# University of Stavanger

## FACULTY OF SCIENCE AND TECHNOLOGY

# MASTER'S THESIS

| Study programme/specialization:<br><br>MSc. Computational Engineering | Spring semester, 2022<br><br>Open |
|---|---|
| Author:<br><br>Ali Tahir | ............................................<br>Ali Tahir |

| Programme coordinator:<br>Øystein Arild<br><br>Supervisor:<br>UiS - Prof. Dan Sui |
|---|

| Title of master's thesis:<br>IMPACT OF DATA PRE-PROCESSING TECHNIQUES ON MACHINE LEARNING MODELS |
|---|

| Credits: 30 |
|---|

| Keywords:<br>Machine Learning, Recurrent neural network, K-Nearest Neighbors Algorithm, Artificial Neural Network. | Number of pages: 87<br><br>+ supplemental material/other: 50<br><br><br>Stavanger, 15th July 2022 |
|---|---|

Ali Tahir

# IMPACT OF DATA PRE-PROCESSING TECHNIQUES ON MACHINE LEARNING MODELS

# Abstract

The Volve dataset, which contains the time series values of different sensors that have been used at the Volve drilling site contains many flaws which make it hard for machine learning models to learn from the dataset and provide useful insights and future predictions. Three flaws have been highlighted including missing data, different frequency rates, and too many attributes (high dimensional data). To solve the issues, present in time series data, a data preprocessing pipeline has been proposed which first removes the noise through the rolling mean. Then applies gap analysis to remove the columns whose gaps can not be filled with data imputation methods. After that gap has been filled by the KNN imputer which imputes the missing values in the data. After that data resampling has been applied to make the sampling rate consistent as the time series prediction model takes a constant sampling rate. For hyper-parameter tuning of the resampling method AIC and BIC value has been created on a grid of hyper-parameters. After resampling, top parameters were selected on basis of Pearson correlation, after which AIC and BIC have been used to select the most relevant 3 parameters. These 3 parameters has then be used to train three models that are: RNN + MLP, LSTM + MLP, and LSTM + RNN + MLP. On basis of mean absolute error (MAE) best model has been selected which is RNN + MLP.

# Acknowledgments

I, first of all, would like to start by being grateful to God Almighty, for giving me all the strength required in the completion of my studies.

This thesis study for the completion owes lot of gratitude and appreciation to number of individuals for their selfless contribution, especially my supervisor Prof. Dan Sui for the guidance and supervising this research by arranging weekly meetings, productive discussion and sharing knowledge.

I would like to show my deepest gratitude to University of Stavanger for giving me platform and all the professors who were part of this journey.

Finally, I would like to thank my family and peers for their complete support.

Ali Tahir

# List of Abbreviations

| | |
|---|---|
| **ANN** | Artificial Neural Network |
| **BHA** | Bottom Hole Assembly |
| **CSV** | Comma-Separated Value |
| **DWOB** | Downhole Weight on Bit |
| **DTQ** | Downhole Torque |
| **ECD** | Equivalen Circulating Density |
| **HPHT** | High Pressure High Temperature |
| **ID** | Inside Diameter |
| **IQR** | Interquartile Range |
| **KNN** | K Nearest Neighbors |
| **LWD** | Logging While Drilling |
| **LSTM** | Long Short Term Memory |
| **MAE** | Mean Absolute Error |
| **ML** | Machine Learning |
| **MLP** | Multi Layer Perceptron |
| **MSE** | Mechanical Specific Energy |
| **MWD** | Measuring While Drilling |
| **NN** | Neural Network |
| **OD** | Outside Diameter |
| **RF** | Random Forest |
| **RNN** | Recurrent Neural Network |
| **ROP** | Rate of Penetration |
| **RPM** | Revolutions per minute |
| **SPP** | Standpipe Pressure |
| **STQ** | Surface Torque |
| **SWOB** | Surface Weight on Bit |
| **WOB** | Weight on Bit |
| **WDP** | Wired Drill Pipe |

# Table of Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1 Data Pre-Processing and Machine Learning

Cleaning of dataset is most of the basic step in pre-processing for a machine learning pipeline. Pre-processing of data include [3]:

- Data cleaning

    - Missing values

    - Noisy Data

- Data integration

    - Data consolidation

    - Data propagation

    - Data virtualization

- Data reduction i.e. dimensionality reduction, feature selection

- Data transformation i.e. feature extraction

Data reduction and Data Transformation are the steps of feature engineering but these two techniques of pre-processing and feature engineering intersect with each other and are thus considered as the same step.

1

The next step in the pipeline formulation is selecting an ML model that needs to be trained on the pre-processed dataset. Selecting the perfect model is also necessary as there are different models that serve the same purpose i.e., classification can be done through Random Forest, Decision tree, and xgBoost algorithms but a specific model can be selected for a specific type of application with a literature study.

In pipeline formation, another thing that needs to be considered is hyper-parameter tuning. Hyper-parameter tuning is to tune parametric values such as learning rate, batch size, and the number of clusters. For example, in the case of the random forest algorithm, a few important parameters to be tuned are maximum depth, the minimum number of leaves, the number of estimators, and minimum number of sample splits, etc. In the case of deep learning models, hyper-parameters can be the number of iterations, learning rate, momentum, and optimizer. There are ranges of values that can be set in hyper-parameters and optimal values can help improve the model significantly.

The final step is evaluating the model, evaluation step tells us how our model performed and applies remedies in case of poor performance. Evaluation methods can tell us where our model is performing poorly i.e., the confusion matrix can tell us that for which class our model is performing poorly (if that is the case).

In petroleum data pre-processing is a need rather than a choice. There are several reasons for this statement. As it can be seen from Figure 1.1 there are several faults in the data. This is the same well, inclination vs measured depth plot. Red is from the time-based log, blue is from depth based log. There are many things wrong here:

- Even though identical variables have been plotted, a shift by 28m was given to the depth for the wells to match.

- 1 – the area where depth-based plot is just missing data

- 2 – depth-based data exists, but at a low frequency

- 3 – time-based data is noisy (can see two distinct paths).  Additionally, depth stops to match

- 4 – discontinuity on time-based data

**Figure 1.1:** Data faults and depth v/s inclination graph

- 5 – artifacts at the start of the well, both datasets.

- Time-based data has "stalactites" hanging from the well path that needs to be filtered out (only thing easy to automate, short median filter)

For using this petroleum data in a time series prediction machine learning pipeline this data needs to be cleaned and processed.

## 1.2 Literature Review

### 1.2.1 Missing Data Imputation

Zhu et al. [4] considered data pre-processing a crucial task before model training and gave a review of different steps in data pre-processing. In this review, missing data imputation was also covered. These are the methods stated in the review for missing value imputation:

- Mean substitution

- Hot-deck substitution

- Regression substitution

- Conditional distribution-based substitution

- Multiple imputation

Zhang et al. [5] give an analysis that existing methods for replacing missing data entries use some deterministic or random imputation methods. Keeping this in mind they proposed clustering-based random imputation (CRI), which makes clusters of values having missing data and values that are complete. Then impute values on basis of the cluster which is near to missing data point. Their method proved effective in missing values imputation tasks.

Somasundaram and Nedunchezhian [6] stated that in real-world datasets missing entries, noise and inconsistencies are a common occurrence and there is a need to correct these problems. They applied three different methods of data imputations (Constant substitution, Mean attribute value substitution, and Random attribute value substitution). They compared the performance of each method using clustering methods. The data set used by them for evaluation was Wisconsin diagnostic breast cancer (WDBC) dataset [**?**].

Huang et al. [7] predicted the cost of an unseen project on basis of existing projects. They analyzed that feature selection affected their model performance significantly. To improve their pre-processing stage of data they applied a three-stage technique based on data imputation, data normalization, and feature selection. Their experiments showed a significant increase in accuracy when Z-Score normalization, kNN imputation, and mutual information-based feature weighting were used.

Xu et al. [8] observed that the existing method of value imputation shows poor performance when the number of missing values increases. They proposed a missing value imputation algorithm based on the evidence chain (MIAEC), that estimates missing data by first mining relevant evidence of missing data and then combining them to build an evidence chain. To further speed up the process map-reduce method was used, not only this the algorithm was made to run on distributed systems. Their method showed elevated performance than existing methods based on naive Bayes and mode imputation methods.

Nelwamondo and Marwala [9] gave a review of missing data imputation methods based on computational intelligence techniques. They highlighted that in literature there is a gap in the use of computational intelligence methods for missing data imputation. The researchers have

mostly used easy methods because computational intelligence methods are complex. They highlighted three patterns of missing data that are, univariate pattern, monotone pattern, and arbitrary pattern. Every missing data pattern requires a different technique for data imputation.

Rja [10] and Rubin [11] showed three types of missing data patterns that affect the choice of the method. These data patterns are missing at random (missing data is unrelated to other missing variables but related to some observable variables), missing completely at random (missing data is unrelated to both other missing and observable variables), and missing not at random (missing data is related to other missing values).

Rimal [12] worked on proving that missing data size matters when applying the data imputation method. In this research, the R programming language was used to carry out the simulations. This research also highlights that it is not possible to find a method of data imputation before analyzing the dataset first.

An extensive study was done by Liu and Hauswirth [13]. They studies 118 missing data imputation methods and highlighted 9 influential factors and 12 selection criteria. They proposed a provenance meta-learning method to select the proper imputation method from the methods they chose.

From this literature, it is quite clear that there are many missing data imputation methods and the selection of a specific method depends highly on missing data patterns.

### 1.2.2   Volve Dataset

The Volve dataset is quite new and there have been only a few works present in the literature. This study is based on a series of works done by Tunkiel et al. [14].

Tunkiel et al. [15] identified that the existing research that claims to have achieved an R2 value of as high as 0.996 for the rate of penetration prediction has no independent datasets that can verify these claims. They worked on providing a benchmark dataset based on Equinor's public Volve dataset. They gathered data from seven wells with nearly 200,000 samples with 12 common attributes. This data can work as the benchmark to test existing studies and verify their claims. Also, this data can work as a base for many new research studies.

Tunkiel et al. [16] highlighted the problem of the machine learning model becoming obso-

lete due to equipment changes on the drilling site. The change happens due to the fast-changing pace of the drilling industry. To solve this problem, they proposed a training while drilling approach. The model is deployed on already working well where data is continuously coming from the sensors. They used recurrent neural networks (RNNs) to capture the dynamic and ever-changing nature of the data. This type of model learns from new features while keeping in mind the patterns learned from previous data.

Tunkiel et al. [17] noticed that the data-driven machine learning models are mostly black boxes and there is no way to find why a model is behaving erratically. TO uncover this mystery sensitivity analysis of data is necessary. Sensitivity analysis of data can help uncover these erratic behaviors which might be caused by overfitting. They used the approach of the one-at-a-time method to cover the hyperspace of potential inputs.

Tunkiel et al. [18] stated that the drilling data is generated continuously from various sensors and this data is huge. But there are several reasons because of which the data gets lost and there are many sensors i.e., gamma and inclination that lags by many meters. They found a solution to this by proposing a novel approach that provides a prediction for lagging data. They combined a trend-based prediction model with traditional artificial intelligence models to predict the lagging data.

Tunkiel et al. [19] worked on exploring Equinor's Volve dataset as the dataset is large and needs pre-processing before performing analysis. The main objective of this study was to overcome the basic obstacles of dealing with this data and to help the new studies happening in the field of oil drilling.

Based on these studies, this study aims to make changes to existing pipelines and apply hyper-parameter tuning to improve the already existing results.

## 1.3   Motivation and Problem Statement

On basis of the fact that petroleum data needs to be cleaned before processing, there is a need for a pipeline that pre-processes the data before feeding it to the machine learning model. This pipeline also needs to be tuned using hyper-parameter optimization to get the most efficient results out of it. To evaluate the results of hyper-parameter optimization, Akaike Information

Criterion (AIC) and Bayesian Information Criterion have been used. Not only pre-processing but model selection is also required in pipeline optimization. For this purpose, multiple time series prediction models have been considered and compared.

## 1.4 Objective

Objectives of this thesis include:

1. Building ML pipeline for time series analysis of petroleum data.

2. Designing pre-processing pipeline and optimizing it via AIC and BIC.

3. Machine learning model selection and optimization via AIC and BIC.

## 1.5 Thesis Organization

The rest of the thesis is organized as follows. Chapter 2 presents the mathematical background of AIC and BIC. It will also give basic ideas behind the steps of pre-processing pipeline and introduces the ML models for time-series prediction. All important definitions are given in this section. Chapter 3 discusses the Recurrent Neural Network (RNN) and Long Short Term Memory (LSTM). This section will also explain the flow diagram of the ML pipeline. Chapter 4 discusses the dataset and simulation environment, and results. It also discusses the outcomes of the results. Chapter 5 concludes the thesis and states future recommendations.

# Chapter 2

# Background

## 2.1  Data Analysis

Data Pre-processing is to transform the raw data into a useful and efficient format. Steps that involve data pre-processing are:

### 2.1.1  Data cleaning

Data has many irrelevant and missing parts, that need to be dealt with. This included missing data handling and noise removal.

**Missing data**

This situation arises when some of the data is missing. There are several reasons that can cause the data to go missing i.e., faulty sensors. Missing data can be dealt with in 2 ways, either remove the missing data entries or even columns or fill the missing data with entries. To fill missing data there are several methods that include: filling entries with mean or most probable value.

**Noisy Data**

The data that can't be interpreted by the machine and is useless is called noisy data. A real-world dataset is affected by several components and one of them is noise [20]. It is an unavoidable

problem and affects the data performance of machine learning models. There are two main causes of noise in data [21]:

- Implicit error (fault of measurement tools).

- Random error (human error by data collectors).

In petroleum datasets, the noise that is present in the data is implicit noise and it is caused by the sensors. It can be handled in several ways that are:

**Binning:** This method works on data that is sorted. This method smoothes that data by dividing the data into several segments and then replacing values in each segment by mean or on basis of boundary values.

**Regression:** In this method data is made smooth by fitting it into a regression model. Both linear (1 independent variable) and multiple (multiple independent variables) regression can be used.

**Clustering:** In this method similar data points fall into the same cluster and outliers are removed.

## 2.1.2 Data Transformation

This step transforms the data and converts it into a suitable form for machine learning models. There are several data transformation methods that include:

**Normalization:** This method scales the data into a defined range i.e., between 0 and 1, or between -1 and 1.

**Attribute Selection:** This method either selects the attributes from within given attributes or constructs new attributes on basis of existing ones.

## 2.2 Data Imputation

Missing data occurs in many statistical analyses. Missing data occurs due to implicit errors where the sensors are unable to provide any value for an entry. There is no concrete method that can help in imputing the missing values and the user needs to select from some existing

methods to achieve this task.  These existing methods apply imputation with various methods
and these methods can sometimes be useful and sometimes that is not helping at all. Selecting
the appropriate methods requires struggle and time.

There are several imputation methods that can help in imputing missing values in the dataset.
There are two basic imputation methods:

- **Single Imputation:** In the single imputation method, each value is estimated separately.
  This method is simple as each value is calculated separately.  This method includes im-
  puting missing values with mean or regression. This imputation method is biased

- **Multiple Imputation:** This method imputes values simultaneously and then calculates
  errors to optimize the prediction.  Sparse matrix completion methods like soft-impute is
  an example of multiple-imputation. This method is unbiased.

## 2.3   Data Resampling

Data sampling is the process of selecting observations at the time of data collection. While data
resampling is to improve the already collected data. Resampling has many advantages in time
series machine learning analysis. The basic reason is that time series prediction models require
data to be evenly sampled. Thus calculating the sampling rate in time series data is a challenge
and requires tuning.

## 2.4   Feature Selection

Feature engineering helps in extracting features from data. It is a crucial task as it affects the
model accuracy directly. Plus there are many methods of feature engineering and selecting the
method according to input data is a difficult task. Feature engineering is also important because
it directly affects the accuracy of the ML model.  If features extracted from the dataset read
the pattern correctly then the model will give good performance.  If the features extracted do
not show any specific pattern then the performance will be poor.  Feature selection also falls
under the domain of feature engineering.  If the dataset is large it causes the model to take

more time and resources and may be due to some poor features model still doesn't show good performance. The pain point is to select which variables to use and which to skip, there are methods to determine that i.e. Z-score or $R^2$, and these are used to apply elimination one variable at a time. This variable selection method is tiresome and consumes a lot of time. In feature engineering feature selection is an important task from a high dimensional space of extracted features. The improvement can be done so the user doesn't have to look at the value and remove variables manually. Plus these techniques are applied when the model has been trained successfully and this causes the use of extra resources and time.

## 2.4.1 AIC & BIC

AIC and BIC are methods of selecting models by scoring them on the basis of their log-likelihood and complexity. AIC and BIC are used to compare statistical models on the basis of the number of free parameters and their values. It takes into account how much a variable contributes to solving a problem. For example, if we have a regression model which is being tested on the basis of four variables A1, A2, A3, and A4, AIC and BIC fit the regression model with different combinations of these variables and calculate the value using the formula:

$$AIC = -2(log\text{-}likelihood) + 2K \tag{2.1}$$

Here, K is the number of parameters used to fit the model, and log-likelihood checks if the model is a good fit or not.

$$BIC = -2(log\text{-}likelihood) + Klog(n) \tag{2.2}$$

Here, n is the number of values in the dataset (sample size), and k is the number of parameters.

The number of combinations for variables can be calculated by $2^N - 1$. Here N is the number of total variables. Combination of these values are shown in table 2.1.

As the number of variables is 4 so number of combinations will be $2^4 - 1 = 15$.

For each combination regression model is trained and AIC and BIC values are calculated using (2.1) and (2.2) respectively. A combination with the minimum value of AIC or BIC is

**Table 2.1:** Possible combination of 4 variables for applying AIC and BIC

| A1 | A2 |
|---|---|
| A3 | A4 |
| A1, A2 | A1, A3 |
| A1, A4 | A2, A3 |
| A2, A4 | A3, A4 |
| A1, A2, A3 | A1, A2, A4 |
| A2, A3, A4 | A1, A3, A4 |
| A1, A2, A3, A4 | |

selected as the best model.

When comparing the two models BIC puts a higher penalty term than AIC.

# 2.5   Time Series Models

## 2.5.1   Time Series Analysis

A specific way of analyzing a sequence of data over time is called time series analysis. In time-series data, data points are at a constant interval in time rather than at random points. Time-series data analysis can be done to see a change in data over time and thus predict future values. Data trends and seasonal changes are also part of the data to see the trends. The time series dataset needs to be extensive in order to ensure the consistency and reliability of the prediction.

Time series datasets are used to understand the trends and patterns over time. Visualization of time series data shows seasonal trends and provides a deeper picture of data. For example, the sale of a specific product increases at a specific time of the year i.e., groceries at the start of each month. Some examples of time-series datasets are:

- Weather Data

- Stock prices data

- Brain Monitoring (EEG)

Time series models can not be generalized to all king of time-series datasets. However deep learning has shown promising results in generalizing the models that can work with most of the datasets. Some important deep learning time series analysis models are:

- Long Short Term Memory (LSTM)

- Recurrent Neural Networks (RNN)

## 2.5.2  Artificial Neural Networks

Artificial Neural Networks are the most basic type of deep learning model that can take tabular data and then classify it. It consists of neurons that try to mimic the behavior of the human brain. Each neuron is connected to other neurons to form a network. Each neuron has weights and biases, whose values change during the optimization of the model. A simple ANN is shown in **??**. It consists of fully connected layers where each neuron is connected to all the neurons in the next layer. At first, each value from independent variables is given as input. Weights and biases are selected at random at the start. On basis of initial values of weights and biases predictions are made. Then error is calculated on basis of predicted values and actual values. On basis of error, weights and biases are updated. This procedure repeats until convergence.

2

## 2.5.3  Long Short Term Memory

LSTM is the neural network that takes previous knowledge into account for future predictions. LSTM takes advantage of both long-term memory (LTM) and short-term memory (STM) to learn. It has two basic mechanisms:

- **Forgetting Mechanism:** Forgets all the information that is not relevant.

- **Saving Mechanism:** Saves the information that is relevant and can help in the future.

To carry out the above tasks LSTM takes advantage of gates. There are 4 kinds of gates:

- **Forget Gate:** LTM uses forget gate to forget unuseful information.

- **Learn gate** Event (current input) and STM are combined so that recently learned information can be applied to the current input.
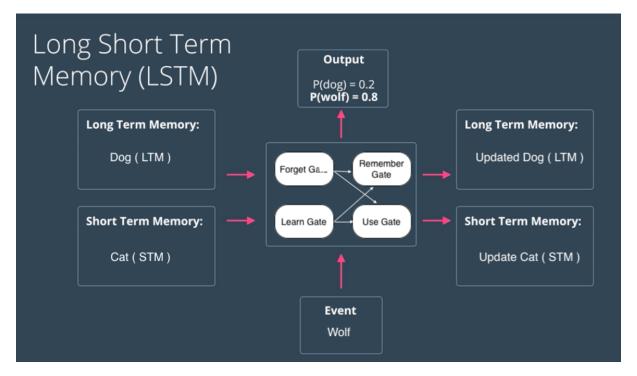
**Figure 2.1:** Operations of LSTM [1]

- **Remember gate:** LTM, STM and event are combined at the remembering gate to update the LTM.

- **Use gate:** LTM, STM, and event are combined to predict the current output and thus it updates STM.

Figure 2.1 shows the working of LTM and STM using the gates and how the values are updated. Figure 2.2 shows the basic architecture of LSTM.

This architecture can be divided for deep understanding. Figure 2.3 shows the architecture of each gate used in the LSTM model.

## 2.5.4  Recurrent Neural Networks

Artificial neural networks and convolution neural networks work best in the case of sequential image datasets respectively but fail for temporal data (dependency over time). Some examples of temporal datasets are speech datasets and stock price datasets etc.

A recurrent neural network (RNN) is similar to ANN with a small addition, which is it also takes previous input into account along with current information. The basic architecture of RNN is shown in Figure 2.4

**Figure 2.2:** Architecture of LSTM [1]

## 2.5.5  Mean Absolute Error

Absolute error is the error in the measurements. Simply speaking it is the difference between the measured values and true values.

$$\delta x = |x_i - x| \tag{2.3}$$

An absolute sign is needed because sometimes the difference is negative.

Mean absolute error is the mean of all the absolute errors.

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |x_i - x| \tag{2.4}$$

**Figure 2.3:** Gates in LSTM

**Figure 2.4:** Architecture of RNN [2]

# Chapter 3

# Mathodology

## 3.1 Machine Learning Pipeline

Building a Machine Learning pipeline is the first step in solving a machine learning problem. Figure 3.1 shows the overall pipeline of the model.
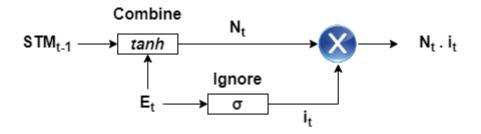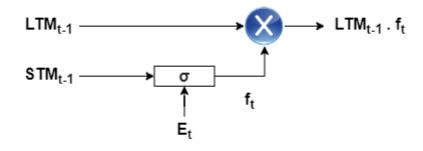
As seen from the Figure 3.1, the first steps are to remove unwanted columns i.e., the indexing column. This also includes the column that has a high correlation with the other independent columns for example ROP and inverse ROP as inversely highly correlated.

The second step is to analyze the data and run the gap analysis on basis of the number of gaps and length of each gap [14]. There are total of 4 scenarios:

1. Few gaps + low overall percentage

   - sensor obstruction
   - temporary sensor failure.

2. Few gaps + high overall percentage

   - equipment setup change
   - permanent sensor failure
   - data corruption

3. Multiple gaps + low overall percentage

**Figure 3.1:** Flowchart for feature selection (left) and Data Resampling hyper-parameter tuning (right) via AIC and BIC

- uneven polling frequency

- sequential sensor use (MWD)

4. Multiple gaps + high overall percentage

- sensors with low polling frequency

- data received occasionally

If the gap sizes are too large then it means that data is actually missing and filling that data is of no use and it is best to remove those columns. To find the col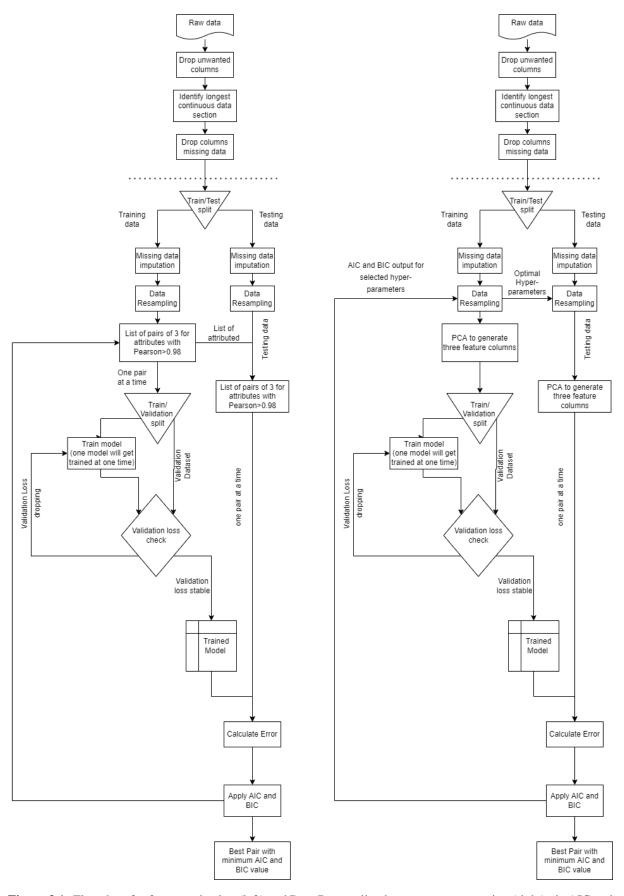umns that need to be removed as they can not be filled, a gap coefficient was used [14]. The gap coefficient is given by:

$$GC_i = \frac{i}{GQ_i \cdot TL} \tag{3.1}$$

where GC is the gap coefficient, i is gap length, $GQ_i$ is the number of gaps of length i, and TL is the total length. The gap coefficient for each attribute is calculated and if GC surpasses a threshold value then it means that those gaps can not be filled and the attributes need to be removed.

After that data is split into training data and testing data. Testing data is the one that needs to be used for final testing after the model has been trained.

### 3.1.1  Missing Data Imputation

The next step in data pre-processing is to impute the missing values. There are several missing data imputation methods but the K-nearest neighbors (KNN) data imputation method has proven to be generally effective. In the KNN data imputation method, k nearest neighbors are identified on basis of Euclidean distance for a specific missing value, and the mean value of the neighbors is selected to be imputed at that particular missing point. Euclidean distance is calculated by Equation 3.2. For data imputation, KNN imputer with 3 nearest neighbors has been used with Euclidean distance.

$$d(p, q) = \sqrt{\sum_{i=1}^{n} (q_i - p_i)^2} \tag{3.2}$$

## 3.1.2　Data Resampling

Data Resampling is an important step as i=time series analysis models require evenly sampled data. Two types of resampling methods have been used and tuned. The methods are:

- Radius Neighbors Regressor

- KNN Regressor

For hyper-parameter optimization of data resampling methods following parameters were considered:

1. Resampling Weights

    - **Uniform:** All points in the neighbourhood are given equal weights.

    - **Distance:** All points are given weights on basis of the distance from the point under consideration. Closer points have more weight then points that are far.

2. Algorithm to compute NN

    - Ball Tree Algorithm

    - KD Tree Algorithm

3. Distance Metric

    - **1:** Manhattan Distance (Equation 3.3)

    - **2:** Euclidean Distance (Equation 3.2)

    - **3:** Minkowski Distance (Equation 3.4)

$$d(p, q) = \sum_{i=1}^{n} |q_i - p_i| \qquad (3.3)$$

$$d(p, q) = \sqrt[r]{\sum_{i=1}^{n} (q_i - p_i)^r} \quad where, \ r > 0 \qquad (3.4)$$

### 3.1.3   Feature Selection

Feature selection is an important step of the pipeline as there are features that support the machine learning model and there are features that degrade the efficiency and accuracy of the model. Feature selection is the process of selecting the relevant features or generating new features for reducing the dimension of the data and fitting the model appropriately. There are several feature selection methods, some of them are:

**Principal Component Analysis (PCA)**

The principal component analysis is a statistical procedure, that summarizes the high dimensional data into lower dimensional data using the principal components of the data points.

PCA finds the lines, planes, and hyper-plane in k-dimensional space and approximates the data on basis of least square approximation. This is done by making the variance of the coordinates as large as possible on the line/plane.

The steps involvind the PCA are:

1. Suppose a matrix X with N rows and K columns. Plot each point on K-dimensional space after scaling to unit variance.

2. Subtract variable averages from the data to get a mean-centering vector, that represents a point in the data.

3. Origin is shifted to the mean-centering point.

4. First principal component (PC1) is found by fitting a line passing through the mean-centering point and is the best fit on basis of the least square. Each point can be projected onto this line to get a new value.

5. Second principal component (PC2) is found by passing a line perpendicular to the first component. This represents the second-largest variation in the data.

PC1 and PC2, combined together represent the whole data. For hyper-parameter tuning of the resampling step, PCA has been used to generate new features.

**Pearson Correlation**

Pearson correlation is the method that finds correlation between independent variables and dependent variables of the data by assigning the value between -1 and 1. Where, -1 means total negative correlation, 0 means no correlation and 1 means total positive correlation. Pearson correlation can be found by:

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \tag{3.5}$$

Pearson correlation can help in removing the independent variables that have a high correlation between them and can also help in identifying the columns that have a high correlation with the dependent variable.

AIC and BIC can be combined with Pearson correlation to find the best features. Top features with high correlation can be used to further reduce the number of training features by using AIC and BIC. After applying Pearson correlation, features are combined in pairs of 3, and then after running the model on a single combination, AIC and BIC value is calculated. Feature pairs with a minimum value of AIC and BIC are the best pair and can be used to get the best results. For the selection of best variables, Pearson with AIC and BIC has been used.

**Time Series Prediction**

Different models has been combined for generating the models that has been used for time series prediction. Time series prediction model has been combined with mylti-layer perceptron (MLP) model. Where, MLP is the fully connected feed forward, artificial neural network. The combinations are as follow:

1. LSTM + MLP (Figure 3.2)

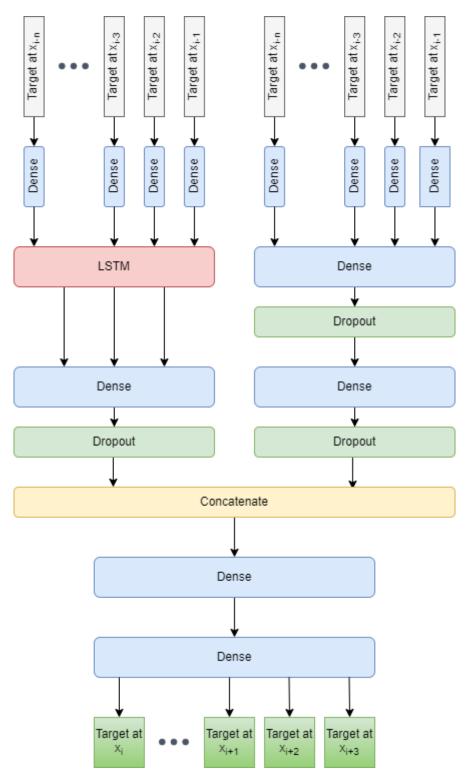2. RNN + MLP (Figure 3.3)

3. LSTM + RNN + MLP (Figure 3.4)
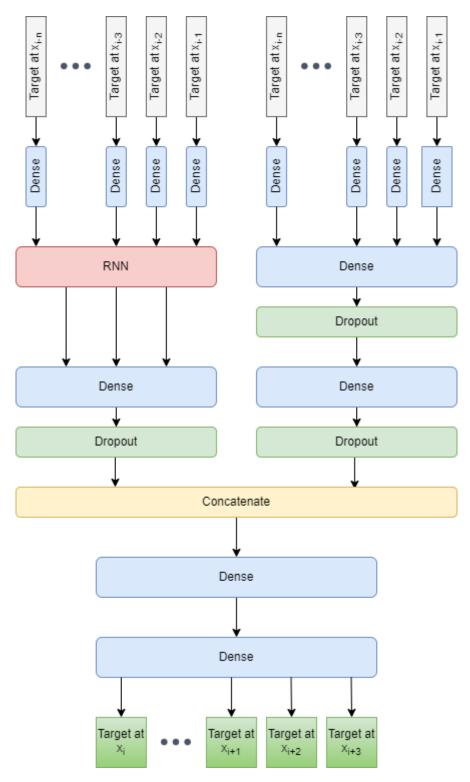
**Figure 3.2:** Model architecture of LSTM + MLP

**Figure 3.3:** Model architecture of RNN + MLP

**Figure 3.4:** Model architecture of LSTM + RNN + MLP

# Chapter 4

# Results and Discussion

## 4.1 Dataset

The Volve dataset Equinor [1] was published in 2018. The fossils here are sandstone of the middle
Jurassic age in the Hugin Formation. The depth of the Volve field is 2700m and 3100m where
the seabed ends at 80m. The field was completely shut down by 2018 and the dataset was
released to enhance the research in the oil and gas industry. During the whole operation, 56000
barrels per day were extracted from the oil field on average. The data that has been released has
40000 files of various kinds. The dataset has 75% and 80% data gaps or empty cells

## 4.2 Simulation Environment

All the simulations took place on Google Colab Jupyter notebooks with 12.75GB RAM and
Tesla P100 16GB and Tesla P4 GPUs.

AT first KNN imputer was applied to the three nearest neighbors. This step does not need
validation as it is a necessary step. Although it can be compared with other imputation models.

After that, the hyper-parameter tuning of the resampling step was done where two methods,
KNN regressor, and Radius Neighbour regressor were compared on basis of their parameter
values. For this stage, the feature selection step was not used and simply PCA was implemented.
After applying PCA, the model used for the time series prediction was LSTM. AIC and BIC

---

[1]https://data.equinor.com/dataset/Volve

have been calculated after fitting the model for each combination of the resampling method and its hyper-parameters.

During the feature selection first, the Pearson correlation was applied. From Pearson correlation, parameters with more than 98% correlation with the output parameter were selected and all possible combinations of 3 parameters were made from them. These combinations of the parameters were used to train the model one by one. AIC and BIC values were calculated after each run and a combination with the minimum value of AIC and BIC was selected as the best. During feature selection resampling method was used which was best in the previous stage, and the model used was LSTM + MLP.

The last step is model selection, in this step different models were tested for the best resampling method and best parameters, In the end, the best model was selected on basis of mean absolute error.

## 4.3   KNN-Imputation

In KNN imputation the nearest neighbors are selected as 3. Each row is imputer separately on basis of the distance from the 3 nearest neighbors. Figure 4.1 shows the plots of some parameters before and after the KNN imputation.

As it can be seen from the Figure 4.1, the graphs are quite predictable when seen, and after imputation, the graphs are as expected.

## 4.4   Resampling Hyperparameter Tuning

As described in the previous section, two methods were compared for the resampling, and for each method, hyper-parameters has been tuned. Table 4.1 shows all the possible combinations that have been tested and their respective AIC, BIC, and mean absolute error value. As it can be seen from the table the minimum value has been obtained for radius neighbor regressor with uniform distance using ball tree algorithm and with Minkowski distance. For this combination, all the output evaluation matrices have minimum values.

**Table 4.1:** Results for resampling method selection and hyper-parameter tuning

| Resampling Method | Weight | Algorithm | P4 | AIC | BIC | MAE |
|---|---|---|---|---|---|---|
| radius | uniform | ball_tree | Manhatten | 1393477.80 | 3716214.07 | 9.17 |
| | | | Euclidean | 1392074.95 | 2267638.43 | 7.57 |
| | | | Minkowski | 1392066.43 | 2267629.91 | 6.42 |
| | | kd_tree | Manhatten | 1392075.74 | 2267639.22 | 7.68 |
| | | | Euclidean | 1392076.66 | 2267640.14 | 7.82 |
| | | | Minkowski | 1392076.05 | 2267639.53 | 7.73 |
| | distance | ball_tree | Manhatten | 1392077.07 | 2267640.54 | 7.88 |
| | | | Euclidean | 1392077.79 | 2267641.27 | 7.99 |
| | | | Minkowski | 1392079.65 | 2267643.13 | 8.28 |
| | | kd_tree | Manhatten | 1392079.22 | 2267642.70 | 8.22 |
| | | | Euclidean | 1392078.76 | 2267642.24 | 8.14 |
| | | | Minkowski | 1392077.07 | 2267640.55 | 7.88 |
| knn | uniform | ball_tree | Manhatten | 1392068.59 | 2267632.07 | 6.69 |
| | | | Euclidean | 1392080.39 | 2267643.87 | 8.40 |
| | | | Minkowski | 1392073.25 | 2267636.73 | 7.32 |
| | | kd_tree | Manhatten | 1392077.83 | 2267641.31 | 8.00 |
| | | | Euclidean | 1392076.82 | 2267640.30 | 7.84 |
| | | | Minkowski | 1392077.79 | 2267641.27 | 7.99 |
| | distance | ball_tree | Manhatten | 1392077.77 | 2267641.25 | 7.99 |
| | | | Euclidean | 1392077.91 | 2267641.38 | 8.01 |
| | | | Minkowski | 1392075.77 | 2267639.25 | 7.68 |
| | | kd_tree | Manhatten | 1392078.59 | 2267642.07 | 8.11 |
| | | | Euclidean | 1392078.09 | 2267641.57 | 8.03 |
| | | | Minkowski | 1392079.04 | 2267642.52 | 8.19 |

**Figure 4.1:** Graphs before (top) and after (bottom) KNN imputation

## 4.5   Feature Selection

AIC and BIC have been used for feature selection too.  Here at first, the Pearson correlation matrix was found which is shown in Figure 4.2.

After Pearson correlation, parameters were extracted having a correlation value greater than 98% with the dependent variable as shown in Table 4.2.

From these parameters, all possible combinations of the three parameters have been made.

**Table 4.2:** Pearson Correlation of parameters with dependent variables

| Parameter | Pearson Correlation |
|---|---|
| Measured Depth m | 0.9813667349941119 |
| Hole depth (MD) m | 0.9813684472500676 |
| Hole Depth (TVD) m | 0.9826332567926808 |
| Corrected Total Hookload kkgf | 0.9827814199492585 |
| Extrapolated Hole TVD m | 0.9828458285022305 |
| Bit Drill Time h | 0.9854056047108368 |
| RGX_RT unitless | 0.9909452831190304 |
| RHX_RT unitless | 0.998696718883108 |

**Figure 4.2:** Heatmap of pearson correlation

**Figure 4.3:** AIC (top) and BIC (bottom) graphs for different parameter combinations

All possible combinations can be seen in Appendix A. These parameters have been then used to train the model and AIC and BIC values have also been calculated for each combination. Figure 4.3 shows the graph of AIC and BIC values for each combination. From the graphs it is clear that the minimum AIC and BIC value has been achieved at combination number 46 which has the following parameters:

- Corrected Total Hookload kkgf

- Extrapolated Hole TVD m

- Bit Drill Time h

AIC and BIC value for this combination is 1393443.28 and 3716179.55 respectively.

**Table 4.3:** Mean Absolute Error after training different models.

| Model | MAE |
|---|---|
| LSTM + MLP | 8.8 |
| RNN + MLP | 8.71 |
| LSTM + RNN + MLP | 8.82 |

# 4.6 Model Selection

All the parameters have been tuned and fixed, models were selected from among the models that have been highlighted in the previous chapter. For each model mean absolute error (MAE) value has been noted and the model has a minimum MAE value that has been selected as the best. The results have been shown in Table 4.3.

For each model, the respective prediction curve can be seen in Figure 4.4. As it can be seen from the figure, RNN is showing the best results, the variance in the prediction line is less as compared to other models.
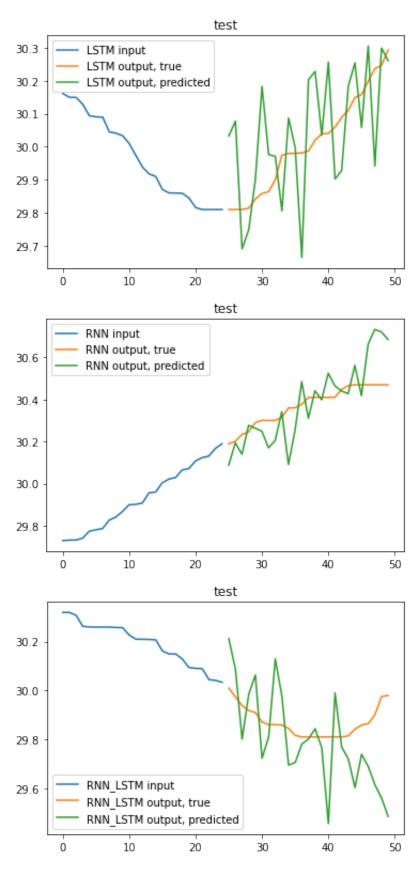
**Figure 4.4:** AIC (top) and BIC (bottom) graphs for different parameter combinations

# Chapter 5

# Conclusion and Future Work

## 5.1   Conclusion

This research aims to build and improve an ML pipeline via hyper-parameter optimization. At first, an already existing pipeline has been studied and work has been done to improve it. This thesis aimed at improving missing data imputation, resampling, feature selection, and in the end model selection. Missing data imputation has been done by KNN imputer which imputed the missing values on basis of nearest neighbors. After that, the resampling method has been selected after tuning their hyper-parameters. For feature selection firstly Pearson correlation has been used to narrow down the top relevant features, after that all possible combinations of these features have been created as a group of three. AIC and BIC values have been used to select the most optimal set of parameters. After feature selection, three models were considered and compared on basis of mean absolute error (MAE). RNN + MLP surpasses all the other models on basis of MAE.

## 5.2   Future Work

In the future, these results can be further improved and tested by joining all the parameter selections as one. As in this research, parameter optimization for each step has been done separately.

# References

[1] Gourav Singh. Understanding architecture of lstm, 2021.

[2] Gourav Singh. Recurrent neural networks for sequence learning, 2020.

[3] Sana Mushtaq. Data preprocessing in detail, 2019.

[4] Jinlin Zhu, Zhiqiang Ge, Zhihuan Song, and Furong Gao. Review and big data perspectives on robust data mining approaches for industrial process modeling with outliers and missing data. *Annual Reviews in Control*, 46:107–133, 2018.

[5] Chengqi Zhang, Yongsong Qin, Xiaofeng Zhu, Jilian Zhang, and Shichao Zhang. Clustering-based missing value imputation for data preprocessing. In *2006 4th IEEE International Conference on Industrial Informatics*, pages 1081–1086. IEEE, 2006.

[6] RS Somasundaram and R Nedunchezhian. Evaluation of three simple imputation methods for enhancing preprocessing of data with missing values. *International Journal of Computer Applications*, 21(10):14–19, 2011.

[7] Jianglin Huang, Yan-Fu Li, Jacky Wai Keung, Yuen Tak Yu, and WK Chan. An empirical analysis of three-stage data-preprocessing for analogy-based software effort estimation on the isbsg data. In *2017 IEEE International Conference on Software Quality, Reliability and Security (QRS)*, pages 442–449. IEEE, 2017.

[8] Xiaolong Xu, Weizhi Chong, Shancang Li, Abdullahi Arabo, and Jianyu Xiao. Miaec: Missing data imputation based on the evidence chain. *IEEE Access*, 6:12983–12992, 2018.

[9] F Nelwamondo and Tshilidzi Marwala. Key issues on computational intelligence techniques for missing data imputation-a review. In *Proc. of world multi conf. on systemics, cybernetics and informatics*, volume 4, pages 35–40, 2008.

[10] Little Rja and Donald B Rubin. Statistical analysis with missing data. *New York*, 1987.

[11] Donald B Rubin. Multiple imputation for survey nonresponse, 1987.

[12] Yagyanath Rimal. Multivariate imputation for missing data handling a case study on small and large data sets. *International Journal of Human Computing Studies*, 2(1):5–11, 2020.

[13] Qian Liu and Manfred Hauswirth. A provenance meta learning framework for missing data handling methods selection. In *2020 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON)*, pages 0349–0358. IEEE, 2020.

[14] Andrzej T. Tunkiel, Dan Sui, and Tomasz Wiktorski. Impact of data pre-processing techniques on recurrent neural network performance in context of real-time drilling logs in an automated prediction framework. *Journal of Petroleum Science and Engineering*, 208:109760, 2022.

[15] Andrzej T. Tunkiel, Dan Sui, and Tomasz Wiktorski. Reference dataset for rate of penetration benchmarking. *Journal of Petroleum Science and Engineering*, 196:108069, 2021.

[16] Andrzej T. Tunkiel, Dan Sui, and Tomasz Wiktorski. Training-while-drilling approach to inclination prediction in directional drilling utilizing recurrent neural networks. *Journal of Petroleum Science and Engineering*, 196:108128, 2021.

[17] Andrzej T. Tunkiel, Dan Sui, and Tomasz Wiktorski. Data-driven sensitivity analysis of complex machine learning models: A case study of directional drilling. *Journal of Petroleum Science and Engineering*, 195:107630, 2020.

[18] Andrzej T Tunkiel, Tomasz Wiktorski, and Dan Sui. Continuous drilling sensor data reconstruction and prediction via recurrent neural networks. In *International Conference on Offshore Mechanics and Arctic Engineering*, volume 84317, page V001T01A002. American Society of Mechanical Engineers, 2020.

[19] Andrzej T Tunkiel, Tomasz Wiktorski, and Dan Sui. Drilling dataset exploration, processing and interpretation using volve field data. In *International Conference on Offshore Mechanics and Arctic Engineering*, volume 84430, page V011T11A076. American Society of Mechanical Engineers, 2020.

[20] Richard Y Wang, Veda C Storey, and Christopher P Firth. A framework for analysis of data quality research. *IEEE transactions on knowledge and data engineering*, 7(4):623–640, 1995.

[21] Xingquan Zhu and Xindong Wu. Class noise vs. attribute noise: A quantitative study. *Artificial intelligence review*, 22(3):177–210, 2004.

# Appendices

# Appendix A

# Parameter Combinations

| Sr. No. | P1 | P2 | P3 |
|---|---|---|---|
| **0** | | | Hole Depth (TVD) m |
| **1** | | | Corrected Total Hookload kkgf |
| **2** | | Hole depth | Extrapolated Hole TVD m |
| **3** | | (MD) m | Bit Drill Time h |
| **4** | | | RGX_RT unitless |
| **5** | | | RHX_RT unitless |
| **6** | | | Corrected Total Hookload kkgf |
| **7** | | | Extrapolated Hole TVD m |
| **8** | | Hole Depth | Bit Drill Time h |
| **9** | | (TVD) m | RGX_RT unitless |
| **10** | Measured | | RHX_RT unitless |
| **11** | Depth m | | Extrapolated Hole TVD m |
| **12** | | Corrected Total | Bit Drill Time h |
| **13** | | Hookload kkgf | RGX_RT unitless |
| **14** | | | RHX_RT unitless |
| **15** | | | Bit Drill Time h |
| **16** | | Extrapolated | RGX_RT unitless |
| **17** | | Hole TVD m | RHX_RT unitless |

| | | | |
|---|---|---|---|
| **18** | | Bit Drill | RGX_RT unitless |
| **19** | | Time h | RHX_RT unitless |
| **20** | | RGX_RT unitless | RHX_RT unitless |
| **21** | | | Corrected Total Hookload kkgf |
| **22** | | Hole Depth (TVD) m | Extrapolated Hole TVD m |
| **23** | | | Bit Drill Time h |
| **24** | | | RGX_RT unitless |
| **25** | | | RHX_RT unitless |
| **26** | | | Extrapolated Hole TVD m |
| **27** | | Corrected Total Hookload kkgf | Bit Drill Time h |
| **28** | Hole depth (MD) m | | RGX_RT unitless |
| **29** | | | RHX_RT unitless |
| **30** | | | Bit Drill Time h |
| **31** | | Extrapolated Hole TVD m | RGX_RT unitless |
| **32** | | | RHX_RT unitless |
| **33** | | Bit Drill | RGX_RT unitless |
| **34** | | Time h | RHX_RT unitless |
| **35** | | RGX_RT unitless | RHX_RT unitless |
| **36** | | | Extrapolated Hole TVD m |
| **37** | | Corrected Total Hookload kkgf | Bit Drill Time h |
| **38** | | | RGX_RT unitless |
| **39** | | | RHX_RT unitless |
| **40** | Hole Depth (TVD) m | | Bit Drill Time h |
| **41** | | Extrapolated Hole TVD m | RGX_RT unitless |
| **42** | | | RHX_RT unitless |
| **43** | | Bit Drill | RGX_RT unitless |
| **44** | | Time h | RHX_RT unitless |

| | | | |
|---|---|---|---|
| **45** | | RGX_RT unitless | RHX_RT unitless |
| **46** | | | Bit Drill Time h |
| **47** | | Extrapolated Hole TVD m | RGX_RT unitless |
| **48** | Corrected Total | | RHX_RT unitless |
| **49** | Hookload kkgf | Bit Drill Time h | RGX_RT unitless |
| **50** | | | RHX_RT unitless |
| **51** | | RGX_RT unitless | RHX_RT unitless |
| **52** | | Bit Drill Time h | RGX_RT unitless |
| **53** | Extrapolated Hole TVD m | | RHX_RT unitless |
| **54** | | RGX_RT unitless | RHX_RT unitless |
| **55** | Bit Drill Time h | RGX_RT unitless | RHX_RT unitless |

# Appendix B

# Python Code

## B.1 Resampling Method Optimization

```python
import fbprophet
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import sys
import random
from sklearn.preprocessing import MinMaxScaler
from numpy.random import seed
import tensorflow as tf
from tensorflow.keras.utils import plot_model
from pca_mod import shift_pca
from pca_mod import shift_notpca
from sklearn.preprocessing import StandardScaler
from statistics_module import stats
import itertools

#%matplotlib qt

df = pd.read_csv('f9ad.csv').iloc[:, 2:]
drops = []

cols = df.columns

for i in range(len(df.T)):
    if str(np.dtype(df.iloc[:, i]))=='object':
        drops.append(cols[i])

print(drops)
```

```
31      df = df.drop(columns = drops)
32
33      index = 'Measured Depth m'
34      target =  'MWD Continuous Inclination dega'
35      smartfill = 0.9
36
37      s, m, per = stats(df)
38
```

**Listing B.1:** Importing data

```
1
2      #%%
3      ## Gap statistics for target
4      #
5      # This chart will show the percentage of dataset occupied by gaps
       of a certain
6      # size. Gaps are normal in drilling logs and nothing to be afraid
       of
7
8      x_label = per[target]['gap_sizes']
9      x = np.arange(0, len(x_label),1)
10
11     y = per[target]['percentage_cells_occupied']
12     plt.figure(figsize=(15, 8))
13     plt.xticks(x, x_label, rotation=90)
14     plt.bar(x,y)
15     plt.title(f'Gap distribution in:\n {target}')
16     plt.xlabel('Gap length')
17     plt.ylabel('Percentage of dataset occupied')
18
19     x_labels = x.tolist()
20     x_labels[0] = 'data'
21     plt.xticks(x, x_labels)
22     plt.grid()
23     plt.show()
24
25     #%%
26
27     ## Outlier detection
28
29     outlier_cutoff = 0.005 #arbitrarily selected
30
31     # calculation that penalizes long, rare, continuous gaps
32     out_coef = per[target]['gap_sizes'] / (per[target]['gap_counts']
       * len(df))
33
34     x = np.arange(0,len(per[target]['gap_sizes']),1)
35     x_label = per[target]['gap_sizes']
36     x = np.arange(0, len(x_label),1)
37     plt.figure(figsize=(15, 8))
38     plt.xticks(x, x_label, rotation=90)
```

```python
39    plt.bar(x, out_coef)
40    plt.ylim(0, 0.01)
41    plt.plot(x, [outlier_cutoff]*len(x), color='red', label='cutoff')
42    plt.legend()
43    x_labels = x.tolist()
44    x_labels[0] = 'data'
45    plt.xticks(x, x_labels)
46    plt.title(f'Gap occupancy = gap size / (relative gap quantity) \n
      in:{target}')
47    plt.xlabel('Gap length')
48    plt.ylabel('Gap occupancy')
49    plt.grid()
50    plt.show()
51
52    #%%
53    ## Automatic proposal of useful area, part 1
54
55    # find the smallest outlier gap
56    # outlier coefficients - True when above outlier cutoff
57    cutoff_map = (out_coef >= outlier_cutoff)
58
59    # Using map created above, what is the smallest outlier?
60    lower_cutoff = np.min(np.asarray(per[target]['gap_sizes'])[
      cutoff_map])
61
62    # This is done to quickly mark gaps bigger than cutoff with zero
      and
63    # other with NaN. This makes a good chart.
64    from functools import partial
65
66    def _cutoff(x, lower_cutoff=0):
67        if x >= lower_cutoff:
68            return 0
69        else:
70            return np.nan
71
72    _cutoff_par = partial(_cutoff, lower_cutoff=lower_cutoff)
73    mapped_outliers = list(map(_cutoff_par, m[target]))
74    plt.figure(figsize=(15, 8))
75    plt.scatter(df[index], df[target], s=1, c='black', label='data')
76    plt.plot(df[index], mapped_outliers, c='red', label='Unusuable
      range')
77    #                                                      # has to be
      plot to avoid index
78    #                                                      #
      discontinuities
79
80    plt.grid()
81    plt.legend()
82    plt.title('Useful range analysis')
83    plt.xlabel(f'{index}')
```

```python
84      plt . ylabel ( f '{ target } ')
85      plt . show ()
86      #%%

88      ## Automatic proposal of useful area , part 2

90      # Simply finds the biggest area with acceptable gaps

92      # TODO - check if the algorithm will detect a stride at the end
        of the dataset
93      #           because I have a feeling it won't !

95      strides = []

97      s_start = -1
98      s_stop = -1

100     for i in range ( len ( df ) ) :
101         if mapped_outliers [ i ] != 0 and s_start == -1:
102             s_start = i
103         elif mapped_outliers [ i ] == 0 and s_start != -1:
104             s_stop = i
105             strides . append ([ s_start , s_stop , s_stop - s_start ])

107             s_start = -1
108             s_stop = -1

110     strides = np . asarray ( strides )
111     strides = strides [ strides [: ,2]. argsort () ][:: -1] # sort by length
        [2]
112                                                           # and reverse
113     print ( f '''Proposed range to use is row { strides [0 ,0]} to row {
        strides [0 ,1]}
114     for total of { strides [0 ,0]} rows
115           ''')
116     print ( f 'All found strides are : [ start , stop , length ]')
117     print ( strides )

119     s_start = strides [0 ,0]
120     s_stop = strides [0 ,1]
121     #%%

123     ## cut the dataframe for the selected stride , redo the stats
124     #  From now on dfs is used ( dataframe stride )
125     margin_percent = 1  # since the edges can be a bit unpredictable ,
        margin is
126                             # removed

128     s_start = s_start + int (( s_stop - s_start ) * 0.01* margin_percent )
129     s_stop = s_stop - int (( s_stop - s_start ) * 0.01* margin_percent )
130     dfs = df . iloc [ s_start : s_stop ] # dfs = DataFrameStride
```

```python
131     s, m, per = stats(dfs)

132
133     #%%

134
135     ## Removing columns that contain big gaps

136
137     from functools import partial

138
139     def _cutoff_inv(x, lower_cutoff=0):
140         if x >= lower_cutoff:
141             return 1
142         else:
143             return 0

144
145     _cutoff_par = partial(_cutoff_inv, lower_cutoff=lower_cutoff)

146
147     killed_cols = []

148
149     for column in list(dfs):
150         mapped_outliers = list(map(_cutoff_par, m[column]))
151         offender_count = np.sum(mapped_outliers)
152         if offender_count > 0:
153             dfs = dfs.drop(columns = [column])
154             killed_cols.append([column, 100*offender_count/len(dfs)])

155
156     killed_cols = pd.DataFrame(killed_cols, columns=['Name','Percent
        offending'])
157     print('Removed following columns due to outlier gap (showing
        under 15% only):')
158     print(killed_cols.sort_values(by='Percent offending')[killed_cols
        ['Percent offending'] < 15])

159
160     #%%

161
162     ## Checking if first derivative of index is stable.

163
164     index_dr = np.diff(dfs[index])

165
166     index_mean = np.mean(index_dr)
167     index_std = np.std(index_dr)
168     index_maxgap = np.max(index_dr)
169     deviation = np.abs(index_dr - index_mean)/index_std

170
171     print(f'Maximum distance from mean is {np.max(deviation):.1f}
        standard deviations')
172     print(f'If this value is above 6, there may be too high sampling
        frequency variation')

173
174     #%%

175
176     ## Counting zeros in the first derivative to see if it should be
```

```
        ffilled
177     ## or linearly interpolated
178
179     ## NOTE: Actual filling will not happen here, but AFTER the data
        split
180
181     fill_method = {}
182
183     for attribute in list(dfs):
184
185         dropna_diff = np.diff(dfs[attribute].dropna())
186         zeros_p = np.count_nonzero(dropna_diff == 0) / len(
        dropna_diff)
187
188         if zeros_p >= smartfill: # Threshold to check?
189             fill_method[attribute] = 'ffill'
190         else:
191             fill_method[attribute] = 'linterp'
192
193
194     #%%
195
196     #%%
197
198     ## Gap filling - but only forward filling. Linear interpolation
        is done later
199
200     # Resampling helps with uneven distribution of data
201     # Timeseries models wants data points to be evenly spaced in time
        domain
202     for attribute in list(dfs):
203         if fill_method[attribute] == 'ffill':
204             dfs[attribute] = dfs[attribute].ffill().bfill()#.rolling
        (5, center=True).mean().ffill().bfill()
205
206
```

**Listing B.2:** Gap Analysis

```
1
2       split = 0.6 #portion of data available
3       future = 0.15 #section after available, for testing
4
5       from fancyimpute import KNN, NuclearNormMinimization, SoftImpute,
        BiScaler
6       from fancyimpute import IterativeImputer
7
8       cols = dfs.columns
9
10      masked = np.nonzero(pd.isnull(dfs.values))
11      xx = masked[0]
12      yy = masked[1]
```

```python
13    missing_mask = np.concatenate((xx[:, None],yy[:, None]), axis=1)
14
15    M_data = np.nan_to_num(dfs)
16
17    masked = np.nonzero(M_data)
18    xx = masked[0]
19    yy = masked[1]
20    observed_mask = np.concatenate((xx[:, None],yy[:, None]), axis=1)
21
22    X_filled_knn = KNN(k=3).fit_transform(dfs)
23
24    dfs1 = pd.DataFrame(X_filled_knn)
25    dfs1.columns = cols
26    dfs1.head()
27
28    X = dfs1.drop(target, axis=1)
29    y = dfs1[target].to_frame()
30
```

**Listing B.3:** KNN Imputation

```python
1
2    from tensorflow.keras import Model, Input
3    from tensorflow.keras.layers import (Dense, Dropout, GRU, Flatten
     ,
4                                          GaussianNoise, concatenate,
     LSTM,
5                                          Bidirectional,
     TimeDistributed)
6    from tensorflow.keras.layers import Conv1D
7    from tensorflow.keras.layers import MaxPool1D
8    from tensorflow.keras.callbacks import EarlyStopping
9    from tensorflow.keras.callbacks import ModelCheckpoint
10   import tensorflow as tf
11
12   from tensorflow.keras.models import load_model
13   import math
14
15   import tensorflow as tf
16   print("Num GPUs Available: ", len(tf.config.list_physical_devices
     ('GPU')))
17
18
```

**Listing B.4:** Importing Neural Networks Libraries

```python
1
2    def models(visible1, visible2, model_name="RNN"):
3      if model_name == 'RNN':
4        x1 = TimeDistributed(Dense(hDense4))(visible1)
5        x1 = GRU(units=hGRU, kernel_initializer = 'glorot_uniform',
6                 recurrent_initializer='orthogonal',
```

```
7                    bias_initializer="zeros", kernel_regularizer='l2',
    recurrent_regularizer=None,
8                    bias_regularizer=None, activity_regularizer=None,
    kernel_constraint=None,
9                    recurrent_constraint=None, bias_constraint=None,
    return_sequences=True,
10                   return_state=False, stateful=False)(x1)
11
12      x1 = Dense(imagination)(x1)
13      x1 = Flatten()(x1)
14      x1 = Dropout(hDrop1)(x1)
15
16
17
18      x2 = TimeDistributed(Dense(hDense5))(visible2)
19      dense2 = Dense(hDense1, activation="linear")(x2)
20      drop2 = Dropout(hDrop2)(dense2)
21      flat2 = Flatten()(drop2)
22      dense2 = Dense(imagination, activation='linear')(flat2)
23      drop2 = Dropout(hDrop3)(flat2)
24
25      model = concatenate([x1, drop2])
26
27    if model_name == 'LSTM':
28      x1 = TimeDistributed(Dense(hDense4))(visible1)
29      x1 = LSTM(units=hGRU, kernel_initializer = 'glorot_uniform',
30                   recurrent_initializer='orthogonal',
31                   bias_initializer="zeros", kernel_regularizer='l2',
    recurrent_regularizer=None,
32                   bias_regularizer=None, activity_regularizer=None,
    kernel_constraint=None,
33                   recurrent_constraint=None, bias_constraint=None,
    return_sequences=True,
34                   return_state=False, stateful=False)(x1)
35
36      x1 = Dense(imagination)(x1)
37      x1 = Flatten()(x1)
38      x1 = Dropout(hDrop1)(x1)
39
40
41
42      x2 = TimeDistributed(Dense(hDense5))(visible2)
43      dense2 = Dense(hDense1, activation="linear")(x2)
44      drop2 = Dropout(hDrop2)(dense2)
45      flat2 = Flatten()(drop2)
46      dense2 = Dense(imagination, activation='linear')(flat2)
47      drop2 = Dropout(hDrop3)(flat2)
48
49      model = concatenate([x1, drop2])
50
51    if model_name=='RNN_LSTM':
```

```
52      x1 = TimeDistributed(Dense(hDense4))(visible1)
53      x1 = GRU(units=hGRU, kernel_initializer = 'glorot_uniform',
54              recurrent_initializer='orthogonal',
55              bias_initializer="zeros", kernel_regularizer='l2',
    recurrent_regularizer=None,
56              bias_regularizer=None, activity_regularizer=None,
    kernel_constraint=None,
57              recurrent_constraint=None, bias_constraint=None,
    return_sequences=True,
58              return_state=False, stateful=False)(x1)
59
60      x1 = Dense(imagination)(x1)
61      x1 = Flatten()(x1)
62      x1 = Dropout(hDrop1)(x1)
63
64      x3 = TimeDistributed(Dense(hDense4))(visible1)
65      x3 = LSTM(units=hGRU, kernel_initializer = 'glorot_uniform',
66              recurrent_initializer='orthogonal',
67              bias_initializer="zeros", kernel_regularizer='l2',
    recurrent_regularizer=None,
68              bias_regularizer=None, activity_regularizer=None,
    kernel_constraint=None,
69              recurrent_constraint=None, bias_constraint=None,
    return_sequences=True,
70              return_state=False, stateful=False)(x3)
71
72      x3 = Dense(imagination)(x3)
73      x3 = Flatten()(x3)
74      x3 = Dropout(hDrop1)(x3)
75
76
77
78      x2 = TimeDistributed(Dense(hDense5))(visible2)
79      dense2 = Dense(hDense1, activation="linear")(x2)
80      drop2 = Dropout(hDrop2)(dense2)
81      flat2 = Flatten()(drop2)
82      dense2 = Dense(imagination, activation='linear')(flat2)
83      drop2 = Dropout(hDrop3)(flat2)
84
85      model = concatenate([x1, x3, drop2])
86
87   return model
88
```

**Listing B.5:** Defining Machine Learning Models

```
1
2    hstep_extension = 5
3    resample='yes'
4    resample_coef = 1
5    #resample_weights='distance'
6    step_length = index_mean * hstep_extension
```

```python
 7
 8    if resample != 'no':
 9      grid = [['radius', 'knn'], ['uniform', 'distance'], ['ball_tree
      ', 'kd_tree'], [1, 2, 3]]
10    elif resample == 'no':
11      grid = [[1], [1], [1], [1]]
12
13    AIC_vals = []
14    BIC_vals = []
15    parameters = []
16    mae = []
17    log_res = []
18
19    for a in range(len(grid[0])):
20      resample=grid[0][a]
21      for b in range(len(grid[1])):
22        resample_weights=grid[1][b]
23        for c in range(len(grid[2])):
24          algorithm_test = grid[2][c]
25          for d in range(len(grid[3])):
26            metric=grid[3][d]
27
28            if split != 1:
29                splitpoint = int(len(dfs)*split)
30                futurepoint = int(len(dfs)*(split+future))
31
32                X_train = X[:splitpoint]
33                y_train = y[:splitpoint]
34                X_test  = X[splitpoint:futurepoint]
35                y_test  = y[splitpoint:futurepoint]
36            else:
37                X_train = X
38                y_train = y
39                X_test  = X
40                y_test  = y
41
42
43            i_train_min = np.min(X_train[index])
44            i_train_max = np.max(X_train[index])
45            i_test_min = np.min(X_test[index])
46            i_test_max = np.max(X_test[index])
47
48            index_train = np.arange(i_train_min, i_train_max,
      step_length).reshape(-1,1)
49            index_test = np.arange(i_test_min, i_test_max,
      step_length).reshape(-1,1)
50
51            parameters.append(np.asarray([resample, resample_weights,
       algorithm_test, metric]))
52            print([resample, resample_weights, algorithm_test, metric
      ])
```

```python
53              if resample != 'no':
54                  if resample == 'radius':
55                      from sklearn.neighbors import
   RadiusNeighborsRegressor
56                      reg = RadiusNeighborsRegressor(radius=
   index_maxgap*resample_coef, weights=resample_weights, algorithm=
   algorithm_test, p=metric)
57                  elif resample == 'knn':
58                      from sklearn.neighbors import KNeighborsRegressor
59                      reg = KNeighborsRegressor(weights=
   resample_weights, n_neighbors=resample_coef, algorithm=
   algorithm_test, p=metric)
60                  else:
61                      sys.exit("Error, incorrect resampling algorithms
   choice")
62
63                  reg.fit(X_train[index].to_numpy().reshape(-1,1),
   y_train[target].to_numpy())
64                  y_train = pd.DataFrame()
65                  y_train[target] = reg.predict(index_train)
66
67                  reg.fit(X_test[index].to_numpy().reshape(-1,1),
   y_test[target].to_numpy())
68                  y_test = pd.DataFrame()
69                  y_test[target] = reg.predict(index_test)
70
71                  X_train_resampled = pd.DataFrame()
72                  for attribute in list(X_train):
73                      reg.fit(X_train[index].to_numpy().reshape(-1,1),
   X_train[attribute].to_numpy())
74                      X_train_resampled[attribute] = reg.predict(
   index_train)
75
76                  X_train = X_train_resampled
77
78                  X_test_resampled = pd.DataFrame()
79                  for attribute in list(X_train):
80                      reg.fit(X_test[index].to_numpy().reshape(-1,1),
   X_test[attribute].to_numpy())
81                      X_test_resampled[attribute] = reg.predict(
   index_test)
82
83                  X_test = X_test_resampled
84
85              elif resample == 'no':
86                  train_q = int((i_train_max - i_train_min)/step_length
   )
87                  test_q  = int((i_test_max - i_test_min)/step_length)
88
89                  #sample_train = np.sort(random.sample(range(len(
   X_train)), train_q))
```

```python
90                #sample_test = np.sort(random.sample(range(len(X_test
    )), test_q))
91
92                sample_train = np.linspace(0, len(X_train)-1, train_q
    , dtype=int)
93                sample_test = np.linspace(0, len(X_test)-1, test_q,
    dtype=int)
94
95                X_train = X_train.iloc[sample_train,:]
96                y_train = y_train.iloc[sample_train,:]
97
98                X_test = X_test.iloc[sample_test,:]
99                y_test = y_test.iloc[sample_test,:]
100
101           else:
102                sys.exit("Error, incorrect resampling choice")
103
104           ## Inclination to delta inclination convertion needed
    here!
105           convert_to_diff = []
106           asel_choice = 'pca'
107           hPcaScaler='mm',
108           #list of parameters that are to be in local coordinate
    system
109
110           for attr in convert_to_diff:
111                if attr == target:
112                    y_train[attr] = y_train[attr].diff().bfill() #
    bfill to kill initial NaN
113                    y_test[attr] = y_test[attr].diff().bfill()
114                else:
115                    X_train[attr] = X_train[attr].diff().bfill()
116                    X_test[attr] = X_test[attr].diff().bfill()
117           #%%
118
119           ## Scaling the data. Note that range is decide on the
    training dataset only!
120
121
122           # Both conditions have same outcome, why?
123           if asel_choice == 'pca' and hPcaScaler == 'ss':
124                scaler_X = MinMaxScaler() # StandardScaler()
125                scaler_y = MinMaxScaler()
126           else:
127                scaler_X = MinMaxScaler()
128                scaler_y = MinMaxScaler()
129
130           pca_allattr = X_train.columns
131           X_train[X_train.columns] = scaler_X.fit_transform(X_train
    [X_train.columns])
132           y_train[y_train.columns] = scaler_y.fit_transform(y_train
```

```
         [ y_train . columns ] )
133
134            ## Test portion is tranformed based on the existing
     scaler
135            X_test [ X_test . columns ] = scaler_X . transform ( X_test [ X_test
     . columns ] )
136            y_test [ y_test . columns ] = scaler_y . transform ( y_test [ y_test
     . columns ] )
137
138            #%%
139
140            ## Dataframe for use in correlation analysis , where
     X_train and y_train is
141            ## together
142            df_train = X_train
143            df_train = df_train . merge ( y_train , how='outer' ,
     left_index =True , right_index =True )
144
145            PCA_n = −1
146            hAttrCount = 3
147
148            asel_choice = 'pca'
149            #asel_choice_1 = 'AIC'
150
151            if asel_choice == 'AIC' :
152
153
154                dfs_corr = df_train . corr (method='pearson' )
155                corr_values = dfs_corr [ target ] . to_numpy ( )
156                corr_index = dfs_corr [ target ] . index . to_numpy ( )
157
158                corr_m = np . column_stack (( corr_values , corr_index ))
159
160                for i in range ( len ( corr_m )) :
161                    if np . isnan ( corr_m [ i ,0 ]) :
162                        corr_m [ i ,0 ] = 0
163                    else :
164                        corr_m [ i ,0 ] = np . abs ( corr_m [ i ,0 ])
165
166                corr_m = corr_m [ corr_m [ : ,0 ] . argsort ( ) ]
167                sns . heatmap ( dfs_corr , linewidth =0.5 )
168                plt . show ( )
169
170                corr_n = corr_m [ corr_m [ : , 0] >0.98][:−1 , : ]
171                keep_columns = corr_m [−1−len ( corr_n ) : −1 ,1 ]
172
173                #X_train = X_train [[ keep_columns [1 ] , keep_columns [3 ] ,
     keep_columns [ 4 ]]]
174                #X_test = X_test [[ keep_columns [1 ] , keep_columns [3 ] ,
     keep_columns [ 4 ]]]
175
```

```python
176                    data = itertools.combinations(keep_columns, 3)
177                    keep_columns_list = list(data)
178
179              elif asel_choice == 'pearson':
180
181                    keep_columns_list = []
182                    dfs_corr = df_train.corr(method='pearson')
183                    corr_values = dfs_corr[target].to_numpy()
184                    corr_index = dfs_corr[target].index.to_numpy()
185
186                    corr_m = np.column_stack((corr_values, corr_index))
187
188                    for i in range(len(corr_m)):
189                        if np.isnan(corr_m[i,0]):
190                            corr_m[i,0] = 0
191                        else:
192                            corr_m[i,0] = np.abs(corr_m[i,0])
193
194                    corr_m = corr_m[corr_m[:,0].argsort()]
195
196                    keep_columns = corr_m[-1-hAttrCount:-1,1]
197
198                    keep_columns_list.append(keep_columns)
199
200                    #X_train = X_train[keep_columns]
201                    #X_test  = X_test[keep_columns]
202
203                    sns.heatmap(dfs_corr, linewidth=0.5)
204                    plt.show()
205
206
207
208              elif asel_choice == 'pca':
209                    # PCA is not implemented here but the asel_choice was
       set to 'pca'. How was the code working then
210                    # Was there no dimensionality reduction?
211                    from sklearn.decomposition import PCA
212                    keep_columns = [] #empty for future code
       compatibility
213
214                    PCA_n = hAttrCount
215                    # applied after sensitivity analysis
216
217              ## ppscore based
218
219              elif asel_choice == 'ppscore':
220                    import ppscore as pps
221                    dfs_corr = pps.predictors(df_train, target, output='
       list')
222
223                    keep_columns_list = []
```

```python
224                    corr_values = []
225                    corr_index = []
226                    for i in dfs_corr:
227                        corr_values.append(i['ppscore'])
228                        corr_index.append(i['x'])
229
230
231                    corr_m = np.column_stack((corr_values, corr_index))
232
233
234                    corr_m = corr_m[corr_m[:,0].argsort()]
235
236                    keep_columns = corr_m[-1-hAttrCount:-1,1]
237                    keep_columns_list.append(keep_columns)
238                    #X_train = X_train[keep_columns]
239                    #X_test  = X_test[keep_columns]
240                    sns.heatmap(dfs_corr, linewidth=0.5)
241                    plt.show()
242
243                else:
244                    sys.exit("Error, incorrect attribute selection choice
    ")
245
246                if PCA_n != -1:
247
248                    keep_columns_list = [1]
249                    scaler_pca = MinMaxScaler() # new scaler here because
    PCA can push
250                                               # variables out of (-1,1)
    bounds
251
252                    pca = PCA(n_components = PCA_n)
253
254                    X_train_new = pca.fit_transform(X_train)
255                    X_train_new = scaler_pca.fit_transform(X_train_new)
256
257                    X_test_new = pca.transform(X_test)
258                    X_test_new = scaler_pca.transform(X_test_new)
259
260                print(keep_columns)
261
262
263                #Model Fitting
264
265                #AIC_vals = []
266                #BIC_vals = []
267                #parameters = []
268                mod = 'LSTM'
269                for k in range(len(keep_columns_list)):
270                  if PCA_n != -1:
271                    keep_columns = []
```

```
272                  else:
273                      keep_columns = keep_columns_list[k]
274                      try:
275                          X_train_new = X_train[keep_columns]
276                          X_test_new  = X_test[keep_columns]
277                          parameters.append(keep_columns)
278                      except:
279                          X_train_new = X_train[list(keep_columns)]
280                          X_test_new  = X_test[list(keep_columns)]
281                          parameters.append(np.asarray(list(keep_columns)))
282              #%%

284              ## Data shaping

286              ## from now on, arrays are being morphed into shapes
     valid for RNN+MLP
287              imagination_meters = 25
288              hMemoryMeters = 25
289              lcs_list = []


292              ## if clean == True:
293              step_length = 1
294              memory = int(hMemoryMeters/step_length)
295              print(memory)
296              imagination = int(imagination_meters/step_length)

298              X_attr = list(X_train_new)

300              try:
301                  X_train_new = X_train_new.to_numpy()
302                  X_test_new = X_test_new.to_numpy()
303                  y_train_new = y_train.to_numpy()
304                  y_test_new = y_test.to_numpy()
305              except:
306                  y_train_new = y_train.to_numpy()
307                  y_test_new = y_test.to_numpy()

309              if split != 1:
310                  X_test_new = np.concatenate([X_train_new[-memory
     +1:,:], X_test_new], axis=0)
311                  y_test_new = np.concatenate([y_train_new[-memory
     +1:], y_test_new], axis=0)
312              #%%

314              # If memopry size is 10 then stacks will be as follow:
315                  #Value 1 to value 10
316                  #Value 2 to value 11
317                  #
318                  #
319                  #
```

```python
320                    #Value N to value N+10
321
322            def prepare(data, start, stop, cut_margin = 0, lcs=
        False):
323                    memory = stop−start
324                    stack = []
325                    for i in range(memory):
326                        stack.append(np.roll(data, −i))
327
328                    stack = np.flip(np.rot90(stack), axis=0)[start:−
        memory+1−cut_margin]
329
330                    if lcs == True:
331                        zero = stack[:,0]
332
333                        for j in range(len(zero)):
334                            stack[j] = stack[j] − zero[j]
335                    return stack
336
337            target_lcs_correction = 1
338
339            if target in lcs_list:
340                X_train_RNN = prepare(np.squeeze(y_train_new), 0,
        memory, cut_margin = imagination, lcs=True)
341                X_test_RNN = prepare(np.squeeze(y_test_new), 0,
        memory, cut_margin = imagination, lcs=True)
342
343                y_train_RNN = prepare(np.squeeze(y_train_new),
        memory, memory+imagination, lcs=True)
344                y_test_RNN = prepare(np.squeeze(y_test_new), memory
        , memory+imagination, lcs=True)
345
346                offset_train = X_train_RNN[:,−1]
347                offset_test = X_test_RNN[:,−1]
348
349                for k in range(len(offset_train)):
350                    y_train_RNN[k] = y_train_RNN[k] + offset_train[
        k]
351
352                for k in range(len(offset_test)):
353                    y_test_RNN[k] = y_test_RNN[k] + offset_test[k]
354
355                target_lcs_correction = 1/np.max(y_train_RNN)
356
357                y_train_RNN = y_train_RNN * target_lcs_correction
358                X_train_RNN = X_train_RNN  * target_lcs_correction
359                y_test_RNN = y_test_RNN * target_lcs_correction
360                X_test_RNN = X_test_RNN * target_lcs_correction
361
362            else:
363                X_train_RNN = prepare(np.squeeze(y_train_new), 0,
```

```
          memory, cut_margin = imagination)
364               X_test_RNN = prepare(np.squeeze(y_test_new), 0,
          memory, cut_margin = imagination)
365               y_train_RNN = prepare(np.squeeze(y_train_new),
          memory, memory+imagination)
366               y_test_RNN = prepare(np.squeeze(y_test_new), memory
          , memory+imagination)
367
368          X_train_MLP = []
369          X_test_MLP = []
370
371          #%%
372          if PCA_n == -1:
373              X_lcs_correction = [1]*len(X_train_new[0])
374
375              for i in range(len(X_train_new[0])):
376                  if keep_columns[i] in lcs_list:
377                      X_train_MLP.append(prepare(X_train_new[:,i
          ],memory,memory+imagination, lcs=True))
378                      X_lcs_correction[i] = 1/np.max(X_train_MLP[
          i])
379                      X_train_MLP[i] = X_train_MLP[i]*
          X_lcs_correction[i]
380                  else:
381                      X_train_MLP.append(prepare(X_train_new[:,i
          ],memory,memory+imagination))
382
383              X_train_MLP = np.asarray(X_train_MLP)
384              X_train_MLP = np.concatenate(X_train_MLP[:,:, np.
          newaxis], axis = 1)
385              X_train_MLP = np.rot90(X_train_MLP, axes=(1,2), k
          =3)
386
387
388
389
390
391              for i in range(len(X_test_new[0])):
392                  if keep_columns[i] in lcs_list:
393                      X_test_MLP.append(prepare(X_test_new[:,i],
          memory,memory+imagination, lcs=True))
394                      X_test_MLP[i] = X_test_MLP[i]*
          X_lcs_correction[i]
395                  else:
396                      X_test_MLP.append(prepare(X_test_new[:,i],
          memory,memory+imagination))
397
398              X_test_MLP = np.asarray(X_test_MLP)
399              X_test_MLP = np.concatenate(X_test_MLP[:,:, np.
          newaxis], axis = 1)
400              X_test_MLP = np.rot90(X_test_MLP, axes=(1,2), k=3)
```

```python
401
402                else:
403                    for i in range(len(X_train_new[0])):
404                        X_train_MLP.append(prepare(X_train_new[:,i],
    memory,memory+imagination))
405
406                    X_train_MLP = np.asarray(X_train_MLP)
407                    X_train_MLP = np.concatenate(X_train_MLP[:,:, np.
    newaxis], axis = 1)
408                    X_train_MLP = np.rot90(X_train_MLP, axes=(1,2), k
    =3)
409
410                    for i in range(len(X_test_new[0])):
411                        X_test_MLP.append(prepare(X_test_new[:,i],
    memory,memory+imagination))
412
413                    X_test_MLP = np.asarray(X_test_MLP)
414                    X_test_MLP = np.concatenate(X_test_MLP[:,:, np.
    newaxis], axis = 1)
415                    X_test_MLP = np.rot90(X_test_MLP, axes=(1,2), k=3)
416                #%%
417
418            X_train_RNN_m = X_train_RNN[:,:,np.newaxis]
419            X_train_m = [X_train_RNN_m, X_train_MLP]
420
421
422            X_test_RNN_m = X_test_RNN[:,:,np.newaxis]
423            X_test_m = [X_test_RNN_m, X_test_MLP]
424            #%%
425
426            #%%
427
428            ## Local coordinate system
429
430            ## Just a cumsum on parameter converted to delta
    earlier
431
432
433            #%%
434
435            ## ML model definition
436
437            hDense4 = 1
438            hGRU = 386
439            hDrop1 = 0
440            hDense5 = 128
441            hDrop2 = 0
442            hDrop3 = 0
443            hDense1 = 1
444            hDense2 = 32
445            hDense3 = 128
```

```
446                    hDense4 = 1
447                    hDense5 = 128
448                    h5prefix=''
449                    verbose=0
450                    sensitivity_analysis = False
451                    plot_samples = True
452
453                    #physical_devices = tf.config.list_physical_devices('
       GPU')
454                    #tf.config.experimental.set_memory_growth(
       physical_devices[0], True)
455
456                    tf.keras.backend.clear_session()
457
458                    visible1 = Input(shape=(memory,1))
459
460
461                    visible2 = Input(shape=((imagination),len(X_train_new
       [0])))
462
463                    combined = models(visible1, visible2, mod)
464
465                    z = Dense(hDense3, activation="relu")(combined)
466                    z = Dense(imagination, activation="linear")(z)
467
468
469                    model = Model(inputs=[visible1, visible2], outputs=z)
470
471
472
473
474                    es = EarlyStopping(monitor='val_loss', mode='min',
       verbose=0, patience=50)
475
476                    mc = ModelCheckpoint(f'{h5prefix}best_model.h5',
       monitor='val_loss',
477                                                    mode='min', save_best_only=
       True, verbose=0)
478
479
480                    model.compile(optimizer='adam',loss='mean_squared_error
       ')
481                    #plot_model(model, to_file='model_plot.png',
       show_shapes=True, show_layer_names=True)
482
483                    ## Training
484                    rowcount = len(y_train_RNN)
485                    val_border = int(rowcount*0.85)
486
487                    X_train_m_a = []
488                    X_train_m_b = []
```

```
489
490             X_train_m_a . append ( X_train_m [ 0 ] [ : val_border ] )
491             X_train_m_a . append ( X_train_m [ 1 ] [ : val_border ] )
492
493             X_train_m_b . append ( X_train_m [ 0 ] [ val_border : ] )
494             X_train_m_b . append ( X_train_m [ 1 ] [ val_border : ] )
495
496
497
498             y_train_RNN_a  =  y_train_RNN [ : val_border ]
499             y_train_RNN_b  =  y_train_RNN [ val_border : ]
500
501
502             history  =  model . fit ( X_train_m_a , y_train_RNN_a ,
    validation_data =( X_train_m_b , y_train_RNN_b ) ,
503                                         epochs =2000, verbose=
    verbose ,  batch_size =32,
504                                         callbacks =[ es ,  mc ] )
505
506             model  =  load_model ( f ' { h5prefix } best_model . h5 ' )
507
508
509             result_test  =  model . evaluate ( X_test_m ,  y_test_RNN ,
    verbose =0)
510
511             if target in convert_to_diff :
512                 truth  =  np . cumsum ( y_test_RNN / target_lcs_correction /
    scaler_y . scale_ ,  axis =1)
513                 pred  =  model . predict ( X_test_m )
514                 pred  =  np . cumsum ( pred / target_lcs_correction /
    scaler_y . scale_ ,  axis =1)
515
516             else :
517                 truth  =  y_test_RNN / target_lcs_correction / scaler_y .
    scale_
518                 pred  =  model . predict ( X_test_m )
519                 pred  =  pred / target_lcs_correction / scaler_y . scale_
520
521             if split == 1 or sensitivity_analysis == True : #
    sensitivity enable !
522
523                 y_test_RNN  =  y_train_RNN
524                 X_test_MLP  =  X_train_MLP
525                 X_test_RNN_m  =  X_train_RNN_m
526
527                 senstable  =  { }
528                 plt . style . use ([ ' science ' , ' no-latex ' ])
529                 singular_sensitivity  =  []
530                 if PCA_n != -1:
531                     for i in range ( pca . n_features_ ) :
532                         X_test_MLP_plus  =  shift_pca ( X_test_MLP ,
```

```
533                                              scaler_pca ,
534                                              pca ,
535                                              channel=i ,
536                                              shift=shift )

538                         X_test_m_plus = [X_test_RNN_m ,
    X_test_MLP_plus ]
539                         results_plus = scaler_y.inverse_transform(
    model.predict(X_test_m_plus))

541                         X_test_MLP_minus = shift_pca(X_test_MLP ,
542                                              scaler_pca ,
543                                              pca ,
544                                              channel=i ,
545                                              shift=-shift )
546                         X_test_m_minus = [X_test_RNN_m ,
    X_test_MLP_minus ]
547                         results_minus = scaler_y.inverse_transform(
    model.predict(X_test_m_minus))

549                         if target in convert_to_diff :
550                             results_plus = np.cumsum(results_plus ,
    axis=1)
551                             results_minus = np.cumsum(results_minus
    , axis=1)

553                         sens = (results_plus - results_minus)/2


556                         ave = np.average(sens , axis=0)
557                         perc5 = np.percentile(sens,5,axis=0)
558                         perc25 = np.percentile(sens,25,axis=0)
559                         perc50 = np.percentile(sens,50,axis=0)
560                         perc75 = np.percentile(sens,75,axis=0)
561                         perc95 = np.percentile(sens,95,axis=0)

563                         plt.figure(figsize=(4,3))
564                         #plt.plot(ave , linewidth=2, color='darkblue
    ')
565                         plt.plot(perc5 , linewidth=1, linestyle=":",
     color='black')
566                         plt.plot(perc25 , linewidth=1, color='black'
    )
567                         plt.plot(perc50 , linewidth=2, color='black'
    )
568                         plt.plot(perc75 , linewidth=1, color='black'
    )
569                         plt.plot(perc95 , linewidth=1, linestyle=":"
    , color='black')
570                         plt.title(pca_allattr[i])
571                         plt.grid()
```

```
572                                    plt.tight_layout()
573
574
575                                    plt.plot([],[], linewidth=1, linestyle=":",
        color='black',
576                                            label='$5^{th}$; $95^{th}$
        percentile')
577
578                                    plt.plot([],[], linewidth=1, color='black',
579                                            label='$25^{th}$; $75^{th}$
        percentile')
580
581                                    plt.plot([],[], linewidth=2, color='black',
582                                            label='$50^{th}$ percentile')
583
584                                    plt.legend()
585                                    plt.ylabel(f'Sensitivity Index\n[{
        pca_allattr[i]}]')
586                                    plt.xlabel('Output\nPrediction distance [m]
        ')
587
588                                    myticks = np.linspace(0,len(perc5),6)
589                                    mylabels = np.linspace(0,imagination_meters
        ,6).astype(int)
590                                    plt.xticks(myticks,  mylabels)
591                                    plt.title(f'Average sensitivity = {np.
        average(sens)}')
592                                    senstable[pca_allattr[i]] = np.average(sens
        )
593
594                                    plt.savefig(f'{pca_allattr[i].replace
        ("/","")}.pdf')
595                                    plt.show()
596                                    singular_sensitivity.append(np.average((
        results_plus - results_minus)/2))
597
598                        print(singular_sensitivity)
599                        print(pca_allattr)
600                    else:
601                        for i in range(len(keep_columns)):
602                            X_test_MLP_plus = shift_notpca(X_test_MLP,
603                                                    channel=i,
604                                                    shift=shift)
605                            X_test_m_plus = [X_test_RNN_m,
        X_test_MLP_plus]
606                            results_plus = scaler_y.inverse_transform(
        model.predict(X_test_m_plus))
607
608                            X_test_MLP_minus = shift_notpca(X_test_MLP,
609                                                    channel=i,
610                                                    shift=-shift)
```

```python
611                              X_test_m_minus = [X_test_RNN_m,
      X_test_MLP_minus]
612                              results_minus = scaler_y.inverse_transform(
      model.predict(X_test_m_minus))
613
614                              if target in convert_to_diff:
615                                  results_plus = np.cumsum(results_plus,
      axis=1)
616                                  results_minus = np.cumsum(results_minus
      , axis=1)
617
618
619                              sens = (results_plus - results_minus)/2
620
621
622                              ave = np.average(sens, axis=0)
623                              perc5 = np.percentile(sens,5,axis=0)
624                              perc25 = np.percentile(sens,25,axis=0)
625                              perc50 = np.percentile(sens,50,axis=0)
626                              perc75 = np.percentile(sens,75,axis=0)
627                              perc95 = np.percentile(sens,95,axis=0)
628
629                              plt.figure(figsize=(4,3))
630                              #plt.plot(ave, linewidth=2, color='darkblue
      ')
631                              plt.plot(perc5, linewidth=1, linestyle=":",
       color='black')
632                              plt.plot(perc25, linewidth=1, color='black'
      )
633                              plt.plot(perc50, linewidth=2, color='black'
      )
634                              plt.plot(perc75, linewidth=1, color='black'
      )
635                              plt.plot(perc95, linewidth=1, linestyle=":"
      , color='black')
636                              plt.title(pca_allattr[i])
637                              plt.grid()
638                              plt.tight_layout()
639
640
641                              plt.plot([],[],linewidth=1, linestyle=":",
      color='black',
642                                          label='$5^{th}$; $95^{th}$
      percentile')
643
644                              plt.plot([],[],linewidth=1, color='black',
645                                          label='$25^{th}$; $75^{th}$
      percentile')
646
647                              plt.plot([],[],linewidth=2, color='black',
648                                          label='$50^{th}$ percentile')
```

```python
                                plt.legend()
                                plt.ylabel(f'Sensitivity Index\n[{
    pca_allattr[i]}]')
                                plt.xlabel('Output\nPrediction distance [m]
    ')

                                myticks = np.linspace(0,len(perc5),6)
                                mylabels = np.linspace(0,imagination_meters
    ,6).astype(int)
                                plt.xticks(myticks, mylabels)
                                plt.title(f'Average sensitivity = {np.
    average(sens)}')
                                senstable[pca_allattr[i]] = np.average(sens
    )
                                plt.savefig(f'{pca_allattr[i].replace
    ("/","")}.pdf')
                                plt.show()
                                singular_sensitivity.append(np.average((
    results_plus - results_minus)/2))
                        print(singular_sensitivity)
                        print(keep_columns)




                ## Sensitivity for RNN input channel

                X_test_m_plus = [X_test_RNN_m + 0.1, X_test_MLP]
                results_plus = scaler_y.inverse_transform(model.
    predict(X_test_m_plus))

                X_test_m_minus = [X_test_RNN_m - 0.1, X_test_MLP]
                results_minus = scaler_y.inverse_transform(model.
    predict(X_test_m_minus))

                # if target in convert_to_diff:
                #     results_plus = np.cumsum(results_plus, axis
    =1)
                #     results_minus = np.cumsum(results_minus, axis
    =1)


                sens = (results_plus - results_minus)/2


                ave = np.average(sens, axis=0)
                perc5 = np.percentile(sens,5, axis=0)
                perc25 = np.percentile(sens,25, axis=0)
                perc50 = np.percentile(sens,50, axis=0)
                perc75 = np.percentile(sens,75, axis=0)
```

```python
689                    perc95 = np.percentile(sens,95,axis=0)
690
691                    plt.figure(figsize=(4,3))
692                    #plt.plot(ave, linewidth=2, color='darkblue')
693                    plt.plot(perc5, linewidth=1, linestyle=":", color='
       black')
694                    plt.plot(perc25, linewidth=1, color='black')
695                    plt.plot(perc50, linewidth=2, color='black')
696                    plt.plot(perc75, linewidth=1, color='black')
697                    plt.plot(perc95, linewidth=1, linestyle=":", color=
       'black')
698                    plt.title("RNN Input sensitivity, full channel")
699                    plt.grid()
700                    plt.tight_layout()
701
702
703                    plt.plot([],[], linewidth=1, linestyle=":", color='
       black',
704                             label='$5^{th}$; $95^{th}$ percentile')
705
706                    plt.plot([],[], linewidth=1, color='black',
707                             label='$25^{th}$; $75^{th}$ percentile')
708
709                    plt.plot([],[], linewidth=2, color='black',
710                             label='$50^{th}$ percentile')
711
712                    plt.legend()
713                    plt.ylabel(f'Sensitivity Index\n[{target}]')
714                    plt.xlabel('Output\nPrediction distance [m]')
715
716                    myticks = np.linspace(0,len(perc5),6)
717                    mylabels = np.linspace(0,hMemoryMeters,6).astype(
       int)
718                    plt.xticks(myticks, mylabels)
719
720                    plt.title(f'Average sensitivity = {np.average(sens)
       }')
721
722                    senstable["RNN"] = np.average(sens)
723                    plt.savefig(f'1.pdf')
724                    plt.show()
725
726
727                    singular_sens_input = []
728                    for i in range(len(X_test_RNN_m[0])):
729                        print(".", end="")
730                        X_test_RNN_m_plus = X_test_RNN_m.copy()
731
732                        localrange = np.abs(np.max(X_test_RNN_m_plus[:,
       i]) - np.min(X_test_RNN_m_plus[:,i]))
733
```

```
734                         X_test_RNN_m_plus[:,i] = X_test_RNN_m_plus[:,i]
       + 0.1*localrange
735                         X_test_m_plus = [X_test_RNN_m_plus, X_test_MLP]
736                         results_plus = scaler_y.inverse_transform(model
       .predict(X_test_m_plus))
737
738                         X_test_RNN_m_minus = X_test_RNN_m.copy()
739                         X_test_RNN_m_minus[:,i] = X_test_RNN_m_minus[:,
       i] - 0.1*localrange
740                         X_test_m_minus = [X_test_RNN_m_minus,
       X_test_MLP]
741                         results_minus = scaler_y.inverse_transform(
       model.predict(X_test_m_minus))
742
743                         sens = (results_plus - results_minus)/2
744
745                         singular_sens_input.append(np.percentile(sens
       ,50, axis=0))
746
747                     plt.figure(figsize=(4,3))
748                     vspread = np.max(np.abs(singular_sens_input))
749                     sns.heatmap(np.rot90(singular_sens_input), vmin = -
       vspread, vmax = vspread,
750                                 cmap="vlag",
751                                 cbar_kws={'label': 'Sensitivity Index'
       })
752
753
754
755
756                     len1 = len(np.rot90(singular_sens_input))
757                     len2 = len(np.rot90(singular_sens_input)[0])
758
759                     plt.xticks(np.linspace(0,len2,6),
760                             np.linspace(-hMemoryMeters,0,6).astype(
       int))
761                     plt.yticks(np.linspace(0,len1,6),
762                             np.linspace(0,imagination_meters,6).
       astype(int))
763                     plt.xlabel('RNN memory location [m]')
764                     plt.ylabel('Prediction distance [m]')
765                     plt.title(f'Average sensitivity = {np.average(sens)
       }')
766                     plt.savefig('2.pdf')
767
768
769                     plt.show()
770                     plt.figure(figsize=(4,3))
771                     plt.plot(np.mean(singular_sens_input, axis=0), c='
       black')
772                     plt.xticks(np.linspace(0,len1,6),
```

```
773                          np.linspace(0,imagination_meters,6).
     astype(int))
774                 plt.xlabel('Prediction distance [m]')
775                 plt.ylabel('Sensitivity Index')
776                 plt.grid()
777                 plt.title(f'Average sensitivity = {np.average(sens)
     }')
778                 plt.savefig('3.pdf')
779                 plt.show()
780
781                 plt.figure(figsize=(4,3))
782                 plt.plot(np.mean(singular_sens_input, axis=1), c='
     black')
783                 plt.ylabel('Sensitivity Index')
784                 plt.xlabel('RNN memory location [m]')
785                 plt.xticks(np.linspace(0,len2,6),
786                         np.linspace(-hMemoryMeters,0,6).astype(
     int))
787                 plt.grid()
788                 plt.title(f'Average sensitivity = {np.average(sens)
     }')
789                 plt.savefig('4.pdf')
790                 plt.show()
791
792
793
794             # Plots
795             if plot_samples==True:
796                 pred = scaler_y.inverse_transform(model.predict(
     X_train_m, verbose=0))
797
798                 if target in convert_to_diff:
799
800                     xtr = np.cumsum(scaler_y.inverse_transform(
     X_train_RNN), axis=1)
801                     off = np.rot90(np.tile(xtr[:,-1], (len(pred[0])
     ,1) ), 3)
802
803                     pred = np.cumsum(pred, axis=1) + off
804                     ytr = np.cumsum(scaler_y.inverse_transform(
     y_train_RNN), axis=1) + off
805                 else:
806                     xtr = scaler_y.inverse_transform(X_train_RNN)
807                     ytr = scaler_y.inverse_transform(y_train_RNN)
808
809
810                 #for i in range(10):
811                 #     s = np.random.randint(0, len(y_train_RNN))
812
813                 #     x = np.arange(0,len(X_train_RNN[0]),1)
814
```

```python
815                    #    plt.title('Train')
816                    #    plt.plot(x, xtr[s], label='RNN input')
817
818
819                    #    x = np.arange(len(X_train_RNN[0]), len(
     X_train_RNN[0]) + len(y_train_RNN[0]),1)
820                    #    plt.plot(x,ytr[s], label='RNN output, true')
821
822
823
824                    #    plt.plot(x,pred[s], label='RNN output,
     predicted')
825                    #    plt.legend()
826
827
828                    #    plt.show()
829
830                    pred = scaler_y.inverse_transform(model.predict(
     X_test_m))
831
832                    if target in convert_to_diff:
833
834                        xts = np.cumsum(scaler_y.inverse_transform(
     X_test_RNN), axis=1)
835                        off = np.rot90(np.tile(xts[:,-1], (len(pred[0])
     ,1) ), 3)
836
837                        pred = np.cumsum(pred, axis=1) + off
838                        yts = np.cumsum(scaler_y.inverse_transform(
     y_test_RNN), axis=1) + off
839                    else:
840                        xts = scaler_y.inverse_transform(X_test_RNN)
841                        yts = scaler_y.inverse_transform(y_test_RNN)
842
843
844                    #for i in range(5):
845                    s = np.random.randint(0, len(y_test_RNN))
846
847                    x = np.arange(0,len(X_test_RNN[0]),1)
848
849                    plt.plot(x, xts[s], label='RNN input')
850
851                    x = np.arange(len(X_test_RNN[0]), len(X_test_RNN
     [0]) + len(y_test_RNN[0]),1)
852                    plt.plot(x,yts[s], label='RNN output, true')
853
854
855                    plt.title('test')
856                    plt.plot(x,pred[s], label='RNN output, predicted')
857                    plt.legend()
858                    plt.show()
```

```python
                if np.isnan(result_test):
                    result_test = 0
                #print(-np.log10(result_test))

                if PCA_n != -1:
                    keep_columns = pca_allattr

                print(f'MAE: {np.average(np.abs(truth-pred))}')

                if split == 1 or sensitivity_analysis == True:
                    print(truth, pred, keep_columns, -np.log10(
    result_test), senstable)

                else:
                    print(truth, pred, keep_columns, -np.log10(
    result_test))

                difference = truth-pred

                #####
                # AIC and BIC
                ssd = np.sum(difference ** 2)
                AIC_wi = 2*model.count_params()+len(difference)*math.
    log(ssd/len(difference))
                BIC_wi = np.log(len(difference))*model.count_params()+
    len(difference)*math.log(ssd/len(difference))
                print(AIC_wi)
                print(BIC_wi)
                AIC_vals.append(AIC_wi)
                BIC_vals.append(BIC_wi)
                mae.append(np.average(np.abs(truth-pred)))
                log_res.append(-np.log10(result_test))

 from google.colab import files
 df_res = pd.concat((pd.DataFrame(parameters), pd.DataFrame(
    AIC_vals), pd.DataFrame(BIC_vals), pd.DataFrame(mae), pd.DataFrame
    (log_res)), axis = 1)
 #if PCA_n ==-1:
 #   df_res.columns = ['P1', 'P2', 'P3', 'AIC', 'BIC']
 #else:
 #   df_res.columns = ['AIC', 'BIC']
 if resample!='no':
    df_res.columns = ['P1', 'P2', 'P3', 'P4', 'AIC', 'BIC', 'MAE',
    'Log Res']
 else:
    df_res.columns = ['AIC', 'BIC']
 df_res.to_csv('IC_results.csv')
 files.download('IC_results.csv')
```

902

**Listing B.6:** Resampling Optimization Loop

## B.2   Parameter Selection

```python
AIC_vals = []
BIC_vals = []
parameters = []
mod = 'LSTM'
for k in range(len(keep_columns_list)):
    if PCA_n != -1:
        keep_columns = []
    else:
        keep_columns = keep_columns_list[k]
        try:
            X_train_new = X_train[keep_columns]
            X_test_new  = X_test[keep_columns]
            parameters.append(keep_columns)
        except:
            X_train_new = X_train[list(keep_columns)]
            X_test_new  = X_test[list(keep_columns)]
            parameters.append(np.asarray(list(keep_columns)))
    #%%

    ## Data shaping

    ## from now on, arrays are being morphed into shapes valid for
    RNN+MLP
    imagination_meters = 25
    hMemoryMeters = 25
    lcs_list = []


    ## if clean == True:
    step_length = 1
    memory = int(hMemoryMeters/step_length)
    print(memory)
    imagination = int(imagination_meters/step_length)

    X_attr = list(X_train_new)

    try:
        X_train_new = X_train_new.to_numpy()
        X_test_new = X_test_new.to_numpy()
        y_train_new = y_train.to_numpy()
        y_test_new = y_test.to_numpy()
    except:
```

```python
43            y_train_new = y_train.to_numpy()
44            y_test_new = y_test.to_numpy()
45
46        if split != 1:
47            X_test_new = np.concatenate([X_train_new[-memory+1:,:],
      X_test_new], axis=0)
48            y_test_new = np.concatenate([y_train_new[-memory+1:],
      y_test_new], axis=0)
49        #%%
50
51        # If memopry size is 10 then stacks will be as follow:
52            #Value 1 to value 10
53            #Value 2 to value 11
54            #
55            #
56            #
57            #Value N to value N+10
58
59        def prepare(data, start, stop, cut_margin = 0, lcs=False):
60            memory = stop-start
61            stack = []
62            for i in range(memory):
63                stack.append(np.roll(data, -i))
64
65            stack = np.flip(np.rot90(stack), axis=0)[start:-memory+1-
      cut_margin]
66
67            if lcs == True:
68                zero = stack[:,0]
69
70                for j in range(len(zero)):
71                    stack[j] = stack[j] - zero[j]
72            return stack
73
74        target_lcs_correction = 1
75
76        if target in lcs_list:
77            X_train_RNN = prepare(np.squeeze(y_train_new), 0, memory,
      cut_margin = imagination, lcs=True)
78            X_test_RNN = prepare(np.squeeze(y_test_new), 0, memory,
      cut_margin = imagination, lcs=True)
79
80            y_train_RNN = prepare(np.squeeze(y_train_new), memory,
      memory+imagination, lcs=True)
81            y_test_RNN = prepare(np.squeeze(y_test_new), memory, memory
      +imagination, lcs=True)
82
83            offset_train = X_train_RNN[:,-1]
84            offset_test = X_test_RNN[:,-1]
85
86            for k in range(len(offset_train)):
```

```
87              y_train_RNN[k] = y_train_RNN[k] + offset_train[k]
88
89          for k in range(len(offset_test)):
90              y_test_RNN[k] = y_test_RNN[k] + offset_test[k]
91
92          target_lcs_correction = 1/np.max(y_train_RNN)
93
94          y_train_RNN = y_train_RNN * target_lcs_correction
95          X_train_RNN = X_train_RNN  * target_lcs_correction
96          y_test_RNN = y_test_RNN * target_lcs_correction
97          X_test_RNN = X_test_RNN * target_lcs_correction
98
99      else:
100         X_train_RNN = prepare(np.squeeze(y_train_new), 0, memory,
    cut_margin = imagination)
101         X_test_RNN = prepare(np.squeeze(y_test_new), 0, memory,
    cut_margin = imagination)
102         y_train_RNN = prepare(np.squeeze(y_train_new), memory,
    memory+imagination)
103         y_test_RNN = prepare(np.squeeze(y_test_new), memory, memory
    +imagination)
104
105     X_train_MLP = []
106     X_test_MLP = []
107
108     #%%
109     if PCA_n == -1:
110         X_lcs_correction = [1]*len(X_train_new[0])
111
112         for i in range(len(X_train_new[0])):
113             if keep_columns[i] in lcs_list:
114                 X_train_MLP.append(prepare(X_train_new[:,i],memory,
    memory+imagination, lcs=True))
115                 X_lcs_correction[i] = 1/np.max(X_train_MLP[i])
116                 X_train_MLP[i] = X_train_MLP[i]*X_lcs_correction[i]
117             else:
118                 X_train_MLP.append(prepare(X_train_new[:,i],memory,
    memory+imagination))
119
120         X_train_MLP = np.asarray(X_train_MLP)
121         X_train_MLP = np.concatenate(X_train_MLP[:,:, np.newaxis],
    axis = 1)
122         X_train_MLP = np.rot90(X_train_MLP, axes=(1,2), k=3)
123
124
125
126
127
128         for i in range(len(X_test_new[0])):
129             if keep_columns[i] in lcs_list:
130                 X_test_MLP.append(prepare(X_test_new[:,i],memory,
```

```python
                memory+imagination, lcs=True))
131                     X_test_MLP[i] = X_test_MLP[i]*X_lcs_correction[i]
132                 else:
133                     X_test_MLP.append(prepare(X_test_new[:,i],memory,
    memory+imagination))
134
135             X_test_MLP = np.asarray(X_test_MLP)
136             X_test_MLP = np.concatenate(X_test_MLP[:,:, np.newaxis],
    axis = 1)
137             X_test_MLP = np.rot90(X_test_MLP, axes=(1,2), k=3)
138
139         else:
140             for i in range(len(X_train_new[0])):
141                 X_train_MLP.append(prepare(X_train_new[:,i],memory,
    memory+imagination))
142
143             X_train_MLP = np.asarray(X_train_MLP)
144             X_train_MLP = np.concatenate(X_train_MLP[:,:, np.newaxis],
    axis = 1)
145             X_train_MLP = np.rot90(X_train_MLP, axes=(1,2), k=3)
146
147             for i in range(len(X_test_new[0])):
148                 X_test_MLP.append(prepare(X_test_new[:,i],memory,memory
    +imagination))
149
150             X_test_MLP = np.asarray(X_test_MLP)
151             X_test_MLP = np.concatenate(X_test_MLP[:,:, np.newaxis],
    axis = 1)
152             X_test_MLP = np.rot90(X_test_MLP, axes=(1,2), k=3)
153     #%%
154
155     X_train_RNN_m = X_train_RNN[:,:,np.newaxis]
156     X_train_m = [X_train_RNN_m, X_train_MLP]
157
158
159     X_test_RNN_m = X_test_RNN[:,:,np.newaxis]
160     X_test_m = [X_test_RNN_m, X_test_MLP]
161     #%%
162
163     #%%
164
165     ## Local coordinate system
166
167     ## Just a cumsum on parameter converted to delta earlier
168
169
170     #%%
171
172     ## ML model definition
173
174     hDense4 = 1
```

```
175     hGRU = 386
176     hDrop1 = 0
177     hDense5 = 128
178     hDrop2 = 0
179     hDrop3 = 0
180     hDense1 = 1
181     hDense2 = 32
182     hDense3 = 128
183     hDense4 = 1
184     hDense5 = 128
185     h5prefix=''
186     verbose=0
187     sensitivity_analysis = False
188     plot_samples = True
189
190     #physical_devices = tf.config.list_physical_devices('GPU')
191     #tf.config.experimental.set_memory_growth(physical_devices[0],
    True)
192
193     tf.keras.backend.clear_session()
194
195     visible1 = Input(shape=(memory,1))
196
197
198     visible2 = Input(shape=((imagination),len(X_train_new[0])))
199
200     combined = models(visible1, visible2, mod)
201
202     z = Dense(hDense3, activation="relu")(combined)
203     z = Dense(imagination, activation="linear")(z)
204
205
206     model = Model(inputs=[visible1, visible2], outputs=z)
207
208
209
210
211     es = EarlyStopping(monitor='val_loss', mode='min', verbose=0,
    patience=50)
212
213     mc = ModelCheckpoint(f'{h5prefix}best_model.h5', monitor='
    val_loss',
214                          mode='min', save_best_only=True,
    verbose=0)
215
216
217     model.compile(optimizer='adam',loss='mean_squared_error')
218     #plot_model(model, to_file='model_plot.png', show_shapes=True,
    show_layer_names=True)
219
220     ## Training
```

```python
221        rowcount = len(y_train_RNN)
222        val_border = int(rowcount*0.85)
223
224        X_train_m_a = []
225        X_train_m_b = []
226
227        X_train_m_a.append(X_train_m[0][:val_border])
228        X_train_m_a.append(X_train_m[1][:val_border])
229
230        X_train_m_b.append(X_train_m[0][val_border:])
231        X_train_m_b.append(X_train_m[1][val_border:])
232
233
234
235        y_train_RNN_a = y_train_RNN[:val_border]
236        y_train_RNN_b = y_train_RNN[val_border:]
237
238
239        history = model.fit(X_train_m_a, y_train_RNN_a, validation_data=(
     X_train_m_b, y_train_RNN_b),
240                                            epochs=2000, verbose=verbose,
     batch_size=32,
241                                            callbacks=[es, mc])
242
243    model = load_model(f'{h5prefix}best_model.h5')
244
245
246    result_test = model.evaluate(X_test_m, y_test_RNN, verbose=0)
247
248    if target in convert_to_diff:
249        truth = np.cumsum(y_test_RNN/target_lcs_correction/scaler_y
     .scale_, axis=1)
250        pred = model.predict(X_test_m)
251        pred = np.cumsum(pred/target_lcs_correction/scaler_y.scale_
     , axis=1)
252
253    else:
254        truth = y_test_RNN/target_lcs_correction/scaler_y.scale_
255        pred = model.predict(X_test_m)
256        pred = pred/target_lcs_correction/scaler_y.scale_
257
258    if split == 1 or sensitivity_analysis == True: # sensitivity
     enable!
259
260        y_test_RNN = y_train_RNN
261        X_test_MLP = X_train_MLP
262        X_test_RNN_m = X_train_RNN_m
263
264        senstable = {}
265        plt.style.use(['science','no-latex'])
266        singular_sensitivity = []
```

```python
          if PCA_n != -1:
              for i in range(pca.n_features_):
                  X_test_MLP_plus = shift_pca(X_test_MLP,
                                              scaler_pca,
                                              pca,
                                              channel=i,
                                              shift=shift)

                  X_test_m_plus = [X_test_RNN_m, X_test_MLP_plus]
                  results_plus = scaler_y.inverse_transform(model.
    predict(X_test_m_plus))

                  X_test_MLP_minus = shift_pca(X_test_MLP,
                                               scaler_pca,
                                               pca,
                                               channel=i,
                                               shift=-shift)
                  X_test_m_minus = [X_test_RNN_m, X_test_MLP_minus]
                  results_minus = scaler_y.inverse_transform(model.
    predict(X_test_m_minus))

                  if target in convert_to_diff:
                      results_plus = np.cumsum(results_plus, axis=1)
                      results_minus = np.cumsum(results_minus, axis
    =1)

                  sens = (results_plus - results_minus)/2


                  ave = np.average(sens, axis=0)
                  perc5 = np.percentile(sens,5,axis=0)
                  perc25 = np.percentile(sens,25,axis=0)
                  perc50 = np.percentile(sens,50,axis=0)
                  perc75 = np.percentile(sens,75,axis=0)
                  perc95 = np.percentile(sens,95,axis=0)

                  plt.figure(figsize=(4,3))
                  #plt.plot(ave, linewidth=2, color='darkblue')
                  plt.plot(perc5, linewidth=1, linestyle=":", color='
    black')
                  plt.plot(perc25, linewidth=1, color='black')
                  plt.plot(perc50, linewidth=2, color='black')
                  plt.plot(perc75, linewidth=1, color='black')
                  plt.plot(perc95, linewidth=1, linestyle=":", color=
    'black')
                  plt.title(pca_allattr[i])
                  plt.grid()
                  plt.tight_layout()


                  plt.plot([],[],linewidth=1, linestyle=":", color='
```

```python
                black',
                                    label='$5^{th}$; $95^{th}$ percentile')

                    plt.plot([],[], linewidth=1, color='black',
                            label='$25^{th}$; $75^{th}$ percentile')

                    plt.plot([],[], linewidth=2, color='black',
                            label='$50^{th}$ percentile')

                    plt.legend()
                    plt.ylabel(f'Sensitivity Index\n[{pca_allattr[i]}]'
        )
                    plt.xlabel('Output\nPrediction distance [m]')

                    myticks = np.linspace(0,len(perc5),6)
                    mylabels = np.linspace(0,imagination_meters,6).
        astype(int)
                    plt.xticks(myticks, mylabels)
                    plt.title(f'Average sensitivity = {np.average(sens)
        }')
                    senstable[pca_allattr[i]] = np.average(sens)

                    plt.savefig(f'{pca_allattr[i].replace("/","")}.pdf'
        )
                    plt.show()
                    singular_sensitivity.append(np.average((
        results_plus - results_minus)/2))

                print(singular_sensitivity)
                print(pca_allattr)
            else:
                for i in range(len(keep_columns)):
                    X_test_MLP_plus = shift_notpca(X_test_MLP,
                                            channel=i,
                                            shift=shift)
                    X_test_m_plus = [X_test_RNN_m, X_test_MLP_plus]
                    results_plus = scaler_y.inverse_transform(model.
        predict(X_test_m_plus))

                    X_test_MLP_minus = shift_notpca(X_test_MLP,
                                            channel=i,
                                            shift=-shift)
                    X_test_m_minus = [X_test_RNN_m, X_test_MLP_minus]
                    results_minus = scaler_y.inverse_transform(model.
        predict(X_test_m_minus))

                    if target in convert_to_diff:
                        results_plus = np.cumsum(results_plus, axis=1)
                        results_minus = np.cumsum(results_minus, axis
        =1)
```

```
355
356                    sens = ( results_plus − results_minus )/2
357
358
359                    ave = np.average(sens, axis=0)
360                    perc5 = np.percentile(sens,5,axis=0)
361                    perc25 = np.percentile(sens,25,axis=0)
362                    perc50 = np.percentile(sens,50,axis=0)
363                    perc75 = np.percentile(sens,75,axis=0)
364                    perc95 = np.percentile(sens,95,axis=0)
365
366                    plt.figure(figsize=(4,3))
367                    #plt.plot(ave, linewidth=2, color='darkblue')
368                    plt.plot(perc5, linewidth=1, linestyle=":", color='
        black')
369                    plt.plot(perc25, linewidth=1, color='black')
370                    plt.plot(perc50, linewidth=2, color='black')
371                    plt.plot(perc75, linewidth=1, color='black')
372                    plt.plot(perc95, linewidth=1, linestyle=":", color=
        'black')
373                    plt.title(pca_allattr[i])
374                    plt.grid()
375                    plt.tight_layout()
376
377
378                    plt.plot([],[],linewidth=1, linestyle=":", color='
        black',
379                            label='$5^{th}$; $95^{th}$ percentile')
380
381                    plt.plot([],[],linewidth=1, color='black',
382                            label='$25^{th}$; $75^{th}$ percentile')
383
384                    plt.plot([],[],linewidth=2, color='black',
385                            label='$50^{th}$ percentile')
386
387                    plt.legend()
388                    plt.ylabel(f'Sensitivity Index\n[{pca_allattr[i]}]'
        )
389                    plt.xlabel('Output\nPrediction distance [m]')
390
391                    myticks = np.linspace(0,len(perc5),6)
392                    mylabels = np.linspace(0,imagination_meters,6).
        astype(int)
393                    plt.xticks(myticks, mylabels)
394                    plt.title(f'Average sensitivity = {np.average(sens)
        }')
395                    senstable[pca_allattr[i]] = np.average(sens)
396                    plt.savefig(f'{pca_allattr[i].replace("/","")}.pdf'
        )
397                    plt.show()
398                    singular_sensitivity.append(np.average((
```

```python
        results_plus - results_minus )/2))
            print( singular_sensitivity )
            print( keep_columns )



        ## Sensitivity for RNN input channel

        X_test_m_plus = [X_test_RNN_m + 0.1, X_test_MLP]
        results_plus = scaler_y.inverse_transform(model.predict(
    X_test_m_plus ))

        X_test_m_minus = [X_test_RNN_m - 0.1, X_test_MLP]
        results_minus = scaler_y.inverse_transform(model.predict(
    X_test_m_minus ))

        # if target in convert_to_diff:
        #       results_plus = np.cumsum(results_plus, axis=1)
        #       results_minus = np.cumsum(results_minus, axis=1)


        sens = (results_plus - results_minus )/2


        ave = np.average(sens, axis=0)
        perc5 = np.percentile(sens,5, axis=0)
        perc25 = np.percentile(sens,25, axis=0)
        perc50 = np.percentile(sens,50, axis=0)
        perc75 = np.percentile(sens,75, axis=0)
        perc95 = np.percentile(sens,95, axis=0)

        plt.figure(figsize=(4,3))
        #plt.plot(ave, linewidth=2, color='darkblue')
        plt.plot(perc5, linewidth=1, linestyle=":", color='black')
        plt.plot(perc25, linewidth=1, color='black')
        plt.plot(perc50, linewidth=2, color='black')
        plt.plot(perc75, linewidth=1, color='black')
        plt.plot(perc95, linewidth=1, linestyle=":", color='black')
        plt.title("RNN Input sensitivity, full channel")
        plt.grid()
        plt.tight_layout()


        plt.plot([],[], linewidth=1, linestyle=":", color='black',
                label='$5^{th}$; $95^{th}$ percentile')

        plt.plot([],[], linewidth=1, color='black',
                label='$25^{th}$; $75^{th}$ percentile')

        plt.plot([],[], linewidth=2, color='black',
```

```
447                        label='$50^{th}$ percentile')

449              plt.legend()
450              plt.ylabel(f'Sensitivity Index\n[{target}]')
451              plt.xlabel('Output\nPrediction distance [m]')

453              myticks = np.linspace(0,len(perc5),6)
454              mylabels = np.linspace(0,hMemoryMeters,6).astype(int)
455              plt.xticks(myticks, mylabels)

457              plt.title(f'Average sensitivity = {np.average(sens)}')

459              senstable["RNN"] = np.average(sens)
460              plt.savefig(f'1.pdf')
461              plt.show()


464              singular_sens_input = []
465              for i in range(len(X_test_RNN_m[0])):
466                  print(".", end="")
467                  X_test_RNN_m_plus = X_test_RNN_m.copy()

469                  localrange = np.abs(np.max(X_test_RNN_m_plus[:,i]) - np
    .min(X_test_RNN_m_plus[:,i]))

471                  X_test_RNN_m_plus[:,i] = X_test_RNN_m_plus[:,i] + 0.1*
    localrange
472                  X_test_m_plus = [X_test_RNN_m_plus, X_test_MLP]
473                  results_plus = scaler_y.inverse_transform(model.predict
    (X_test_m_plus))

475                  X_test_RNN_m_minus = X_test_RNN_m.copy()
476                  X_test_RNN_m_minus[:,i] = X_test_RNN_m_minus[:,i] -
    0.1*localrange
477                  X_test_m_minus = [X_test_RNN_m_minus, X_test_MLP]
478                  results_minus = scaler_y.inverse_transform(model.
    predict(X_test_m_minus))

480                  sens = (results_plus - results_minus)/2

482                  singular_sens_input.append(np.percentile(sens,50, axis
    =0))

484              plt.figure(figsize=(4,3))
485              vspread = np.max(np.abs(singular_sens_input))
486              sns.heatmap(np.rot90(singular_sens_input), vmin = -vspread,
     vmax = vspread,
487                          cmap="vlag",
488                          cbar_kws={'label': 'Sensitivity Index'})

490
```

```
491

492

493            len1 = len(np.rot90(singular_sens_input))
494            len2 = len(np.rot90(singular_sens_input)[0])
495

496            plt.xticks(np.linspace(0,len2,6),
497                       np.linspace(-hMemoryMeters,0,6).astype(int))
498            plt.yticks(np.linspace(0,len1,6),
499                       np.linspace(0,imagination_meters,6).astype(int))
500            plt.xlabel('RNN memory location [m]')
501            plt.ylabel('Prediction distance [m]')
502            plt.title(f'Average sensitivity = {np.average(sens)}')
503            plt.savefig('2.pdf')
504

505

506            plt.show()
507            plt.figure(figsize=(4,3))
508            plt.plot(np.mean(singular_sens_input, axis=0), c='black')
509            plt.xticks(np.linspace(0,len1,6),
510                       np.linspace(0,imagination_meters,6).astype(int))
511            plt.xlabel('Prediction distance [m]')
512            plt.ylabel('Sensitivity Index')
513            plt.grid()
514            plt.title(f'Average sensitivity = {np.average(sens)}')
515            plt.savefig('3.pdf')
516            plt.show()
517

518            plt.figure(figsize=(4,3))
519            plt.plot(np.mean(singular_sens_input, axis=1), c='black')
520            plt.ylabel('Sensitivity Index')
521            plt.xlabel('RNN memory location [m]')
522            plt.xticks(np.linspace(0,len2,6),
523                       np.linspace(-hMemoryMeters,0,6).astype(int))
524            plt.grid()
525            plt.title(f'Average sensitivity = {np.average(sens)}')
526            plt.savefig('4.pdf')
527            plt.show()
528

529

530

531        # Plots
532        if plot_samples==True:
533            pred = scaler_y.inverse_transform(model.predict(X_train_m,
    verbose=0))
534

535            if target in convert_to_diff:
536

537                xtr = np.cumsum(scaler_y.inverse_transform(X_train_RNN)
    , axis=1)
538                off = np.rot90(np.tile(xtr[:,-1], (len(pred[0]),1) ),
    3)
```

```
539
540                    pred = np.cumsum(pred, axis=1) + off
541                    ytr = np.cumsum(scaler_y.inverse_transform(y_train_RNN)
        , axis=1) + off
542                else:
543                    xtr = scaler_y.inverse_transform(X_train_RNN)
544                    ytr = scaler_y.inverse_transform(y_train_RNN)
545
546
547            #for i in range(10):
548            #    s = np.random.randint(0, len(y_train_RNN))
549
550            #    x = np.arange(0,len(X_train_RNN[0]),1)
551
552            #    plt.title('Train')
553            #    plt.plot(x, xtr[s], label='RNN input')
554
555
556            #    x = np.arange(len(X_train_RNN[0]), len(X_train_RNN[0])
        + len(y_train_RNN[0]),1)
557            #    plt.plot(x,ytr[s], label='RNN output, true')
558
559
560
561            #    plt.plot(x,pred[s], label='RNN output, predicted')
562            #    plt.legend()
563
564
565            #    plt.show()
566
567            pred = scaler_y.inverse_transform(model.predict(X_test_m))
568
569            if target in convert_to_diff:
570
571                xts = np.cumsum(scaler_y.inverse_transform(X_test_RNN),
        axis=1)
572                off = np.rot90(np.tile(xts[:,-1], (len(pred[0]),1) ),
        3)
573
574                pred = np.cumsum(pred, axis=1) + off
575                yts = np.cumsum(scaler_y.inverse_transform(y_test_RNN),
        axis=1) + off
576            else:
577                xts = scaler_y.inverse_transform(X_test_RNN)
578                yts = scaler_y.inverse_transform(y_test_RNN)
579
580
581            #for i in range(5):
582            s = np.random.randint(0, len(y_test_RNN))
583
584            x = np.arange(0,len(X_test_RNN[0]),1)
```

```python
          plt.plot(x, xts[s], label='RNN input')

          x = np.arange(len(X_test_RNN[0]), len(X_test_RNN[0]) + len(
     y_test_RNN[0]),1)
          plt.plot(x,yts[s], label='RNN output, true')


          plt.title('test')
          plt.plot(x,pred[s], label='RNN output, predicted')
          plt.legend()
          plt.show()


     if np.isnan(result_test):
          result_test = 0
     #print(-np.log10(result_test))

     if PCA_n != -1:
          keep_columns = pca_allattr

     print(f'MAE: {np.average(np.abs(truth-pred))}')

     if split == 1 or sensitivity_analysis == True:
          print(truth, pred, keep_columns, -np.log10(result_test),
     senstable)

     else:
          print(truth, pred, keep_columns, -np.log10(result_test))

     difference = truth-pred

     #####
     # AIC and BIC
     ssd = np.sum(difference ** 2)
     AIC_wi = 2*model.count_params()+len(difference)*math.log(ssd/
     len(difference))
     BIC_wi = np.log(len(difference))*model.count_params()+len(
     difference)*math.log(ssd/len(difference))
     print(AIC_wi)
     print(BIC_wi)
     AIC_vals.append(AIC_wi)
     BIC_vals.append(BIC_wi)

 from google.colab import files
 df_res = pd.concat((pd.DataFrame(parameters), pd.DataFrame(
     AIC_vals), pd.DataFrame(BIC_vals)), axis = 1)
 df_res.columns = ['P1', 'P2', 'P3', 'AIC', 'BIC']
 df_res.to_csv('IC_results.csv')
 files.download('IC_results.csv')
```

630

**Listing B.7:** Parameter Selection Loop