

The manuscript describes the new python package AI4Water, intended as a modelling tool for hydrological predictions. It incorporates the basic steps of data-driven analysis and modelling - preprocessing, choosing one or several modelling approaches, post-processing including error analysis and visualization - and makes extensive use of existing python libraries. The focus is strongly on machine learning approaches and the package embraces many of the currently discussed and newly developed algorithms.

The paper is easy to read and, while surely leaving many details out and thus not a manual for the ambitious user, summarizes the fundamental steps in the modelling process in a concise manner. The authors do not develop their own routines or approaches, but the collection of state-of-the-art modelling utilities and approaches is impressive.

Response: We thank the reviewer for their valuable comments and suggestions. Your suggestions have substantially improved the quality of the manuscript. Our point-by-point responses to the reviewers' comments are provided below. We have updated the code and archived its latest release at Zenodo (<https://zenodo.org/record/5595680>). We have updated the documentation of the code that is available at <https://ai4water.readthedocs.io/en/latest/>. The code to reproduce figures shown in the manuscript is given in the “examples/paper” folder in the code repository.

it would be nice if the authors could address three major issues relevant for anybody intending to analyze their own data:

Response: We have revised the manuscript. Accordingly, the code of *AI4Water* has been assigned to it so that others can analyze their own data using this package. Please find the detailed responses below each comment.

1. It is not obvious how users might get their own data (time series) into the package. The access to two existing databases is implemented, CAMELS and LamaH, and the authors rightly remark that the different data formats, conventions etc. are an obstacle slowing down the analysis process. How generic are the input options to accomodate own data in different formats (text of Excel files, spatially extended time series in netcdf files, and the like)?

Response: *AI4Water* contains a *DataHandler* class that pre-processes the input data and prepares the training, validation, and test data. This class can read data from various files as long as the data are in the correct format in those files. *AI4Water* is designed for tabular data; therefore, *DataHandler* expects the data to be in tabular form in a given file. For example, in a csv or excel file, if the data are arranged in a tabular form, one column represents one input or output feature, and each row indicates one example. In this case, the *DataHandler* class reads data from the given file. However, the user must specify the names of the input and output features, which must correspond to the names of columns in the files. Internally, *DataHandler* reads the input file and converts it into a *pandas*' *DataFrame* object. A *DataFrame* object is a data model of *pandas* for tabular data (McKinney, 2011). *DataHandler* can read tabular data from the most commonly used file systems, such as a comma separated file (.csv), Microsoft Excel (.xlsx), network common data form (netCDF), feather, parquet (Vohra, 2016), npz, and mat files. If the given file is in netCDF format, it is read using the *xarray* package (Hoyer and Hamman, 2017) and then converted into

pandas *DataFrame*. We have added an ipython notebook named “input_data_file_types.ipynb” to demonstrate this. This notebook shows how users can bring data from .csv, .xlsx, and netcdf files into the *model* using *DataHandler*. The *DataHandler* can save the processed data into an HDF5 file, which can be used by the user for inspection. The processed data comprise training, validation, and test data. In reference to this, we have added the following lines to the manuscript.

Lines 228 - 233: The *DataHandler* class prepares the input data for the machine learning model and acts as an intermediate between the *Model* class and other preprocessing classes such as *Imputation* and *Transformation* classes. The *DataHandler* can read data from various files as long as the data are in a tabular format in those files. The complete list of allowed file types and their accepted file extensions is provided in Table S5. Internally, the *DataHandler* class stores data as a *pandas DataFrame* object, which is a data model of pandas for tabular data (Mckinney, 2011). *DataHandler* can also save processed data as an HDF5 file, which can be used to inspect the processed input data.

Table S5. File types and their extensions accepted by *AI4Water*.

File extension	File type
.csv	Comma separated file
.xlsx	Microsoft Excel
.npz	Numpy zipped file
.parquet	Parquet
.feather	Feather
.nc	netCDF5
.mat	MATLAB

2. Expandibility: it might well be that for the specific data at hand or the particular user, other methods than the ones already provided might be desirable. An example would be gap-filling, but also others. The part of the manuscript describing that (chapter 3) is very vague and general, please be more specific.

Response: We agree with the reviewer that the ability to customize a certain functionality will be advantageous to the user. Therefore, we used the object-oriented programming (OOP) paradigm for writing the code of *AI4Water*. This paradigm makes it easier to expand or enhance a specific functionality of *AI4Water*, such as customizing the training loop and loss function, or adding an extra pre-processing step to the input data before feeding it to the model. We have also provided examples of codes that customize the loss function, training step, and training loop. These notebooks are available under the examples/paper folder in the code repository. We have added these details in the manuscript in the following lines.

Lines 381 - 386: *AI4Water* was built using the object-oriented programming (OOP) paradigm. Its core logic was implemented by the *Model* class. The use of OOP allows a user to customize any steps of model building, training, or testing by sub-classing the *Model* class. This may include the implementation of a custom training loop or a customized loss function. Similarly, the pre-processing and data preparation steps implemented in the *Model* class can also be overwritten for specific usages. For example, if users want to implement another transformation on the training data, they can subclass the *Model* class and overwrite the “*training_data*” method. Similarly, the user can customize the loss function by overwriting the “*loss*” method of *Model* class.

3. Interpretation: it would be wonderful if the package could produce a comprehensive interpretation of the results achieved with the chosen model approaches. Interpretation also

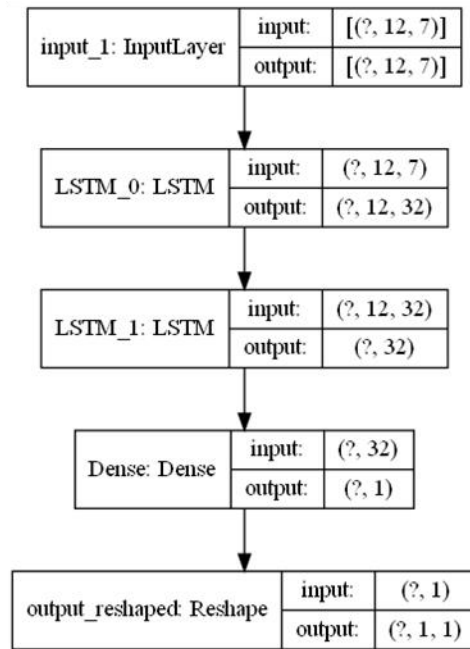
implies making connections to existing hydrological knowledge (process understanding) as well as local conditions (metadata) available for the site, its peculiarities. However, in this context, interpretation is merely a visualization of the model architecture (e.g. the weights in the case of NNs). A more modest phrasing, e.g. "Model Visualization" instead of "Interpret" as the class name, seems to be more appropriate.

Response: We have added a separate sub-module named “*visualize*” to view the model. This sub-module exposes a class named *Visualize* to the user. This class plots the decision tree learned by the tree-based machine learning model. For neural network-based deep learning models, this class can plot the outputs of intermediate layers, weights, gradients of layer outputs, and gradients of weights. This sub-module is separate from the “*Interpret*” and “*Explain*” sub-modules. The *Interpret* sub-module is used to interpret the behavior of attention-based deep learning models, such as DA-LSTM (Qin et al., 2017) or temporal fusion transformers (Lim et al., 2021). The purpose of “*Explain*” sub-module is to explain the output of the machine learning model by considering it as black-box. This sub-module comprises two classes: *ShapExplainer* and *LimeExplainer*, which explain the model using the SHAP (Lundberg and Lee, 2017) and LIME (Ribeiro et al., 2016) methods, respectively. All of these sub-modules are part of the postprocessing sub-module of *AI4water*. We have added figures (S2–S5) using the *Visualize* sub-module of an LSTM model, which show the outputs and weights of LSTM along with the gradients of LSTM outputs and gradients of weights of LSTM. We have added separate sub-sections about interpretability (2.10.1) and visualization (2.10.2) in the manuscript.

Lines 298 - 305: The “*visualize*” sub-module, consisting of a *Visualize* class, is used to examine inside the machine learning model. When the model comprises several layers of neural networks,

this class plots the outputs of the intermediate layers, gradients of these outputs, weights and biases of intermediate layers, and gradients of these weights. Thus, this class helps to visualize the working of neural networks and can be used to plot the decision tree learned by the tree-based machine learning model. We demonstrate the use of this class by building a four-layer neural network to predict streamflow using the CAMELS dataset (Fowler et al., 2021). The four-layered neural network comprises an input layer, two layers of LSTM, and one output layer (Fig. S1). Figures S2–S5 show the outputs of the first LSTM layer and its gradients along with the weights of the first LSTM layer, and the gradients of those weights.

a)



b)

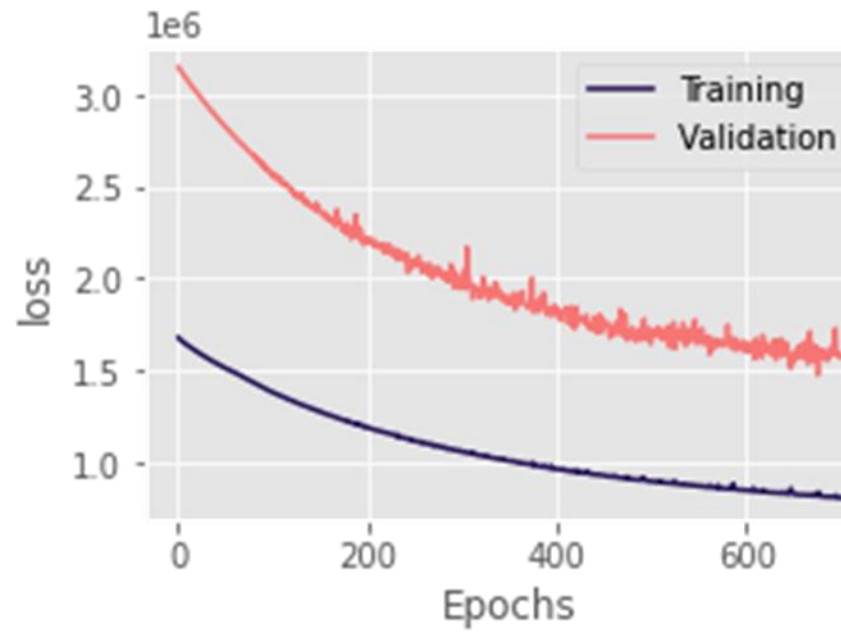


Figure S1: Architecture of a four-layer neural network used for prediction of streamflow at catchment number 224206 of CAMELS-AUS. Seven climate variables were used and 12 days of historical data was used for training the model. A) The model consisted of 2 LSTM layers followed by a Dense layer as output layer. The output was finally reshaped into 3d array. B) Training and validation loss curves during model training. The model was trained for 700 epochs.

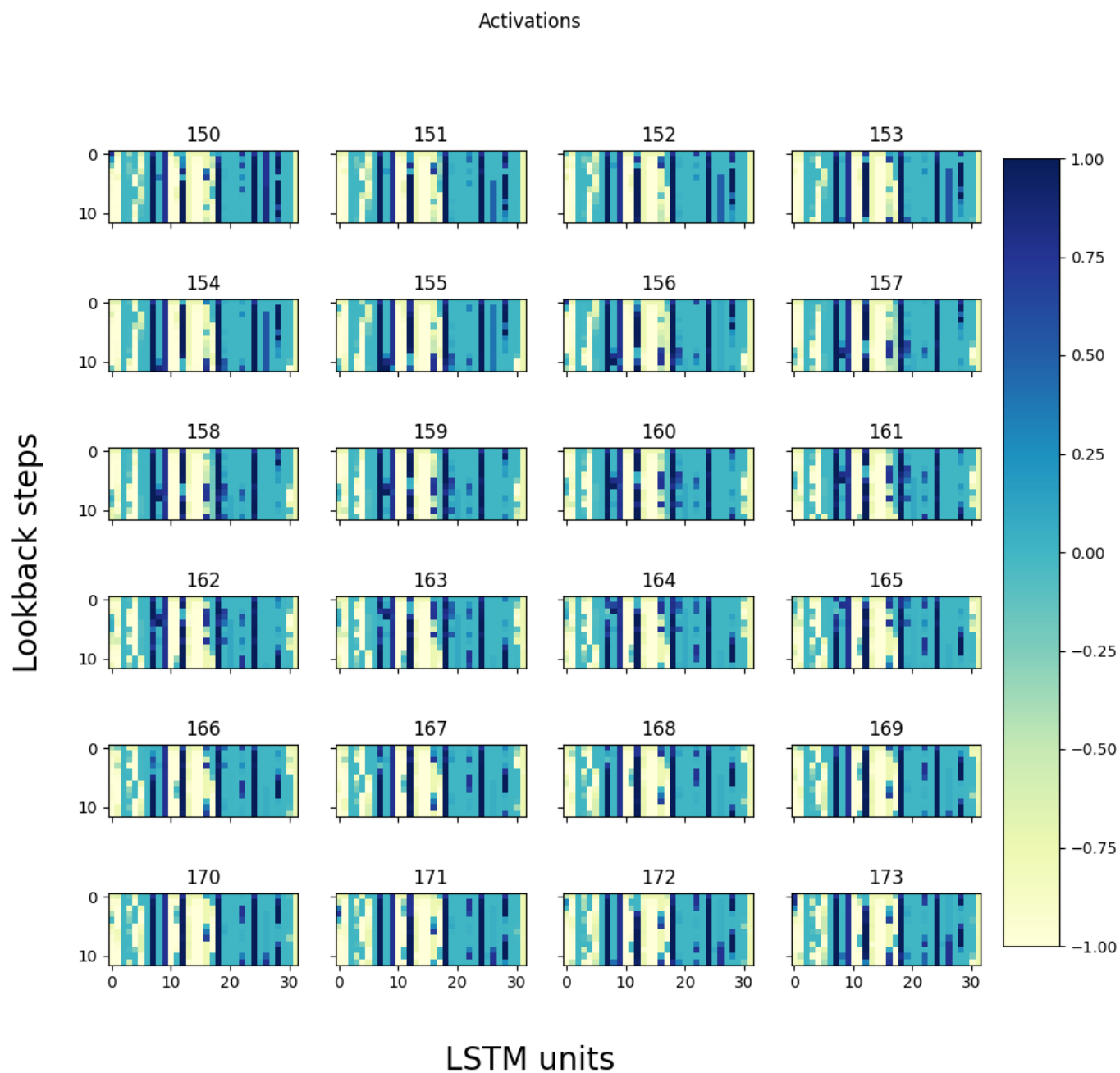


Figure S2: Output of first LSTM for 24 days. The model consisted of two LSTM layers with 32 units for each LSTM. The lookback steps indicate the number of historical days used by the model to predict value for next day. The titles for each subplot indicate Julian day for the year 2000.

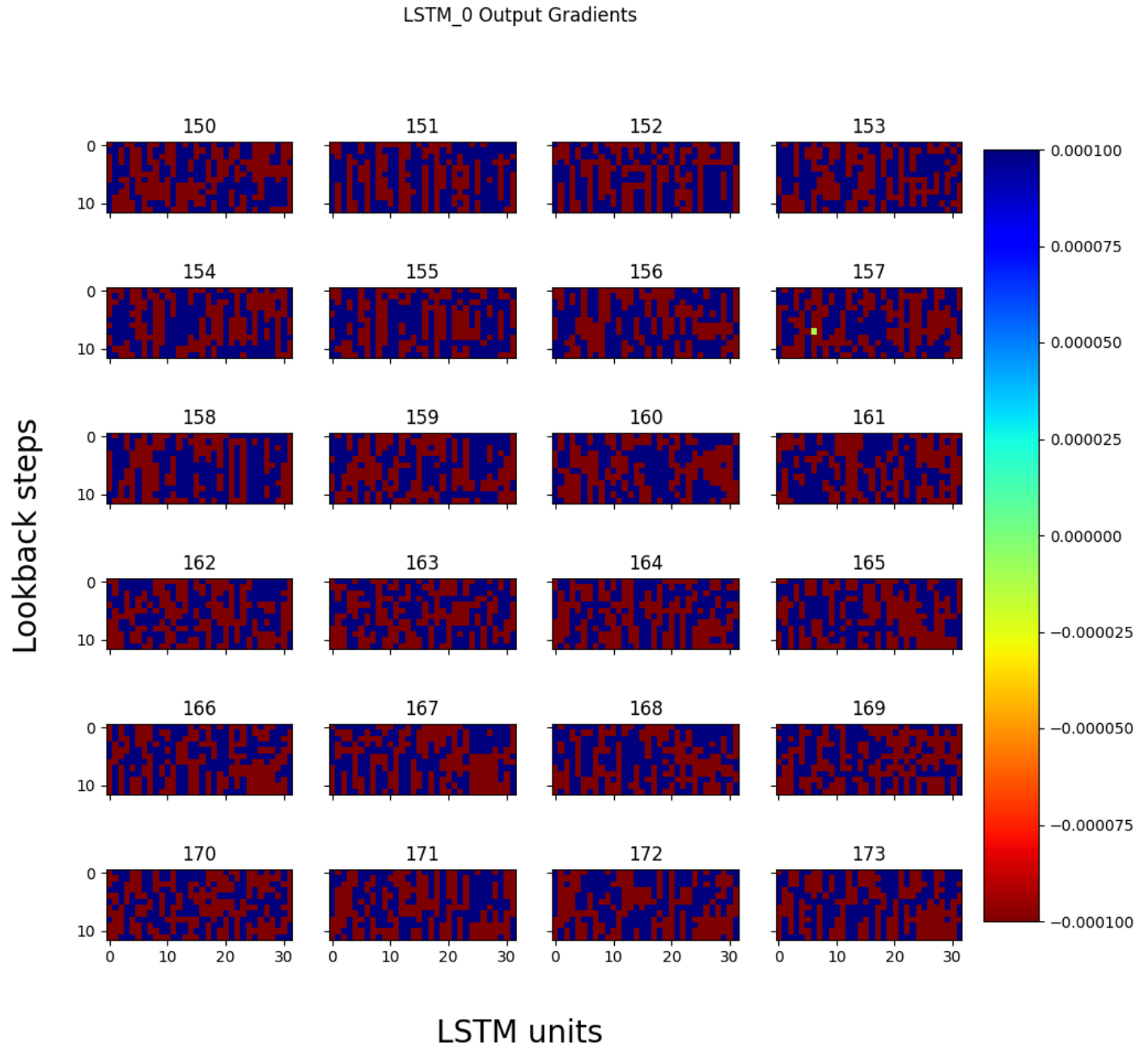


Figure S3: Gradients of outputs of first LSTM for 24 days. The model consisted of two LSTM layers with 32 units for each LSTM. The lookback steps indicate the number of historical days used by the model to predict value for next day. The titles for each subplot indicate Julian day for the year 2000.

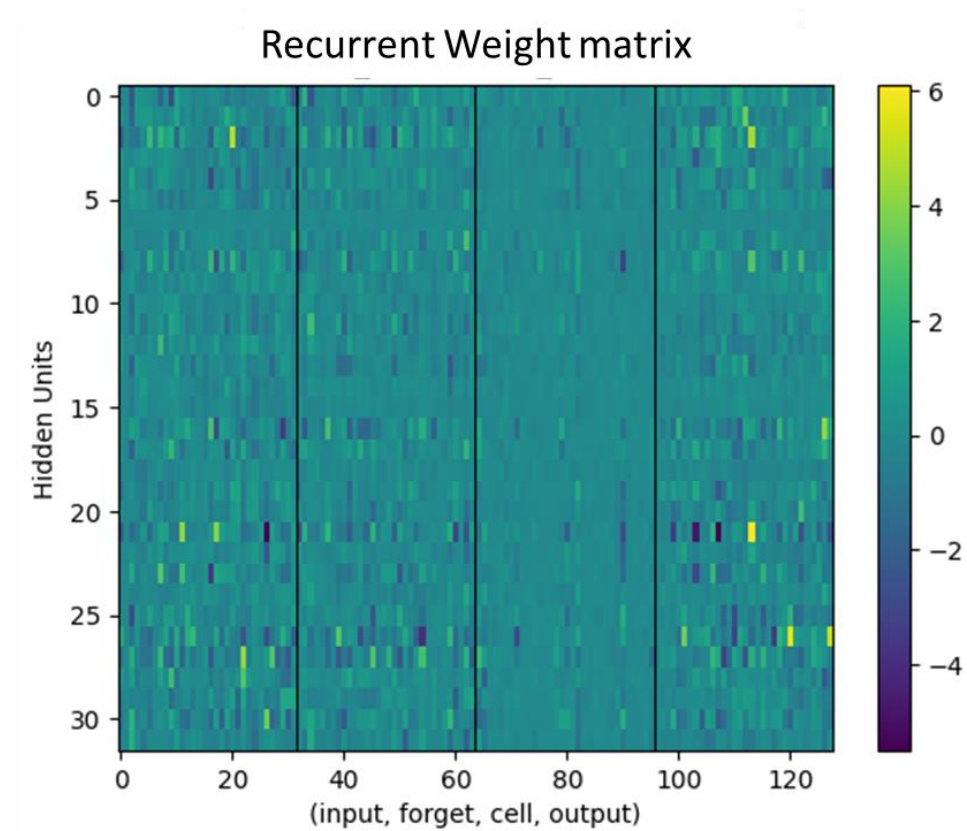
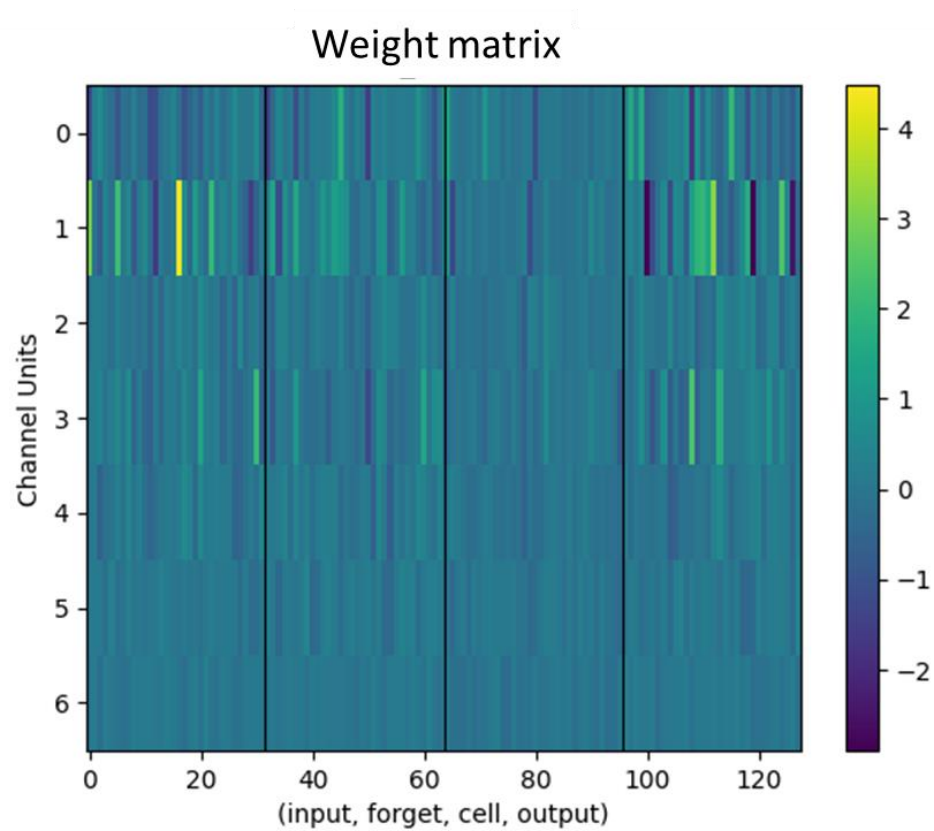


Figure S4: Weight matrices of LSTM layer. The LSTM layer consists of two weight matrices. The portion of weight matrices responsible for input gate, forget gate, output gate and cell state are highlighted by lack lines.

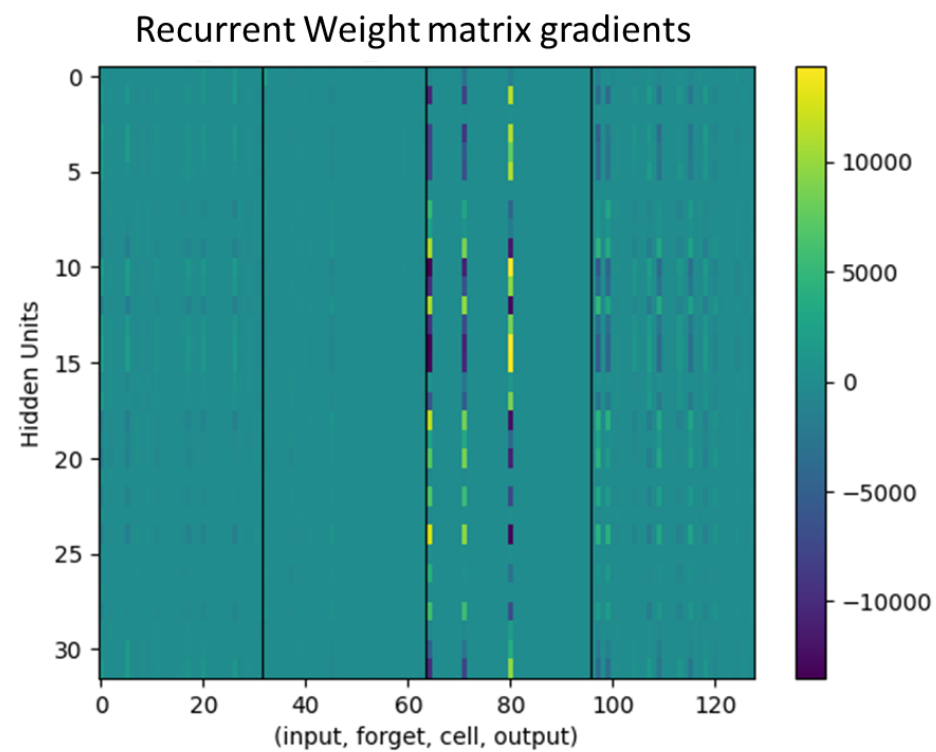
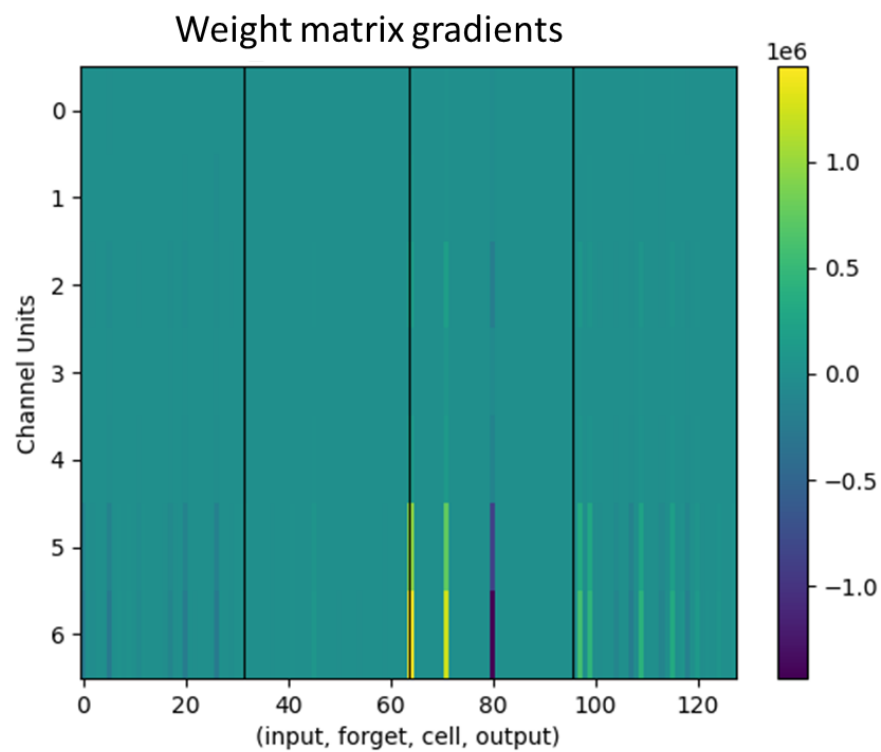


Figure S5: Gradients of weight matrices of LSTM layer. The LSTM layer consists of two weight matrices. The portion of weight matrices responsible for input gate, forget gate, output gate and cell state are highlighted by lack lines.

Lines 327 – 355: Several model-agnostic methods have recently been developed to explain black-box machine learning models, such as local independent model explanations (LIME) (Ribeiro et al., 2016) and Shapely Additive Explanations (SHAP) (Lundberg and Lee, 2017). These methods explain the behavior of complex machine-learning models (such as black-box) using a simplified but interpretable model. However, using these methods in high-stake decision-making has been criticized (Rudin, 2019). The explanations of these methods can be local or global. A local explanation describes the behavior of the model for a single example, whereas a global explanation can describe the model's behavior for all examples. The LIME method is only relevant for local explanations, whereas SHAP also provides explanations for approximating the global importance of a feature. *AI4Water* consists of *LimeExplainer* and *ShapExplainer* classes to explain its behavior using the LIME and SHAP methods.

We built an XGBoost (Chen and Guestrin, 2016) model for the prediction of *E. coli* in a Laotian catchment (Boithias et al., 2021). Fig. S10 shows the output of the *LimeExplainer* class, whereas Fig. S11 shows the output of the *ShapExplainer* class. In Fig. S10, a large horizontal bar for a given feature indicate that this feature strongly affected the model's prediction. A positive value indicate that the given feature caused increase in model's prediction. On the other hand, the negative value indicate that it caused decrease in model's prediction. Thus, large negative value for solar radiation in example 41 indicate that the solar radiation causes large reduction in model's prediction. Large positive values for water level in examples 42 to 46 indicate that the water level in these cases strongly increased model's prediction. The numerical values of features along y-axis indicate which value of feature was responsible for the aforementioned behaviour. Thus, more precisely, the water level above 147.8 causes very large increase in model's prediction. Therefore,

we can verify that the *E. coli* prediction during flood events are more strongly affected by water level.

The SHAP module provides more detailed explanation about local as well as global importance of input features on model's prediction. Fig. S11a and Fig. S11b show the local explanation summary of model in the form of SHAP of each input feature for each example (Lundberg et al., 2020). Fig. S11a shows that the examples with large SHAP values of water level and suspended matter resulted in large *E. coli* prediction. The $f(x)$ in Fig. S11a indicate model's prediction. The examples in Fig. S11a are clustered in such a way that examples with similar explanations are grouped together. Fig S11b indicate that the large values of water level and suspended particulate matter results in increase in *E. coli*. On the other, large values of solar radiation resulted in negative SHAP values. This shows that large solar radiation causes reduction in *E. coli* prediction. Fig S11c shows the global importance of input features for *E. coli* prediction. This global importance is obtained by calculating mean of SHAP value of a feature for all examples (Lundberg and Lee, 2016). The explanations from Fig. S11 correlate with our background understanding of *E. coli* behavior. Several studies have shown that *E. coli* in surface water is strongly affected by suspended solids, water level and solar radiations (Nakhle et al., 2021; Pandey and Soupir, 2013).

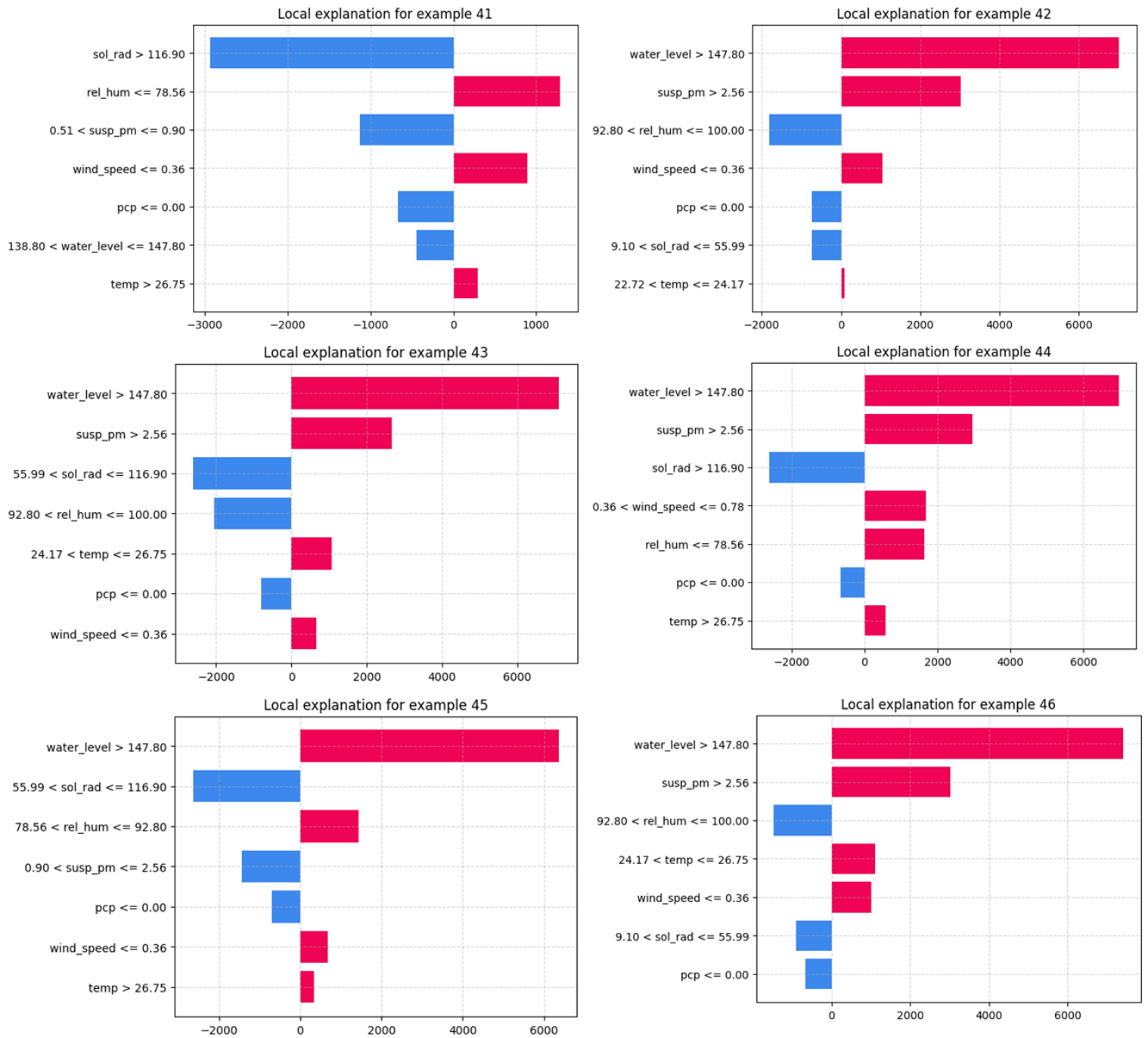


Figure S10. Explanation of XGBoost model for *E. coli* prediction using LIME method for six selected examples from test data. The explanations show the importance for each input feature by the model.

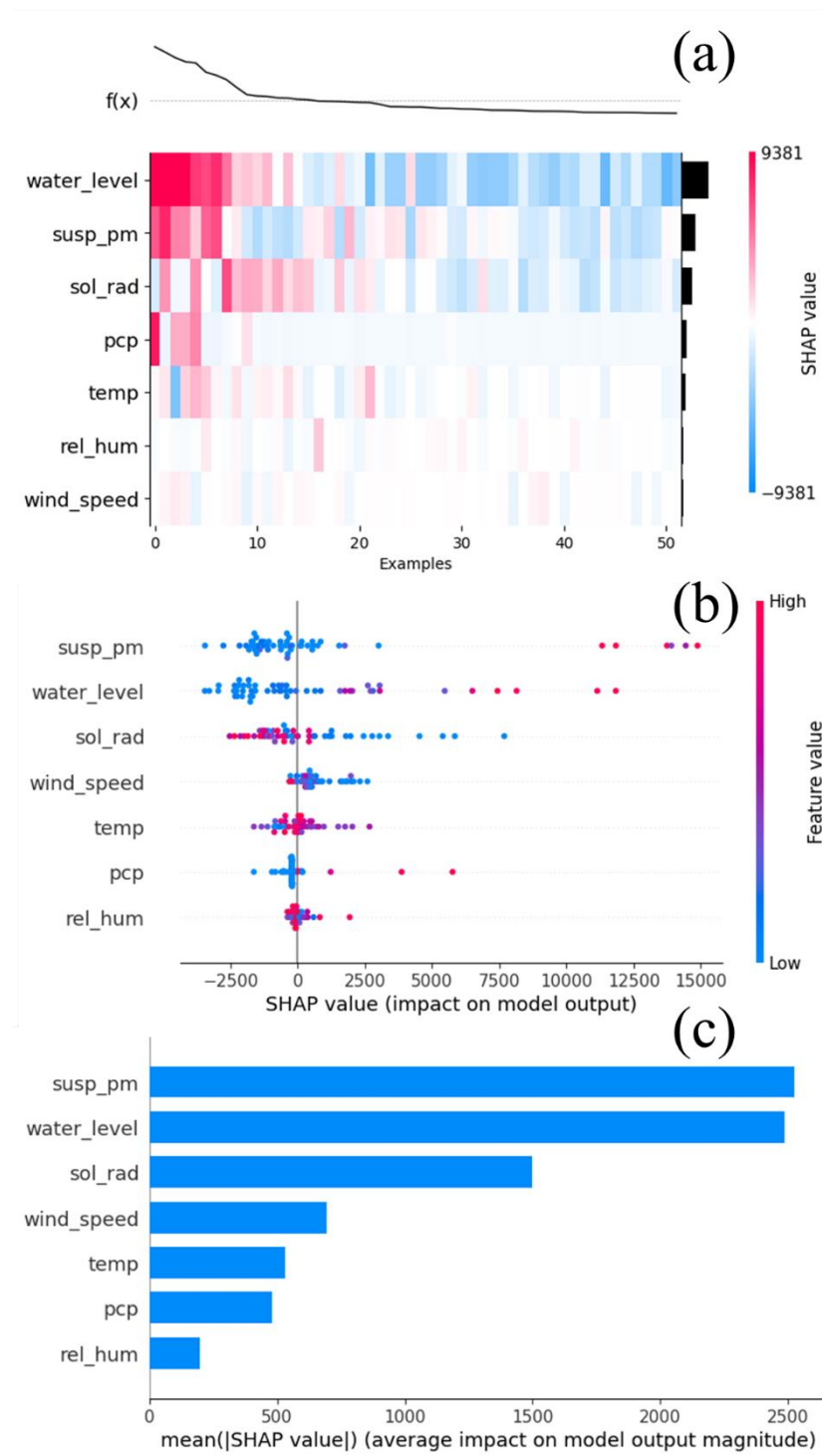


Figure S11. Explanation of XGBoost model for *E. coli* prediction using SHAP method. The explanations show the importance for each input feature by the model. (a) SHAP values as

heatmap (b) SHAP values for individual examples in test data, (c) global feature importance based upon SHAP values.

The language quality is good to very good with very few typos etc. Some specific comments and corrections:

Response: Please find the responses to the specific comments attached below.

l. 78: "time series errors": do you rather mean performance measures rather than errors?

Response: Yes, we have replaced the word, "errors" with "performance metrics"

Lines 78-79: The *SeqMetrics* sub-module calculates several time-series performance metrics for regression and classification problems.

l. 112: "Fig. 3 shows examples of the three configuration files" -> "Fig. 3 shows three examples for configuration files"

Response: We have corrected the sentence.

Lines 375 - 375: Fig. 4 shows three examples of configuration files.

l. 118: "obtain large and diverse data" - no, this cannot be guaranteed, and the hope is that modelling is also possible when there is only a limited amount of data from a given catchment, as is often the case!

Response: As suggested by the reviewer, we have removed the term “large and diverse” from the sentence.

Line 138: The first step in building a data-driven hydrological model is to obtain the data.

l. 139: what is the difference between "scaling" and "transforming the data onto a different scale" ?

Response: We agree that the terms “scaling” and “transforming the data onto a different scale” are similar. Therefore, we have removed the word “scaling” from this sentence.

Line 159: Data transformation includes standardizing and transforming the data onto a different scale.

l. 142: EMD is a decomposition, not a transformation method, much like PCA. Of course, using IMFs as input rather than the original variables does change the model setup and has an impact on performance etc. as is correctly stated further down.

Response: We have modified the sentence to highlight that EMD is a decomposition method.

Lines 161-163: Additionally, several decomposition methods such as empirical mode transformation (EMD), ensemble EMD (EEMD), wavelet transform (Sang, 2013), and fast

Fourier transform (Sang et al., 2009) were found to improve the performance of hydrological models.

l. 146 "were" -> "are"

Response: We have replaced the word “were” with “are”.

Line 162: The predicted features are transformed back after the prediction.

l. 153 "(McKinney, 2011) scikit" -> "(McKinney, 2011), 2) scikit"

Response: We have corrected the sentence.

Lines 172 - 173: These include using either the 1) *pandas* library (McKinney, 2011), 2) scikit-learn library-based methods, or 3) dedicated algorithms to fill the missing input data.

ch. 2.4 Missing labels: it should be mentioned that this refers to a classification task only, not to regression.

Response: We apologize for the confusion in this section. The absence of target data is common in regression tasks. *AI4Water* can handle these situations for both regression and classification problems.

Also, what is the difference between "exclude examples" (l. 170) and "skip these examples" (l. 173)?

Response: The words “exclude” and “skip” mean the same in these lines. Thus, we now use only the term “exclude”.

Lines 192-194: However, the user can also opt to exclude these examples, although this can reduce the number of examples in water quality problems where the number of samples is already very small.

l. 179 "later" -> "latter"

Response: We have replaced “later” with “latter.”

Line 199: The latter can be achieved by setting the “*input_steps*” argument to a value >1.

l. 198 "time series weather data" -> "time series of weather data"

Response: We have replaced “time-series weather data” with “time-series of weather data.”

Line 217-218: *AI4Water* contains a sub-module *MakeHRUs*, which helps in distributing the time-series of weather data into HRUs using different HRU definitions.

l. 205 how does the user provide HRUs / soil types, land use classes etc. ? Through shapefiles if available?

Response: Yes, the module requires shapefiles of soil types, land use classes, slope, and sub-basins to make the HRUs according to a given definition. We have also specified this in the manuscript.

Line 225 - 226: The *MakeHRUs* sub-module requires shapefiles of land use, soil and slope to make the HRU according to a given definition.

l. 212 "large" -> "many"

Response: We have replaced the word “large” with “many.”

Line 238 - 239: These include complex methods such as Penman–Monteith (Allen et al., 1998), which require many input variables.

l. 272 "all possible results" -> "many different results"

Response: We have corrected the sentence.

Line 311-312: The *Interpret* class takes the trained model of *AI4Water* as input and plots numerous results, which help to explain the behavior of the model.

l. 294 "cannot be defined" - why not?

Response: We have mentioned that scale-independent error metrics cannot be defined for some cases. This is true for percentage errors, such as mean absolute percentage error (MAPE), where one or more values in the observed array can be equal to zero. In such cases, the MAPE calculation yields infinity as result. For cases where one or more values are close to zero, the calculated MAPE

values are extremely skewed. This has been emphasized in the literature, such as Hyndman (2006) and Prestwich et al. (2014). We have elaborated this in the manuscript as well.

Line 362-363: However, certain scale-independent error metrics cannot be defined when one or more observed values are zero, such as percentage errors or relative errors (Hyndman, 2006).

l. 335 delete the first occurrence of "training" in this line

Response: We have deleted the first occurrence of ‘training.’

Lines 415 - 417: Modeling hydrological processes by machine learning requires the development of pipelines that encompasses data retrieval, feature extraction, visualization, building, training, and testing the machine learning model, along with visualization and interpretation of its results

l. 346 Christine, 2014 does not seem to be in the reference list

Response: We have corrected this and added a reference for MKDocs in the reference list.

Lines 425-426: The user manual is built into the source code *Docstring* and compiled into a “read the docs” web page (<https://ai4water.readthedocs.io/en/latest/>) using the MKDocs (Christie, 2014) software.

If these comments are taken into account by the authors, the paper should be published by GMD.

References

Boithias, L., Auda, Y., Audry, S., Bricquet, J. p., Chanhphengxay, A., Chaplot, V., de Rouw, A., Henry des Tureaux, T., Huon, S., and Janeau, J. I.: The Multiscale TROPical CatchmentS critical zone observatory M-TROPICS dataset II: land use, hydrology and sediment production monitoring in Houay Pano, northern Lao PDR, *Hydrological Processes*, 35, e14126, 2021.

Chen, T. and Guestrin, C.: Xgboost: A scalable tree boosting system, *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, 785-794, <https://doi.org/10.1145/2939672.2939785>,

MkDocs. Project documentation with Markdown.: <https://www.mkdocs.org/>, last

Fowler, K. J., Acharya, S. C., Addor, N., Chou, C., and Peel, M. C.: CAMELS-AUS: Hydrometeorological time series and landscape attributes for 222 catchments in Australia, *Earth System Science Data*, 13, 3847-3867, <https://doi.org/10.5194/essd-13-3847-2021>, 2021.

Hoyer, S. and Hamman, J.: xarray: ND labeled arrays and datasets in Python, *Journal of Open Research Software*, 5, 2017.

Hyndman, R. J.: Another look at forecast-accuracy metrics for intermittent demand, *Foresight: The International Journal of Applied Forecasting*, 4, 43-46, 2006.

Lim, B., Arık, S. Ö., Loeff, N., and Pfister, T.: Temporal fusion transformers for interpretable multi-horizon time series forecasting, *International Journal of Forecasting*, 2021.

Lundberg, S. and Lee, S.-I.: An unexpected unity among methods for interpreting model predictions, *arXiv preprint arXiv:1611.07478*, 2016.

Lundberg, S. M. and Lee, S.-I.: A unified approach to interpreting model predictions, *Proceedings of the 31st international conference on neural information processing systems*, 4768-4777,

Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., and Lee, S.-I.: From local explanations to global understanding with explainable AI for trees, *Nature machine intelligence*, 2, 56-67, <https://doi.org/10.1038/s42256-019-0138-9>, 2020.

McKinney, W.: pandas: a foundational Python library for data analysis and statistics, *Python for high performance and scientific computing*, 14, 1-9, 2011.

Nakhle, P., Ribolzi, O., Boithias, L., Rattanavong, S., Auda, Y., Sayavong, S., Zimmermann, R., Souleleuth, B., Pando, A., and Thammahacksa, C.: Effects of hydrological regime and land use on in-stream *Escherichia coli* concentration in the Mekong basin, Lao PDR, *Scientific reports*, 11, 1-17, 2021.

Pandey, P. K. and Soupir, M. L.: Assessing the impacts of *E. coli* laden streambed sediment on *E. coli* loads over a range of flows and sediment characteristics, *JAWRA Journal of the American Water Resources Association*, 49, 1261-1269, <https://doi.org/10.1038/s41598-017-12853-y>, 2013.

Prestwich, S., Rossi, R., Armagan Tarim, S., and Hnich, B.: Mean-based error measures for intermittent demand forecasting, *International Journal of Production Research*, 52, 6782-6791, 2014.

Qin, Y., Song, D., Chen, H., Cheng, W., Jiang, G., and Cottrell, G.: A dual-stage attention-based recurrent neural network for time series prediction, *arXiv preprint arXiv:1704.02971*, 2017.

Ribeiro, M. T., Singh, S., and Guestrin, C.: " Why should i trust you?" Explaining the predictions of any classifier, *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135-1144,

Rudin, C.: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead, *Nature Machine Intelligence*, 1, 206-215, <https://doi.org/10.1038/s42256-019-0048-x>, 2019.

Sang, Y.-F.: A review on the applications of wavelet transform in hydrology time series analysis, *Atmospheric research*, 122, 8-15, 2013.

Sang, Y.-F., Wang, D., Wu, J.-C., Zhu, Q.-P., and Wang, L.: The relation between periods' identification and noises in hydrologic series data, *Journal of Hydrology*, 368, 165-177, 2009.

Vohra, D.: Apache parquet, in: *Practical Hadoop Ecosystem*, Springer, 325-335, 2016.

