

# Automated Data-driven Approach for Gap Filling in the Time Series using Evolutionary Learning

Mikhail Sarafanov<sup>1</sup>, Nikolay O. Nikitin<sup>1</sup>, and Anna V. Kalyuzhnaya<sup>1</sup>

ITMO University, Saint Petersburg, Russia

**Abstract.** In the paper, we propose an adaptive data-driven model-based approach for filling the gaps in time series. The approach is based on the automated evolutionary identification of the optimal structure for a composite data-driven model. It allows adapting the model for the effective gap-filling in a specific dataset without the involvement of the data scientist. As a case study, both synthetic and real datasets from different fields (environmental, economic, etc) are used. The experiments confirm that the proposed approach allows achieving the higher quality of the gap restoration and improve the effectiveness of forecasting models.

**Keywords:** time series forecasting, gap filling, machine learning, AutoML

## 1 Introduction

Time series is a common way to represent real-world time process data. As an example of the widely-known applications, different cases of time series can be noted: weather stations, financial stocks, industrial sensors, etc. Due to failures of the sensors themselves or the connection issues, the time series may have gaps. The presence of gaps can be a significant problem for time series forecasting since the vast majority of forecasting models can not proceed with gaps in training data.

The existing open-source solutions (e.g. SSGP-Toolbox [12]) provide only relatively simple methods for the gap-filling. But simple methods can be very inaccurate if the size of the gaps is large. In the paper, we consider an approach for filling gaps based on composite models [9]. The composite models consist of several machine learning models and can be generated using automated machine learning methods (AutoML). The use of composite models could potentially decrease the error of the forecast. Since such models with several levels can approximate more complex dependencies in the data than single models do.

The restored time series can be used in time series forecasting tasks in the future. Since the restored parts may differ from the original time series, it becomes hard to approximate relationships between elements in the time series. In this case, the forecast based on incorrectly reconstructed historical time series can be

inaccurate [11]. To investigate this problem, research has also been conducted within the paper.

The paper is organized as follows. The problem statement is described in Sec. 2. The existing approaches and methods for gap filling in the time series are analyzed in Sec. 3. Sec. 4 contains the description of the proposed adaptive approach to the gap filling. The results of the experimental evaluation of the proposed approach for different datasets are described in Sec. 5. The main conclusions are proved in Sec. 6.

## 2 Problem Statement

The analysis of time series is a difficult task if there are a big amount of gaps. If the number of gaps is too large, then trying to restore them will affect the quality of the data-driven forecasting model that is fitted using this data. In practical engineering tasks, it is not the common practice to use time series with the percent of gaps that exceed a certain threshold [13].

A gap filling problem can be seen as an interpolation problem. Classic definition of interpolation function says that values of interpolation function  $\tilde{f}$  should be equal to the values of original function  $f$  on a set of primary points  $x_{\{i\}}$ ,  $i \in \{1, N\}$ .

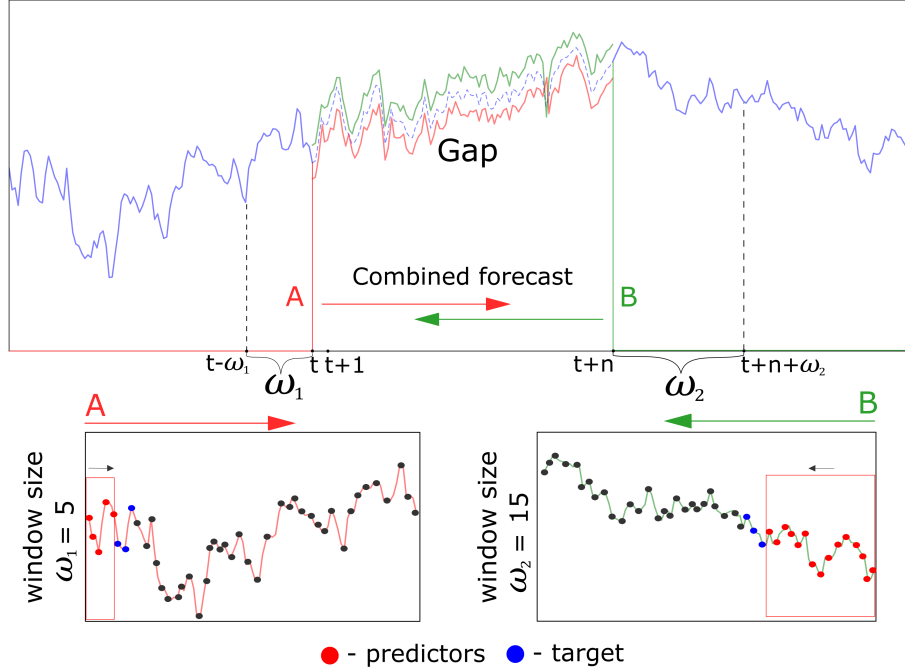
There is a set of points with coordinates on a time axis  $x_{\{G\}}$ ,  $G \in \{t, t + n\}$ , where  $n$  is the number of elements in the gap. In these points we would like to restore values of the time series using information about behavior of the time series before the gap  $x_{\{G\}}$  on the set of points  $x_{\{H\}}$ ,  $H \in \{t - w_1, t - 1\}$  and after – on the set of points  $x_{\{N\}}$ ,  $N \in \{t + n + 1, t + n + w_2\}$ . Here  $w_1$  and  $w_2$  are the time lags that describe time intervals that contain significant information about behavior of the time series on the interval  $x_{\{G\}}$ .

In the current paper, we suggest an automated data-driven approach that allows working with non-stationary time-series and non-linear predictive models for better gap filling. The main idea of this approach is shown in Fig. 1. In this paper, we aim to study the both efficiency of the AutoML-based models for time series restoration and their impact on the error of the forecasting models that are build using this data.

## 3 Related Work

For time series forecasting tasks there are existing models, for example, convolutional and recurrent neural networks, which can eliminate gaps in data during training [1]. However, such models are difficult to set up and effectively train. When building neural network models, the amount of data for training is also an important factor - for short time series, it is hard to train complex deep neural networks. Therefore, there are implemented methods that fill in the gaps before applying predictive models.

There are both simple methods, such as linear interpolation and moving average, and more complex ones, for example, spline interpolation and Kalman



**Fig. 1.** Scheme of the gap-filling approach based on a bi-directional forecasting. A (red curve) - prediction based on the source part of the time series, B (green curve) - prediction based on the inverted part of the time series.  $t$  - time step-index,  $w$  - size of the historical time window for forecasting,  $n$  - forecast length (less or equal to the gap size).

filters [3]. The disadvantages of simple methods are their inaccuracy for long-term omissions. Partly with this problem can cope more complex algorithms include Radial Basis Functions, Moving Least Squares, Adaptive Inverse Distance Weighted, which are better recovering long skips [2]. Especially difficult cases can be considered attempts to restore data in time series, where the percentage of gaps exceeds 30%.

The problem of filling in gaps can also be solved using time series forecasting algorithms. A sequential time-series forecast can be used to fill in the gaps [4]. For such a classical approach autoregressive models, such as AR and ARIMA, can be used. On the other hand, this approach does not take into account the specifics of the gap-filling task. Therefore, a potentially more accurate modification of this approach may be an algorithm that uses both the pre-history and past-history parts of the skip. This approach can be called the forward and backward imputation method [7]. For more efficient time series forecasting, evolutionary algorithms can be used as auxiliary tools for selecting important hyperparameters [5]. Therefore, in the paper, we relied on evolutionary computing to form an effective algorithm for forward and backward imputation.

## 4 Evolutionary design of model-based gap filling approach

In this paper combine the existing approaches and automate forecasting-based gap-filling application to increase the quality of the gap-filling.

### 4.1 Model-based gap filling with ML and AutoML

The main idea of the proposed approach is the following. The problem of gap filling in the time series can be reduced to the well-studied problem of time series forecasting. That is, use only the data before the gap (pre-history) to configure the model, and then apply it to get a sequence of values of the same size as the length of the gap. To build the model, AutoML approaches can be used.

However, in this case, the specifics of the task are not taken into account and the part of the time series after the gap (post-history) is not used. To resolve this issue, we have used the bi-directional approach to restore values (the pseudo-code of the underlying algorithm provided in Alg. 1).

In this case, all possible information available in the time series is used to fill in the gaps. This approach is potentially more accurate but requires more computing resources.

### 4.2 Composite modelling

The data-driven bi-directional model for gap-filling can have a complex structure. Models can be combined into ensembles or stacked into multi-level pipelines, where predictions from one level of models can be predictors for the next level. These structures are naming composite models and can be effectively generated using evolutionary structural learning [9].

The proposed self-adapting gap-filling algorithm is implemented on a basis of the open-source automated modeling framework FEDOT<sup>1</sup>. It allows building the data-driven and hybrid models consist of several atomic blocks [10]. Regression models on lagged features (e.g. lasso regression or K-nearest neighbors) can be used as such blocks.

The evolutionary model design is implemented on a basis of the custom graph-based evolutionary approach. The common pipeline of adaptive evolutionary-based gap filling is presented in Fig. 2.

The Fig. 2 shows that for each time series with gaps (corrupted time series), an algorithm for automatic model design is started to restore the values. During this algorithm's execution, the population with models is initialized. And then the crossover and mutation operators are applied to it. As a result of the structure of the composite model search and tuning the hyperparameters, the final model (Gap filling algorithm) is ready for use.

We implement both basic methods of gap-filling and the proposed model-based approach as parts of the framework. These implementations and different examples of its applications are available as open-source code.

<sup>1</sup> <https://github.com/nccr-itmo/FEDOT>

---

**Algorithm 1:** Pseudocode of the algorithm for the bi-directional gap filling based at composite model, obtained using evolutionary optimisation. The full description of evolutionary part is provided in Fig. 2

---

```

Data: time_series_with_gaps;
 $w \leftarrow$  moving window size;
Result: time_series_without_gaps
gaps  $\leftarrow$  all gap-induced segments in time_series_with_gaps
correct_data  $\leftarrow$  gap-free part of time_series_with_gaps
time_series_without_gaps = time_series_with_gaps
for gap in gaps do
    gap_id  $\leftarrow$  index of gap segment start
    pre_window  $\leftarrow$  subsec(gap, - $w$ )
    post_window  $\leftarrow$  subsec(gap,  $w$ )
    models  $\leftarrow$  generate_population
    for generation in amount_of_generations do
        models  $\leftarrow$  apply_mutation(models)
        models  $\leftarrow$  apply_crossover(models)
        best_models  $\leftarrow$  select_fittest(models, pre_window, post_window)
        models  $\leftarrow$  apply_reproduction(best_models)
    end
    gap_filling_model  $\leftarrow$  select_best_model(correct_data)
    forward_prediction  $\leftarrow$  gap_filling_model.predict(pre_window,)
    backward_prediction  $\leftarrow$  gap_filling_model.predict(post_window)
    prediction  $\leftarrow$  ensemble_model(forward_prediction,
        backward_prediction)
    time_series_without_gaps[gap_id] = prediction
end

```

---

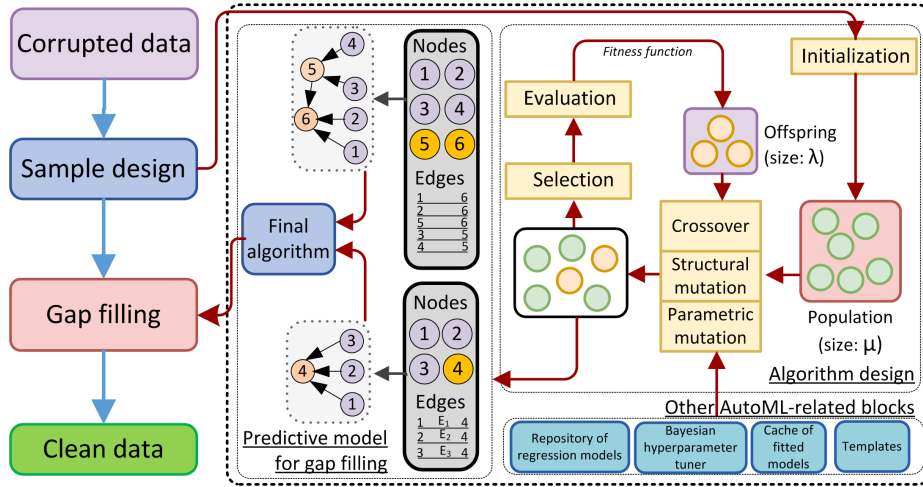
## 5 Experimental Study

To validate the proposed approach, we run a set of experiments using several datasets with different properties. We used an artificially generated (synthetic) time series, a time series of sea surface height with hourly and daily discreteness, time series with air temperature and economic time series.

### 5.1 Datasets

Five time series of different nature were prepared: a synthetic time series (1) sea level (2) time series obtained from the reanalysis grid of satellite altimetry, a sea level (3) time series obtained from simulating the sea surface height in the Arctic Ocean, time series of the air temperature (4), economic time series (5).

However, the synthetic data with desired properties can be obtained with an equation-based model that may be approximately restored from the real-date using algebraic terms approach [6] (if necessary). The length of each obtained



**Fig. 2.** The pipeline of the evolutionary design of gap filling algorithm for specific problem using AutoML-based techniques. The combination of the modelling blocks and data preprocessing blocks allow obtaining the composite predictive model that is used in the algorithm.

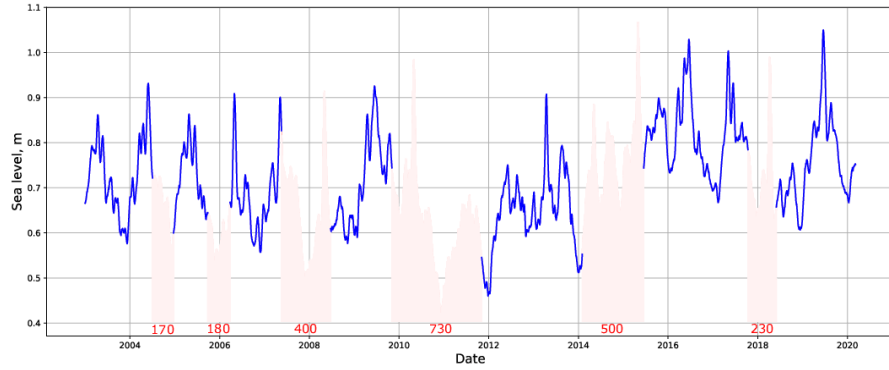
time series is 6276 elements. Synthetic gaps were generated with a total size of 30% of the length of the time series (an example of time series with synthetic gaps can be seen in Fig. 3). Also, the long gap (1500 steps) was generated in the central part of the time series to assess the errors of data recovering for long gaps.

For the reconstructed sections of the series, the predicted values were compared with the actual using Mean Absolute Percentage Error (MAPE) and symmetric MAPE (SMAPE). After gap-filling procedure the reconstructed time series were used to train predictive models, and then a forecast was given. For validation, MAPE measures on the sequence of 20% of the length of the time series were used. The more the error of the forecast increases, the worse the reconstructed time series is suitable for model training.

## 5.2 Experimental setup

There are different algorithms that were applied to the datasets. The list of algorithms, involved in the experimental studies is the following:

- Linear interpolation - as baseline;
- Local approximation by polynomial functions (Savitzky–Golay filters), Batch approximation by polynomial functions - as simple approximation-based approaches. The polynomial batch approximation differs from the local one in that the batch approximation use one polynomial function for the entire interval with gaps, no matter how many elements are missed in it. And with



**Fig. 3.** Examples of generated gaps. The total length of gaps is 2210 elements. Red labels indicate the number of time series elements in the gap.

local approximation, a polynomial is constructed for each gap element (one polynomial function per one element);

- Kalman filter, moving average, spline interpolation - as methods from widely-used library "imputeTS" [8];
- Non-linear time series forecasting model, which iteratively predict missing parts in time series. The pipeline was identified manually and comprise of a lagged transformation and a decision tree model;
- Composite model identified by AutoML - as an approach based on the framework FEDOT.

To make the research more valuable and reproducible, we implemented the evolutionary approach described in the paper as a part of the FEDOT AutoML framework functionality.

### 5.3 Results of the experiments

The average values for each time series are shown for several cases shown in Table 1. The table also contains information about forecasting errors when the restored time series was used.

As can be seen, the best results were obtained based on the composite AutoML model (Fig. 4). The composite model has well reconstructed the phases of fluctuation of the height of the sea surface (the convergence of the structural learning of the composite model is demonstrated in Fig. 5). According to SMAPE measure, the following list ranked from best to the worst was obtained: proposed approach (16.4%), linear interpolation (48.4%), moving average (55.4%), batch approximation by polynomials (64.0%), local approximation by polynomials (110.8%), Kalman filter (115.9%), spline interpolation (231.6%).

The composite model has well reconstructed the phases of fluctuation of the height of the sea surface (the convergence of the structural learning of the composite model is demonstrated in Fig. 5).

**Table 1.** The results for estimating the error of gap-filling algorithms using MAPE measure for time series with different parameters. The results of the error of the predictive model using the reconstructed time series are shown

Algorithm	MAPE					MAPE of forecast on restored series
	Synthetic	Sea level, hourly	Sea level, daily	Temperature	Economic	
Linear interpolation	16.0	20.3	14.4	3.5	193.1	11.6
Local approximation by polynomials	74.4	156.3	181.6	9.5	141.8	13.6
Batch approximation by polynomials	28.5	50.6	57.2	3.6	174.1	12.4
Kalman filter	41.4	105.5	151.4	22.2	247.9	17.3
Moving average	18.4	14.1	25.1	4.0	205.9	11.8
Spline interpolation	223.0	394.6	312.6	31.3	146.5	13.6
Non linear	13.6	17.9	15.1	3.3	136.4	12.0
<b>Proposed approach</b>	<b>11.9</b>	<b>16.7</b>	<b>14.9</b>	<b>2.7</b>	<b>32.4</b>	<b>11.2</b>

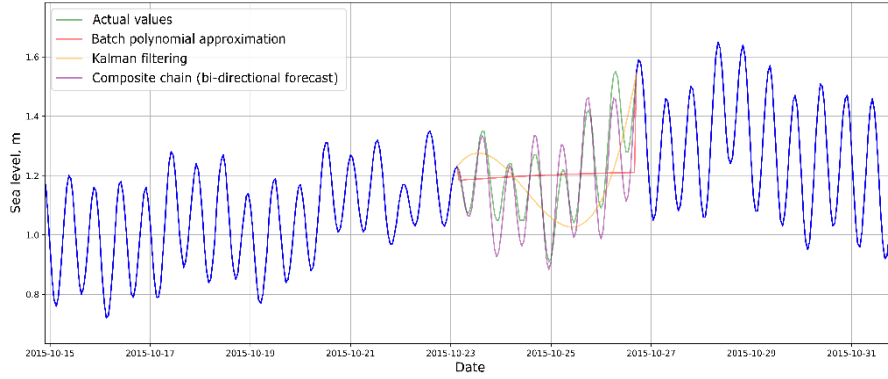
## 6 Conclusions

An approach to fill gaps in time series, using automatic machine learning methods was presented and validated. It was found that the most effective approach was bi-directional forecasting using the evolutionary algorithm with automatic identification of the model. The obtained AutoML-based solution got averaged MAPE of 15.7% and SMAPE of 16.4% for the gap-filling task, while the competitive algorithms could not have error less than 49.5% MAPE (48.4% SMAPE). For a synthetic time series with a breakpoint (when the periodicity of one of the component components changed), the bi-directional approach also can be considered as the better solution.

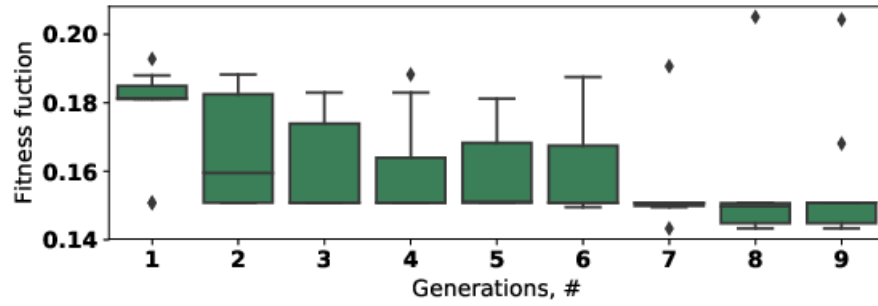
The reconstructed time series were used to fit forecasting models and predict new values of the time series. The fitting forecasting model on time series reconstructed by the proposed approach distorted the forecasts less than using the others series. It confirms that the proposed data-driven approach in conjunction with AutoML techniques allows efficiently recovering gaps in time series.

It is planned to explore the possibilities of using the AutoML approach to restore multivariate time series. It is also planned to improve the performance of the proposed approach.





**Fig. 4.** Examples of time series restoration by different gap-filling approaches for the sea surface height dataset with hourly discreteness.



**Fig. 5.** The convergence of the RMSE-based fitness function during the evolutionary optimization of the composite model structure for the gap-filling task (sea surface height case).

## Code and data availability

All implemented approaches are available in repository <https://github.com/nccr-itmo/FEDOT> as a part of the open-source FEDOT framework and can be used for practical purposes. Data and scripts used to conduct the experiments in the paper are available in the additional repository<sup>2</sup>.

## Acknowledgements

This research is financially supported by The Russian Science Foundation, Agreement #17-71-30029 with cofinancing of Bank Saint Petersburg.

<sup>2</sup> [https://github.com/ITMO-NSS-team/FEDOT-benchmarks/tree/master/experiments/gap\\_filling](https://github.com/ITMO-NSS-team/FEDOT-benchmarks/tree/master/experiments/gap_filling)

## References

1. Che, Z., Purushotham, S., Cho, K., Sontag, D., Liu, Y.: Recurrent neural networks for multivariate time series with missing values. *Scientific reports* 8(1), 1–12 (2018)
2. Ding, Z., Mei, G., Cuomo, S., Li, Y., Xu, N.: Comparison of estimating missing values in iot time series data using different interpolation algorithms. *International Journal of Parallel Programming* 48(3), 534–548 (2020)
3. Lepot, M., Aubin, J.B., Clemens, F.H.: Interpolation in time series: An introductory overview of existing methods, their performance criteria and uncertainty assessment. *Water* 9(10), 796 (2017)
4. Loukopoulos, P., Sampath, S., Pilidis, P., Zolkiewski, G., Bennett, I., Duan, F., Mba, D.: Dealing with missing data for prognostic purposes. In: *2016 Prognostics and System Health Management Conference (PHM-Chengdu)*. pp. 1–5 (2016)
5. Lukoseviciute, K., Ragulskis, M.: Evolutionary algorithms for the selection of time lags for time series forecasting by fuzzy inference systems. *Neurocomputing* 73(10–12), 2077–2088 (2010)
6. Merezhnikov, M., Hvatov, A.: Closed-form algebraic expressions discovery using combined evolutionary optimization and sparse regression approach. *Procedia Computer Science* 178, 424–433 (2020)
7. Moahmed, T.A., El Gayar, N., Atiya, A.F.: Forward and backward forecasting ensembles for the estimation of time series missing data. In: El Gayar, N., Schwenker, F., Suen, C. (eds.) *Artificial Neural Networks in Pattern Recognition*. pp. 93–104. Springer International Publishing, Cham (2014)
8. Moritz, S., Bartz-Beielstein, T.: imputeTS: Time Series Missing Value Imputation in R. *The R Journal* 9(1), 207–218 (2017), <https://doi.org/10.32614/RJ-2017-009>
9. Nikitin, N.O., Polonskaia, I.S., Vychuzhanin, P., Barabanova, I.V., Kalyuzhnaya, A.V.: Structural evolutionary learning for composite classification models. *Procedia Computer Science* 178, 414–423 (2020)
10. Nikitin, N.O., Vychuzhanin, P., Sarafanov, M., Polonskaia, I.S., Revin, I., Barabanova, I.V., Maximov, G., Kalyuzhnaya, A.V., Boukhanovsky, A.: Automated evolutionary approach for the design of composite machine learning pipelines (2021)
11. Saad, M., Chaudhary, M., Karray, F., Gaudet, V.: Machine learning based approaches for imputation in time series data and their impact on forecasting. In: *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. pp. 2621–2627 (2020)
12. Sarafanov, M., Kazakov, E., Nikitin, N.O., Kalyuzhnaya, A.V.: A machine learning approach for remote sensing data gap-filling with open-source implementation: An example regarding land surface temperature, surface albedo and ndvi. *Remote Sensing* 12(23), 3865 (2020)
13. Weigend, A.S.: *Time series prediction: forecasting the future and understanding the past*. Routledge (2018)