
THE DEEPCAR METHOD: FORECASTING TIME-SERIES DATA THAT HAVE CHANGE POINTS

Ayla Jungbluth
Department of Mathematics
Ruhr-University Bochum
Bochum 44801, Germany
ayla.jungbluth@rub.de

Johannes Lederer
Department of Mathematics
Ruhr-University Bochum
Bochum 44801, Germany
johannes.lederer@rub.de

February 23, 2023

ABSTRACT

Many methods for time-series forecasting are known in classical statistics, such as autoregression, moving averages, and exponential smoothing. The DeepAR framework is a novel, recent approach for time-series forecasting based on deep learning. DeepAR has shown very promising results already. However, time series often have change points, which can degrade the DeepAR’s prediction performance substantially. This paper extends the DeepAR framework by detecting and including those change points. We show that our method performs as well as standard DeepAR when there are no change points and considerably better when there are change points. More generally, we show that the batch size provides an effective and surprisingly simple way to deal with change points in DeepAR, Transformers, and other modern forecasting models.

1 Introduction

1.1 The DeepAR Forecasting Algorithm

Time-series forecasting has become textbook material [Box et al., 1974, De Gooijer & Hyndman, 2006], including classical approaches such as moving averages and exponential smoothing [Hyndman et al., 2008]. A much newer development is the use of recurrent neural networks [Graves, 2013, Sutskever et al., 2014]. An architecture particularly suited for this are LSTMs Gers et al. [2001]. A prominent example of time-series forecasting with LSTM-based neural networks is the DeepAR algorithm Alexandrov et al. [2019], Salinas et al. [2020].

One of the main features of DeepAR is that it can combine related time series—rather than considering each time series by itself. The model produces probabilistic forecasts based on Gaussian (or negative binomial) likelihoods [Dagan et al., 2022, Snyder et al., 2012, Zhang et al., 1998]. A key aspect of the DeepAR framework is the selection of the batches in the training process. Multiple training instances are created for each time series by selecting windows with different starting points from the original time series. It is ensured that the entire prediction range is covered when choosing these windows; in particular, the starting point could lie before the beginning of the time series, in which case the unobserved targets are filled with zeros.

At each time step, the next step is to be predicted. The network receives the previous observations together with a set of covariates as input. The information is passed through the hidden layer to the likelihood function. The error is calculated during training using the current parameterization of the likelihood function. Thus, when performing backpropagation, the weights are updated and the values optimized.

After the weights of the network are trained, the forward propagation is performed using the previous input to obtain the distribution parameters μ (mean) and σ (standard deviation) of the Gaussian likelihood.

1.2 Limitation: Change Points

Time-series data often have change points. Frequently, these change points are the result of manual interventions, such as resetting a machine after maintenance or replacing players or coaches in a football team. But sometimes, these change points can also have more systematic reasons, such as the end of a football season. The simplest way to deal with these change points is to ignore the time periods around them, but this approach potentially dismisses large amounts of relevant data. Our question is, therefore, if we can account for change points without losing predictive quality—or even improving prediction.

An example of a time series of machine data with change points can be seen in Figure 2. It represents the hydraulic oil pollution of a shredder machine. The points, where the time series drops from higher to lower values indicate replacements of the machine’s pollution filter. The change points themselves do not have a pattern and, therefore, are not necessarily predictable. The vanilla DeepAR model would try to learn the change points in the time series: indeed,

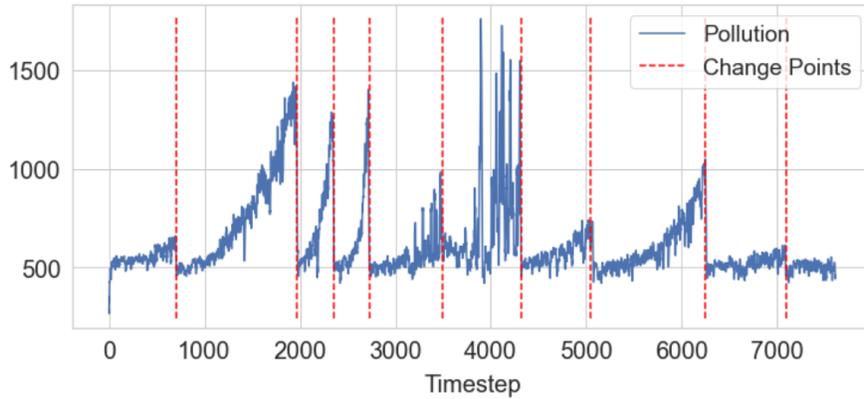


Figure 1: Example of the pollution time series (blue) with manual change points (red). The change points indicate machine maintenance. The behavior of the pollution data shows that it was replaced at times with change points. The change points are not at regular intervals and do not show seasonal behavior since the machine is not used consistently and has irregular deterioration. The change points are located at the time indices 700, 1970, 2350, 2730, 3500, 4320, 5050, 6250, 7100.

the algorithm takes batches from the time series, no matter if there is a change point in it or not. This can deteriorate the forecasts and predict change points that do not exist.

To conclude, DeepAR requires little manual feature engineering and recognizes seasonality, but it does not incorporate change points correctly. More generally, we are not aware of any deep-learning-based method for time-series forecasting that integrates change points.

1.3 Research Question and Solution

While, as discussed earlier, change points are common in time-series data, there is no accepted method for handling them in forecasting. Our goal is, therefore, to interweave change-point detection with DeepAR and other forecasting methods. By focusing on the training batches, we will be able to do so while (i) preserving the original forecasting methods’ characteristics and (ii) preserving the implementations.

Besides those two features, a strength of our approach is that (iii) it allows for the inclusion of any method for change-point detection. This provides access to a vast pool of existing methods and potential future methods.

2 Methods

Since we refer to the DeepAR algorithm for our method, we will introduce it briefly.

2.1 Mathematical Description of the DeepAR Model

The following presents the autoregressive recurrent network architecture of the DeepAR model and the training process Salinas et al. [2020].

2.1.1 Model

The value of time series i at time t is denoted by $z_{i,t}$. With given past data

$$[z_{i,1}, \dots, z_{i,t_0-2}, z_{i,t_0-1}] := \mathbf{z}_{i,1:t_0-1},$$

the conditional distribution

$$P(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T})$$

of the future of each time series

$$[z_{i,t_0}, z_{i,t_0+1}, \dots, z_{i,T}] := \mathbf{z}_{i,t_0:T},$$

is to be modeled. Here t_0 denotes the time point from which we assume that $z_{i,t}$ is unknown at prediction time. The covariates $\mathbf{x}_{i,1:T}$ are assumed to be known for all time points. For the past, we write the time ranges $[1, t_0 - 1]$, called the conditioning range. For the future, we write $[t_0, T]$ as prediction range, respectively. We assume that the model distribution

$$Q_{\Theta}(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T})$$

consists of a product of likelihood factors

$$Q_{\Theta}(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T}) = \prod_{t=t_0}^T Q_{\Theta}(z_{i,t} | \mathbf{z}_{i,1:t-1}, \mathbf{x}_{i,1:T}) = \prod_{t=t_0}^T \ell(z_{i,t} | \theta(\mathbf{h}_{i,t}, \Theta)),$$

parameterized by the output $\mathbf{h}_{i,t}$ of an autoregressive recurrent network

$$\mathbf{h}_{i,t} = h(\mathbf{h}_{i,t-1}, z_{i,t-1}, \mathbf{x}_{i,t}, \Theta).$$

Here h denotes a function implemented by a multilayer recurrent neural network with LSTM cells. We fed the network output $\mathbf{h}_{i,t}$ into a function $\theta(\mathbf{h}_{i,t}, \Theta)$, which then builds the parameters of the fixed distribution $\ell(z_{i,t} | \theta(\mathbf{h}_{i,t}))$.

2.1.2 Training

A time-series data set $\{\mathbf{z}_{i,1:T}\}_{i=1,\dots,N}$ and associated covariates $\mathbf{x}_{i,1:T}$ are obtained by choosing a time range such that $z_{i,t}$ is known in the prediction range. The parameters Θ of the model consist of the parameters of the RNN $h(\cdot)$ and the parameters of $\theta(\cdot)$. They can be learned by maximizing the log-likelihood

$$\mathcal{L} = \sum_{i=1}^N \sum_{t=t_0}^T \log \ell(z_{i,t} | \theta(\mathbf{h}_{i,t})). \quad (1)$$

Since $\mathbf{h}_{i,t}$ is a deterministic function of the input, all quantities needed to compute (1) are observed, so (1) can be optimized directly via stochastic gradient descent. In this process, the gradients will be computed with respect to Θ .

For a given dataset, we create multiple training instances for each time series by selecting windows with different starting points from the original time series.

The total length T and the relative lengths of the conditioning and prediction regions are fixed for all training examples. Selecting these windows ensures that the entire prediction range is always covered by the available baseline data. If the start point of the time series is chosen before $t = 1$, the unobserved targets are set to zero. In this way, the model can learn the behavior of the new time series given all other available features. Enriching the data with this windowing procedure ensures that information about absolute time is available to the model only through covariates, not through the relative position of $z_{i,t}$ in the time series.

The detailed training is as follows. A network with parameters Θ has three inputs: the covariates $x_{i,t}$, the target value of the previous time step $z_{i,t-1}$, and the network output of the previous time step $\mathbf{h}_{i,t-1}$. The network output

$$\mathbf{h}_{i,t} = h(\mathbf{h}_{i,t-1}, z_{i,t-1}, \mathbf{x}_{i,t}, \Theta)$$

is then used to calculate the parameters

$$\theta_{i,t} = \theta(\mathbf{h}_{i,t}, \Theta)$$

of the likelihood $\ell(z|\theta)$, which is used to train the model parameters. For prediction, the history of the time series $z_{i,t}$ is fed into the network for $t < t_0$, then a sample is drawn in the prediction domain for $t \geq t_0$ and fed back for the next point until the end of the prediction range $t = t_0 + T$. Repeating this prediction process results in multiple traces that represent the joint predicted distribution.

2.2 The DeepCAR Method: Incorporate Change Points

Our *DeepCAR* extends the vanilla DeepAR method by incorporating change points. Change points are common: in our pollution data, for example, there are unpredictable breaks when the machine is being repaired. These change points disrupt the normal behavior of the time series; if tolerated, they can harm predictions because the algorithm misinterprets them. Change points can be identified in various ways.

As mentioned before, a known method to deal with change points is to remove the data or replace it with linear interpolation methods. However, in our case and also for many real datasets, this approach manipulates the dataset and disrupts the time series behavior. Our DeepCAR method does not change the data but adapts the training process as follows.

The vanilla DeepAR learns by taking a random batch. A start index is chosen; then a batch with a fixed length is built. For each step, a new start index is chosen. In our DeepCAR method, we instead successively select windows from the time series, which are chosen to not contain a change point, being able to use the time series as a whole.

In detail, the DeepCAR method searches for the change points before the training starts. When selecting the batches, we define a new function that indicates whether the change point lies in the batch; if the change point lies in the batch, the search is repeated until we find a batch that does not contain a change point.

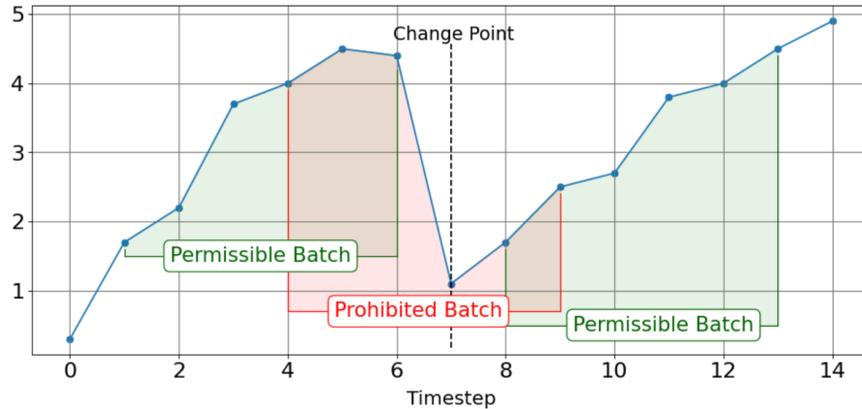


Figure 2: Visualization of the change point method *BatchCP*. It shows an example of a batch selection around a change point. The blue points represent time-series data. The first green area shows a batch created by choosing $t = 1$ as the start index, and with a selected batch size of $s = 6$, the batch contains the points $t = 1$ to $t = 6$. The change point is located at $t = 7$ and is therefore outside the detected batch. The batch is therefore permissible. The red area next to it shows a prohibited batch. The start index of this batch is $t = 4$ and the end index at $t = 9$, so the change point is located inside the batch and is not allowed in training. A new batch must be found. The green area to the right indicates another permissible batch, since the change point lies outside of it.

We want to go into more detail about the training process of the DeepCAR. Just as in vanilla DeepAR, the parameters Θ of the model are learned by optimizing the log-likelihood (1) using stochastic gradient descent by computing the gradients with respect to Θ . In the vanilla DeepAR, the batches are selected from the entire time series with a fixed batch size and different start points. Before we look at the batch generation, we need to mention the change point density and the size of the batches.

It is relevant how large the batch size is selected. It should be chosen depending on the number of change points and the size of the training data set. If there are many change points or if they lie close together due to less training data, the batch size should be chosen accordingly small, so that many batches can be found around the change points and no information gets lost. With small data sets, the batch size should also not be chosen too large, so that enough different batches can be generated.

In detail, the batch size s is one of the most important control parameters of our method. We have to choose s in such a way that it could cover at most one change point. If s is larger than the distance between two change points c_i, c_j , the part of the time series between time index i and j would never be learned. The batches containing this part of the time series would always contain at least one change point and would, therefore, not be allowed. To avoid this, we develop a method that automatically selects the maximum batch size based on the change points. In this method we calculate the distance between all detected change points. We then set s_{\max} to be half of the smallest distance between

Algorithm 1 find change points and maximum batch size

Input: training data d , change point detection method A
Output: maximal batch size s_{\max} , change points c_0, \dots, c_{k-1}
 $c_0, \dots, c_{k-1} \leftarrow A(d)$ # determine change points
for $i, j \in \{0, \dots, k-1\}, i \neq j$ **do**
 $\text{diff}_{i,j} \leftarrow \|c_i - c_j\|_2$
end for
 $s_{\max} \leftarrow \lceil \min_{i,j} \{\text{diff}_{i,j}\} / 2 \rceil$
return $c_0, \dots, c_{k-1}, s_{\max}$

the change points. We take half of this distance, because we increase the probability of selecting permissible batches between the change points. In addition, the range of the batches is more varied, and we have more diverse training input. This can be seen in Algorithm 1. The number we get for s_{\max} is the maximum batch size that should be selected. Our experiments have shown that the prediction improves with the choice of a batch size $s \leq s_{\max}$.

Now, we consider the generation of the batches based on this batch size. In the DeepCAR method, we pass a batch size $s \leq s_{\max}$ and the detected change points c_0, \dots, c_{k-1} to the algorithm. Our goal is to find the start and end points of the batches that do not contain any change points. To do this, we select indices that lie between 0 and $n - s - 1$. For each change point, we check whether the index of the respective change point lies between the previously selected start index and end index of the batch. If the change point lies in the batch, we repeat the method and select a new start index, for which the same procedure is performed—see Algorithm 2.

We go through the possible cases in detail: We have found a set of change points

$$C := \{c_0, \dots, c_{k-1}\}$$

and have chosen the batch size $s \leq s_{\max}$ in Algorithm 1. We select a random start index i_{start} which lies in the range $[0, n - s - 1]$ of the time series. This gives us a batch with start point i_{start} and end point

$$i_{\text{end}} := i_{\text{start}} + s - 1.$$

Assume that one of the change points lies in the batch, thus without loss of generality $\exists j \in \{0, \dots, k-1\}$ such that $c_j \in [i_{\text{start}}, i_{\text{end}}]$. Then again a random start index $i_{\text{start}'}$ is chosen and

$$i_{\text{end}'} := i_{\text{start}'} + s - 1.$$

If now $\forall j \in \{0, \dots, k-1\}$ holds

$$c_j \notin [i_{\text{start}'}, i_{\text{end}'}],$$

then in the following the batch $[i_{\text{start}'}, i_{\text{end}'}]$ is a valid batch.

The choice of the batches is somewhat random. We have ensured that the batches are either completely away from a change point or at the beginning or end of the batch bordering on a change point. Like the DeepAR, we also use a Gaussian likelihood in the model and apply LSTM layer in addition to the Gaussian layer. In other aspects, the training process follows the DeepAR.

We will illustrate later in Section 4 that the very same method of including change points also applies to other forecasting methods, such as Transformers or TFTs. We thus refer to the general method as *BatchCP*.

2.3 Methods for Change-Point Detection

Time series can have special points such as anomalies or outliers generated by measurement errors, for example. These points are usually isolated phenomena and, therefore, relatively straightforward to deal with in terms of prediction: usually, these points are simply deleted. But there is also another, more structural type of special points: change points. There are changes in mean, which are the most common change points. There is also change in variance or change in periodicity. A rarer case and also more difficult to detect is change in pattern. Other complex topics are for example changes in multidimensional time series, where changes in correlation are considered. Some change points can hardly be recognized visually. Therefore, there are numerous change point detection methods Aminikhanghahi & Cook [2017], Dehling et al. [2022], Eichinger & Kirch [2018], Li et al. [2022], Maidstone et al. [2017], Yamin et al. [2022]. Among others, there is the MOSUM method [Meier et al., 2021], which is particularly well suited for the detection of the change points of the first type, the change in mean.

We want to talk about the change point detection method MOSUM, that we use in the experiments. For the DeepCAR model, it is especially important to identify the change points as precisely as possible. If this cannot be guaranteed, the

Algorithm 2 find valid batch

Input: training data d , batch size $s \leq s_{\max}$,
change points c_0, \dots, c_{k-1}
Output: start and end points of the batch
 $n \leftarrow \text{len}(d)$ # number of samples
repeat
 $i_{\text{start}} \leftarrow$ random number in $\{0, \dots, n - s - 1\}$
valid \leftarrow True
for $j \in \{0, \dots, k - 1\}$ **do**
if $c_j \geq i_{\text{start}}$ and $c_j \leq i_{\text{start}} + s - 1$ **then**
valid \leftarrow False
break for
end if
end for
until valid = True
 $i_{\text{end}} \leftarrow i_{\text{start}} + s - 1$
return $i_{\text{start}}, i_{\text{end}}$

DeepCAR method would create batches that contain the actual change points and omit the change points that were found incorrectly. In reality, however, it is often not possible to determine change points manually, because either the data sets are too large, or the determination of the change points is too complex. The idea behind the MOSUM method is the detection of several change points in the mean. In MOSUM, when the bandwidth G is small, the η criterion is used, and when G is large, the ε criterion is used. This is because a small value of G is best for highly fluctuating time series to detect small changes, while a large value for G is best for detecting large changes. The disadvantage is that one must specify the ε criterion for large values of G so that MOSUM is not distracted by neighboring small changes and misses large changes.

3 Experiments

This section demonstrates the practical performance of our DeepCAR model. We consider three real-world data sets: data from a recycling company (the company already uses our pipeline on a daily basis), football data, and treasury data. We also consider four different scenarios for each data set: (I) baseline model naïve; (II) ignore all change points, which amounts to vanilla DeepAR, as a reference to show that our model performs as well as DeepAR in any case; (III) include hand-picked change points; and (IV) include change points detected by MOSUM to show that our method can largely outperform DeepAR in practice.

Code examples are given in <https://github.com/LedererLab/DeepCAR>.

3.1 Pollution Dataset (Multivariate)

Our idea is to detect the change points in the time series. We use the change point information when training the model. We consider time series data about oil filter pollution in a shredding machine. Changing the filter changes the data-generating process. We have time series data for one specific machine. The machine data consists of a timestamp with a value for the pollution filter and 10 additional features, like certain temperatures and pressures, which are dependent on the pollution. These data were prepared in advance to fit the time-series format of the DeepAR model. The goal is to predict the future values for the pollution in order to detect too high values early and to avoid failures and damage to the machine. We start with a baseline model.

The entire dataset has 10 200 rows, each row corresponds to the machine status in a given hour. We start with Scenario (I), the naïve method—see Table 1. The naïve method is very simple. It takes the current known value and predicts that it will be the same in the future. The errors for test and training samples are quite high with this method.

In the next step, we ignore the change points in the data and run the vanilla DeepAR model Arrigoni [2018]. The batch size is calculated with Algorithm 1. Since the minimum difference is between the change points $c_1 = 1970$ and $c_2 = 2350$, we get the difference $\text{diff}_{1,2} = 380$. Thus, the batch size should be at most

$$s_{\max} = \left\lceil \frac{380}{2} \right\rceil = 190.$$

We set $s = 30$. The loss is the Gaussian likelihood function with parameters μ and σ . The model consists of 3 layers: a LSTM layer with 4 units, a dense layer with 3 units, and a Gaussian layer. ReLU is used as the activation function Lederer [2021]. We do not attempt to optimize the hyperparameters. We split the data set into 60% training samples, 20% validation samples, and 20% test samples. The previous settings will be the same for all methods. We calculate the errors as root mean squared error (RMSE).

For our method in Scenario (II), that is, for vanilla DeepAR, we get an RMSE of 98.51 on the training data and 147.83 on the test data—see the first line in Table 1.

The change points were trained 265 times in the training process. In Scenario (III), we hand-picked the change points of the pollution time series. We know when oil-filter damage occurred to the machine and when the filter was replaced. We picked change points at the following time indices: 700, 1970, 2350, 2730, 3500, 4320, 5050, 6250, 7100.

We pass this list of change points in the training process. When the method that generates the random batches is called, we check whether a change point from our list is in the generated batch. If this is the case, we search for a new valid batch. For our method in Scenario (III), we get an RMSE of 81.27 on the training data and 112.67 on the test data—see the second line in Table 1. Compared to the first line, both the training and test score have improved significantly. We note that both the vanilla DeepAR and DeepCAR methods predict the distribution of the predicted values. This can be seen in Figure 3, where the prediction for one batch from the test data is shown.

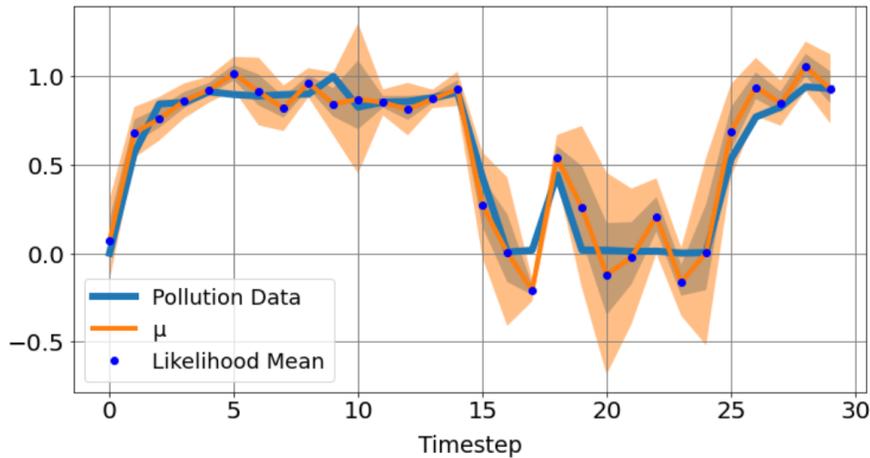


Figure 3: One Batch predicted by the DeepCAR Model. The shown batch has a length of $s = 6$, therefore 6 time steps are shown. The pollution data is blue, the likelihood mean is marked by dark blue data points and μ corresponds to the orange curve.

Now, we include the well-known change-point-detection method MOSUM Meier et al. [2021], which consists of procedures for the multiple mean change problem using moving sum statistics—see Section 2.3 . We choose a bandwidth of 0.2, and a value for $\eta = 0.1$ and let the MOSUM method detect the change points of the pollution time series. This can be seen in Figure 4.

The resulting list of change points is given to the DeepCAR model and the further procedure is analogous to the previous experiment. For this method (Scenario (IV)), we get a training error of 83.42 and a test error of 117.75—see the second line in Table 1. The error has worsened slightly, which can be explained by the fact that the MOSUM change point detection is not 100% correct but also forgot some points. The error is still smaller than with the DeepAR without taking the change points into account.

All results of the experiments can be seen in Table 1. The model is clearly performing better than the vanilla deepAR on the same dataset. Our DeepCAR method (Scenario (III)) is more than three times better than the naïve baseline model (Scenario (I)) and improves the vanilla DeepAR (Scenario (II)) by 23,78 %. Even with automatic change point detection (Scenario (IV)), which could be further optimized, we are 20, 35 % better than the vanilla DeepAR.

3.2 Football Dataset (Univariate)

Next, we consider a football dataset Bundesliga. The dataset consists of a date—the match day—and the difference in goals scored and received by a team. After each season, the statistics are reset, resulting in the change points in the

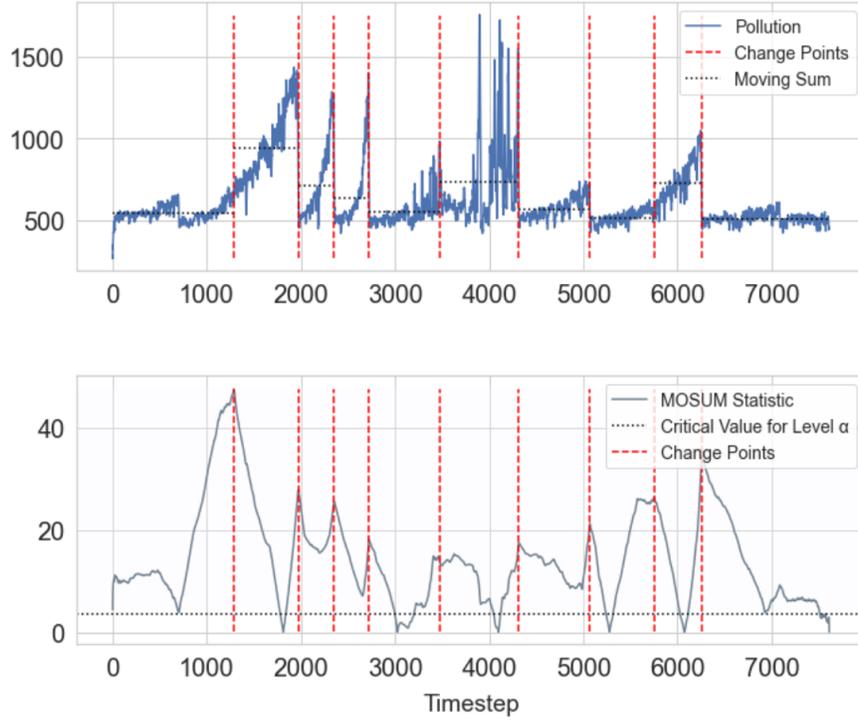


Figure 4: Upper diagram: The pollution time series (blue) in millibar with moving sum (gray) and change points (red) detected with MOSUM. Lower diagram: The MOSUM statistic, which identifies the changes in the mean, the threshold level (dotted line), and the corresponding change points (red). Most of the detected change points equal the actual change points.

Table 1: Prediction error of the different methods on the pollution data. Reported are test and training RMSE of (I) naïve, (II) DeepAR and our DeepCAR equipped with (III) manual and (IV) MOSUM change-point detection.

SCENARIO	TRAIN RMSE	TEST RMSE
(I) BASELINE NAÏVE	198.55	382.08
(II) DEEPAR-NO CPD	98.51	147.83
(III) DEEPCAR-MANUAL	81.27	112.67
(IV) DEEPCAR-MOSUM	83.42	117.75

dataset. The dataset is univariate since there are no features. The dataset with change points can be seen in Figure 5. The change points are at time indices 31, 65, 99, 133, 157, 174, 191, 208. We again use the naïve method (Scenario I) as a baseline model—see Table 2.

As already mentioned, it is relevant how large the batch size is selected. In the case of the football dataset, we only have a small amount of training data, so the batch size should be correspondingly smaller. We again calculate the batch size with Algorithm 1. The minimum difference is between the change points $c_4 = 157$ and $c_5 = 174$, we get the difference $\text{diff}_{4,5} = 17$. Thus, the batch size should be at most

$$s_{\max} = \left\lfloor \frac{17}{2} \right\rfloor = 9.$$

We apply this batch size.

We now start with our method in Scenario (II), that is, for vanilla DeepAR, we get a RMSE of 15.41 on the training data and 51.78 on the test data—see the first line in Table 2.

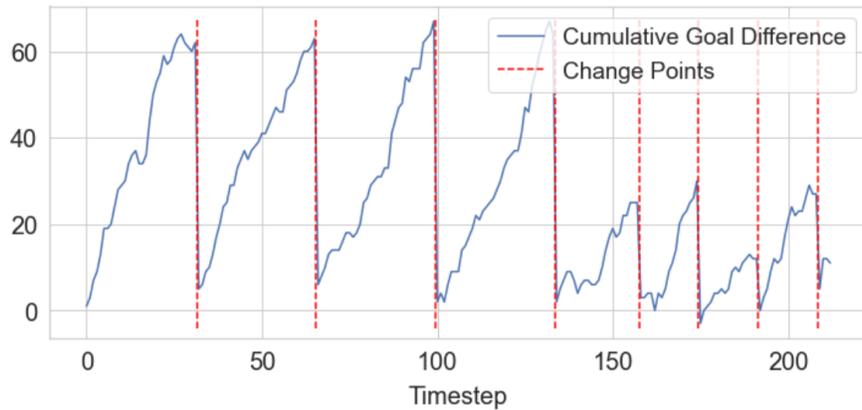


Figure 5: Football data (blue) with manual change points (red) located at time indices 31, 65, 99, 133, 157, 174, 191, 208.

Next, we use our method in Scenario (III) with the change points we hand-picked manually before. The change points are passed into the training. The procedure is the same as in the pollution experiment, except that this time, we have a univariate time series. We get an RMSE of 8.20 on the training data and 14.89 on the test data—see the second line in Table 2. Compared to the first Scenario, both the training and test score have improved significantly.

Now we again use MOSUM to detect the change points of the football dataset—see Figure 6.

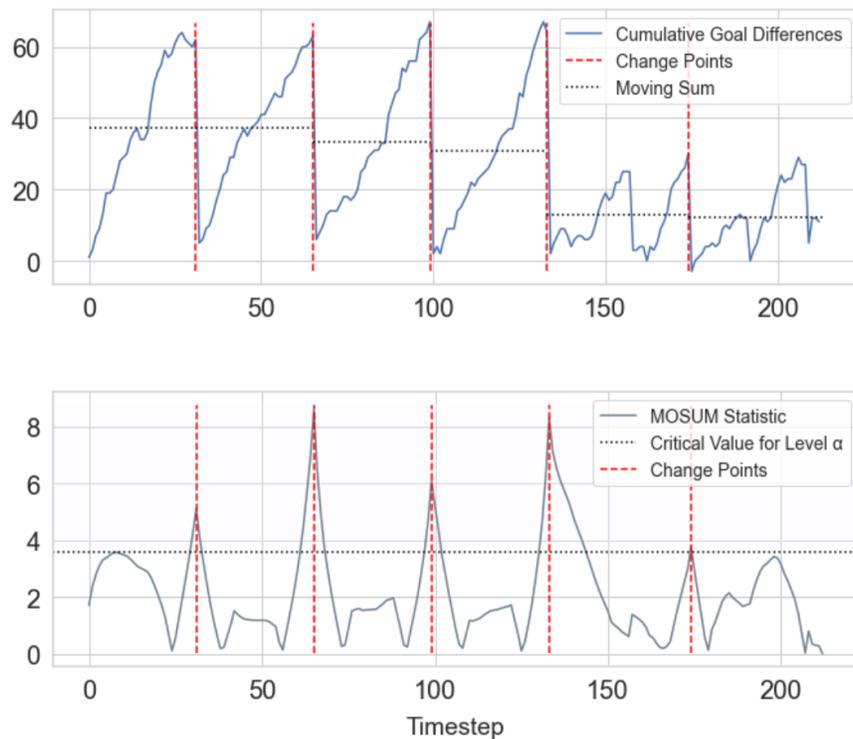


Figure 6: Upper diagram: The football data (blue) with change points (red) detected with MOSUM, and the corresponding moving sum (grey). Lower diagram: The MOSUM statistic, which identifies the changes in the mean, and the corresponded change points in red. The defined threshold value is also plotted in grey.

Table 2: Prediction error of the different methods on the football data. Reported are test and training RMSE of (I) naïve, (II) DeepAR and our DeepCAR equipped with (III) manual and (IV) MOSUM change-point detection.

SCENARIO	TRAIN RMSE	TEST RMSE
(I) BASELINE NAÏVE	25.29	51.78
(II) DEEPAR-NO CPD	15.41	25.62
(III) DEEPCAR-MANUAL	8.20	14.89
(IV) DEEPCAR-MOSUM	10.44	18.63

The detection is very good: most manually-set change points are detected. We pass this list of change points into the training of the DeepCAR. Thus for Scenario (IV) we get a training error of 10.44 and a test error of 18.63.

We again compare all results of the football data experiments in Table 2. We again use the test error for a comparison. The results with our DeepCAR method are again significantly better than with naïve and with vanilla DeepAR. The test error with DeepCAR (Scenario (III)) is 41.88 % better than the vanilla DeepAR. With MOSUM (Scenario (IV)) we are 27.09 % better than the vanilla DeepAR.

3.3 Treasury Rate (Univariate)

We now consider univariate data of the Federal Reserve Board One Year Treasury Rate Dataset. The samples are daily averages of yields of several treasury securities, all adjusted to the equivalent of a one-year maturity.

With our baseline model (I) we get an RMSE of 3.98 for training samples and 8.17 for test samples.

The change points are not immediately recognizable here. But we know that recessions have occurred at time points 1992, 3155, 4544, 5105, 7065, 9710, 11480, 14543—see Figure 7. MOSUM recognizes these change points relatively accurately again (for detailed MOSUM results of the treasury dataset see Appendix A). Algorithm 1 yields the maximal batch size

$$s_{\max} = \left\lceil \frac{4544 - 3155}{2} \right\rceil = \left\lceil \frac{1389}{2} \right\rceil = 695.$$

We choose $s = 50$.

The results are summarized in Table 3. Once more, our DeepCAR outmatches DeepAR: the test errors of DeepCAR are about 9 % (Scenario (III)) and 5 % (Scenario (IV)) better than for DeepAR and also much better than for naïve.

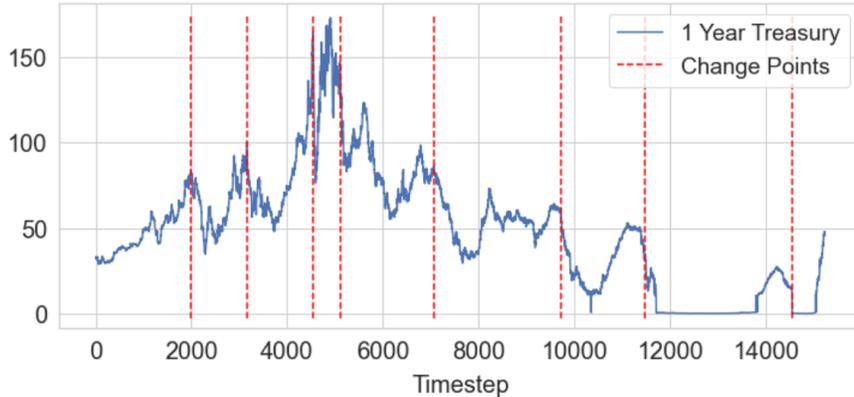


Figure 7: Treasury dataset (blue) with manual change points (red) located at time indices 1992, 3155, 4544, 5105, 7065, 9710, 11480, 14543.

3.4 Synthetic data

We substantiate the empirical results further by considering synthetic data, where we know and control the data-generating process. The data comprises 3000 samples with 13 change points—see Figure 8. The naïve method (I) yields 3.57 RMSE for the training samples and 9.33 for the test samples.

Table 3: Prediction error of the different methods on the Treasury data. Reported are test and training RMSE of (I) naïve, (II) DeepAR, and our DeepCAR equipped with (III) manual and (IV) MOSUM change-point detection.

SCENARIO	TRAIN RMSE	TEST RMSE
(I) BASELINE NAÏVE	3.98	8.17
(II) DEEPAR-NO CPD	4.13	5.93
(III) DEEPCAR-MANUAL	3.40	5.37
(IV) DEEPCAR-MOSUM	3.69	5.61

Table 4: Prediction error of the different methods on the synthetic data. Reported are test and training RMSE of (I) naïve, (II) DeepAR and our DeepCAR equipped with (III) manual and (IV) MOSUM change-point detection.

SCENARIO	TRAIN RMSE	TEST RMSE
(I) BASELINE NAÏVE	3.57	9.33
(II) DEEPAR-NO CPD	2.81	5.77
(III) DEEPCAR-MANUAL	2.42	5.40
(IV) DEEPCAR-MOSUM	2.51	5.62

Algorithm 1 allows for a batch size of $s = 50$. The results for Scenarios (I), (II), (III), and (IV) are in Table 4. DeepCAR, both with manually and automatically selected change points, outmatches vanilla DeepAR and naïve once more.

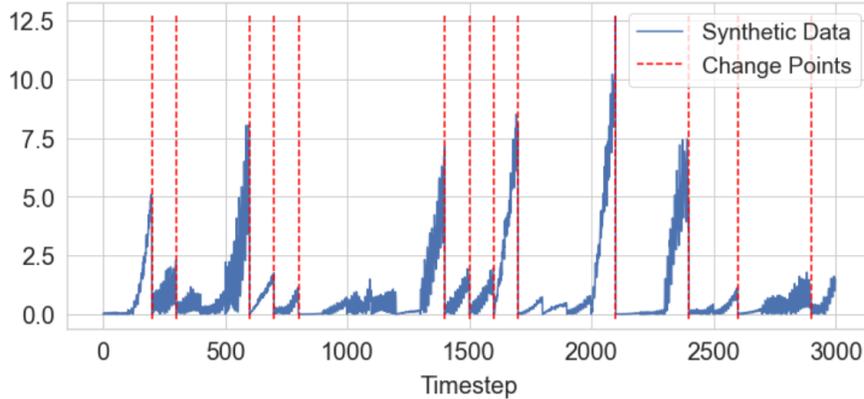


Figure 8: Synthetic data (blue) with manual change points (red) located at time indices 200, 300, 600, 700, 800, 1400, 1500, 1600, 1700, 2100, 2400, 2600, 2900.

4 Extensions Beyond DeepAR

So far, we have extended DeepAR. But importantly, our concepts apply to batch-trained methods much more generally. Examples are the attention-based Transformer or Temporal Fusion Transformer (TFT) [Lim et al., 2021, Vaswani et al., 2017]. To illustrate this, we implement Transformer via Gluon Time Series (GluonTS) Alexandrov et al. [2019]. The batches are generated according to Algorithm 2.

We test the pipeline on the pollution data of Section 3.1. We consider the following three scenarios: (A) ignore all change points; (B) include hand-picked change points (the same ones as in Section 3.1); and (C) include change points detected by MOSUM (the same ones as in Section 3.1). The results are summarized in Table 5. We observe that the Transformer model generally seems to work better than the DeepAR model in this example—compare to Table 1—but this is not necessarily the case in other examples. Important for us is that we can improve the Transformer noticeably by our BatchCP method for including change points.

Table 5: Prediction errors on the pollution data. Reported are test and training RMSE of (A) Transformer without considering change points and Transformer equipped with (B) manual, and (C) MOSUM change-point detection.

SCENARIO	TRAIN RMSE	TEST RMSE
(A) TRANSFORMER-NO CPD	75.26	96.23
(B) TRANSFORMER-MANUAL	67.85	91.02
(C) TRANSFORMER-MOSUM	70.22	93.12

5 Conclusion

Our generalization of DeepAR, called DeepCAR, improves on the vanilla version substantially when there are change points and equals that version if no change points are detected. Thus, DeepCAR can generally be used as a drop-in replacement for DeepAR.

But equally importantly, our work identifies the batch size as a way to account for change points in time-series analyses more generally. As an illustration, we have shown that our general method to account for change points, called DeepCAR, can improve Transformers as well.

6 Acknowledgments

We sincerely thank Kata Vuk and David Salinas for their helpful input.

References

- Alexandrov, A., Benidis, K., Bohlke-Schneider, M., Flunkert, V., Gasthaus, J., Januschowski, T., Maddix, D., Rangapuram, S., Salinas, D., Schulz, J., Stella, L., Türkmen, A., and Wang, Y. GluonTS: Probabilistic Time Series Modeling in Python. *CoRR, arXiv:1906.05264*, abs/1906.05264, 2019.
- Aminikhanghahi, S. and Cook, D. A survey of methods for time series change point detection. *Knowledge and Information Systems*, 51:339–367, 2017.
- Arrigoni, A. Paper review & code: Amazon DeepAR, 2018. URL <https://medium.com/@albertoarrigoni/paper-review-code-amazon-deepar-809938a319d9>. Online, accessed: 2022-12-20.
- Box, G., Jenkins, G., and MacGregor, J. Some recent advances in forecasting and control. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 23(2):158–179, 1974.
- Bundesliga. Football-data.co.uk. URL <https://www.football-data.co.uk/germanym.php>. Online, accessed: 2022-12-20.
- Dagan, Y., Kandiros, V., and Daskalakis, C. EM’s convergence in Gaussian latent tree models. *arXiv:2211.11904*, 2022.
- De Gooijer, J. and Hyndman, R. 25 years of time series forecasting. *International Journal of Forecasting*, 22(3): 443–473, 2006.
- Dehling, H., Vuk, K., and Wendler, M. Change-point detection based on weighted two-sample U-statistics. *Electronic Journal of Statistics, arXiv:2003.12573*, 16(1):862–891, 2022.
- Eichinger, B. and Kirch, C. A MOSUM procedure for the estimation of multiple random change points. *Bernoulli*, 24(1):526–564, 2018.
- Gers, F., Eck, D., and Schmidhuber, J. Applying LSTM to time series predictable through time-window approaches. *Artificial Neural Networks — ICANN 2001, Springer Berlin Heidelberg*, pp. 669–676, 2001.
- Graves, A. Generating sequences with recurrent neural networks. *CoRR, arXiv:1308.0850*, 2013.
- Hyndman, R., Koehler, A., Ord, K., and Snyder, R. (eds.). *Forecasting with Exponential Smoothing: The State Space Approach. Springer Series in Statistics*. Springer Berlin, Heidelberg, 2008.
- Lederer, J. Activation functions in artificial neural networks: A systematic overview. *arXiv:2101.09957*, 2021.

- Li, J., Fearnhead, P., Fryzlewicz, P., and Wang, T. Automatic change-point detection in time series via deep learning. *arXiv:2211.03860*, 2022.
- Lim, B., Arik, S., Loeff, N., and Pfister, T. Temporal Fusion Transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4):1748–1764, 2021.
- Maidstone, R., Hocking, T., Rigaiil, G., and Fearnhead, P. On optimal multiple changepoint algorithms for large data. *Statistics and Computing*, 27(2):519–533, 2017.
- Meier, A., Kirch, C., and Cho, H. mosum: A package for moving sums in change-point analysis. *Journal of Statistical Software*, 97(8):1–42, 2021.
- One Year Treasury Rate Dataset. URL <https://www.macrotrends.net/2492/1-year-treasury-rate-yield-chart>. Online, accessed: 2022-12-20.
- Salinas, D., Flunkert, V., Gasthaus, J., and Januschowski, T. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, *arXiv:1704.04110*, 36(3):1181–1191, 2020.
- Snyder, R., Ord, K., and Beaumont, A. Forecasting the intermittent demand for slow-moving inventories: A modelling approach. *International Journal of Forecasting*, 28(2):485–496, 2012.
- Sutskever, I., Vinyals, O., and Le, Q. Sequence to sequence learning with neural networks. *In Advances in Neural Information Processing Systems*, 4:3104–3112, 2014.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. Attention is all you need. *Advances in Neural Information Processing Systems*, *Curran Associates, Inc.*, 30:5998–6008, 2017.
- Yamin, K., Wang, H., Montreuil, B., and Xie, Y. Online detection of supply chain network disruptions using sequential change-point detection for Hawkes processes. *arXiv:2211.12091*, 2022.
- Zhang, G., Patuwo, E., and Hu, M. Forecasting with artificial neural networks: The state of the art. *International Journal of Forecasting*, 14(1):35–62, 1998.

A Appendix

Figure 9 shows the treasury data with change points detected by MOSUM.

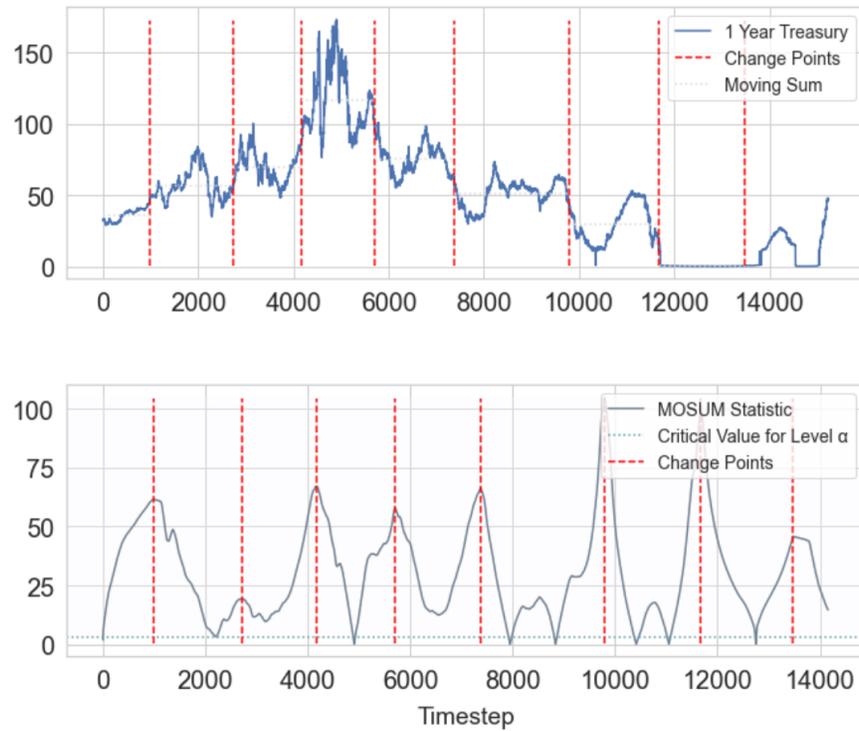


Figure 9: Upper diagram: The one-year-treasury time series (blue) with change points (red) detected with MOSUM, and the corresponding moving sum (grey). Lower diagram: The MOSUM statistic, which identifies the changes in the mean, and correspondingly the detected change points in red. The defined threshold value is plotted in grey.