



National and Kapodistrian
University of Athens

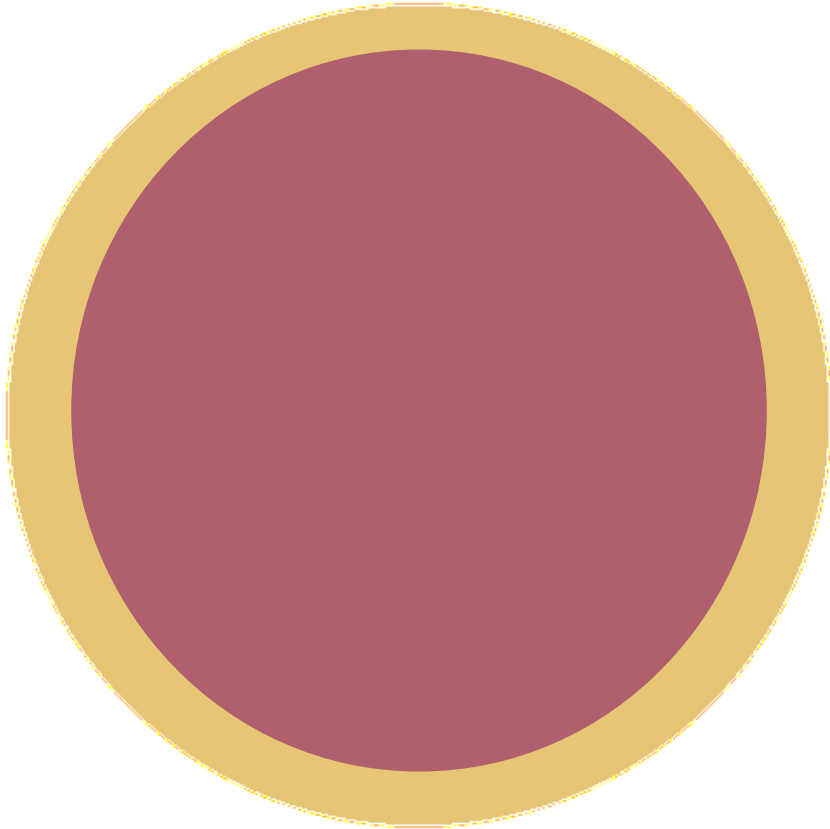
A Spatio-Temporal Data Imputation Model for Supporting Analytics at the Edge

Kostas Kolomvatsos, Panagiota Papadopoulou,
Christos Anagnostopoulos, Stathes Hadjiefthymiades

The 18th IFIP Conference on E-Business, E-Services
and E-Society (I3E 2019)
Trondheim, Norway
18-20 September 2019



Outline



01 Introduction

Current state of the art, our contribution and novelty of our work.

02 Problem Description

Description of the envisioned setting.

03 The Proposed Model

Our approach for data imputation at the edge of the network.

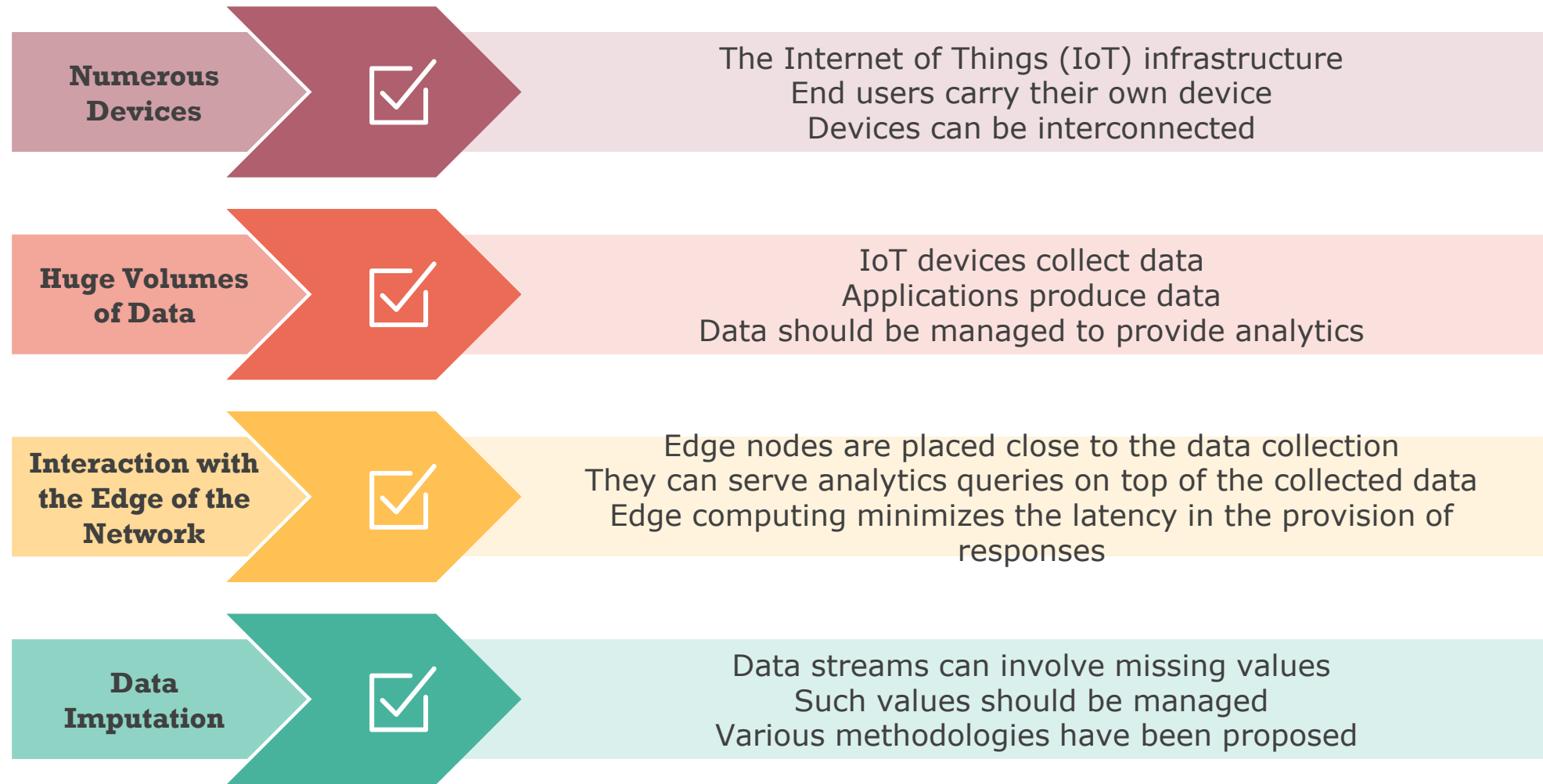
04 Experimental Evaluation

Description of our experiments and the delivered outcomes.

05 Conclusions

Our conclusions and vision for future work.

Edge Computing



Our Contribution

State of the Art

- Legacy techniques mainly focus on the statistics of data
- They try to find the best value to replace the missing data
- The detection of the exact distribution of data can be difficult

Our focus

- We incorporate the dynamics of the environment where IoT devices act
- We take into consideration the group of nodes reporting data
- We adopt a spatio-temporal approach

We propose

- We combine the reports of multiple IoT devices
- We consider the proximity of:
 - the location
 - the reported data
- Our two-layer clustering scheme combined with a consensus model assists in the missing values replacement

The Envisioned Setting



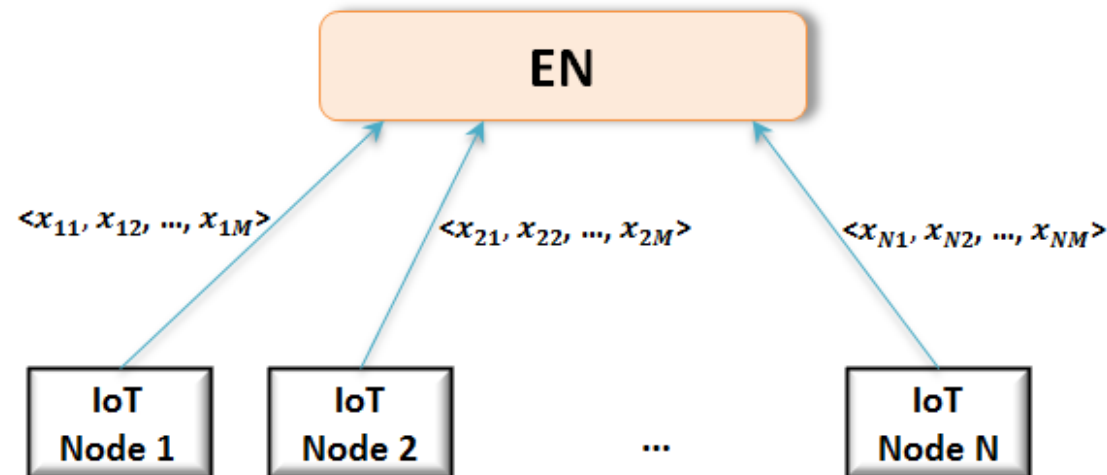
IoT devices collect multivariate data from their environment

A set of Edge Nodes (ENs) receive data from a set of IoT devices

Every EN should detect if missing values are present and, then, apply the proposed mechanism

Missing can refer in the entire multivariate vector or specific dimensions

We consider a sliding window approach and take into consideration the location of each device



Clustering and Correlation



1st level of clustering



2nd level of clustering



We focus on the spatial proximity of the devices



We focus on the data proximity of the IoT devices



We adopt the k-Means algorithm



The clustering is applied for a 1st level cluster



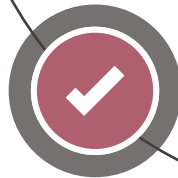
We create a set of clusters with devices in close distance



We create a set of clusters with devices reporting 'similar' data



Clusters are produced at pre-defined intervals being aligned with the mobility of the devices



We adopt the k-Means algorithm



Data Proximity

We get the IDs of the IoT devices and focus at each time step

Every cluster represents a 'transaction' where a (sub-)set of IDs are involved

For each t, we provide the delivered clusters

The presence of IDs in a cluster at t depicts the data correlation between the corresponding IoT devices

When a missing value is present, we get the intersection of clusters where the device ID with missing value is present

We adopt the Pearson Correlation Coefficient (PCC) in multiple data dimensions

| | Node 1 | Node 2 | ... | Node N |
|-------|---|---|-----|---|
| t = 1 | $\langle x_{11}^1, x_{12}^1, \dots, x_{1M}^1 \rangle$ | $\langle x_{21}^1, x_{22}^1, \dots, x_{2M}^1 \rangle$ | ... | $\langle x_{N1}^1, x_{N2}^1, \dots, x_{NM}^1 \rangle$ |
| ... | ... | | | |
| t = W | $\langle x_{11}^W, x_{12}^W, \dots, x_{1M}^W \rangle$ | $\langle x_{21}^W, x_{22}^W, \dots, x_{2M}^W \rangle$ | ... | $\langle x_{N1}^W, x_{N2}^W, \dots, x_{NM}^W \rangle$ |



Data Imputation

Devices

We rely on the delivered clusters

We focus on the dimension where missing values are observed

When multiple dimensions are involved, we adopt an iterative approach

Aggregation

We adopt the linear opinion pool model

It is standard approach for aggregating multiple experts opinion

Only devices with strong correlation are involved

Imputation

We consider the weighted average of data at the same dimension

The weight of a device is high when a high correlation is observed

Weights are calculated on the correlation of all dimensions to avoid random events

Experimental Setup



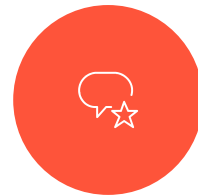
We report on the performance of the model

We aim at detecting if the replacements are close to the real values

We adopt widely known performance metrics



Mean Absolute Error (MAE)



Datasets:

- Unmanned Surface Vehicles Sensor Data¹
- Intel Berkeley Research Lab Dataset²
- the Iris dataset³



Root Mean Squared Error (RMSE)



We compare our model with:

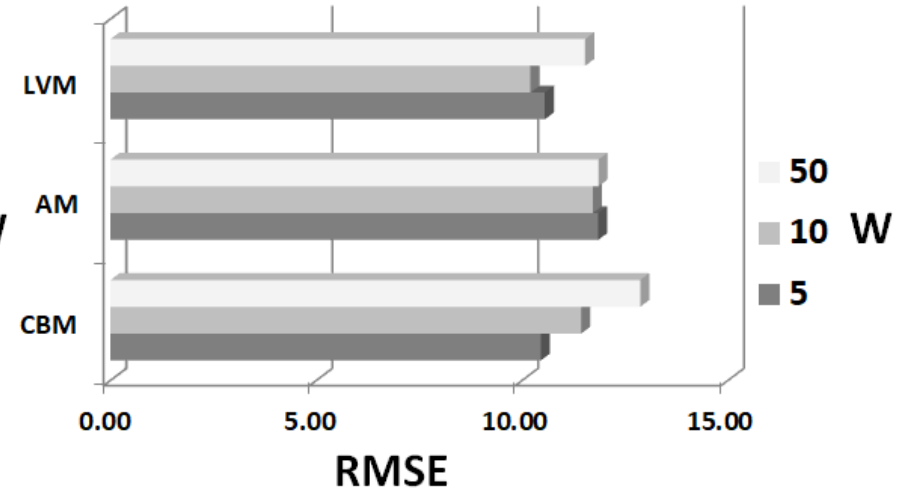
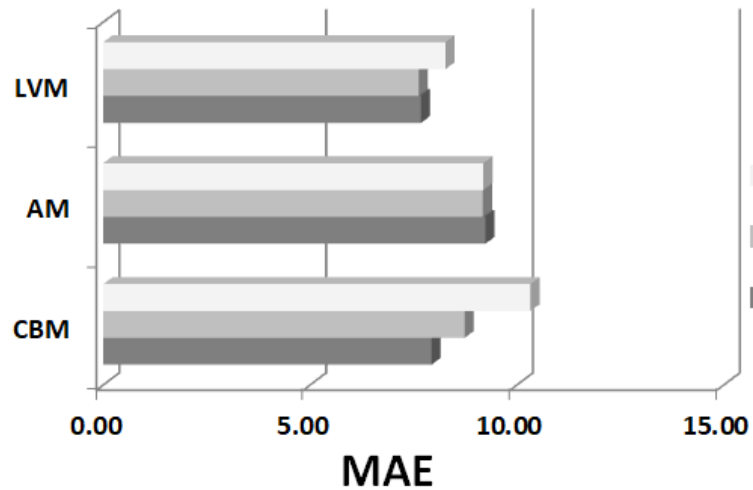
- an Averaging Mechanism
- the Last Value Mechanism

¹ Harth, N., Anagnostopoulos, C., 'Edge-centric Efficient Regression Analytics', IEEE EDGE, 2018

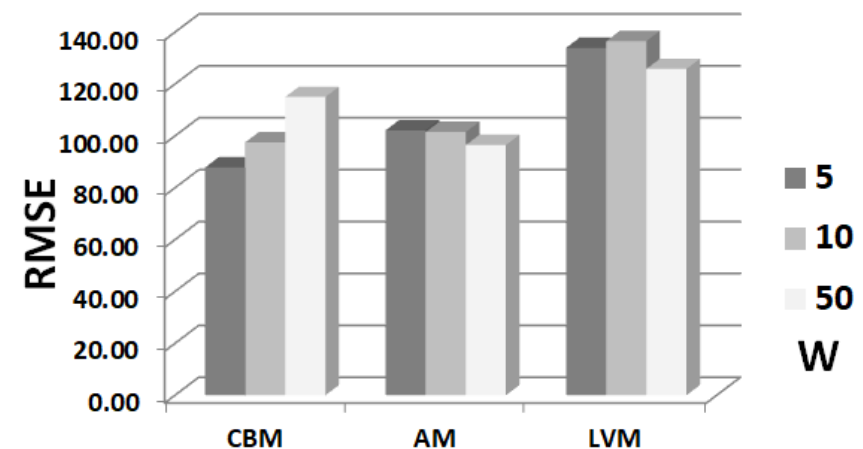
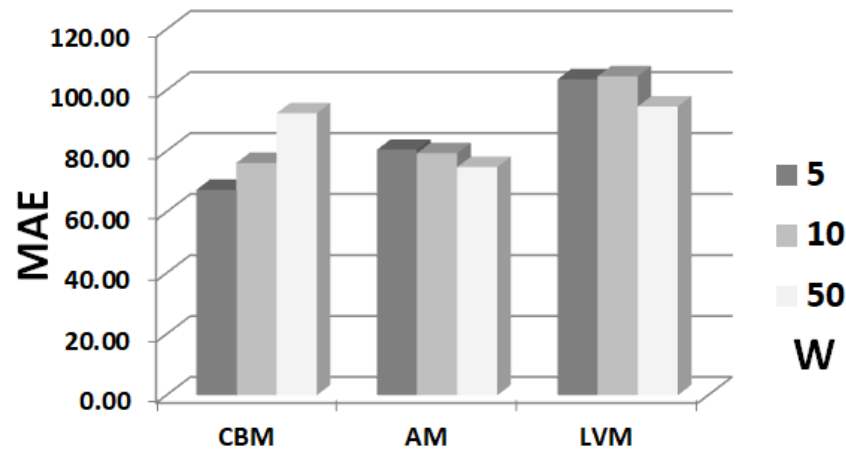
² <http://db.csail.mit.edu/labdata/labdata.html>

³ <http://archive.ics.uci.edu/ml/datasets/iris>

Results

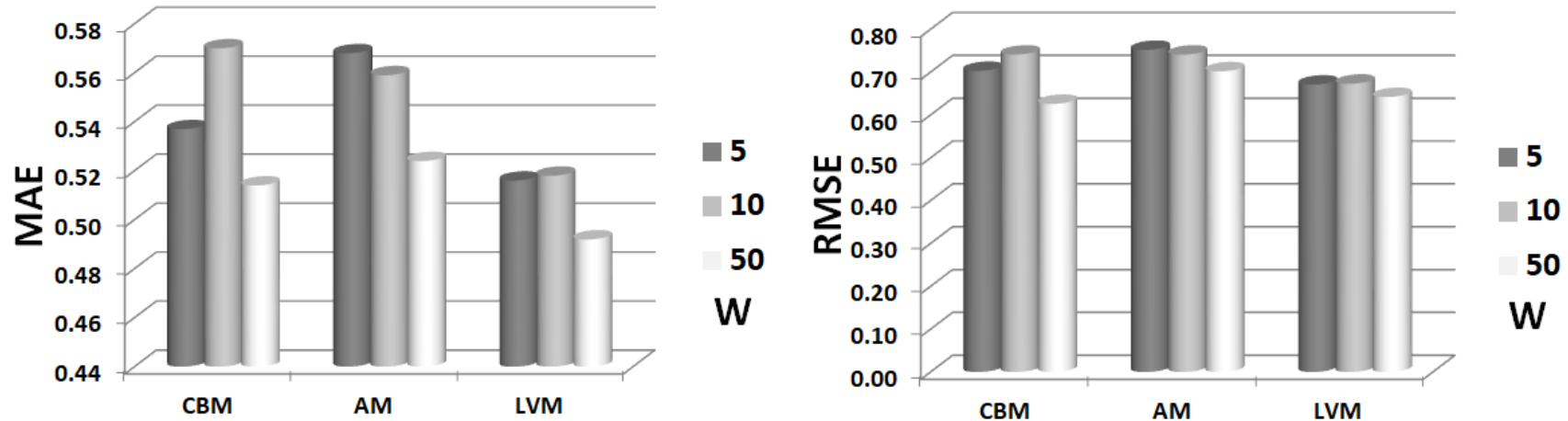


MAE and RMSE for the Unmanned vehicles dataset (different window values)

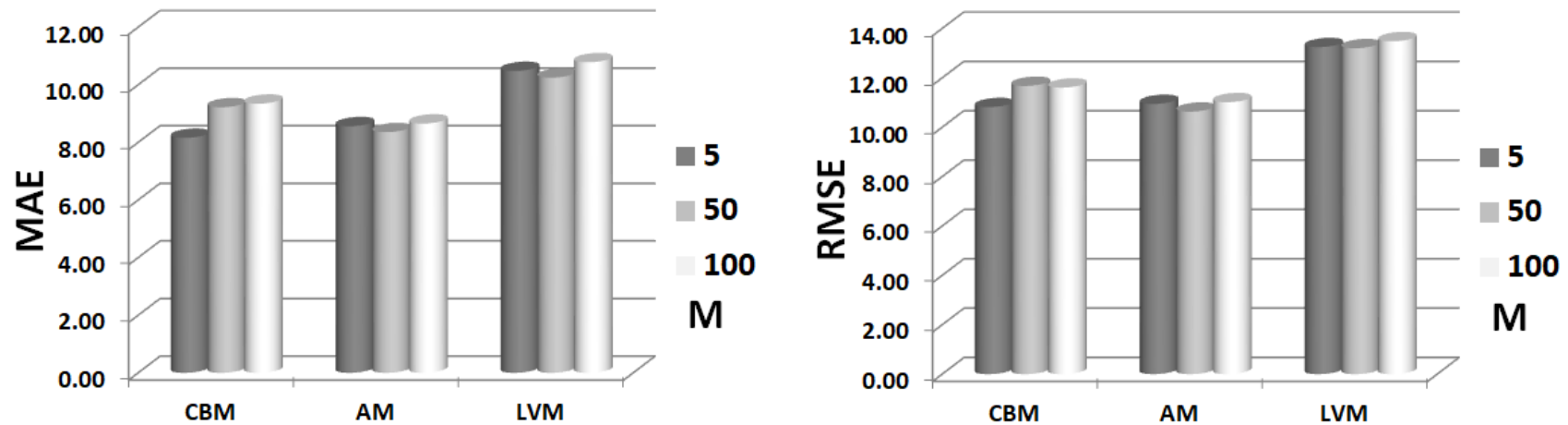


MAE and RMSE for the Intel dataset (different window values)

Results



MAE and RMSE for the Iris dataset (different window values)



MAE and RMSE for the Unmanned vehicles dataset (different number of dimensions)

Conclusions & Future Work

Efficiency

The proposed scheme can efficiently replace the missing values
Our mechanism outperforms in the comparative assessment for the majority of the experimental scenarios

Uncertainty

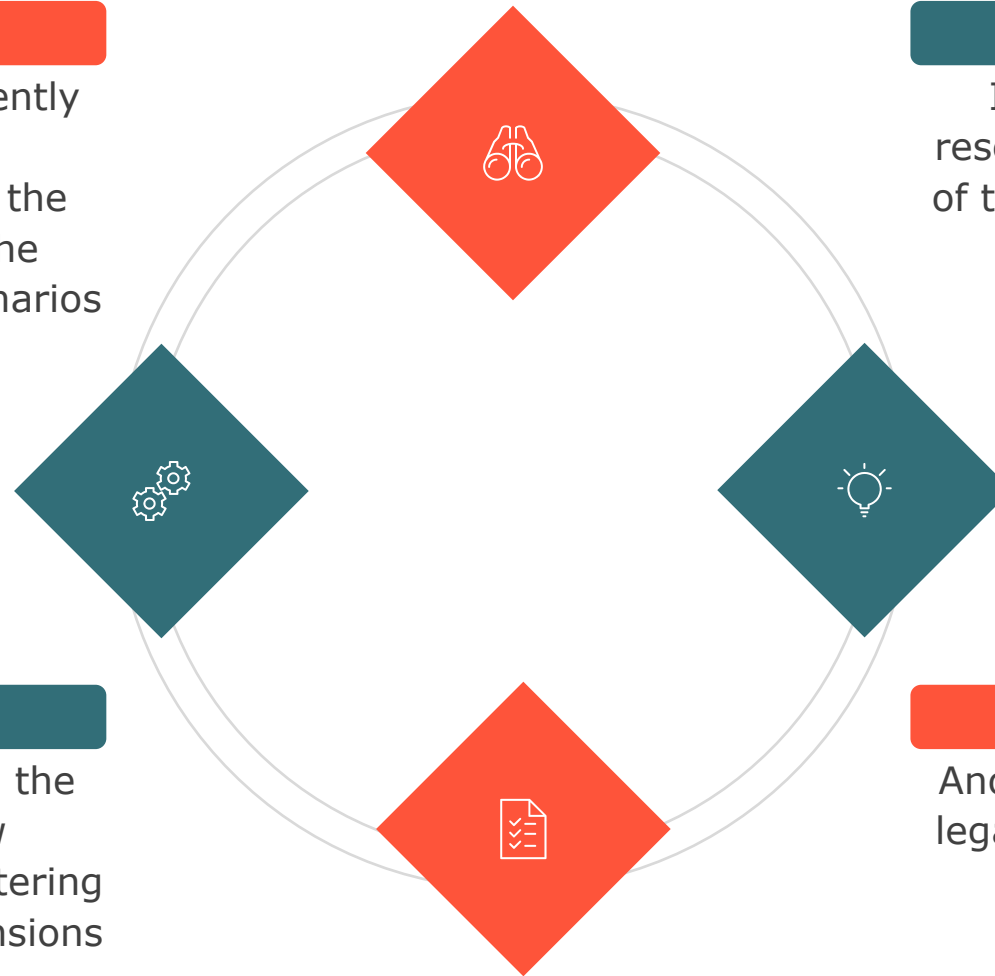
In the first places of our future research agenda is the management of the uncertainty in the aggregation process

Realization

The performance is better when the number of dimensions is low
The model is affected by the clustering process and the number of dimensions involved in the calculations

Statistics

Another research plan is to combine legacy techniques with the proposed model



Thank You!

Questions?

