
Monash Time Series Forecasting Archive

Rakshitha Godahewa
Monash University
Melbourne, Australia
rakshitha.godahewa@monash.edu

Christoph Bergmeir
Monash University
Melbourne, Australia
christoph.bergmeir@monash.edu

Geoffrey I. Webb
Monash University
Melbourne, Australia
geoff.webb@monash.edu

Rob J. Hyndman
Monash University
Melbourne, Australia
rob.hyndman@monash.edu

Pablo Montero-Manso
University of Sydney
Australia
pmontm@gmail.com

Abstract

1 Many businesses nowadays rely on large quantities of time series data making
2 time series forecasting an important research area. Global forecasting models and
3 multivariate models that are trained across sets of time series have shown huge
4 potential in providing accurate forecasts compared with the traditional univariate
5 forecasting models that work on isolated series. However, there are currently no
6 comprehensive time series forecasting archives that contain datasets of time series
7 from similar sources available for researchers to evaluate the performance of new
8 global or multivariate forecasting algorithms over varied datasets. In this paper, we
9 present such a comprehensive forecasting archive containing 20 publicly available
10 time series datasets from varied domains, with different characteristics in terms of
11 frequency, series lengths, and inclusion of missing values. We also characterise
12 the datasets, and identify similarities and differences among them, by conducting
13 a feature analysis. Furthermore, we present the performance of a set of standard
14 baseline forecasting methods over all datasets across ten error metrics, for the
15 benefit of researchers using the archive to benchmark their forecasting algorithms.

16 1 Introduction

17 Accurate time series forecasting is important for many businesses and industries to make decisions,
18 and consequently, time series forecasting is a popular research area. The field of forecasting has
19 traditionally been advanced by influential forecasting competitions. The most popular forecasting
20 competition series is the M-competition series [1–5]. Other well-known forecasting competitions
21 include the NN3 and NN5 Neural Network competitions [6], and Kaggle competitions such as the
22 Wikipedia web traffic competition [7].

23 The winning approaches of many of the most recent competitions such as the winning method of the
24 M4 by Smyl [8] and the winning method of the M5 forecasting competition [5], consist of global
25 forecasting models [9] which train a single model across all series that need to be forecast. Compared
26 with local models, global forecasting models have the ability to learn cross-series information during
27 model training and can control model complexity and overfitting on a global level [10].

28 This can be seen as a paradigm shift in forecasting. Over decades, single time series were seen as a
29 dataset that should be studied and modelled in isolation. Nowadays, we are oftentimes interested in
30 models built on sets of series from similar sources, such as series which are all product sales from a
31 particular store, or series which are all smart meter readings in a particular city. Here, time series

32 are seen as an instance in a dataset of many time series, to be studied and modelled together. Global
33 (univariate) models and (local) multivariate models are the methods of choice here, the difference
34 being that global models train across series, but predict each series in isolation, with no need for
35 the time series to have the same length or to be aligned in time, whereas multivariate models train
36 and test over time series that are all the same length and all aligned in time, so that dependencies at
37 certain time points and for the forecasts can be modelled. Thus, global models are applicable more
38 broadly than multivariate models.

39 Both global and multivariate models get attention lately in machine learning (especially deep learning),
40 with Li et al. [11], Rangapuram et al. [12], Wen et al. [13] presenting global models and Salinas
41 et al. [14], Sen et al. [15], Yu et al. [16], Zhou et al. [17] discussing novel approaches for multivariate
42 modelling. However, when it comes to benchmarking, these recent works use a mere two [13] to
43 seven [11] datasets to evaluate the performance of the new algorithms and the chosen datasets are
44 different in each work. The datasets mainly belong to the energy, transport, and sales domains, and
45 they do not include datasets from other domains such as banking, healthcare, or nature.

46 In contrast, other areas of machine learning, such as general classification and regression, or time
47 series classification, have greatly benefitted from established benchmark dataset archives, which
48 allow a much broader and more standardised evaluation. The University of California Irvine (UCI)
49 repository [18] is the most common and well-known benchmarking archive used in general machine
50 learning, with currently 507 datasets from various domains. In time series classification, the dataset
51 archives from the University of California Riverside (UCR) [19] and from the University of East
52 Anglia (UEA) [20], contain 128 sets of univariate time series, and 30 datasets with multivariate time
53 series, respectively, allowing routinely for much broader and more standardised evaluations of the
54 methods, and therewith enabling more streamlined, robust, and reliable progress in the field.

55 The time series classification datasets, though they contain time series, do usually not resemble
56 meaningful forecasting problems, so they cannot be used for the evaluation of forecasting methods.
57 Also in the time series forecasting space there are a number of benchmarking archives, but they
58 follow the paradigm of single series as datasets, and consequently contain mostly unrelated single
59 time series such as the Time Series Data Library [21] and ForeDeCk [22].

60 There are currently no comprehensive time series forecasting benchmarking archives, to the best of
61 our knowledge, that focus on sets of time series to evaluate the performance of global and multivariate
62 forecasting algorithms. We introduce such an archive, available at [https://forecastingdata.
63 org/](https://forecastingdata.org/). The archive contains 20 publicly available time series datasets covering varied domains, with
64 both equal and variable lengths time series. Many datasets have different versions based on the
65 frequency and the inclusion of missing values, resulting in a total of 50 dataset variations.

66 We also introduce a new format to store time series data, based on the Weka ARFF file format [23],
67 to overcome some of the shortcomings we observe in the .ts format used in the sktime time series
68 repository [24]. We use a .tsf extension for this new format. This format stores the meta-information
69 about a particular time series dataset such as dataset name, frequency, and inclusion of missing values
70 as well as the series specific information such as starting timestamps, in a non-redundant way. The
71 format is very flexible and capable of including any other attributes related to time series as preferred
72 by the users.

73 Furthermore, we analyse the characteristics of different series to identify the similarities and differ-
74 ences among them. For that, we conduct a feature analysis using tsfeatures [25] and catch22 features
75 [26] extracted from all series of all datasets. The extracted features are publicly available for further
76 research use. The performance of a set of baseline forecasting models including both traditional
77 univariate forecasting models and global forecasting models are also evaluated over all datasets across
78 ten error metrics. The forecasts and evaluation results of the baseline methods are publicly available
79 for the benefits of researchers that use the repository to benchmark their forecasting algorithms.

80 **2 Data records**

81 This section details the datasets in our time series forecasting archive. The current archive contains
82 20 time series datasets. Furthermore, the archive contains in addition 6 single very long time series.
83 As a large amount of data oftentimes renders machine learning methods feasible compared with
84 traditional statistical modelling, and we are not aware of good and systematic benchmark data in this

85 space either, these series are included in our repository as well. A summary of all primary datasets
 86 included in the repository is shown in Table 1.

87 A total of 50 datasets have been derived from these 26 primary datasets. Nine datasets contain
 88 time series belonging to different frequencies and the archive contains a separate dataset per each
 89 frequency. Seven of the datasets have series with missing values. The archive contains 2 versions of
 90 each of these, one with and one without missing values. In the latter case, the missing values have
 91 been replaced by using an appropriate imputation technique.

Table 1: Datasets in the current time series forecasting archive

	Dataset	Domain	No: of Series	Min. Length	Max. Length	No: of Freq.	Missing	Competition
1	M1	Multiple	1001	15	150	3	No	Yes
2	M3	Multiple	3003	20	144	4	No	Yes
3	M4	Multiple	100000	19	9933	6	No	Yes
4	Tourism	Tourism	1311	11	333	3	No	Yes
5	NN5	Banking	111	791	791	2	Yes	Yes
6	CIF 2016	Banking	72	34	120	1	No	Yes
7	Web Traffic	Web	145063	803	803	2	Yes	Yes
8	Solar	Energy	137	52560	52560	2	No	No
9	Electricity	Energy	321	26304	26304	2	No	No
10	London Smart Meters	Energy	5560	288	39648	1	Yes	No
11	Wind Farms	Energy	339	6345	527040	1	Yes	No
12	Car Parts	Sales	2674	51	51	1	Yes	No
13	Dominick	Sales	115704	28	393	1	No	No
14	FRED-MD	Economic	107	728	728	1	No	No
15	San Francisco Traffic	Transport	862	17544	17544	2	No	No
16	Pedestrian Counts	Transport	66	576	96424	1	No	No
17	Hospital	Health	767	84	84	1	No	No
18	COVID Deaths	Nature	266	212	212	1	No	No
19	KDD Cup	Nature	270	9504	10920	1	Yes	Yes
20	Weather	Nature	3010	1332	65981	1	No	No
21	Sunspot	Nature	1	73931	73931	1	Yes	No
22	Saugeen River Flow	Nature	1	23741	23741	1	No	No
23	US Births	Nature	1	7305	7305	1	No	No
24	Electricity Demand	Energy	1	17520	17520	1	No	No
25	Solar Power	Energy	1	7397222	7397222	1	No	No
26	Wind Power	Energy	1	7397147	7397147	1	No	No

92 Out of the 26 datasets, 8 originate from competition platforms, 3 from research conducted by Lai
 93 et al. [27], 6 are taken from R packages, 1 is from the Kaggle platform [28], and 1 is taken from a
 94 Johns Hopkins repository [29]. The other datasets have been extracted from corresponding domain
 95 specific platforms. The datasets mainly belong to 9 different domains: tourism, banking, web, energy,
 96 sales, economics, transportation, health, and nature. Three datasets, the M1 [1], M3 [2], and M4
 97 [3, 4] datasets, contain series belonging to multiple domains. All datasets are explained in detail in
 98 the Appendix (supplementary materials).

99 2.1 Data format

100 We introduce a new format to store time series data, based on the Weka ARFF file format [23]. We
 101 use the file extension .tsf and it is comparable with the .ts format used in the sktime time series
 102 repository [24], but we deem it more streamlined and more flexible. The basic idea of the file format
 103 is that each data file can contain 1) attributes that are constant throughout the whole dataset (e.g., the
 104 forecasting horizon, whether the dataset contains missing values or not), 2) attributes that are constant
 105 throughout a time series (e.g., its name, its position in a hierarchy, product information for product
 106 sales time series), and 3) attributes that are particular to each data point (the value of the series, or
 107 timestamps for non-equally spaced series). An example of series in this format is shown in Figure 1.

108 The original Weka ARFF file format already deals well with the first two types of such attributes.
 109 Using this file format, in our format, each time series file contains tags describing the meta-information

```

# Dataset Information
# This dataset was used in the NN5 forecasting competition.
# It contains 111 daily time series from the banking domain.
# The goal is predicting the daily cash withdrawals from ATMs in UK.
#
# For more details, please refer to
# Taieb, S.B., Bontempi, G., Atiya, A.F., Sorjamaa, A., 2012.
# A review and comparison of strategies for multi-step ahead time series forecasting based on
# the nn5 forecasting competition. Expert Systems with Applications 39(8), 7067 - 7083
#
# Neural Forecasting Competitions, 2008.
# NN5 forecasting competition for artificial neural networks and computational intelligence.
# Accessed: 2020-05-10. URL http://www.neural-forecasting-competition.com/NN5/
#
@relation NN5
@attribute series_name string
@attribute start_timestamp date
@frequency daily
@horizon 56
@missing true
@equallength true
@data
T1:1996-03-18 00-00-00:13.4070294784581,14.7250566893424,20.5640589569161,34.7080498866213,26
T2:1996-03-18 00-00-00:11.5504535147392,13.5912698412698,15.0368480725624,21.5702947845805,19
T3:1996-03-18 00-00-00:5.640589569161,14.3990929705215,24.4189342403628,28.7840136054422,20.6
T4:1996-03-18 00-00-00:13.1802721088435,8.44671201814059,19.515306122449,28.8832199546485,19.
T5:1996-03-18 00-00-00:9.77891156462585,10.8134920634921,21.6128117913832,38.5204081632653,24
T6:1996-03-18 00-00-00:9.24036281179138,11.6354875283447,12.1031746031746,21.4143990929705,24
T7:1996-03-18 00-00-00:14.937641723356,16.2840136054422,16.6666666666667,23.5685941043084,26.
T8:1996-03-18 00-00-00:2.89115646258503,12.3582766439909,16.3832199546485,30.1587301587302,31
T9:1996-03-18 00-00-00:7.34126984126984,9.15532879818594,10.5867346938776,12.5,7.157029478458
T10:1996-03-18 00-00-00:10.2891156462585,12.7125850340136,14.4416099773243,19.4019274376417,2.

```

Figure 1: An example of the file format for the NN5 daily dataset.

of the corresponding dataset such as *@frequency* (seasonality), *@horizon* (expected forecast horizon), *@missing* (whether the series contain missing values) and *@equallength* (whether the series have equal lengths). We note that in principle these attributes can be freely defined by the user and the file format does not need any of these values to be defined in a certain way, though the file readers reading the files may rely on existence of attributes with certain names and assume certain meanings. Next, there are attributes in each dataset which describe series-wise properties, where the tag *@attribute* is followed by the name and type. Examples are *series_name* (the unique identifier of a given series) and *start_timestamp* (the start timestamp of a given series). Again, the format has the flexibility to include any additional series-wise attributes as preferred by users.

Following the ARFF file format, the data are then listed under the *@data* tag after defining attributes and meta-headers, and attribute values are separated by colons. The only extension that our format has compared with the original ARFF file format, is that the time series then are appended to their attribute vector as a comma-separated variable-length vector. As this vector can have a different length for each instance, this cannot be represented in the original ARFF file format. In particular, a time series with m number of attributes and n number of values can be shown as:

$$\langle attribute_1 \rangle : \langle attribute_2 \rangle : \dots : \langle attribute_m \rangle : \langle s_1, s_2, \dots, s_n \rangle \quad (1)$$

The missing values in the series are indicated using the “?” symbol. Code to load datasets in this format into R and Python is available in our github repository at <https://github.com/rakshitha123/TSForecasting>.

3 Methods

This section details the feature analysis and baseline evaluation we conducted on the datasets in our repository.

3.1 Feature analysis

We characterise the datasets in our archive to analyse the similarities and differences between them, to gain a better understanding on where gaps in the repository may be and what type of

134 data are prevalent in applications. This may also help to select suitable forecasting methods for
135 different types of datasets. We analyse the characteristics of the datasets using the *tsfeatures* [25] and
136 *catch22* [26] feature extraction methods. All extracted features are publicly available in our website
137 <https://forecastingdata.org/> for further research use. Due to their large size, we have not
138 been able to extract features from the London smart meters, wind farms, solar power, and wind power
139 datasets, which is why we exclude them from this analysis.

140 We extract 42 features using the *tsfeatures* function in the R package *tsfeatures* [25] in-
141 cluding mean, variance, autocorrelation features, seasonal features, entropy, crossing points,
142 flat spots, lumpiness, non-linearity, stability, Holt-parameters, and features related to the
143 Kwiatkowski–Phillips–Schmidt–Shin (KPSS) test [30] and the Phillips–Perron (PP) test [31]. For
144 all series that have a frequency greater than daily, we consider multi-seasonal frequencies when
145 computing features. Therefore, the amount of features extracted is higher for multi-seasonal datasets
146 as the seasonal features are individually calculated for each season presented in the series. Further-
147 more, if a series is short and does not contain two full seasonal cycles, we calculate the features
148 assuming a non-seasonal series (i.e., setting its frequency to “one” for the feature extraction). We use
149 the *catch22_all* function in the R package *catch22* [32] to extract the *catch22* features from a given
150 time series. The features are a subset of 22 features from the *hctsa* package [33] which includes the
151 implementations of over 7000 time series features. The computational cost of the *catch22* features is
152 low compared with all features implemented in the *hctsa* package.

153 For the feature analysis, we consider 5 features, as suggested by Bojer and Meldgaard [34]: first
154 order autocorrelation (ACF1), trend, entropy, seasonal strength, and the Box-Cox transformation
155 parameter, lambda. The *BoxCox.lambda* function in the R package *forecast* [35] is used to extract the
156 Box-Cox transformation parameter from each series, with default parameters. The other 4 features are
157 extracted using *tsfeatures*. Since this feature space contains 5 dimensions, to compare and visualise
158 the features across multiple datasets, we reduce the feature dimensionality to 2 using Principal
159 Component Analysis [PCA, 36].

160 The numbers of series in each dataset are significantly different, e.g., the CIF 2016 monthly dataset and
161 M4 monthly dataset contain 72 and 48,000 series, respectively. Hence, if all series were considered
162 to calculate the PCA components, those components would be dominated by datasets that have
163 large amounts of series. Therefore, for datasets that contain more than 300 series, we randomly
164 take a sample of 300 series, before constructing the PCA components across all datasets. Once the
165 components are calculated, we map all series of all datasets into the resulting PCA feature space. We
166 note that we use PCA for dimensionality reduction over other advanced dimensionality reduction
167 algorithms such as t-Distributed Stochastic Neighbor Embedding [t-SNE, 37] due to this capability
168 of constructing the basis of the feature space with a reduced sample of series with the possibility to
169 then map all series into the space afterwards.

170 3.2 Baseline forecasting models

171 In the forecasting space, benchmarking against simple benchmarks is vital [38] as even simple
172 benchmarks can oftentimes be surprisingly competitive. However, many works in the machine
173 learning space are notoriously weak when it comes to proper benchmarking for time series forecasting
174 [39]. To fill this gap, we evaluate the performance of 11 different baseline forecasting models over the
175 datasets in our repository using a fixed origin evaluation scheme, so that researchers that use the data
176 in our repository can directly benchmark their forecasting algorithms against these baselines. The
177 baseline models include 6 traditional univariate forecasting models: Exponential Smoothing [ETS,
178 40], Auto-Regressive Integrated Moving Average [ARIMA, 41], Simple Exponential Smoothing
179 (SES), Theta [42], Trigonometric Box-Cox ARMA Trend Seasonal [TBATS, 43] and Dynamic
180 Harmonic Regression ARIMA [DHR-ARIMA, 44], and 5 global forecasting models: a linear Pooled
181 Regression model [PR, 45], Feed-Forward Neural Network [FFNN, 46], CatBoost [47], DeepAR
182 [48] and N-BEATS [49].

183 Again, we do not consider the London smart meters, wind farms, solar power, and wind power
184 datasets for both univariate and global model evaluations, and the Kaggle web traffic daily dataset for
185 the global model evaluations, as the computational cost of running these models was not feasible in
186 our experimental environment.

187 We use the R packages *forecast* [50], *glmnet* [51], *catboost* [47] and *nnet* [52] to implement the 6
188 traditional univariate forecasting methods, the globally trained PR method, CatBoost and FFNN,
189 respectively. We use the implementations of DeepAR and N-BEATS available from the Python
190 package *GLuonTS* [53] with their default hyperparameters.

191 The Theta, SES, and PR methods are evaluated for all datasets. ETS and ARIMA are evaluated
192 for yearly, quarterly, monthly, and daily datasets. We consider the datasets with small frequencies,
193 namely, 10 minutely, half hourly, and hourly as multi-seasonal and hence, TBATS and DHR-ARIMA
194 are evaluated for those datasets instead of ETS and ARIMA due to their capability of dealing with
195 multiple seasonalities [54]. TBATS and DHR-ARIMA are also evaluated for weekly datasets due to
196 their capability of dealing with long non-integer seasonal cycles present in weekly data [55].

197 Forecast horizons are chosen for each dataset to evaluate the model performance. For all competition
198 datasets, we use the forecast horizons originally employed in the competitions. For the remaining
199 datasets, 12 months ahead forecasts are obtained for monthly datasets, 8 weeks ahead forecasts
200 are obtained for weekly datasets, except the solar weekly dataset, and 30 days ahead forecasts are
201 obtained for daily datasets. For the solar weekly dataset, we use a horizon of 5 as the series in this
202 dataset are relatively short compared with other weekly datasets. For half-hourly, hourly and other
203 low frequency datasets, we set the forecasting horizon to one week, e.g., 168 is used as the horizon
204 for hourly datasets.

205 The number of lagged values used in the PR models are determined based on a heuristic suggested
206 in prior work [56]. Generally, the number of lagged values is chosen as the seasonality multiplied
207 with 1.25. If the datasets contain short series and it is impossible to use the above defined number
208 of lags, for example in the Dominick and solar weekly datasets, then the number of lagged values
209 is chosen as the forecast horizon multiplied with 1.25, assuming that the horizon is not arbitrarily
210 chosen but based on certain characteristics of the time series structure. When defining the number of
211 lagged values for multi-seasonal datasets, we consider the corresponding weekly seasonality value,
212 e.g., 168 for hourly datasets. If it is impossible to use the number of lagged values obtained with the
213 weekly seasonality due to high memory and computational requirements, for example with the traffic
214 hourly and electricity hourly datasets, then we use the corresponding daily seasonality value to define
215 the number of lags, e.g., 24 for hourly datasets. In particular, due to high memory and computational
216 requirements, the number of lagged values is chosen as 50 for the solar 10 minutely dataset which is
217 less than the above mentioned heuristics based on seasonality and forecasting horizon suggest.

218 4 Results & discussion

219 This section details the results of feature analysis and baseline evaluation together with a discussion
220 of the results.

221 4.1 Feature analysis results

222 Figure 2 shows hexbin plots of the normalised density values of the low-dimensional feature space
223 generated by PCA across ACF1, trend, entropy, seasonal strength and Box-Cox lambda for 20
224 datasets. The figure highlights the characteristics among different datasets. For the M competition
225 datasets, the feature space is highly populated on the left-hand side and hence, denoting high trend
226 and ACF1 levels in the series. The tourism yearly dataset also shows high trend and ACF1 levels. In
227 contrast, the car parts, hospital, and Kaggle web traffic datasets show high density levels towards
228 the right-hand side, indicating a higher degree of entropy. The presence of intermittent series can
229 be considered as the major reason for the higher degree of entropy in the Kaggle web traffic and car
230 parts datasets. The plots confirm the claims of prior similar studies [34, 57] that the M competition
231 datasets are significantly different from the Kaggle web traffic dataset.

232 The monthly datasets generally show high seasonal strengths compared with datasets of other
233 frequencies. Quarterly datasets also demonstrate high seasonal strengths except for the M4 quarterly
234 dataset. In contrast, the datasets with high frequencies such as weekly, daily, and hourly show low
235 seasonal strengths except for the NN5 weekly and NN5 daily datasets.

236 Related to the shapes of the feature space, the 3 yearly datasets: M3, M4, and tourism show very
237 similar shapes and density populations indicating they have similar characteristics. The M4 quarterly
238 dataset also shows a similar shape as the yearly datasets, even though it has a different frequency. The

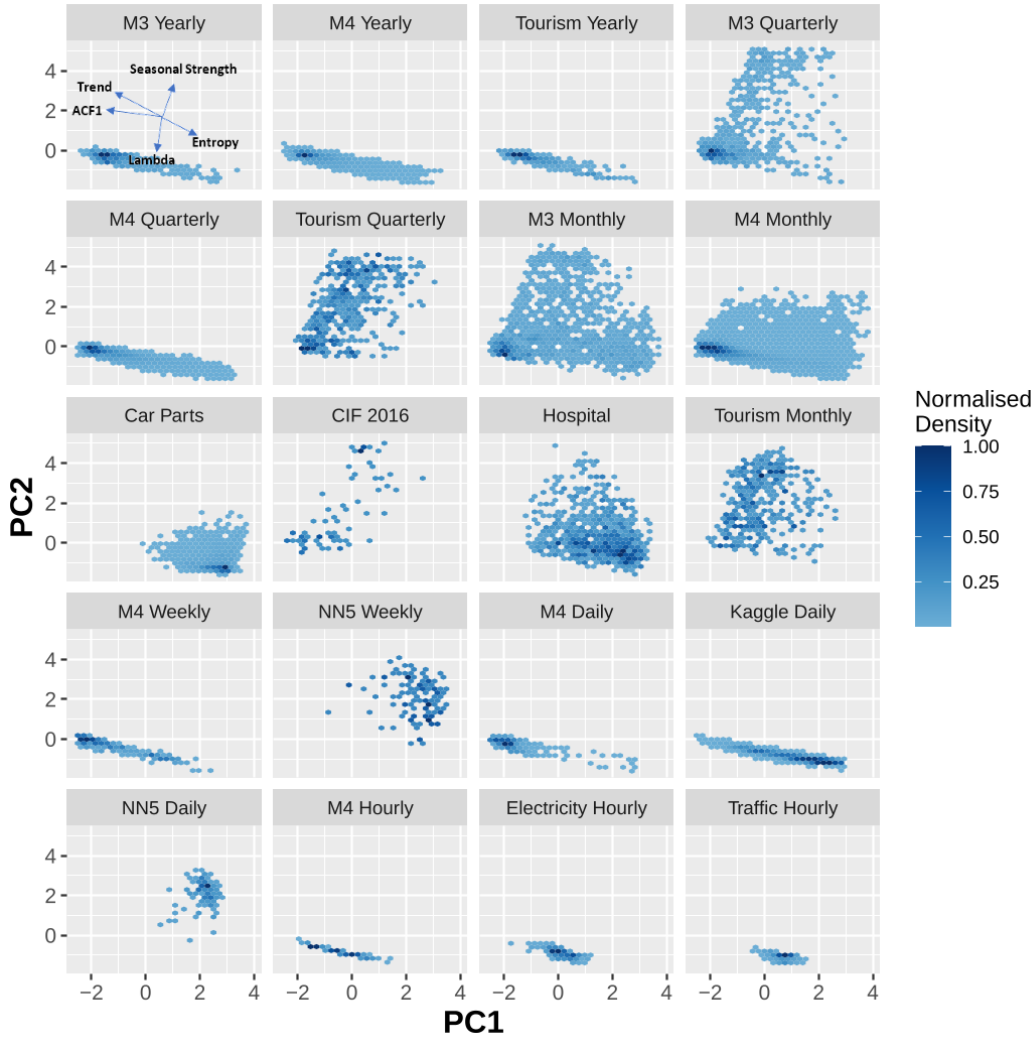


Figure 2: Hexbin plots showing the normalised density values of the low-dimensional feature space generated by PCA across ACF1, trend, entropy, seasonal strength, and Box-Cox lambda for 20 datasets. The dark and light hexbins denote the high and low density areas, respectively. The M3 Yearly facet shows the directions of the 5 features, which are the same across all facets.

239 other 2 quarterly datasets M3 and tourism are different, but similar to each other. The M3 and M4
 240 monthly datasets are similar to each other in terms of both shape and density population. Furthermore,
 241 the electricity hourly and traffic hourly datasets have similar shapes and density populations, whereas
 242 the M4 hourly dataset has a slightly different shape compared with them. The daily datasets show
 243 different shapes and density populations, where the NN5 daily dataset is considerably different from
 244 the other 2 daily datasets: M4 and Kaggle web traffic, in terms of shape and all 3 daily datasets are
 245 considerably different from each other in terms of density population. The weekly datasets also show
 246 different shapes and density populations compared with each other.

247 PCA plots showing the normalized density values of all datasets corresponding with both tsfeatures
 248 and catch22 features are available in the Appendix (supplementary materials).

249 4.2 Baseline evaluation results

250 It is very difficult to define error measures for forecasting that perform well under all situations
 251 [58], in the sense that it is difficult to define a scale-free measure that works for any type of non-

stationarity in the time series. Thus, how to best evaluate forecasts is still an active area of research, and (especially in the machine learning area) researchers often use ad-hoc, non-adequate measures. For example, usage of the Mean Absolute Percentage Error (MAPE) for normalised data between 0 and 1 may result in undefined or heavily skewed measures, or errors using the mean of a series like the Root Relative Squared Error [RSE, 27] will not work properly for series where the mean is essentially meaningless, such as series with steep trends. We use four error metrics that – while having their own problems – are common for evaluation in forecasting, namely the Mean Absolute Scaled Error [MASE, 59], symmetric MAPE (sMAPE), Mean Absolute Error [MAE, 60], and Root Mean Squared Error (RMSE) to evaluate the performance of the seven baseline forecasting models explained in Section 3.2. For datasets containing zeros, calculating the sMAPE error measure may lead to divisions by zero. Hence, we also consider the variant of the sMAPE proposed by Suilin [61] which overcomes the problems with small values and divisions by zero of the original sMAPE. We report the original sMAPE only for datasets where divisions by zero do not occur. Equations 2, 3, 4, 5, and 6, respectively, show the formulas of MASE, sMAPE, modified sMAPE, MAE, and RMSE, where M is the number of data points in the training series, S is the seasonality of the dataset, h is the forecast horizon, F_k are the generated forecasts and Y_k are the actual values. We set the parameter ϵ in Equation 4 to its proposed default of 0.1.

$$MASE = \frac{\sum_{k=M+1}^{M+h} |F_k - Y_k|}{\frac{h}{M-S} \sum_{k=S+1}^M |Y_k - Y_{k-S}|} \quad (2)$$

$$sMAPE = \frac{100\%}{h} \sum_{k=1}^h \frac{|F_k - Y_k|}{(|Y_k| + |F_k|)/2} \quad (3)$$

$$msMAPE = \frac{100\%}{h} \sum_{k=1}^h \frac{|F_k - Y_k|}{\max(|Y_k| + |F_k| + \epsilon, 0.5 + \epsilon)/2} \quad (4)$$

$$MAE = \frac{\sum_{k=1}^h |F_k - Y_k|}{h} \quad (5)$$

$$RMSE = \sqrt{\frac{\sum_{k=1}^h |F_k - Y_k|^2}{h}} \quad (6)$$

The MASE measures the performance of a model compared with the in-sample average performance of a one-step-ahead naïve or seasonal naïve (snaïve) benchmark. For multi-seasonal datasets, we use the length of the shortest seasonality to calculate the MASE. For the datasets where all series contain at least one full seasonal cycle of data points, we consider the series to be seasonal and calculate MASE values using the snaïve benchmark. Otherwise, we calculate the MASE using the naïve benchmark, effectively treating the series as non-seasonal.

The error metrics are defined for each series individually. We further calculate the mean and median values of the error metrics over the datasets to evaluate the model performance and hence, each model is evaluated using 10 error metrics for a particular dataset: mean MASE, median MASE, mean sMAPE, median sMAPE, mean msMAPE, median msMAPE, mean MAE, median MAE, mean RMSE and median RMSE. Table 2 shows the mean MASE of the SES, Theta, ETS, ARIMA, TBATS, DHR-ARIMA, and PR models on the same 20 datasets we considered for the feature analysis. The results of all baselines across all datasets on all 10 error metrics are available in the Appendix.

Overall, SES shows the worst performance and Theta shows the second-worst performance across all error metrics. ETS and ARIMA show a mixed performance on the yearly, monthly, quarterly, and daily datasets but both outperform SES and Theta. TBATS generally shows a better performance than DHR-ARIMA on the high frequency datasets. For our experiments, we always set the maximum order of Fourier terms used with DHR-ARIMA to $k = 1$. Based on the characteristics of the datasets, k can be tuned as a hyperparameter and it may lead to better results compared with our results. Compared with SES and Theta, both TBATS and DHR-ARIMA show superior performance.

Table 2: Mean MASE results. The best model across each dataset is highlighted in boldface.

Dataset	SES	Theta	ETS	ARIMA	TBATS	DHR-ARIMA	PR	CatBoost	FFNN	DeepAR	N-BEATS
NN5 Daily	1.521	0.885	0.865	1.013	-	-	1.263	0.973	1.409	-	-
NN5 Weekly	0.903	0.885	-	-	0.872	0.887	0.854	0.853	1.009	-	-
CIF 2016	1.291	0.997	0.841	0.929	-	-	1.019	1.175	2.434	-	-
Kaggle Daily	0.924	0.928	1.231	0.890	-	-	-	-	-	-	-
Tourism Yearly	3.253	3.015	3.395	3.775	-	-	3.516	3.553	7.352	-	-
Tourism Quarterly	3.210	1.661	1.592	1.782	-	-	1.643	1.793	6.424	-	-
Tourism Monthly	3.306	1.649	1.526	1.589	-	-	1.678	1.699	5.159	-	-
Traffic Hourly	1.922	1.922	-	-	2.482	2.535	1.281	-	-	-	-
Electricity Hourly	4.544	4.545	-	-	3.690	4.602	2.912	-	-	-	-
M3 Yearly	3.167	2.774	2.860	3.417	-	-	3.223	3.788	7.938	-	-
M3 Quarterly	1.417	1.117	1.170	1.240	-	-	1.248	1.441	4.212	-	-
M3 Monthly	1.091	0.864	0.865	0.873	-	-	1.010	1.065	2.215	-	-
M4 Yearly	3.981	3.375	3.444	3.876	-	-	3.625	-	-	-	-
M4 Quarterly	1.417	1.231	1.161	1.228	-	-	1.316	-	-	-	-
M4 Monthly	1.150	0.970	0.948	0.962	-	-	1.080	-	-	-	-
M4 Weekly	0.587	0.546	-	-	0.504	0.550	0.481	0.615	5.266	-	-
M4 Daily	1.154	1.153	1.239	1.179	-	-	1.162	1.593	-	-	-
M4 Hourly	11.607	11.524	-	-	2.663	13.557	1.662	1.771	-	-	-
Carparts	0.897	0.914	0.925	0.926	-	-	0.755	-	-	-	-
Hospital	0.813	0.761	0.765	0.787	-	-	0.782	0.798	0.986	-	-

289 The globally trained PR models show a mixed performance compared with the traditional univariate
 290 forecasting models. The performance of the PR models is considerably affected by the number of
 291 lags used during model training, performing better as the number of lags is increased. The number
 292 of lags we use during model training is quite high with the high-frequency datasets such as hourly,
 293 compared with the other datasets and hence, PR models generally show a better performance than the
 294 traditional univariate forecasting models on all error metrics across those datasets. But on the other
 295 hand, the memory and computational requirements are also increased when training PR models with
 296 larger numbers of lags. Furthermore, the PR models show better performance across intermittent
 297 datasets such as car parts, compared with the traditional univariate forecasting models.

298 We note that the MASE values of the baselines are generally high on multi-seasonal datasets. For
 299 multi-seasonal datasets, we consider longer forecasting horizons corresponding to one week unless
 300 they are competition datasets. As benchmark in the MASE calculations, we use a seasonal naïve
 301 forecast for the daily seasonality. As therewith the MASE compares the forecasts of longer horizons
 302 (up to one week) with the in-sample naïve forecasts obtained with shorter horizons (one day),
 303 the MASE values of multi-seasonal datasets are considerably greater than one across all baselines.
 304 Furthermore, the error measures are not directly comparable across datasets as we consider different
 305 forecasting horizons with different datasets.

306 5 Conclusion

307 Recently, global forecasting models and multivariate models have shown huge potential in providing
 308 accurate forecasts for collections of time series compared with the traditional univariate benchmarks.
 309 However, there are currently no comprehensive time series forecasting benchmark data archives
 310 available that contain datasets to facilitate the evaluation of these new forecasting algorithms. In
 311 this paper, we have presented the details of an archive that contains 20 publicly available time series
 312 datasets with different frequencies from varied domains. We have also characterised the datasets
 313 and have identified the similarities and differences among them by conducting a feature analysis
 314 exercise using tsfeatures and catch22 features extracted from each series. Finally, we have evaluated
 315 the performance of seven baseline forecasting models over all datasets across ten error metrics to
 316 enable other researchers to benchmark their own forecasting algorithms directly against those.

317 Acknowledgements

318 This research was supported by the Australian Research Council under grant DE190100045, a Face-
 319 book Statistics for Improving Insights and Decisions research award, Monash University Graduate
 320 Research funding and the MASSIVE High performance computing facility, Australia.

References

- [1] S. Makridakis, A. Andersen, R. F. Carbone, R. Fildes, M. Hibon, R. Lewandowski, J. Newton, E. Parzen, and R. L. Winkler. The accuracy of extrapolation (time series) methods: results of a forecasting competition. *Journal of Forecasting*, 1(2):111–153, 1982.
- [2] S. Makridakis and M. Hibon. The M3-competition: results, conclusions and implications. *International Journal of Forecasting*, 16(4):451–476, 2000.
- [3] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. The M4 competition: results, findings, conclusion and way forward. *International Journal of Forecasting*, 34(4):802–808, 2018.
- [4] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. The M4 competition: 100,000 time series and 61 forecasting methods. *International Journal of Forecasting*, 36(1):54 – 74, 2020. ISSN 0169-2070.
- [5] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. The M5 accuracy competition: results, findings and conclusions, 2020.
- [6] S. Ben Taieb, G. Bontempi, A. F. Atiya, and A. Sorjamaa. A review and comparison of strategies for multi-step ahead time series forecasting based on the NN5 forecasting competition. *Expert Systems with Applications*, 39(8):7067 – 7083, 2012.
- [7] Google. Web traffic time series forecasting, 2017. URL <https://www.kaggle.com/c/web-traffic-time-series-forecasting>.
- [8] S. Smyl. A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting. *International Journal of Forecasting*, 36(1):75–85, 2020.
- [9] T. Januschowski, J. Gasthaus, Y. Wang, D. Salinas, V. Flunkert, M. Bohlke-Schneider, and L. Callot. Criteria for classifying forecasting methods. *International Journal of Forecasting*, 36(1):167–177, 2020.
- [10] P. Montero-Manso and R. J. Hyndman. Principles and algorithms for forecasting groups of time series: locality and globality. *International Journal of Forecasting*, 2021. URL <https://arxiv.org/abs/2008.00444>. to appear.
- [11] S. Li, X. Jin, Y. Xuan, X. Zhou, W. Chen, Y.-X. Wang, and X. Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In *Advances in Neural Information Processing Systems*, page 5243–5253, 2019.
- [12] S. S. Rangapuram, M. W. Seeger, J. Gasthaus, L. Stella, Y. Wang, and T. Januschowski. Deep state space models for time series forecasting. In *Advances in Neural Information Processing Systems*, page 7785–7794, 2018.
- [13] R. Wen, K. Torkkola, B. Narayanaswamy, and D. Madeka. A multi-horizon quantile recurrent forecaster. In *31st Conference on Neural Information Processing Systems, Time Series Workshop*, 2017.
- [14] D. Salinas, M. Bohlke-Schneider, L. Callot, R. Medico, and J. Gasthaus. High-dimensional multivariate forecasting with low-rank Gaussian copula processes. In *Advances in Neural Information Processing Systems*, page 6827–6837, 2019.
- [15] R. Sen, H.-F. Yu, and I. S. Dhillon. Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting. In *Advances in Neural Information Processing Systems*, page 4837–4846, 2019.
- [16] H.-F. Yu, N. Rao, and I. S. Dhillon. Temporal regularized matrix factorization for high-dimensional time series prediction. In *Advances in Neural Information Processing Systems*, page 847–855, 2016.
- [17] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang. Informer: beyond efficient transformer for long sequence time-series forecasting. In *AAAI Conference on Artificial Intelligence*, 2020.

- 368 [18] D. Dua and C. Graff. UCI machine learning repository. <https://archive.ics.uci.edu/>,
369 2017.
- 370 [19] H. A. Dau, A. Bagnall, K. Kamgar, C.-C. M. Yeh, Y. Zhu, S. Gharghabi, C. A. Ratanamahatana,
371 and E. Keogh. The UCR time series archive. *IEEE/CAA Journal of Automatica Sinica*, 6(6):
372 1293–1305, 2019.
- 373 [20] A. Bagnall, H. A. Dau, J. Lines, M. Flynn, J. Large, A. Bostrom, P. Southam, and E. Keogh.
374 The UEA multivariate time series classification archive, 2018, 2018.
- 375 [21] R. J. Hyndman and Y. Yang. tsdl: time series data library. v0.1.0. [https://pkg.
376 yangzhuoranyang.com/tsdl/](https://pkg.yangzhuoranyang.com/tsdl/), 2018.
- 377 [22] Forecasting & Strategy Unit. Foredeck. <http://fsudataset.com/>, 2019.
- 378 [23] G. Paynter, L. Trigg, and E. Frank. Attribute-relation file format (arff). [https://www.cs.
379 waikato.ac.nz/ml/weka/arff.html](https://www.cs.waikato.ac.nz/ml/weka/arff.html), 2008.
- 380 [24] M. Löning, A. Bagnall, S. Ganesh, V. Kazakov, J. Lines, and F. J. Király. sktime: a unified
381 interface for machine learning with time series. In *Workshop on Systems for ML at NeurIPS
382 2019*, 2019.
- 383 [25] R. J. Hyndman, Y. Kang, P. Montero-Manso, T. Talagala, E. Wang, Y. Yang, and M. O’Hara-
384 Wild. *tsfeatures: time Series Feature Extraction*, 2020. URL [https://pkg.robjhyndman.
385 com/tsfeatures/](https://pkg.robjhyndman.com/tsfeatures/).
- 386 [26] Carl H Lubba, Sarab S Sethi, Philip Knaute, Simon R Schultz, Ben D Fulcher, and Nick S Jones.
387 catch22: Canonical time-series characteristics. *Data Mining and Knowledge Discovery*, 33(6):
388 1821–1852, 2019.
- 389 [27] G. Lai, W.-C. Chang, Y. Yang, and H. Liu. Modeling long- and short-term temporal patterns
390 with deep neural networks. In *The 41st International ACM SIGIR Conference on Research and
391 Development in Information Retrieval*, page 95–104, New York, NY, USA, 2018.
- 392 [28] Kaggle. <https://www.kaggle.com/>, 2019.
- 393 [29] Center for Systems Science and Engineering at Johns Hopkins University. COVID-19 data
394 repository. <https://github.com/CSSEGISandData/COVID-19>, 2020.
- 395 [30] C. Baum. KPSS: stata module to compute Kwiatkowski-Phillips-Schmidt-Shin test for station-
396 arity, 2018. URL <https://EconPapers.repec.org/RePEc:boc:bocode:s410401>.
- 397 [31] P. Phillips and P. Perron. Testing for a unit root in time series regression. *Cowles Foundation,
398 Yale University, Cowles Foundation Discussion Papers*, 75, January 1986.
- 399 [32] C. H. Lubba. *catch22: subset of hctsa-features*, 2018.
- 400 [33] B. D. Fulcher and N. S. Jones. hctsa: a computational framework for automated time-series
401 phenotyping using massive feature extraction. *Cell Sys.*, 5:527, 2017.
- 402 [34] C. S. Bojer and J. P. Meldgaard. Kaggle forecasting competitions: an overlooked learning
403 opportunity. *International Journal of Forecasting*, 2020. ISSN 0169-2070.
- 404 [35] R. J. Hyndman and Y. Khandakar. Automatic time series forecasting: the forecast package for
405 R. *Journal of Statistical Software*, 27(3):1–22, 2008. URL [http://www.jstatsoft.org/
406 v27/i03](http://www.jstatsoft.org/v27/i03).
- 407 [36] I. Jolliffe. *Principal Component Analysis*, pages 1094–1096. Springer, 2011. ISBN 978-3-642-
408 04898-2.
- 409 [37] L. J. P. van der Maaten and G. E. Hinton. Visualizing high-dimensional data using t-SNE.
410 *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- 411 [38] J. S. Armstrong. *Evaluating Forecasting Methods*, pages 443–472. Springer US, Boston, MA,
412 2001. ISBN 978-0-306-47630-3.

- 413 [39] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. Statistical and machine learning forecasting
414 methods: concerns and ways forward. *PLOS ONE*, 13(3):1–26, March 2018.
- 415 [40] R. J. Hyndman, A. B. Koehler, J. K. Ord, and R. D. Snyder. *Forecasting with exponential
416 smoothing: the state space approach*. Springer, 2008.
- 417 [41] G. E. P. Box and G. Jenkins. *Time Series Analysis, Forecasting and Control*. HoldenDay Inc.,
418 1990.
- 419 [42] V. Assimakopoulos and K. Nikolopoulos. The theta model: a decomposition approach to
420 forecasting. *International Journal of Forecasting*, 16(4):521–530, 2000.
- 421 [43] A. M. De Livera, R. J. Hyndman, and R. D. Snyder. Forecasting time series with complex
422 seasonal patterns using exponential smoothing. *Journal of the American Statistical Association*,
423 106(496):1513–1527, 2011.
- 424 [44] R. J. Hyndman. *fpp2: data for "Forecasting: Principles and Practice" (2nd Edition)*, 2018.
425 URL <https://CRAN.R-project.org/package=fpp2>.
- 426 [45] J. R. Trapero, N. Kourentzes, and R. Fildes. On the identification of sales forecasting models
427 in the presence of promotions. *Journal of the Operational Research Society*, 66(2):299 – 307,
428 2015.
- 429 [46] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- 431 [47] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, and A. Gulin. Catboost: unbiased
432 boosting with categorical features. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
433 N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
434 volume 31. Curran Associates, Inc., 2018.
- 435 [48] V. Flunkert, D. Salinas, and J. Gasthaus. DeepAR: probabilistic forecasting with autoregressive
436 recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2017.
- 437 [49] B. N. Oreshkin, D. Carpov, N. Chapados, and Y. Bengio. N-BEATS: neural basis expansion
438 analysis for interpretable time series forecasting. <https://arxiv.org/abs/1905.10437>,
439 2019.
- 440 [50] R. J. Hyndman, G. Athanasopoulos, C. Bergmeir, G. Caceres, L. Chhay, M. O’Hara-Wild,
441 F. Petropoulos, S. Razbash, E. Wang, F. Yasmeeen, R Core Team, R. Ihaka, D. Reid, D. Shaub,
442 Y. Tang, and Z. Zhou. *forecast: Forecasting Functions for Time Series and Linear Models*, 2021.
443 URL <https://CRAN.R-project.org/package=forecast>. R package version 8.14.
- 444 [51] J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models
445 via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2010. URL [http://www.
446 jstatsoft.org/v33/i01/](http://www.jstatsoft.org/v33/i01/).
- 447 [52] W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer, New York, fourth
448 edition, 2002. URL <http://www.stats.ox.ac.uk/pub/MASS4>.
- 449 [53] A. Alexandrov, K. Benidis, M. Bohlke-Schneider, V. Flunkert, J. Gasthaus, T. Januschowski,
450 D. C. Maddix, S. Rangapuram, D. Salinas, J. Schulz, L. Stella, A. C. TÅ¼rkmen, and Y. Wang.
451 GluonTS: probabilistic and neural time series modeling in python. *Journal of Machine Learning
452 Research*, 21(116):1–6, 2020.
- 453 [54] K. Bandara, C. Bergmeir, and H. Hewamalage. LSTM-MSNet: leveraging forecasts on sets of
454 related time series with multiple seasonal patterns. *IEEE Transactions on Neural Networks and
455 Learning Systems*, 2019.
- 456 [55] R. Godahewa, C. Bergmeir, G. I. Webb, and P. Montero-Manso. A strong baseline for weekly
457 time series forecasting. <https://arxiv.org/abs/2010.08158>, 2020.
- 458 [56] H. Hewamalage, C. Bergmeir, and K. Bandara. Recurrent neural networks for time series
459 forecasting: current status and future directions. *International Journal of Forecasting*, 2021.

- 460 [57] C. Fry and M. Brundage. The M4 forecasting competition - a practitioner's view. *International*
461 *Journal of Forecasting*, 36(1):156 – 160, 2020. ISSN 0169-2070.
- 462 [58] R. J. Hyndman and G. Athanasopoulos. *Forecasting: principles and practice*. OTexts, 3rd
463 edition, 2021. URL <http://OTexts.com/fpp3>.
- 464 [59] R. J. Hyndman and A. B. Koehler. Another look at measures of forecast accuracy. *International*
465 *Journal of Forecasting*, 22(4):679–688, 2006.
- 466 [60] Mean absolute error. In C. Sammut and G. I. Webb, editors, *Encyclopedia of Machine Learning*,
467 pages 652–652. Springer US, Boston, MA, 2010. ISBN 978-0-387-30164-8.
- 468 [61] A. Suilin. `kaggle-web-traffic`. <https://github.com/Arturus/kaggle-web-traffic>,
469 2017.

470 Checklist

- 471 1. For all authors...
- 472 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
473 contributions and scope? [Yes] See Sections 2, 3 and 4.
- 474 (b) Did you describe the limitations of your work? [Yes] See Sections 3 and 4. We mention
475 that we could not consider all datasets for feature analysis and benchmark evaluation
476 due to the high computational requirements of some datasets.
- 477 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Not
478 applicable to our work.
- 479 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
480 them? [Yes]
- 481 2. If you are including theoretical results...
- 482 (a) Did you state the full set of assumptions of all theoretical results? [N/A] Not applicable
483 to our work.
- 484 (b) Did you include complete proofs of all theoretical results? [N/A] Not applicable to our
485 work.
- 486 3. If you ran experiments (e.g. for benchmarks)...
- 487 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
488 mental results (either in the supplemental material or as a URL)? [Yes] See Appendix
489 A of supplementary materials.
- 490 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
491 were chosen)? [Yes] See Section 3.2.
- 492 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
493 ments multiple times)? [Yes] See Section 4.2 and Appendix C of supplementary
494 materials.
- 495 (d) Did you include the total amount of compute and the type of resources used (e.g., type
496 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D of supplementary
497 materials.
- 498 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 499 (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix A.6
500 and A.7 of supplementary materials.
- 501 (b) Did you mention the license of the assets? [Yes] See Appendix A.3 and A.5 of
502 supplementary materials.
- 503 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
504 See Appendix A of supplementary materials.
- 505 (d) Did you discuss whether and how consent was obtained from people whose data you're
506 using/curating? [Yes] See Appendix A.5 of supplementary materials.
- 507 (e) Did you discuss whether the data you are using/curating contains personally identifiable
508 information or offensive content? [N/A] Not applicable to our work.

- 509 5. If you used crowdsourcing or conducted research with human subjects...
- 510 (a) Did you include the full text of instructions given to participants and screenshots, if
- 511 applicable? [N/A] Not applicable to our work.
- 512 (b) Did you describe any potential participant risks, with links to Institutional Review
- 513 Board (IRB) approvals, if applicable? [N/A] Not applicable to our work.
- 514 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 515 spent on participant compensation? [N/A] Not applicable to our work.