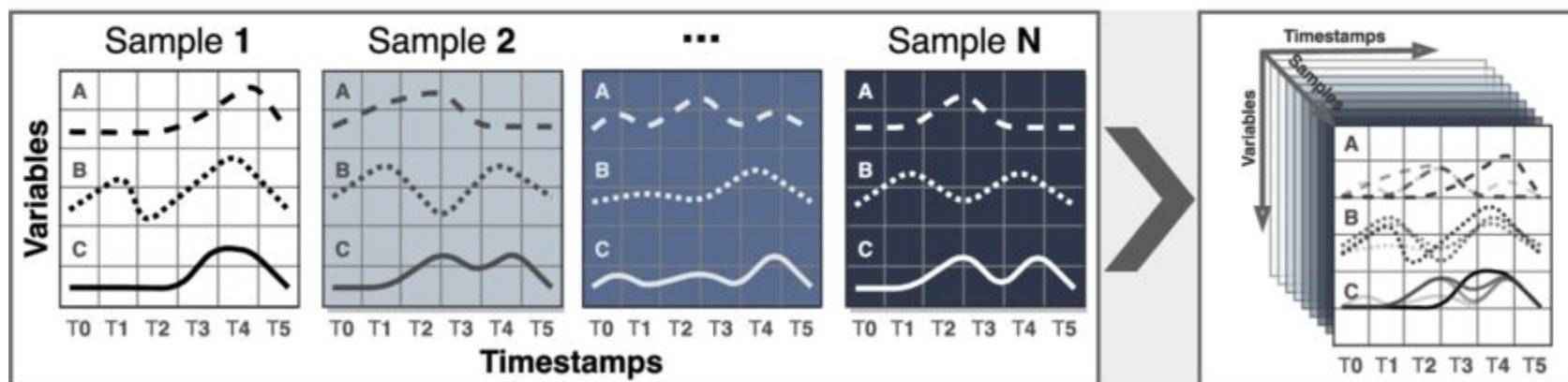# Brief overview of time series data

❏ Time series data:
  ❏ describes how something changes over time,
  ❏ is a sequential set of data points,
  ❏ can present temporal patterns.

❏ Understanding time series data enables us to make predictions, identify patterns, and make informed decisions.

# What is a multivariate time series?

❏ Multivariate time series data:
  ❏ Evolution of several variables over time
  ❏ Different variables might influence each other
  ❏ We might find more complex patterns, such as spatio-temporal patterns

# Challenges in analyzing time series data (1)

- ❏ Complexity
    - ❏ Complex systems
    - ❏ Complex relationships between features

- ❏ Data quality
    - ❏ Missing data
    - ❏ Outliers
    - ❏ Noise
    - ❏ …

- ❏ Non-stationary time series

- ❏ High-Dimensionality

# Challenges: Complexity



- ❏ Supply and demand
- ❏ Geopolitical events
- ❏ Financial crises
- ❏ Pandemic
- ❏ Weather conditions
- ❏ Seasonal patterns
- ❏ Cyclic patterns
- ❏ Currency fluctuations

# Challenges: Data Quality

❏ The quality of data significantly impacts the reliability of data analysis, modeling, and decision-making

❏ Problems commonly found:
  ❏ Missing data
  ❏ Anomalies
  ❏ Inconsistent sampling rates
  ❏ Duplicate data
  ❏ Lack of consistency in units
  ❏ Data drift
  ❏ Noise
  ❏ Data quality degradation over time

# Use cases

- ❏ Finances
    - ❏ Stock Market Analysis
    - ❏ Forecasting stock prices

- ❏ Healthcare
    - ❏ Patient monitoring
    - ❏ Detecting diseases

- ❏ Urban mobility
    - ❏ Traffic Management
    - ❏ Real-time traffic forecasting
    - ❏ Optimize traffic signal timing

My research aims to **forecast traffic metrics** even in the presence of **high ratios of missing data**.
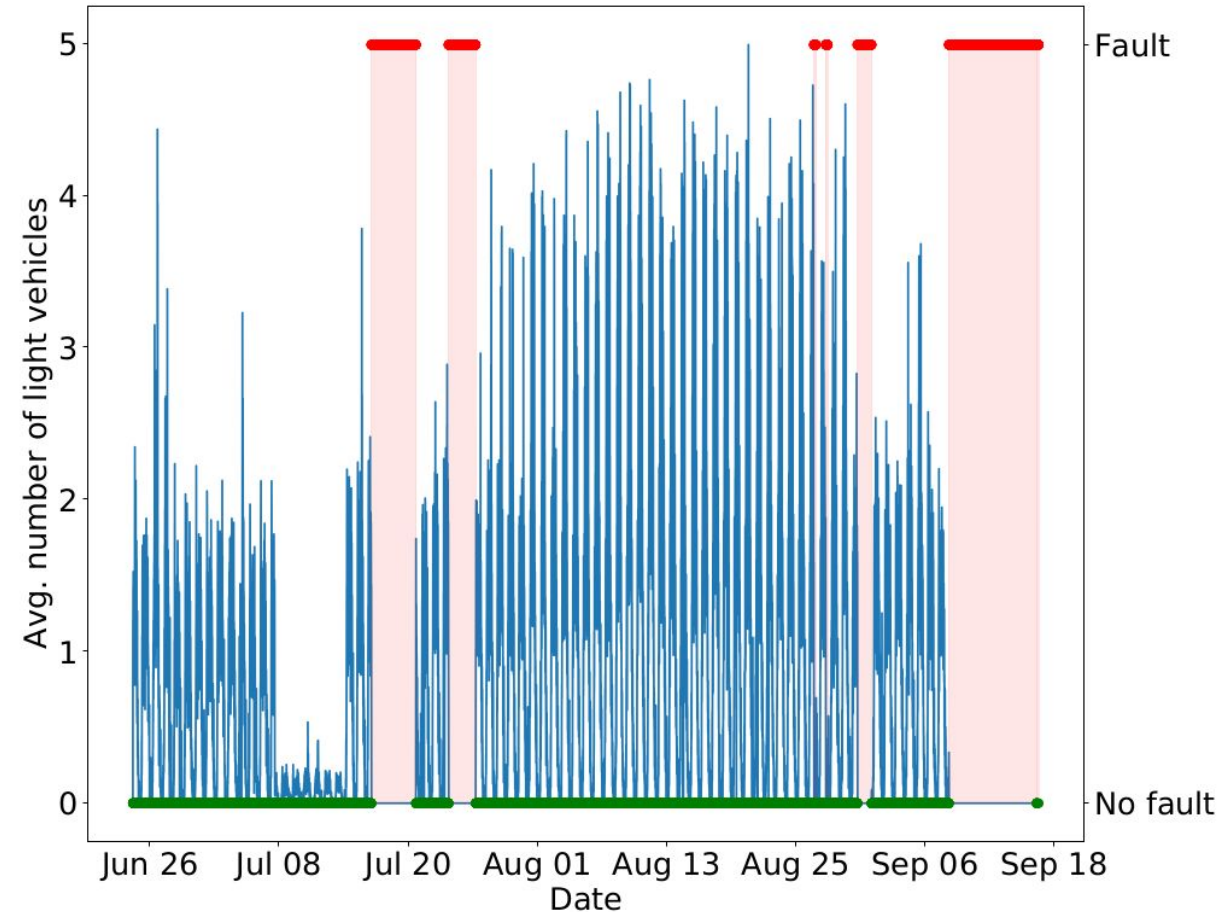
# Dealing with Missing Data

# Time series missing data

❑ Missing data in time series can hide existing patterns and trends

❑ Simple imputation methods can fail in the presence of long intervals of missing data

❑ Data analysis and algorithms can be affected by the existence of missing data

# How can we notice missing data?

- ❏ '*Nan*' values
- ❏ Masked missing data
  - ❏ Default values out of range/ that are not possible

> The number of people in Hamburg yesterday was -1

  - ❏ Default values inside the range/ are possible
    - ❏ This can be dangerous depending on the context
    - ❏ This can make it difficult for us to identify missing data

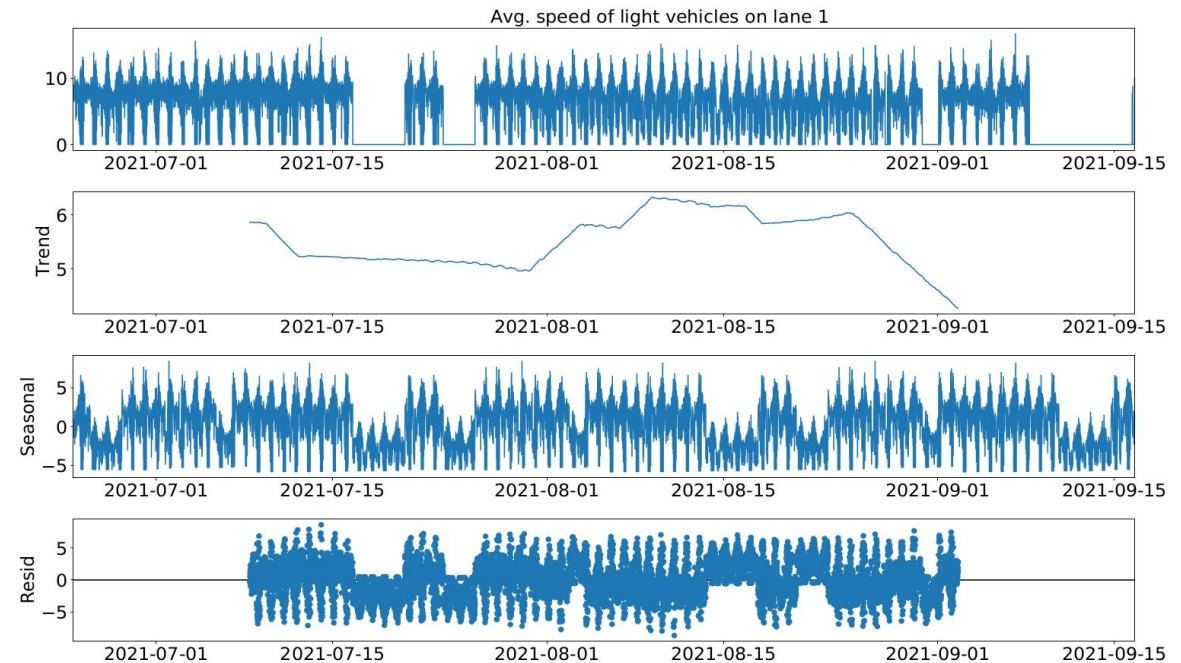> The number of people in Hamburg yesterday was 0

> The speed of cars between 9 a.m. and 10 a.m. yesterday was 0 km/h

# What is the impact of missing data?

❏ We can have **different percentages** of missing data
❏ Missing data can **affect one or more features**
❏ We can have problems that are more tolerant to missing data than others
❏ We can have models that deal better with lower rates of missing data, while other may deal better with higher rates of missing data
❏ We can have different scenarios of missing data

**Missing data can affect data analysis, models, and algorithms. Having a negative impact on its applications in business and research.**



Avg. speed of light vehicles on lane 1

# How can we handle missing data?

❏ Before it happen
  ❏ We can prevent it from happen by **monitoring our system and preventing conditions that lead to missing data**
  ❏ However, this is not always possible!!!!
❏ After it happen
  ❏ Ignore missing data (not a very good solution)
  ❏ Delete observations with missing data (not adequate for time series, even if we have few observations with missing data)
  ❏ Replace missing data for a value
    ❏ 0 or values out of range (however, this can bring additional problems)
    ❏ Mean, median, mode…
    ❏ Simple univariate imputation techniques
      ❏ e.g., Moving Average
  ❏ Multivariate imputation techniques

universidade de aveiro
theoria poiesis praxis

NAP

deti universidade de aveiro
departamento de eletrónica,
telecomunicações e informática

instituto de
telecomunicações

# Case study: OpenWeather dataset

❏ Dataset provided by OpenWeather
❏ Data from sensors such as temperature, pressure, humidity, wind speed, wind direction, wind gust, and cloudiness
❏ Data from 20 cities
❏ Hourly data from 2022

# Workflow

# Generation of synthetic missing data



Overlap

Disjoint

# Generation of synthetic missing data

# Algorithms

- ❏ We selected the top 3 from 20 statistical methods as baseline to evaluate our algorithm:
    - ❏ Replaced missing data using a specific value:
        - ❏ Mean, median, last value, previous value, the nearest value, zero
    - ❏ Interpolation techniques:
        - ❏ Barycentric, pchip, splines, polynomial functions, piecewise polynomial, akima
- ❏ Experimented with the KNN imputer
- ❏ Propose the Focalized KNN algorithm:
    - ❏ Based on KNN imputer

universidade de aveiro
theoria poiesis praxis

NAP

deti
universidade de aveiro
departamento de eletrónica,
telecomunicações e informática

instituto de
telecomunicações

# k-NN Algorithm

# k-NN imputer

❏ Based on k-NN
❏ It can be used for imputation by using similar points to guess the missing data
❏ However, k-NN has some problems…
  ❏ Suffers from the curse of high dimensionality
  ❏ Stores the complete dataset in memory


❏ KNN imputer can be used for time series datasets; however, it does not take advantage of time series properties!!!

# Correlation between features



Spatio-temporal patterns!

# Correlation between temporal lags

# Focalized KNN

- ❏ Select the most correlated features
- ❏ Select the most relevant temporal lags
- ❏ Select the column with less missing data
    - ❏ Apply the KNN imputer for the matrix composed with:
        - ❏ the column c with missing data + correlated features + relevant temporal lags
    - ❏ Replace column c with the column with the imputed data
    - ❏ Repeat until there is no more missing data

# Evaluation metrics

❏ Mean Squared Error (MSE)

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$

❏ R²-Score

$$R^2 - Score = 1 - \frac{MSE}{MSE_{baseline}}$$

# Overlap versus Disjoint missing pattern



Overlap

Disjoint

# Discussion

❏ Pros
  ❏ We do not need to train our algorithm
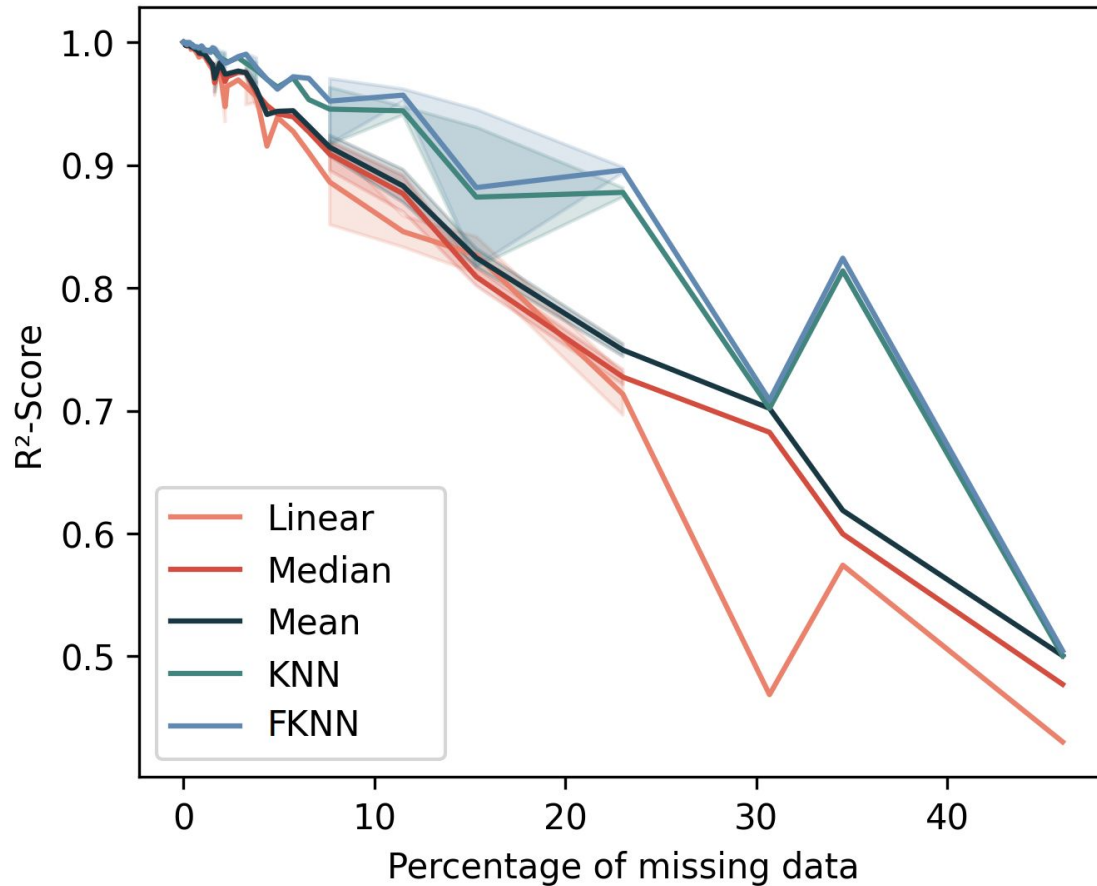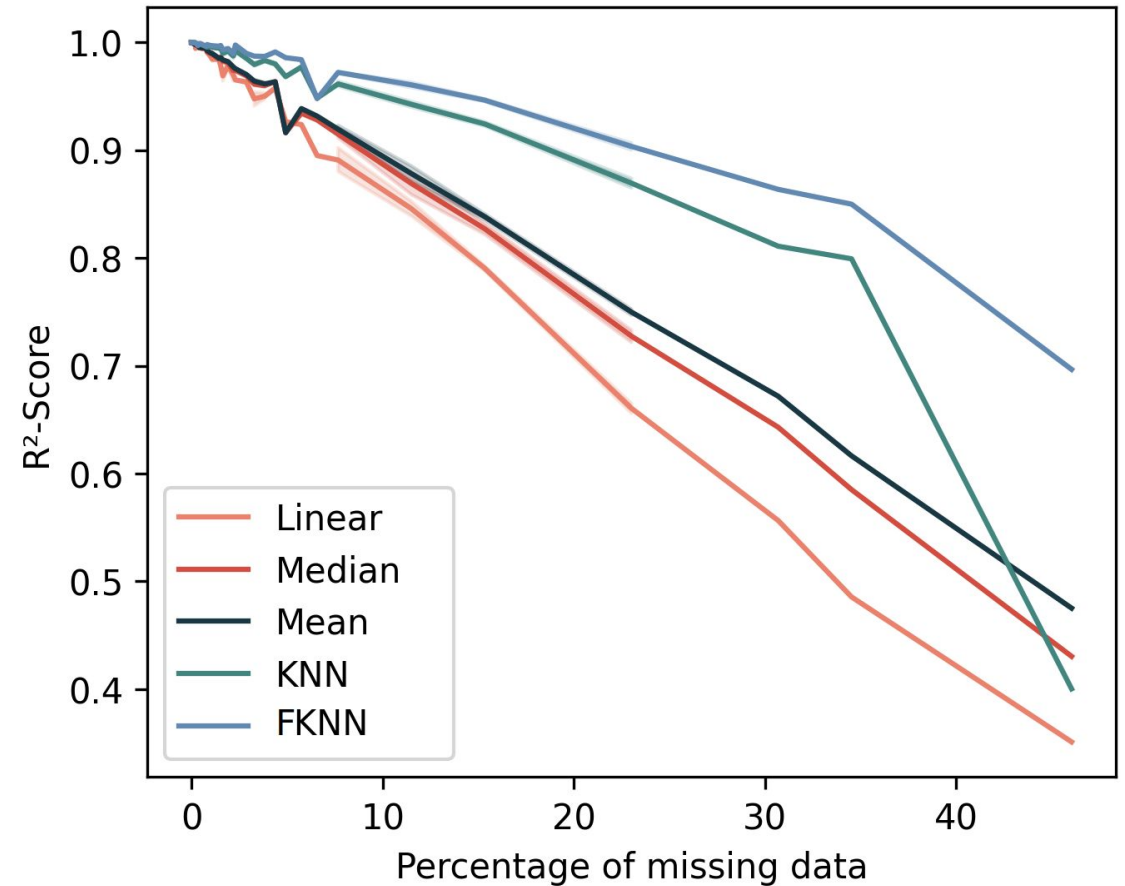  ❏ Good with disjoint patterns
  ❏ Our solution helps with the curse of high dimensionality
❏ And Cons
  ❏ Usually, it takes more time than regular KNN, especially if we have missing data affecting several columns
  ❏ Not very good with overlap patterns

In the future, we would like to:

❏ Develop other methods to perform imputation in time series, such as methods based on Autoencoders
❏ Create more patterns of missing data to evaluate our models

SPRINGER LINK

Find a journal   Publish with us   🔍 Search

Iberian Conference on Pattern Recognition and Image Analysis
↳ IbPRIA 2023: **Pattern Recognition and Image Analysis** pp 28–39 | Cite as

Home > Pattern Recognition and Image Analysis > Conference paper

## Time Series Imputation in Faulty Systems

Ana Almeida ✉, Susana Brás, Susana Sargento & Filipe Cabral Pinto

Conference paper | First Online: 25 June 2023

**348** Accesses

Part of the Lecture Notes in Computer Science book series (LNCS,volume 14062)

### Abstract

Time series data has a crucial role in business. It reveals temporal trends and patterns, making it possible for decision-makers to make informed decisions and mitigate problems even before they happen. The existence of missing values in time series can bring difficulties in the analysis and lead to inaccurate conclusions. Thus, there is the need to solve this issue by performing missing data imputation on time series.

In this work, we propose a Focalize KNN that takes advantage of time series properties to perform missing data imputation. The approach is tested with different methods, combinations of parameters and features. The results of the proposed approach, with overlap and disjoint missing patterns, show Focalize KNN is very beneficial in scenarios with disjoint missing patterns.

### Keywords

Missing data imputation   K-Nearest Neighbors   Time series   Overlap missing data   Disjoint missing data

Access via your institution →

Chapter   EUR 29.95
Price includes VAT (Germany)
• Available as PDF
• Read on any device
• Instant download
• Own it forever

Buy Chapter

eBook   EUR 82.38
Softcover Book   EUR 104.85

Tax calculation will be finalised at checkout
**Purchases are for personal use only**
Learn about institutional subscriptions

Sections   References
Abstract
Notes
References
Author information
Editor information
Rights and permissions
Copyright information
About this paper

This work is supported by FEDER, through POR LISBOA 2020 and COMPETE 2020 of the

universidade de aveiro
theoria poiesis praxis

NAP

deti universidade de aveiro
departamento de eletrónica,
telecomunicações e informática

instituto de
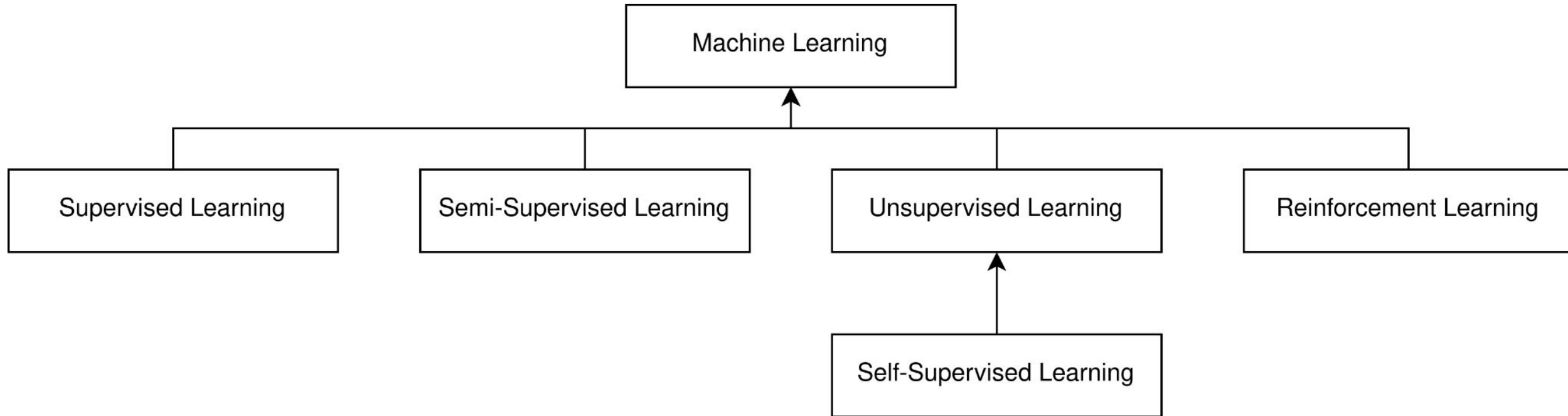telecomunicações

# Forecasting

# Brief overview of forecasting

❏ Forecasting is the process of making predictions about future events.
❏ Forecasting can be very beneficial in decision-making, planning, and risk management.
❏ There are different time horizons for forecasting, such as short, medium, and long-term.
❏ Before applying forecasting techniques, we should analyze the time series.

# Forecasting methods

- ❏ Naïve methods
  - ❏ Use the last value to forecast the next one
- ❏ Statistical methods
  - ❏ AutoRegressive (AR)
  - ❏ Moving Average (MA)
  - ❏ Autoregressive Integrated Moving Average (ARIMA)
  - ❏ Seasonal ARIMA (SARIMA)
- ❏ Machine Learning methods
  - ❏ SVM
  - ❏ kNN
  - ❏ LightGBM
- ❏ Deep Learning methods
  - ❏ FNN
  - ❏ GRU
  - ❏ LSTM
  - ❏ CNN

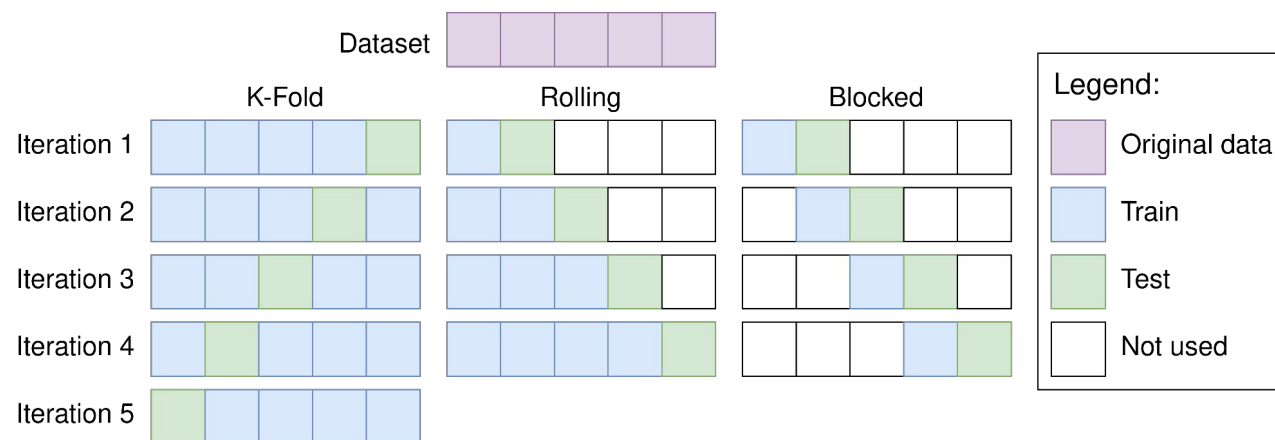# What type of learning problem is forecasting?

# Choosing the best model for forecasting methods

❏ Most used evaluation metrics in the literature:
  ❏ Mean Absolute Error (MAE)
  ❏ Mean Absolute Percentage Error (MAPE)
  ❏ Root Mean Squared Error (RMSE)
  ❏ $R^2$-Score

> Metrics based on comparing the **observed value** with the **predicted value**

❏ Cross-Validation for Time-series

# Case study: Forecasting traffic flow

Dataset:

- ❏ Traffic counters deployed in Oporto, Portugal
- ❏ Data from September 30 to November 3 of 2019
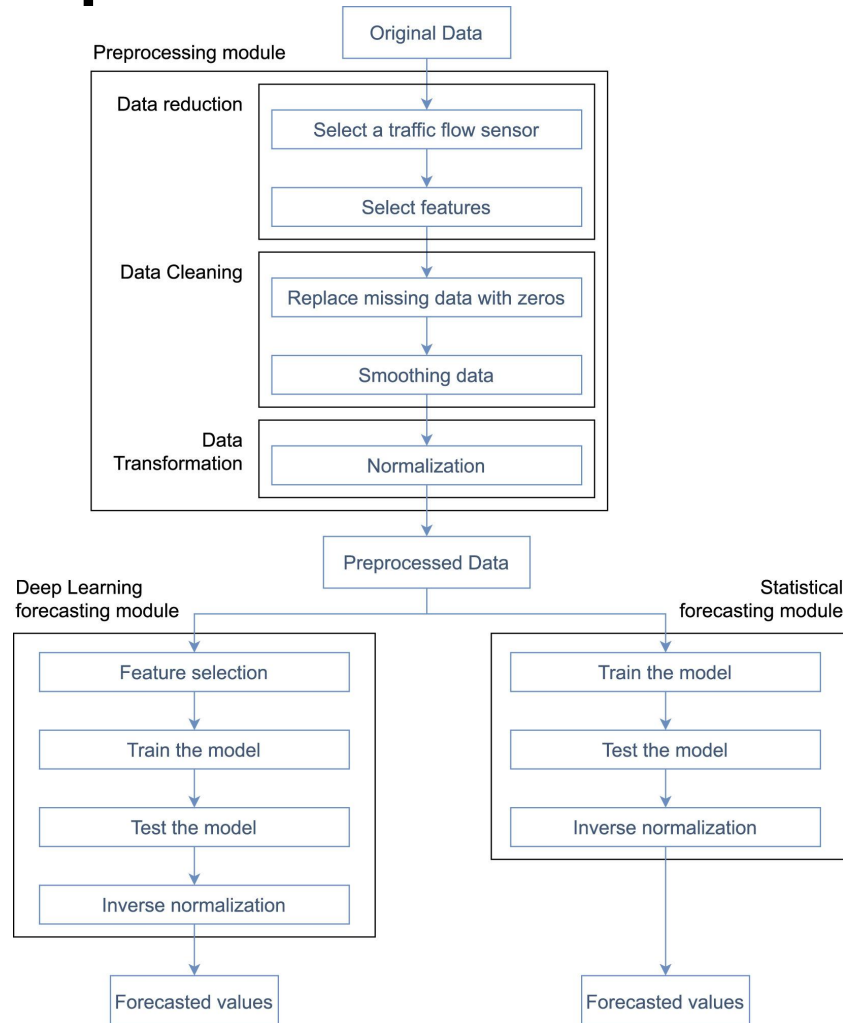- ❏ 5 minutes interval
- ❏ More than 100 sensors

# Approach and Methods

- ❏ SARIMA
- ❏ FNN models
- ❏ LSTM-based models
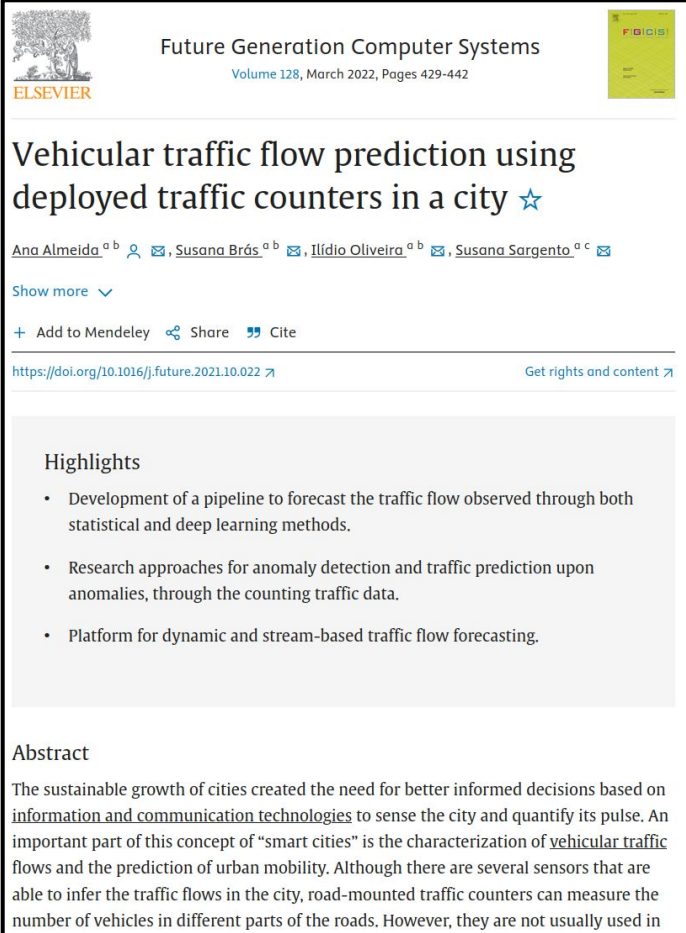- ❏ CNN-based models
- ❏ Hybrid LSTM-CNN models

# Forecasting Pipeline

# Results and Discussion

❏ Short-term forecasting
  ❏ SARIMA achieves a good performance
  ❏ Computationally light
  ❏ (We can also use deep learning strategies)
❏ Long-term forecasting
  ❏ SARIMA is not suitable for long-term forecasting
  ❏ Deep learning strategies achieve good performance
  ❏ The best model was based on CNNs
  ❏ LSTMs take more time to train than CNNs or FNNs

Vehicular traffic flow prediction using deployed traffic counters in a city ☆

Ana Almeida ᵃ ᵇ, Susana Brás ᵃ ᵇ, Ilídio Oliveira ᵃ ᵇ, Susana Sargento ᵃ ᶜ
Show more ∨

+ Add to Mendeley   ⬀ Share   🙶 Cite

Highlights

- Development of a pipeline to forecast the traffic flow observed through both statistical and deep learning methods.
- Research approaches for anomaly detection and traffic prediction upon anomalies, through the counting traffic data.
- Platform for dynamic and stream-based traffic flow forecasting.

Abstract

The sustainable growth of cities created the need for better informed decisions based on information and communication technologies to sense the city and quantify its pulse. An important part of this concept of "smart cities" is the characterization of vehicular traffic flows and the prediction of urban mobility. Although there are several sensors that are able to infer the traffic flows in the city, road-mounted traffic counters can measure the number of vehicles in different parts of the roads. However, they are not usually used in

universidade de aveiro
theoria poiesis praxis

NAP

deti universidade de aveiro
departamento de eletrónica,
telecomunicações e informática

instituto de
telecomunicações

# Future Work

universidade de aveiro
theoria poiesis praxis

NAP

deti universidade de aveiro
departamento de eletrónica,
telecomunicações e informática

instituto de
telecomunicações

# Future Work

❏ Develop a model able predict future values even in the presence of missing data
  ❏ Imputation + Forecasting

# Q&A

Thank you!

anaa@ua.pt