

DATA ANALYSIS AND VISUALIZATION OF MULTIVARIATE TIME SERIES SOFTWARE

Victor C, Flor De Luz

2024-05-20

Contents

1	Introduction	2
2	Goals	2
3	State of the art	2
4	Pre-implementation analysis	2
4.1	Multilayer design	2
4.2	Frontend specifics	4
4.3	Backend specifics	5
4.4	Detected issues	5
5	DONE Implementation	6
5.1	Activities summary	6
5.2	DONE Activities execution	6
5.2.1	Change log	6
5.2.2	Backlog	6
5.2.3	DONE Backend data structures	6
5.2.4	DONE Backend new specifications	9
5.2.5	DONE Backend new data structure and frontend integration	10
5.2.6	Data repository and cache	11
6	DONE Requirements	11
6.1	DONE Software	11
6.2	DONE Hardware	12
7	DONE Visual Analytics Guidance Development	12
7.1	DONE Spiral diagram analysis improvement	12
7.2	DONE Independent section for statistics	12
7.3	DONE Flow diagram improvement	12
8	DONE Format Multivariate Time Series Data	13
9	References	13

10 Appendix	14
10.1 Time Estimation Plan	14
10.2 Testing functions	14
10.3 Dictionary	15
10.4 Network coloring	16
10.5 Change Log	17
10.6 Backlog	17

1 Introduction

The project explores, analyzes, and integrates Vue.JS and Python code with different models for clean and data completion.

2 Goals

- Complete the code by integrating different Python models into the software project.
- Data structure improvement to adapt new ML models for cleaning data and completion.
- Technical documentation of the software.

3 State of the art

- Source code of the project(1), with documentation to install and deploy the software. The software currently reads data and graphs time series. It has options for data completion through rolling mean and kNN and partially cleaning data alternatives. It tracks for changes while navigating through the Diagram Operator interface graphically. Integrates Radial Chart for time series cycles.
- Self-documented notebook with multiple machine-learning techniques and their variants for data completion like Rolling Mean, Decision Trees, Stochastic Grading Boosting, Locally Weighted Regression, Legendre Polynomials Regression, Random Forest Regressor, k-nearest Neighbors. Includes removing features with no data at all. Computes a Dicky Fuller Stationarity Test. Automatically computes Weighted MAPE and R-Score + RMSE to detect and suggest which model fits better. Computes Autocorrelation, Fourier, and Hodrick Prescott to detect Cyclicity.

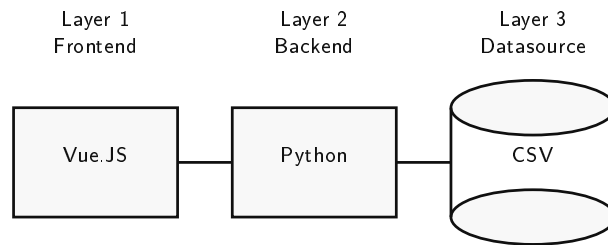
4 Pre-implementation analysis

4.1 Multilayer design

The critical aspects of the software are:

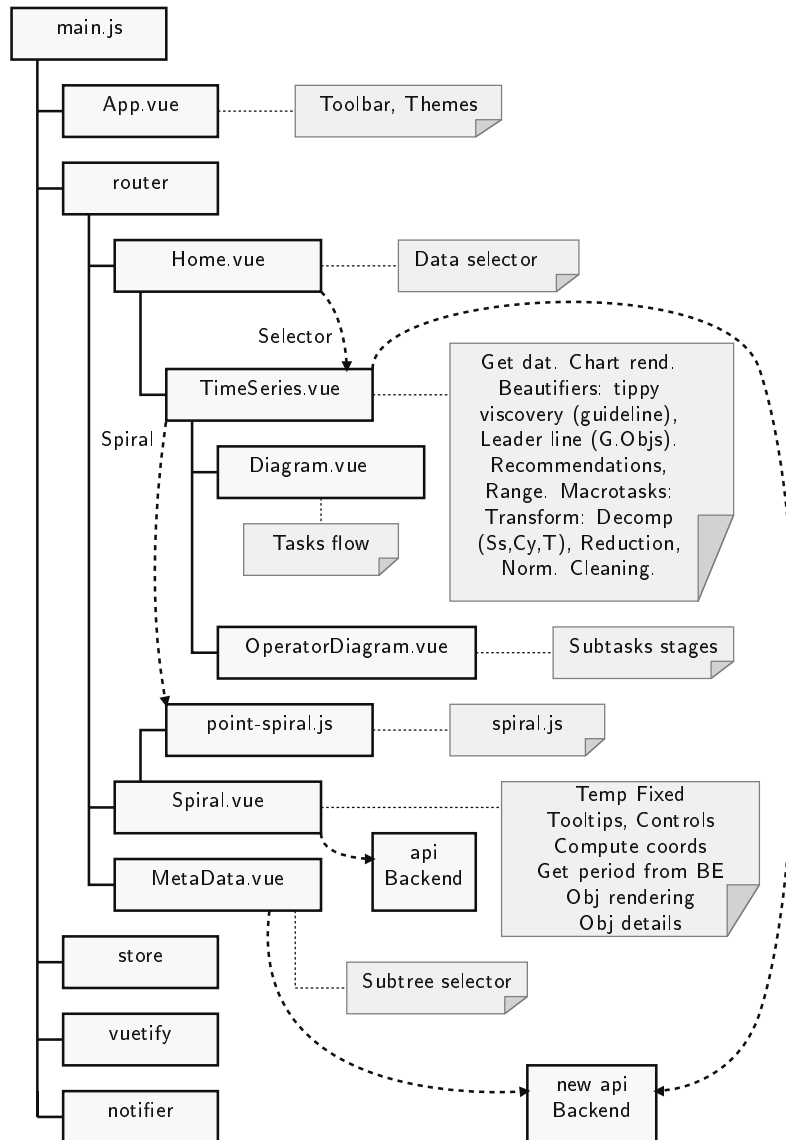
- Frontend developed in Vue.JS

- Backend developed in Python
- The data source layer primarily are CSV files ¹

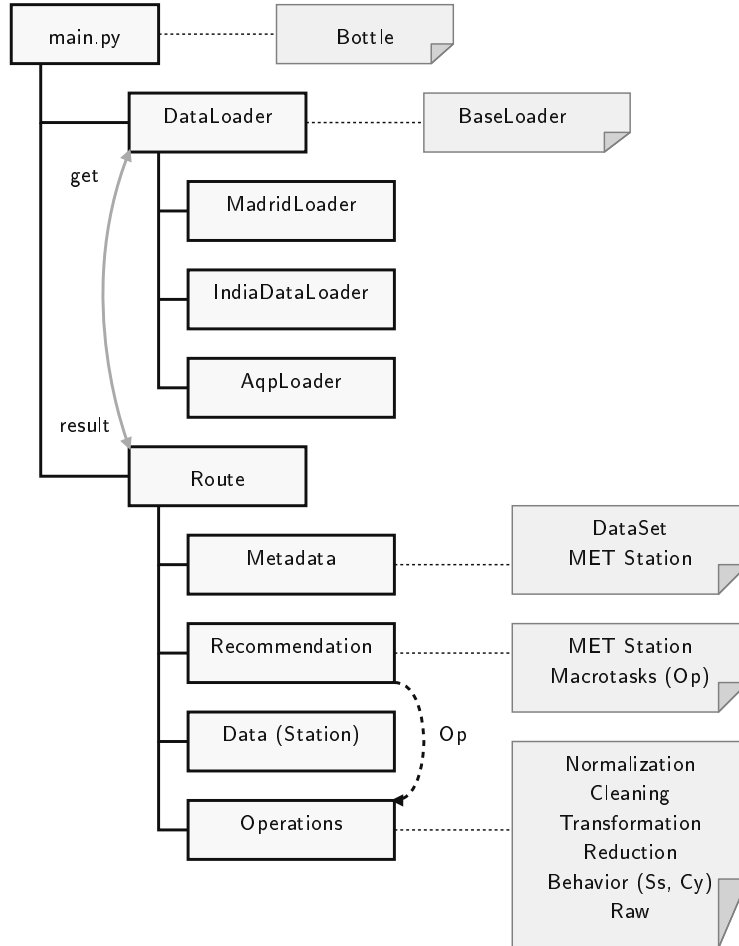


¹HDF5 detected

4.2 Frontend specifics



4.3 Backend specifics



4.4 Detected issues

- Data structures. It requires adapting the current data structures of the program to the new dictionary of the tested and trained models.
- Exceptional cases. There are some treatments in the frontend code for specific data sources ². Those treatments will be done in the backend.
- The Radial Diagram component ³ hardcode for temperature and precipitation. We will rewrite the code to accept different datasets and features. Additionally, it is currently consuming the old backend⁴. The old backend has to be re-implemented in the new backend⁵.

²TimeSeries.vue, temperature, and precipitation, line 308

³Spiral.vue

⁴api/main.py

⁵new_api/main.py

5 DONE Implementation

5.1 Activities summary

- Adapt the Backend Data Structure.
- Data Structure Integration. It will include saving a dictionary and allowing notebook compatibility.
- Adapt the Frontend Data Structure.
- Implement new ML models in the back end.
- Adapt the front end for new ML models.
- Frontend UI additional improvements.
- Format multivariate time series data.

5.2 DONE Activities execution

5.2.1 Change log

Tracking changes on the program will be done by using diff to create a patch component that allows the creation of a checkpoint and registering the changes by the size of modify or created code⁶. The folder structure is:

Table 1: Change log folder structure

A	File System Structure	Description
drwx	VisWeb-AlgoritmosLimpieza.orig	Original Source Code
drwx	VisWeb-AlgoritmosLimpieza.incr	Incremental Checkpoint
drwx	VisWeb-AlgoritmosLimpieza.diff	Diff/Patch Repository
drwx	VisWeb-AlgoritmosLimpieza	Development Folder
-r-	checkpoint.lisp	The program for recording changes

5.2.2 Backlog

The pending activities are in the appendix backlog section. These activities correspond to changes or reviews that depend on multiple program files around the software, and their resolutions will come on the project's timeline. Bugs will gradually fixed.

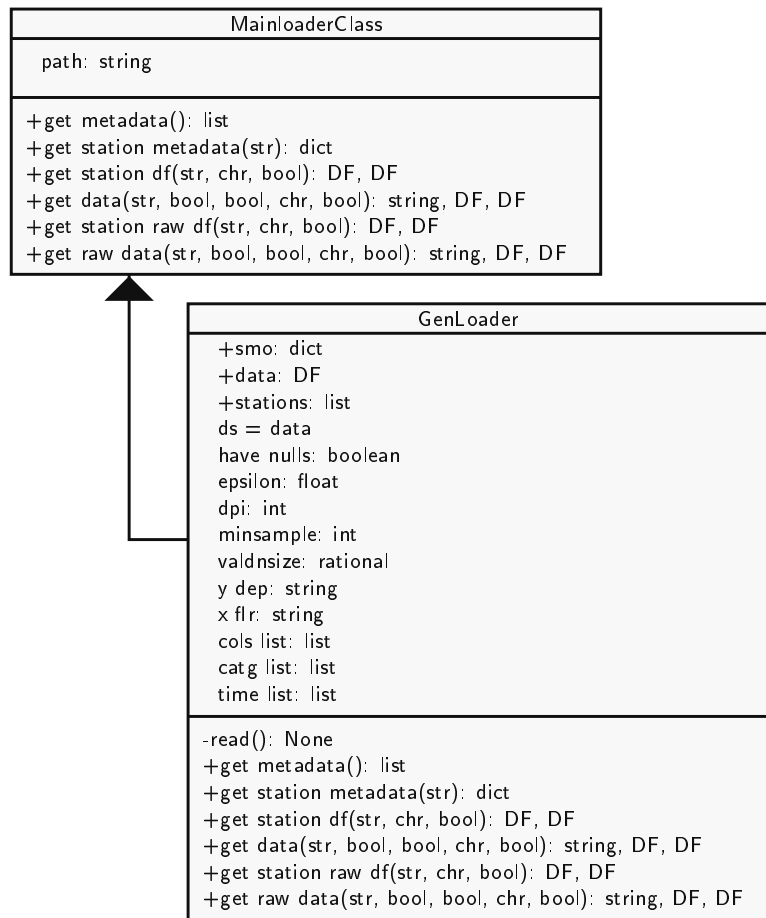
5.2.3 DONE Backend data structures

Programs adaptation in the new_api.

- Creating a new MainloaderClass and its derived GenLoader class as a generic data loader that extends MainloaderClass⁷. Methods' names remain unchanged to keep compatibility with the original code. **Created.**

⁶Tracking log in the appendix

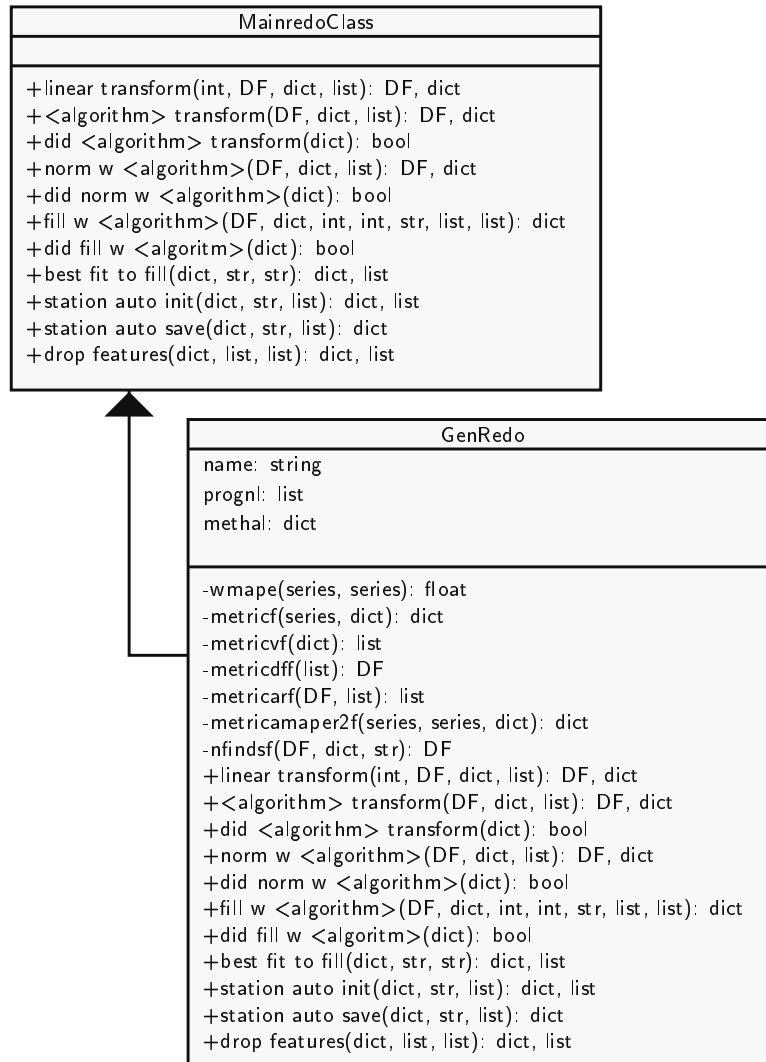
⁷The new Class defines the notebook's data structure



- The derived class **GenLoader** improves:
 - The differences in the data structures returned by some Data Loaders. Those differences crash the Python kernel. **Fixed.**
 - Some old Classes' methods trunk data to 1000 rows for interprocess communication. **Fixed.**
- Normalization will use **MaxAbsScaler**. Max scaler detected without considering negative values in some normalizations through the *main.py* in the new_api code. **Fixed.**
- The **GenLoader**-derived class allows different datasets to load. It returns a new loader object with the needed structure. It loads the data from the source when the constructor creates an instance of the class. The main program will map correctly the different datasets when the front end requires it. **Implemented.**
- The front end becomes slow with too much data. A resample by day solves the issue. Internally, it will keep the original data, and it will transform

the data without resampling. The resamples are just for visualization purposes. **Fixed.**

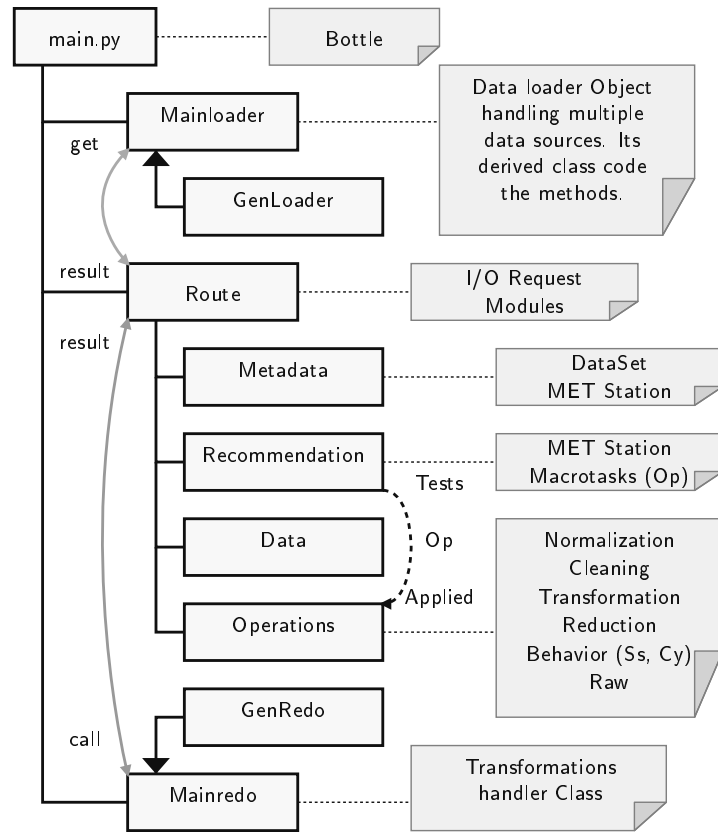
- Code refactoring is mandatory⁸. Some functions have no real input for data observations and use a random dataset created inside the function. Other functions presented incomplete treatment. Dimensionality reduction rewritten functions (2). **Implemented.**
- Duplicated and non-relevant functions have been removed. **Fixed.**
- The transformations will be packed in the MainredoClass and its derived class GenRedo, where the methods will be different techniques to clean and complete data. **Created.**



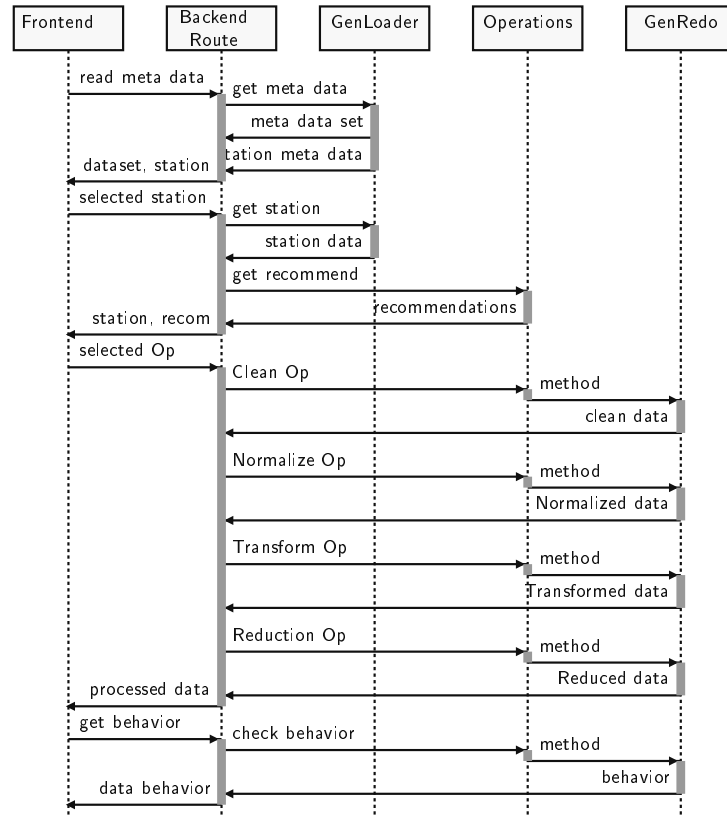
⁸Testing functions in the appendix

5.2.4 DONE Backend new specifications

1. DONE Classes and modules



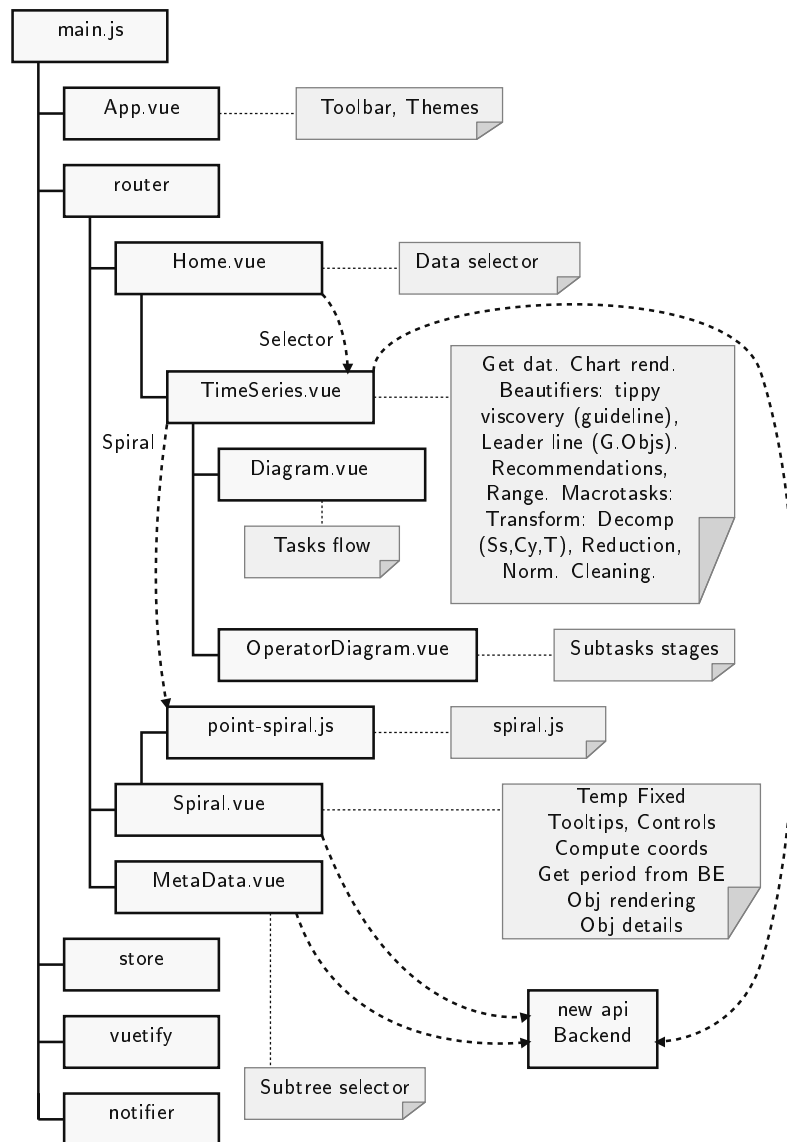
2. DONE Sequence diagram



5.2.5 DONE Backend new data structure and frontend integration

The initial state of the front end had two APIs. All the methods to grab data were migrated and transformed in the new API, leaving just one centralized processing Python API.

The main structure of the front end has not changed too much. The most relevant changes were in coding because the program was created for Peruvian temperature and precipitation datasets.



5.2.6 Data repository and cache

The backend will support dataset uploading. The structure has been defined on folders as an ID inside a .data repository. It will allow us to keep the data treatment progress in folders with the same ID inside a .cache repository.

6 DONE Requirements

6.1 DONE Software

- Python 3.12.7

- NodeJS 18.19.1
- GNU/Linux distribution with kernel 5.15.19 or superior
- Windows has not been tested, but it may work

6.2 DONE Hardware

- Processor AMD64 or x86₆₄ architecture
- 16GB RAM
- 32GB Swap
- 64GB SSD (128GB SSD Recommended)
- GPU (Optional)

7 DONE Visual Analytics Guidance Development

7.1 DONE Spiral diagram analysis improvement

- Compute segments and timespan for the Spiral. **Done.**
- Multiple time series integration. **Done.**
- Fix polygon coordinates when there are more than three dimensions. **Done.**

7.2 DONE Independent section for statistics

- Statistics of the time series (size, nulls, type of columns, or dimensions). **Done.**
- Distribution type of the time series (e.g., normal). **Done.**
- Outliers. **Done.**
- Correlation matrix. **Done.**

7.3 DONE Flow diagram improvement

- Data cleaning, Normalization and transformation, and Time-series behavior. **Done.**
- Stoppers and prereqs control between stages of the Guidance flow. **Done.**
- Network graph controller. **Done.**

8 DONE Format Multivariate Time Series Data

The program makes data preprocessing by inspecting data and giving recommendations through a hierarchical network graph navigation (3). The steps (5) related to this are:

- Load data. Data from multiple sources can include data synchronization techniques like aggregation, disaggregation, or data alignment in terms of date, creating additional missing values.
- Define appropriate data types and structures for time series.
- Inspect data for missing values, addressing them through interpolation or imputation methods (4).
- Stationarity testing with Augmented Dickey-Fuller test.
- Look for characteristics that change the model assumptions, like exponential growth, periodicities, or cyclicity patterns.
- Transformation like a logarithmic or square root to stabilize variance and reduce skewness.
- Data normalization by scaling data to enhance comparability.
- Decomposition into a trend, seasonality, and residuals to isolate cyclic behavior.

9 References

1. <https://github.com/flordeluz/VisWeb>
2. Mahmood Al-khassaweneh, Mark Bronakowski, Esraa Al-Sharoa (2023). Multivariate and Dimensionality-Reduction-Based Machine Learning Techniques for Tumor Classification of RNA-Seq Data. Engineering, Computing and Mathematical Sciences, Lewis University, Romeoville, USA. Computer Engineering Department, Yarmouk University, Jordan. Electrical Engineering Department, Jordan University of Science and Technology, Jordan.
3. Stefan Gladisch, Heidrun Schumann, Christian Tominski. Navigation Recommendations for Exploring Hierarchical Graphs. Institute for Computer Science, University of Rostock, Germany.
4. <https://ch.mathworks.com/help/econ/multivariate-time-series-data-structures.html>
5. Can Zhou, Masami Fujiwara, William E. Grant (2016). Finding regulation among seemingly unregulated populations: a practical framework for analyzing multivariate population time series for their interactions. Springer Science Business Media New York. https://www.researchgate.net/figure/Steps-in-analyzing-multivariate-time-series-data-a-After-plotting-each-indivi-fig2_288179902

10 Appendix

10.1 Time Estimation Plan

- 2 months +4 backup weeks

Table 2: Activities in weeks

Activities	1	2	3	4	5	6	7	8	9	10	11	12	13
Backend DS													
DS Integration													
Frontend DS													
New ML Backend models													
Frontend new ML models													
Frontend UI changes													

10.2 Testing functions

Testing functions are executed to give advice when the system is processing the signal. In the case of nulls, it uses WMAPE to prioritize the imputation algorithms, the rest are selected by the majority because each algorithm evaluates the same. When testing, more than 50% evaluate true and false otherwise. Trend algorithms do not detect, just use decomposition to break the signal in its main components.

Table 3: Testing function modules

Module	Functions
Nulls imputation	null_values_data(DF): bool fill_w_meanmedian(*dict) fill_w_decisiontree(*dict) fill_w_gradientboosting(*dict) fill_w_locallyweighted(*dict) fill_w_legendre(*dict) fill_w_randomforest(*dict) fill_w_kneighbors(*dict)
Cleaning	obtener_ruido_de(DF, int): bool obtener_ruido_cv(DF, float): bool obtener_outlier_iqr(DF): bool obtener_outlier_zscore(DF, int): bool obtener_outlier_grubbs(DF, float): bool
Normalization	obtener_no_patrones_estacionalidad(DF, int): bool obtener_distribucion_conocida(DF): bool
Transformation	obtener_no_estacionariedad_adf(DF, float): bool obtener_no_estacionariedad_kpss(DF, float): bool obtener_comportamiento_persistente_hurst(DF): bool
Reduction	verificar_correlacion_pearson(DF, float): bool verificar_correlacion_spearman(DF, float): bool verificar_correlacion_kendall(DF, float): bool

Continued on next page

Continued from previous page

Module	Functions
	check_multicollinearity(DF, int): bool
	check_dimensionality_reduction_pca(DF, float): bool
	check_dimensionality_reduction_fa(DF, int, float): bool
Decomposition	seasonality_detection(Series, Array): LOB
	trend_detection(Series, Array): LOB
	noise_detection(Series, Array): LOB

10.3 Dictionary

The dictionary is dynamic, it grows at any algorithm or action made. Data is mostly held on Dataframes. Classes can retain information on the object like testing functions described in the previous subsection. The function contains the signal.

The following is an excerpt from the dictionary. It repeats for each loader class.

```
loaders (dict): Class type, one per dataset.
|__ aqp: GenLoader
|__ brasil: GenLoader
|__ btc: GenLoader
|__ chiguata: GenLoader
|__ india: GenLoader
|__ madrid: GenLoader

aqp: Main dataset.
|__ RM (dict): Imputation algorithm applied.
|   |__ MAJES (dict): Group or substation.
|       |__ PPT: function
|       |__ TNM: function
|       |__ TXM: function
|__ cache (dict)
|__ full: DataFrame: Full group or substation-treated data.
|__ iqr: DataFrame: An entry of data for each executed algorithm.
|__ raw: DataFrame: Full group or substations unchanged data.
|__ redo (dict)
|   |__ fill (dict)
|       |__ decisiontree: bool
|       |__ gradientboosting: bool
|       |__ kneighbors: bool
|       |__ legendre: bool
|       |__ locallyweighted: bool
|       |__ meanmedian: bool
|       |__ randomforest: bool
|   |__ norm (dict)
|       |__ maxabs: bool
|       |__ minmax: bool
|       |__ robust: bool
|       |__ standard: bool
```

```

|__ outliers (dict)
|   |__ iqr: bool
|   |__ sdv: bool
|__ transform (dict)
|   |__ diff: bool
|   |__ linear: bool
|   |__ log: bool
|   |__ quadratic: bool
|   |__ sqrt: bool

```

10.4 Network coloring

Color attributes are applied with information that comes from the backend as arrays.

The coloring algorithm works as each node is gray by default, then each node name that comes in the subprocess array overwrites as red then adjacent nodes on the right are colored green, and left or parents are colored red. If an array exception shows up it replaces the red node with orange. The activities array carries information about when a dialog will be launched if the user picks the activity. The array made-path carries the executed actions and allows the algorithm to overwrite the colors of the nodes as blue, the coloring process follows the graph adjacency theory.

Trace:

Stage 1:

```

[ Algorithms priority ]:
Array [ "Rolling Mean" ]

```

```

[ Subprocesses ]:
Array [ "Clean", "Nulls" ]

```

Stage 2:

```

[ Made Path ]:
Array(3) [ "Clean", "Nulls", "Rolling Mean" ]

```

```

[ Subprocesses ]:
Array [ "Clean", "Outliers" ]

```

Stage 3:

```

[ Made Path ]:
Array(6) [ "Clean", "Nulls", "Rolling Mean", "Clean", "Outliers",
"Interquartile Range" ]

```

```

[ Subprocesses ]:
Array [ "DimRed" ]

```

```

[ Activities ]:
Array(3) [ "Multicollinearity Dim.Reduction=['PPT']", "PCA
Dim.Reduction", "FA Dim.Reduction=['PPT']" ]

```


Stage 4:

```
[ Made Path ]:  
Array(8) [ "Clean", "Nulls", "Rolling Mean", "Clean", "Outliers",  
"Interquartile Range", "DimRed", "Factor Analysis" ]  
  
[ Subprocesses ]:  
Array [ "Analysis" ]
```

10.5 Change Log

Table 4: List of patches by timestamp

A	Size	M	D	H	Patch
-rw-r--r--	918794	May	22	01:38	20240522-013818.diff
-rw-r--r--	32116	May	24	01:39	20240524-013931.diff
-rw-r--r--	55590	May	24	18:17	20240524-181739.diff
-rw-r--r--	40553	May	27	01:09	20240527-010932.diff
-rw-r--r--	64068	May	28	01:45	20240528-014458.diff
-rw-r--r--	160383	Jun	2	03:04	20240602-030410.diff
-rw-r--r--	110627	Jun	4	02:45	20240604-024500.diff
-rw-r--r--	63591	Jun	5	00:46	20240605-004621.diff
-rw-r--r--	68193	Jun	6	02:31	20240606-023122.diff
-rw-r--r--	33434	Jun	7	03:04	20240607-030453.diff
-rw-r--r--	91512	Jun	8	01:59	20240608-015902.diff
-rw-r--r--	76847	Jun	10	18:28	20240610-182855.diff
-rw-r--r--	78312	Jun	16	01:54	20240616-015442.diff
-rw-r--r--	24940	Jun	17	02:10	20240617-021030.diff
-rw-r--r--	73776	Jun	18	01:05	20240618-010542.diff
-rw-r--r--	27856	Jun	18	10:28	20240618-102832.diff
-rw-r--r--	138469	Jun	25	02:45	20240625-024532.diff
-rw-r--r--	65443	Jun	26	02:59	20240626-025901.diff
-rw-r--r--	56919	Jun	27	00:56	20240627-005611.diff
-rw-r--r--	24307	Jul	1	00:37	20240701-003738.diff
-rw-r--r--	30011	Jul	3	16:51	20240703-165116.diff
-rw-r--r--	91786	Jul	7	01:55	20240707-015505.diff
-rw-r--r--	46989	Jul	9	01:33	20240709-013313.diff
-rw-r--r--	96392	Jul	12	01:51	20240712-015153.diff
-rw-r--r--	231634	Jul	17	01:38	20240717-013822.diff
-rw-rw-r--	1025353	Jan	5	16:25	20250105-162546.diff
-rw-rw-r--	52880	Jan	20	12:08	20250120-120815.diff
-rw-rw-r--	108172	Jan	21	23:19	20250121-231951.diff
-rw-rw-r--	48118	Jan	29	01:22	20250129-012158.diff
-rw-rw-r--	12685	Jan	30	10:32	20250130-103221.diff
-rw-rw-r--	36912	Feb	2	00:08	20250202-000816.diff
-rw-rw-r--	115526	Feb	10	09:59	20250210-095953.diff

10.6 Backlog

Table 5: List of Backlog activities

Activity	Status
Some parts of the code are setting -1 to complete null values. It needs a review of the data with negative values like temperature. Resolution: Used for initial visualization purposes.	<i>Closed</i>
View action buttons have to be reviewed in the front end. There are buttons in the tree view requesting labels instead of the key data used to look for. Resolution: Remove action buttons on non-relevant leaf labels.	<i>Closed</i>
Integrate data and cache repositories with the front end. Resolution: Implemented.	<i>Closed</i>
Does get-raw-data() revert .ds to .smo["raw"] when invoked in the frontend? if so, uncomment the referred line in the get-raw-data(). Resolution: We are using different structures inside .smo.	<i>Closed</i>
For large datasets the prediction time of null values pre-evaluation goes like: decisiontree : too slow, > 2 mins kneighbors : too slow, > 2 mins gradientboosting: slow, > 1 min < 2 mins randomforest : medium, > 20 secs < 1 min meanmedian : regular, <= 20 secs locallyweighted : fast, <= 6 secs legendre : very fast, <= 2 secs Resolution: For large datasets predict with medium to fast algos. Author's note: While testing large dataset stations of Madrid's data, KNN gave better predictions. It is too risky to exclude slower algos in the early stage of nulls' pre-evaluation.	<i>Unsolved</i>