# pyWATTS: Python Workflow Automation Tool for Time Series

B. Heidrich<sup>a</sup>, A. Bartschat<sup>a</sup>, M. Turowski<sup>a</sup>, O. Neumann<sup>a</sup>, K. Phipps<sup>a</sup>, S. Meisenbacher<sup>a</sup>, K. Schmieder<sup>a</sup>, N. Ludwig<sup>ab</sup>, R. Mikut<sup>a</sup>, V. Hagenmeyer<sup>a</sup>

<sup>a</sup>Institute for Automation and Applied Informatics (IAI), Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany; <sup>b</sup>Cluster of Excellence "Machine Learning: New Perspectives for Science", University of Tübingen, Germany

#### ABSTRACT

Time series data are fundamental for a variety of applications, ranging from financial markets to energy systems. Due to their importance, the number and complexity of tools and methods used for time series analysis is constantly increasing. However, due to unclear APIs and a lack of documentation, researchers struggle to integrate them into their research projects and replicate results. Additionally, in time series analysis there exist many repetitive tasks, which are often re-implemented for each project, unnecessarily costing time. To solve these problems we present pyWATTS, an open-source Python-based package that is a non-sequential workflow automation tool for the analysis of time series data. pyWATTS includes modules with clearly defined interfaces to enable seamless integration of new or existing methods, subpipelining to easily reproduce repetitive tasks, load and save functionality to simply replicate results, and native support for key Python machine learning libraries such as scikit-learn, PyTorch, and Keras.

#### **KEYWORDS**

Time Series Analysis; Python; Workflow Automation; Machine Learning; Pipeline

#### 1. Introduction

In many areas, time series data are the most prominent form of data collected. In contrast to other sequential data such as speech data, time series data are not only ordered, but the time stamp associated with the observation might also have explicit information. For example, looking at energy time series, the demand at a specific time step depends on calendar-based information such as the day of the week or the season. Generally, time series analysis uses various algorithms from statistics to deep learning to answer questions about time-dependent systems.

Although more and more code from researchers focusing on time-dependent data is publicly available, there is still a need for respective tools. These tools should allow automating the workflow in time series analysis and an easy integration of new research approaches with third-party code. Automating the workflow is necessary, since many preprocessing tasks are repetitive, such as accounting for seasonality, adding calendar-based features, or detecting and imputing missing values. As a result of the lacking tools, researchers often re-implement these repetitive tasks at the unnecessary expense of time. Moreover, it is challenging to integrate new or alternative approaches into existing code workflows and, although the push towards open science increases the importance of reproducibility, it is often difficult to replicate earlier experimental results. Thus, any tool to aid researchers in automated time series analysis needs to

focus on two features: re-usability and reproducibility of existing and new code.

Several factors currently prevent good integration and re-usability of publicly available code for time series analysis. For example, most authors only publish their proposed new algorithm or method, excluding any steps necessary to prepare the data. Using their code then entails re-writing the required preprocessing method. Additionally, interfaces are hardly ever defined and basic unit-testing is often non existent, which regularly leads to re-implementation being the only quick and attainable solution. Regarding reproducibility, some of the issues include platform-dependent code, no information on parameter settings, or insufficient description on the order in which function or scripts need to be executed, making it almost impossible to reproduce results.

A remedy to the issues mentioned is workflow automation using pipelines and modules. In a pipeline, one can define the workflow, i. e. the exact order in which several modules, each including a method, are run to achieve a specific result. No matter if we wish to use or reproduce code with pipelines, the steps needed to reach a specific result can be non-sequential. For example, we might wish to run parts of the code in parallel, have branching and merging pathways in the workflow, or even condition-dependent paths. Furthermore, the modules have a clearly defined structure which allows simple integration of new or alternative methods into an existing workflow.

Current Python tools which allow the realisation of pipelines are, for example, scikit-pipeline<sup>1</sup> [8] and river<sup>2</sup> [6]. While scikit-pipeline is part of the package scikit-learn [8], river is a merger of CremeML and Sk-multiflow. However, both tools only allow linear execution of modules, where neither parallel nor conditional execution is possible. Only the package baikal<sup>3</sup> provides non-sequential pipelines inherited from scikit-learn. It is based on wrappers for scikit-learn modules, where each module has to be wrapped individually. Therefore, it is rather tedious. Furthermore, it does not allow integration of other libraries such as PyTorch [7] and Keras [1], which are useful for deep learning-based time series approaches. Additionally, baikal aims to combine several scikit-learn modules such that they work as one module and thus focuses on the model creation only.

In the present paper, we introduce pyWATTS, an open-source Python-based package that provides a non-sequential workflow automation tool for the analysis of time series data. In contrast to baikal, pyWATTS includes generic wrappers for libraries such as scikit-learn, PyTorch, and Keras, allows the pipelines to have conditions, and is able to visualise intermediate results. Summarising the key features, pyWATTS

- is a platform-independent solution to implement workflows from start to finish using pipelines. Thereby, time series experiments can be performed in an organised manner and in any environment that supports pyWATTS.
- enables re-usability through subpipelining. Any useful part of a time series experiment, e.g. preprocessing, can be defined as a subpipeline and integrated into other pipelines without further adaption and independently of the original experiment.
- allows saving and loading of any given pipeline configuration to reproduce results at a later date.
- enables simple integration of new research approaches through a plug-and-play style environment where modules implemented in pyWATTS can be exchanged

<sup>&</sup>lt;sup>1</sup>https://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html

<sup>&</sup>lt;sup>2</sup>https://github.com/online-ml/river

<sup>&</sup>lt;sup>3</sup>https://github.com/alegonz/baikal/

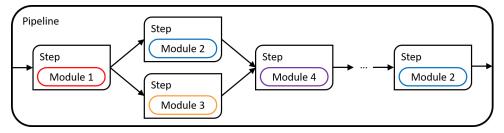


Figure 1.: pyWATTS uses the three classes *pipeline*, *step*, and *module* to realise non-sequential workflows.

seamlessly between pipelines through a modular architecture with data handling through xarray [4].

- includes a clear API of the modules, i.e. transform and fit methods, ensuring that pipelines within pyWATTS are adaptable and that modules can easily run on multiple data sets, at different points in a pipeline, and in various pipelines.
- allows using different modules for the same part in the pipeline such that a condition mechanism decides which module is executed depending on the characteristics of the applied data.
- takes pandas DataFrame[5, 9] or xarray Dataset as input, allowing users to flexibly read the data from any source (file, database, website) with their method of choice.
- is able to use callbacks, e.g. for visualising, analysing, and writing the intermediate results of modules.

To the best of our knowledge, pyWATTS is the first tool to automate time series analysis workflows using non-sequential pipelines in this form. The remainder of the paper is structured as follows. We first introduce the implementation and architecture of pyWATTS before providing an overview on the availability and re-usability. We conclude with an outlook on ongoing and future projects that use pyWATTS.

#### 2. Implementation and architecture

Implementing the features mentioned above requires a careful design of the architecture. In this section, we, therefore, describe pyWATTS' architecture and implementation.

The pyWATTS package is written in the programming language Python<sup>4</sup>. To realise non-sequential workflows, it uses three classes as illustrated in Figure 1; a *pipeline*, representing the workflow as a graph, a *step* which represents a node in the graph referencing its dependencies with edges, and a *module*, representing the algorithm running in a step. We introduce these three classes in more detail in the following.

Every algorithm used in pyWATTS is implemented with the module class. pyWATTS distinguishes between algorithms requiring training and algorithms that can be applied without training. For both types, the module class's transform method must be implemented to apply the algorithm. For algorithms requiring training, we additionally have to implement the module class's fit method that defines the training of the machine learning model. Modules must also include methods to save and load all information necessary for executing the module.

In general, the implementation of modules follows the concepts introduced by

<sup>4</sup>https://www.python.org/

Table 1.: The library of pyWATTS contains several utile algorithms when dealing with time series.

Module name	Description
Calendar extraction	Extracts or extends a time-series with calendar information such as weekdays or holidays
Change Direction	Extracts if the change is positive or negative for each time point in a time series
Clock Shift	Shifts the data with a certain offset
Differentiate	Calculates the n-th order difference of a time series
Linear Interpolator	Creates a linear interpolation
Missing Value Detector	Detects missing values such as "NaN"
Resampler	Reduces or increases the temporal resolution of a given time series
Rolling Mean	Calculates a rolling mean over a specific window size
RMSE Calculator	Calculates the Root Mean Squared Error (RMSE)
Sampler	Creates samples with a specified sample size
Trend Extraction	Extracts a trend specified by a period and a length
Sklearn Wrapper	Wraps machine learning modules from the scikit-learn library
Keras Wrapper	Wraps deep learning neural networks implemented in the ${\tt Keras}$ library
PyTorch Wrapper	Wraps deep learning neural networks implemented in the ${\tt PyTorch}$ library

scikit-learn [8]. pyWATTS itself provides a comprehensive library of algorithms, as listed in Table 1. The currently available modules implement utile algorithms for time series analytics and serve as a guideline to implement further modules. pyWATTS also provides special modules, called wrappers, to seamlessly integrate existing algorithms and models from scikit-learn [8] or deep learning models implemented in Keras [1], or PyTorch [7],

Given a module, pyWATTS creates one or multiple steps<sup>5</sup>. The step class organises the execution of the pipeline. A step collects and merges the results of its dependencies, calls the fit and transform method of its module, and provides its output to the pipeline for the subsequent steps. Moreover, a step can execute callbacks defined by the user, e.g. for visualising, analysing, and writing the module's intermediate result. Furthermore, a step controls the execution of the module based on conditions defined by the user.

The pipeline class organises the steps in nodes and creates a graph, where every step input is represented as an edge. This graph supports branching and merging of paths and is used to define the execution order of the steps. This way, all previous steps represented as dependencies have to be successfully executed before the current step itself is executed. The pipeline also serves as the interface to the user and provides control commands. These commands include training and executing a pipeline, as well as saving and loading the whole pipeline.

Based on the mentioned three classes, pyWATTS implements the following three functionalities for an easy structuring and flexible application of pipelines:

<sup>&</sup>lt;sup>5</sup>A step contains zero or one module. A module can be used in multiple steps.

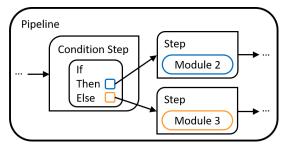


Figure 2.: pyWATTS uses condition steps with "if-then-else" for conditional branching in pipelines.

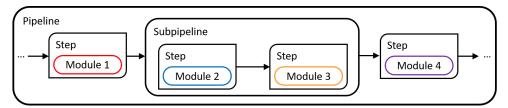


Figure 3.: pyWATTS makes use of subpipelines to easily structure and name parts of the workflow.

**Batch/online learning:** By specifying that the pipeline processes only one time step at a time, the pipeline can be executed iteratively.

Conditional branching: Depending on the applied data, condition steps with "ifthen-else" can be used to select different paths of the pipeline for execution (see Figure 2).

**Subpipelines:** Grouping steps of the pipeline in subpipelines allows an easy structuring and naming of certain parts of the workflow (see Figure 3).

#### 3. Quality control

For quality control, we apply comprehensive testing to pyWATTS, use programming guidelines and code reviews, as well as provide up-to-date documentation and examples. For automated testing, we use GitHub Actions. All core classes implementing the main functionality of the pipeline and the steps are tested with unit tests. Furthermore, unit tests cover the wrappers and the modules of the library. To ensure the correct interaction of steps and modules as well as saving and loading of whole workflows, exemplary pipelines are implemented as integration tests.

Furthermore, we follow programming guidelines to ensure high-quality source code. These guidelines define the naming of branches and code conventions as well as prescribe the use of linting software such as pylint<sup>6</sup>. Additionally, the guidelines require developers to implement tests for each module and to use loggers with appropriate logging messages. Finally, the guidelines demand developers to use type annotations for arguments, variables, and return values. In the GitHub Actions, we automatically check whether code conventions are met using flake8<sup>7</sup>. To ensure compliance with all guidelines, we additionally perform manual code reviews on pull requests. Maintainers

<sup>6</sup>https://www.pylint.org/

<sup>&</sup>lt;sup>7</sup>https://gitlab.com/pycqa/flake8

review pull requests concerning their correct operation, coverage through tests, and compliance to programming guidelines and code conventions before the pull requests are merged into the master branch.

Lastly, we maintain an up-to-date documentation. Based on the annotated source code and restructured text files, the documentation<sup>8</sup> of pyWATTS is automatically generated using sphinx<sup>9</sup> and readthedoc<sup>10</sup>.

Besides serving as integration tests, the provided examples introduce new users to pyWATTS and its features and support them in creating working pipelines in pyWATTS. In the following, we briefly describe the provided examples, which are detailed in the documentation.

- To prevent fundamental errors during the creation of a pipeline, a simple example explains how one can create a pipeline for electrical load forecasting and how one can add modules such as the Calendar Extraction to the pipeline.
- To test the functionality of the condition mechanism, we provide an example that changes the method for electrical load forecasting depending on day-time and night-time.
- Advanced examples aim to avoid mistakes in the application of deep learning frameworks in pyWATTS. In the examples using Keras [1] or PyTorch [7], the pipelines train simple deep learning models.

### 4. Availability

### Operating system

Platform independent

#### Programming language

Python

#### Additional system requirements

pyWATTS is designed to perform various time series analysis tasks on data sets of arbitrary size. Therefore, hardware requirements depend on the size of the data set and the task being performed.

## Dependencies

The core pyWATTS dependencies are the following:

- $\bullet$  scikit-learn -0.23.2
- cloudpickle 1.6.0
- xarray 0.16.1
- numpy 1.19.2
- pandas 1.1.5

<sup>&</sup>lt;sup>8</sup>The pyWATTS Documentation is available at https://pywatts.readthedocs.io/en/latest/.

<sup>9</sup>https://www.sphinx-doc.org/en/master/

<sup>10</sup>https://readthedocs.org/.

- matplotlib 3.3.2
- tensorflow 2.3.1
- $\bullet$  workalendar 12.0.0

Dependencies required for development purposes comprise the following:

- pytest 6.1.1
- sphinx 3.2.1
- pylint 2.6.0
- pytest-cov 2.10.1

### Software location:

#### Archive

Name: Zenodo

Persistent identifier: https://doi.org/10.5281/zenodo.4637197

Licence: MIT Licence<sup>11</sup>
Publisher: Zenodo

Version published: 0.1.0 Date published: 25.03.2021 Code repository GitHub

Name: pyWATTS

Persistent identifier: https://github.com/KIT-IAI/pyWATTS

Licence: MIT Licence

**Date published:** 25/09/2020

#### Language

English

### 5. Reuse potential

Due to the architecture and the modular structure of pyWATTS, anyone who wishes to analyse time series can use pyWATTS out-of-the-box. It enables the users to easily select the modules and determine the pipeline structure relevant for their specific use case, such as forecasting. Additionally, the possibility to save and load pipelines together with the platform-independence of pyWATTS, allows easy reproduction of research results. Moreover, common Python-based machine learning libraries can be used within pyWATTS. For example, we provide wrapper modules for scikit-learn [8], Keras [1], and PyTorch [7] to allow the inclusion of the available functions.

Moreover, pyWATTS' users are supported by comprehensive documentation for its core structure and the individual modules as well as detailed examples. In case of questions, the core developer team can also be contacted with the help of GitHub issues or the pyWATTS contact email address. The generous MIT license<sup>11</sup> allows research, commercial and non-commercial use, and development of the package as either an anonymous user, private developer or publicly contributing developer. All users can stick to the existing modules and pipelines, extend them based on known or unknown

<sup>&</sup>lt;sup>11</sup>https://opensource.org/licenses/MIT

issues, or create new modules and pipelines. Whether any changes to the modules are made locally or through the public repository is up to the user to decide.

The developer team, for example, wants to use pyWATTS in various research applications in the future. For preprocessing, we plan to extend pyWATTS with the Copy Paste Imputation of missing values for energy time series as described in [10]. We also plan to use pyWATTS for time series forecasting, e.g. by using Profile Neural Networks [3]. Furthermore, we intend to extend pyWATTS for the insertion of typical anomalies in energy time series to have data sets with ground truth for anomaly detection and anomaly handling. To generate realistic synthetic energy time series, we also aim to use pyWATTS. An interface for pipeline tuning and selection will further assist in automating the iterative design process. Lastly, we want to deploy pyWATTS as an execution environment in the research infrastructure Energy Lab 2.0 [2].

Putting it all in a nutshell, pyWATTS provides an extendable framework for automating time series analysis workflows. It uses comprehensible pipelines and is able to integrate established statistical, machine learning, and deep learning frameworks. Thus, pyWATTS makes it easy to develop, adapt, and reproduce pipeline-based experiments for energy time series analysis.

#### Acknowledgements

We thank Simon Waczowicz for the valuable input on the concept of pyWATTS.

### Funding statement

This project is funded by the Helmholtz Association's Initiative and Networking Fund through Helmholtz AI, the Helmholtz Association under the Program "Energy System Design", the Joint Initiative "Energy System Design - A Contribution of the Research Field Energy", the Helmholtz Metadata Collaboration, and the German Research Foundation (DFG) as part of the Research Training Group 2153 "Energy Status Data: Informatics Methods for its Collection, Analysis and Exploitation" and under Germany's Excellence Strategy – EXC number 2064/1 – Project number 390727645.

## References

- [1] François Chollet et al. Keras. https://keras.io. 2015.
- [2] Veit Hagenmeyer, Hüseyin Kemal Çakmak, Clemens Düpmeier, Timm Faulwasser, Jörg Isele, Hubert B. Keller, Peter Kohlhepp, Uwe Kühnapfel, Uwe Stucky, Simon Waczowicz, and Ralf Mikut. "Information and Communication Technology in Energy Lab 2.0: Smart Energies System Simulation and Control Center with an Open-Street-Map-Based Power Flow Simulation Example". In: Energy Technology 4.1 (2016), pp. 145–162. DOI: 10.1002/ente.201500304.
- [3] Benedikt Heidrich, Marian Turowski, Nicole Ludwig, Ralf Mikut, and Veit Hagenmeyer. "Forecasting Energy Time Series with Profile Neural Networks". In: Proceedings of the Eleventh ACM International Conference on Future Energy Systems. Association for Computing Machinery, 2020, pp. 220–230. DOI: 10.1145/3396851.3397683.
- [4] Stephan Hoyer and Joseph J. Hamman. "xarray: N-D labeled Arrays and Datasets in Python". In: *Journal of Open Research Software* 5 (2017). DOI: 10.5334/jors.148.
- [5] Wes McKinney. "Data Structures for Statistical Computing in Python". In: Proceedings of the 9th Python in Science Conference. Ed. by Stéfan van der Walt and Jarrod Millman. 2010, pp. 56–61. DOI: 10.25080/Majora-92bf1922-00a.
- [6] Jacob Montiel, Max Halford, Saulo Martiello Mastelini, Geoffrey Bolmier, Raphael Sourty, Robin Vaysse, Adil Zouitine, Heitor Murilo Gomes, Jesse Read, Talel Abdessalem, and Albert Bifet. "River: machine learning for streaming data in Python". In: arXiv:2012.04740 (2020).
- [7] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. "PyTorch: An Imperative Style, High-Performance Deep Learning Library". In: Advances in Neural Information Processing Systems. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc., 2019, pp. 8026–8037.
- [8] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. "Scikitlearn: Machine Learning in Python". In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [9] The pandas development team. pandas-dev/pandas: Pandas. Version latest. Feb. 2020. DOI: 10.5281/zenodo.3509134.
- [10] Moritz Weber, Marian Turowski, Hüseyin Kemal Çakmak, Ralf Mikut, Uwe Kühnapfel, and Veit Hagenmeyer. "Data-Driven Copy-Paste Imputation for Energy Time Series". In: arXiv:2101.01423 (2021).