# PYRREGULAR: A UNIFIED FRAMEWORK FOR IRREGULAR TIME SERIES, WITH CLASSIFICATION BENCHMARKS

**Anonymous authors** 

Paper under double-blind review

#### **ABSTRACT**

Irregular temporal data, characterized by varying recording frequencies, differing observation durations, and missing values, presents significant challenges across fields like mobility, healthcare, and environmental science. Existing research communities often overlook or address these challenges in isolation, leading to fragmented tools and methods. To bridge this gap, we introduce a unified framework, and the first standardized dataset repository for irregular time series classification, built on a common array format to enhance interoperability. This repository comprises 34 datasets on which we benchmark 12 classifier models from diverse domains and communities. This work aims to centralize research efforts and enable a more robust evaluation of irregular temporal data analysis methods.

# 1 Introduction

High-dimensional temporal data is increasingly accessible to decision-makers, domain experts, and researchers (Shumway et al., 2000). It is vital in fields like mobility, healthcare, and environmental science to capture dynamic changes over time. Yet, variations in recording frequencies, durations across sensors, and occasional failures lead to signals with unequal lengths, gaps, and missing values (Harvey et al., 1998). These traits make real-world temporal data irregular and hard to manage.

Several research communities address the challenge of irregular temporal data from different perspectives, as its analysis depends heavily on the task, application setting, and modeling approach. As a result, the problem spans multiple fields, including mobility analytics (da Silva et al., 2019), irregular time series classification (Kidger et al., 2020), forecasting (Weerakody et al., 2021), and imputation (Luo et al., 2018; Li & Marlin, 2020), to name a few. Due to this vast amount of tasks, and despite some shared challenges, communities working on irregular temporal data tend to be separated, each relying on its own set of techniques, such as traditional statistical or data mining models (Hamilton, 2020), neural networks (Wang et al., 2024), or differential equations (Rubanova et al., 2019), often resulting in domain-specific tools and libraries. This is not inherently a drawback, but can lead to fragmented research efforts. The challenges of irregular temporal data are amplified in supervised learning, where standardized benchmarks are notably lacking. While repositories exist for regular time series classification (Dau et al., 2019), truly irregular datasets, capturing real-world missingness and variability, remain scarce. Researchers often resort to artificially manipulated datasets (Weerakody et al., 2021), introducing assumptions that overlook structural missingness tied to data collection (Mitra et al., 2023). As a result, and given that many studies rely on a narrow range of datasets, the generalizability of their methods often remains untested.

We bridge this gap by proposing pyrregular, a unified framework for irregular time series. (1) We introduce a taxonomy of irregularities and a dataset structure in a common array format that improves interoperability across libraries while supporting the handling, visualization, and modeling of irregular time series using existing analysis methods. (2) We introduce the first standardized dataset repository for irregular time series classification, and (3) we leverage this repository to propose the first generalized benchmark for state-of-the-art classifiers from different research domains, in an effort to centralize research on this topic. Specifically, we curate 34 irregular time series datasets and evaluate 12 time series classifiers. Our goal is to empower users to seamlessly explore and evaluate a wide range of libraries to address the challenges of irregular temporal data.

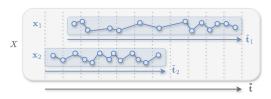


Figure 1: An example of an irregular time series, X, comprising two signals  $\mathbf{x}_1, \mathbf{x}_2$  with indices  $\tilde{\mathbf{t}}_1, \tilde{\mathbf{t}}_2$ , and the combined shared index  $\tilde{\tilde{\mathbf{t}}}$ .

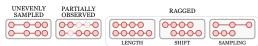


Figure 2: Different kinds of irregularity shown on a multivariate time series with 2 signals and containing up to 5 timestamps. Missing values are depicted as faded red if they were expected to be recorded, while they are omitted if they are caused by raggedness.

#### 2 Organizing Irregularity

As our first contribution, we propose a systematic taxonomy that clearly distinguishes among different forms of irregularity. We begin by defining a time series signal.

**Definition 2.1** (Time Series Signal). A signal (or channel) is a sequence of  $\tau$  observations, each associated to a timestamp, i.e.,  $\mathbf{x} = [(x_1, t_1), \dots, (x_{\tau}, t_{\tau})] = [x_{t_1}, \dots, x_{t_{\tau}}] \in \mathbb{R}^{\tau}$ .

A single signal can be *irregular* for two reasons: *uneven sampling*, when at least one interval  $t_{k+1} - t_k$  differs from a constant  $\Delta t$ , and *partially observed*, when expected values are missing and marked as NaN. The set of real numbers extended with the NaN symbol is here represented as  $\mathbb{R}$ . We denote with  $\tilde{\mathbf{t}} = [t_1, \dots, t_{\tau}] \in \mathbb{R}^{\tau}$ , the sorted collection of all timestamps where an observation of signal  $\mathbf{x}$  was, or should have been recorded, and with  $\tau = |\tilde{\mathbf{t}}|$  the number of observations.

**Definition 2.2** (Time Series). A time series is a collection of d signals,  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_d\} \in \mathbb{R}^{d \times T}$ .

Time series timestamps are the sorted union of all signal timestamps, i.e.,  $\tilde{\mathbf{t}} = \bigcup_{j=1}^d \tilde{\mathbf{t}}_j \in \mathbb{R}^T$ , with  $T = |\tilde{\mathbf{t}}|$ , as shown in Figure 1. In addition to these intrinsic irregularities, tensor representations introduce a third, structural type: raggedness, that is the necessity of padding due to length, sampling, or alignment mismatches between signals. Hence, there are three independent irregularity causes:  $uneven\ sampling$ ,  $partial\ observation$ , and raggedness, as depicted in Figure 2. While these categories have appeared informally in prior literature, here we show that they are independent: none implies the others. Unevenly sampled time series do not necessarily imply the presence of partially observed data, as seen in Figure 2 (left). This commonly happens in trajectory data, where the timestamps are usually highly uneven, but shared across the latitude and longitude signals. Vice versa, the presence of unobserved data does not imply uneven timestamps, as an observation may be accidentally missing from an overall constant sampling. Finally, neither unevenly sampled nor partially observed data imply raggedness. In particular, the two leftmost time series shown in Figure 2 could be stored in  $2 \times 4$  and  $2 \times 5$  matrices, respectively, without requiring any padding.

Raggedness arises because of different issues created when storing a multivariate time series in an array-like structure. As so, a single, univariate signal cannot be ragged by itself. In general, raggedness arises when at least two signals, a and b, do not share the same timestamps, i.e.,  $\mathbf{t}_a \neq \mathbf{t}_b$ . We identify three independent fundamental reasons for why this can happen. The first is ragged length, when a and b have a different number of observations:  $\tau_a \neq \tau_b$ . The second is shift, where at least one signal starts and ends before another:  $(t_{a,1} < t_{b,1}) \land (t_{a,\tau_a} < t_{b,\tau_b})$ . The third is ragged sampling, when at least one element of the sampling intervals differs between two signals, i.e.,  $\Delta t_{a,k} \neq \Delta t_{b,k}$  for some k, where  $\Delta t_{a,k} = t_{a,k+1} - t_{a,k}$  and  $\Delta t_{b,k} = t_{b,k+1} - t_{b,k}$ . Again, none of these, by itself, implies the other, as shown in Figure 2, and, in more detail, in Appendix B. Combinations of these issues yield highly irregular data, where NaN can indicate either a missing value in a partially observed time series or padding due to raggedness. Moreover, raggedness can exist also in a time series dataset, i.e., a collection of n time series,  $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_n\} \in \mathbb{R}^{n \times d \times \mathcal{T}}$ , as all instances share the same sorted timestamps,  $\mathbf{t} = \bigcup_{i=1}^n \tilde{\mathbf{t}}_i \in \mathbb{R}^{\mathcal{T}}$ , with  $\mathcal{T} = |\mathbf{t}|$ . The timestamp index for the whole dataset is denoted as  $\mathbf{k} = [1, \dots, \mathcal{T}]$ .

Associated with time series datasets are often *static attributes*, which refer to information linked to individual instances that remain independent of the time dimension. These attributes can also serve as targets in supervised tasks. Specifically, we focus on classification, i.e., targets are categorical.

# 3 RELATED WORK

108

109 110

111

112

113

114

115 116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

147

148

149

150

151

152

153 154 155

156 157

158

159

160

161

**Datasets and Benchmarks.** There is a significant divide in the literature in the availability of datasets and benchmarking efforts, between regular and irregular time series data. Supervised learning for regular time series data is extensively addressed in the literature, with numerous "bake-offs" (Bagnall et al., 2017; Ruiz et al., 2021; Middlehurst et al., 2024b) benchmarking state-of-the-art classifiers on hundreds of standard datasets from the UEA and UCR repositories (Dau et al., 2019; Bagnall et al., 2018). On the contrary, the benchmarking literature on *irregular* time series remains limited. While secondary sources, such as (Weerakody et al., 2021; Wang et al., 2024), offer surveys on specific tasks like ITS imputation, comprehensive benchmarks for downstream tasks like classification are largely confined to primary studies (Kidger et al., 2020; Shukla & Marlin, 2021; Du et al., 2023). Even within these studies, evaluations are often performed on a small number of datasets. Moreover, benchmark datasets are not always inherently irregular; instead, they are commonly derived from regular datasets through simulation, i.e., dropping valid observations (Weerakody et al., 2021). Although this strategy can create ITS, introducing missingness is a non-trivial process requiring careful decisions about the type of missingness to simulate (Rubin, 1976). Adding to these challenges, a recent study (Mitra et al., 2023) highlighted that most research neglects structural missingness, referring to non-random, multivariate patterns of missingness within datasets. Such patterns can be faithfully preserved only by maintaining the original data with minimal modifications, which is the central focus of this proposal.

Libraries. Regarding regular time series data, Python libraries such as sktime (Löning et al., 2019), aeon (Middlehurst et al., 2024a), and tslearn (Tavenard et al., 2020) provide a wide range of classifier implementations, along with access to the UEA and UCR repositories, enabling systematic and reproducible evaluations. Although some of these datasets contain irregularities, the typical approach involves imputing missing values and discarding timestamps during downstream tasks. The most prominent Python library for irregular time series analysis is pypots (Du, 2023). pypots offers several classifiers, a few partially observed time series datasets, and provides an interface for adding missingness in regular datasets. A limitation of pypots is that it overlooks irregularity from uneven sampling, ignoring timestamps. It also operates within its own ecosystem, lacking interfaces for cross-library comparisons. This makes using ITS with libraries like aeon and sktime difficult, due to incompatible data formats and requirements, hindering standardization efforts. The primary reason for these challenges is the difficulty in managing ITS due to high dimensionality, missing values, and timestamps. Most libraries for time series prediction require dense 3D tensors to represent time series, signals, and identifiers (IDs), often demanding extensive padding and increased memory usage. To mitigate this, special arrays to represent missing values or variable-length instances are often used. For example, numpy masked arrays (Harris et al., 2020) indicate valid entries with masks but are memory-inefficient since they store both data and masks. Alternatives include awkward arrays (Pivarski et al., 2020), jagged pytorch arrays (Paszke et al., 2017), ragged tensorflow arrays (Abadi et al., 2015), zarr, pyarrow, or sparse arrays (Abbasi, 2018). Although efficient in managing varied-sized data, these structures cannot inherently handle timestamps. Forecasting libraries like nixtla or gluonTS (Alexandrov et al., 2020) typically use a long format, representing data as tuples (i, j, t, x) with instance and signal IDs, timestamps, and observed values. While efficient for forecasting, this format requires pivoting for classification tasks, and static variables are either duplicated or stored separately, causing inefficiencies. Lastly, xarray (Hoyer & Hamman, 2017) supports timestamped multi-dimensional arrays but lacks native support for sparse ITS.

In summary, to the best of our knowledge, no existing array format is capable of representing *ITS* data in all their nuances. To address this limitation, we propose a framework that serves as a compatibility layer based on a unified array format, facilitating comprehensive benchmarking across a wide range of datasets and methods from diverse time series communities.

#### 4 A Unified Framework for Irregular Time Series

This work addresses the gap in the literature on irregular time series by introducing an efficient container specifically designed for such data. This facilitates the integration of methods and datasets from various research communities into a unified framework. We outline key aspects of this solution. (i) Ease of Use: the framework supports several stages of the data science workflow, including visualization, preprocessing with classical and temporal slicing, and seamless conversion to dense arrays used in leading machine learning libraries. (ii) Robustness: the implementation leverages established

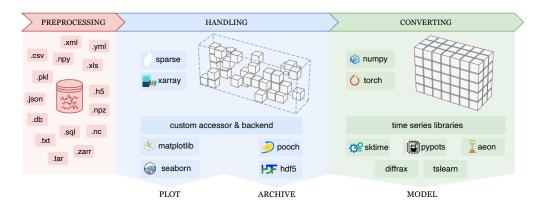


Figure 3: A simplified schema of our framework. (left) Data from different sources is preprocessed and represented in our proposed array container (center), which combines xarray with an underlying sparse tensor via a custom accessor and backend. This container can be easily manipulated, plotted, and stored. (right) Finally, it can also be converted into a more common dense representation, which can be used for downstream tasks with any standard time series library.

and well-maintained libraries, as there is no point in reinventing the wheel. (iii) Flexibility: the container supports several types of time series irregularities. (iv) Replicability: to ensure comparable results, preprocessing is standardized, addressing the variability in ITS. A depiction of the three steps of pyrregular is shown in Figure 3: preprocessing, where the original ITS is transformed into our proposed container; handling, where the data can be explored, manipulated, and stored; and converting, where the data is prepared for downstream tasks. <sup>1</sup>

**Preprocessing.** The first step in our framework involves transforming *ITS* datasets into the proposed representation. *ITS* can be found in a wide variety of sources and formats (Figure 3, left), presenting unique challenges in terms of preprocessing. Regardless of the original data structure, our framework requires only a function capable of yielding the data in the standardized *long format*. In this representation, each row captures the time series ID, signal ID, timestamp, and observed value: (i, j, t, x). The core intuition behind our approach is that the long format closely resembles the sparse coordinate (COO) representation (Duff et al., 2017).

The COO format, as implemented by sparse (Abbasi, 2018), can efficiently encode sparse 3D tensors, by using indices for the time series, signal, and timestamp, accompanied by an observed value entry, formally (i, j, k, x). The key distinction between the long format and the COO representation lies in the handling of the timestamps: while the COO format requires discrete timestamp indices, k, the long format uses real-valued timestamps, t. An example is reported in Figure 4 (left). This difference, however, can be easily bridged by mapping the timestamps, t, to discrete positions within the COO array, k. Formally, given the timestamps vector  $\mathbf{t} = [t_1, \dots, t_T]$ , each timestamp can be mapped to its corresponding position (index), in the COO format as  $\mathbf{k} = [1, \dots, T]$  (and vice-versa), as depicted in Figure 4 (center). With this mapping, converting between the long format and the COO representation can be easily accomplished, as the time series dataset is read once to construct the mapping and a second time to incrementally build the COO matrix by yielding each row as it is generated (Figure 4, right). Practitioners need only to define a custom function that, given their own data, incrementally produces rows in the long format. Even when the initial dataset is not organized in this manner, the conversion to the long format is typically straightforward. This process ensures uniformity across input formats and transparency, as the preprocessing steps are explicitly documented in this function, and can be reproduced at any time. Though it may be runtime-intensive, this step needs to be performed only once, after which the library streamlines all subsequent transformations and processing. The output after preprocessing is a sparse tensor, denoted as  $\mathbf{X} \in \mathbb{R}^{n \times d \times T}$ .

**Handling.** The COO representation offers advantages over the classical long format. First, it supports array-like operations with reasonable performance, including reshaping and slicing. Moreover, it allows for rapid conversion to task-specific array structures, such as other sparse formats like GCXS

<sup>&</sup>lt;sup>1</sup>The code is provided in the Supplementary Materials. Examples are available in Appendix H.

(Shaikh & Hasan, 2015). Compared to classical dense arrays, its primary advantage lies in memory efficiency, as only the recorded observations are stored. All padding is represented by a *fill value* and remains implicit, meaning it is not directly stored but is generated only when the sparse array is transformed into a dense form. We propose setting such value to NaN to capture raggedness. Further, the COO format naturally accommodates partially observed data by explicitly storing a fill value. This allows for distinguishing between the two types of missing data previously discussed. Specifically, an explicitly stored fill value, i.e., a row (i, j, k, NaN), can indicate a missing entry that should be present, while implicit NaNs reflect missingness due to data raggedness. In this sense, the COO tensor by itself is enough to represent both ragged and partially observed time series.

However, to capture an unevenly sampled time series, it is also essential to store the timestamps. To achieve this, we leverage the timestamp to COO (t to k) mapping using xarray (Figure 3, center). In particular, we use xarray (Hoyer & Hamman, 2017) to store the timestamps and extend it to utilize an underlying sparse COO tensor. These functionalities are possible through our custom backend and accessor, which extend the xarray library, to support sparse arrays. Further, xarray naturally facilitates the storage of static attributes linked to any dataset dimension, such as class labels in classification tasks. Overall, this approach offers significant

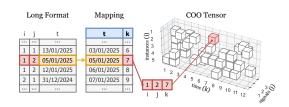


Figure 4: Long format to COO tensor conversion process. Each row of the long format is processed to retrieve the absolute position k of a given timestamp t. The triplet, instance ID (i=1), signal ID (j=2), and timestamp index (k=7), is used to populate the sparse COO tensor.

storage efficiency, particularly given the typically high data sparsity (see Appendix E), and ensures ease of use by supporting all existing xarray functions like timestamp range queries. Further, our accessor enables plotting, while our backend allows direct saving and loading to a hierarchical data format, locally or online, eliminating the need to perform the preprocessing step again.

Converting. Despite its advantages, xarray is not directly supported by most libraries for supervised learning tasks. Therefore, it is crucial to demonstrate how this array structure can be efficiently prepared for such applications. Specifically, for classification tasks,  $\mathbf{X} \in \mathbb{R}^{n \times d \times T}$  should be transformed into a dense tensor that minimizes raggedness while preserving the inherent missingness from partially observed time series and maintaining the order of observations within the same time series. This conversion is important because, in classification tasks, raggedness is typically irrelevant to the target and would otherwise result in vast dense arrays filled predominantly with NaNs. For instance, the specific starting dates of time series, such as a beginning on January 23rd and b on January 30th, are typically uninformative with respect to the output class, so we generally want to avoid introducing 7 leading NaNs in time series b to account for the shift. For a COO array, this transformation corresponds to a dense ranking operation on the timestamp index, k, performed time series-wise. Formally, for each COO entry (i, j, k, x), we produce  $(i, j, rank_i(k), x)$ , where:

$$rank_i(k) = 1 + |\{k' \in [1, T_i] : k' < k\}|.$$

This process shifts the timestamp indices within each time series,  $X_i$ , into a consecutive sequence ranging from 1 to its length,  $T_i$ . As a result, the tensor  $\mathbf{X} \in \mathbb{R}^{n \times d \times T}$  can be densified into a more compact,  $\mathbf{X}' \in \mathbb{R}^{n \times d \times T}$ , where  $T = \max_i^n(T_i)$ . This ensures minimal raggedness, with the timestamp dimension set to the maximum number of timestamps in any time series.  $\mathbf{X}'$  can be used by downstream libraries such as sktime (Löning et al., 2019), aeon (Middlehurst et al., 2024a), tslearn (Tavenard et al., 2020), pypots (Du, 2023) and diffrax (Kidger, 2021).

### 5 CLASSIFICATION BENCHMARKS

We present a comprehensive benchmark enabled by pyrregular, in which we evaluate 12 classifiers from a variety of time series libraries on a curated collection of 34 *ITS* datasets. We assess model performance from multiple perspectives, including dataset characteristics, robustness across irregularity types, and the potential for performance improvement through fine-tuning.

271

272273274275276277278279

281

284

287

289

291292293

295296297298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318 319

320

321

322

323

Table 1: Datasets used for our benchmarks, divided by irregularity type: unevenly sampled (US), partially observed (PO), unequal length (UL), shift (SH), ragged sampling (RS).

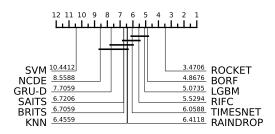
	h	eal	th			hui	nai	n a	ctiv	ity	re	cog	nii	tion	ı					тc	bil	ity				se	ense	or		0	the	r		synth
	MI3	P12	P19	CJ	GM1	GM2	GM3	GP1	GP2	Ϋ́S	ĞΥ	ZS	LPA	PAM	PGZ	SgZ	AN	AOC	APT	ARC	SS	Æ	SE	TA	ΛE	8	DG	MQ	MI	Ŋ	PGE	PL	SAD	ABF
US	/	/	/	Х	Х	Х	Х	Х	Х	Х	Х	Х	/	/	Х	Х	/	Х	Х	Х	1	Х	/	/	/	Х	Х	Х	Х	Х	/	X	Х	/
PO	1	1	1	X	X	X	X	X	X	X	X	X	X	1	X	X	X	X	X	X	X	X	X	X	X	1	1	1	X	X	X	X	X	X
UL	1	1	1	/	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	X	X	X	1	1	1	1	1	X
SH	1	1	1	X	X	X	X	X	X	X	X	X	1	1	X	X	X	X	X	X	1	X	1	1	X	X	X	X	X	X	1	X	X	X
RS	✓	1	✓	X	X	X	X	X	X	X	X	X	1	1	X	X	1	X	X	1	✓	X	1	1	1	X	X	X	X	X	✓	X	X	X

Table 2: Summary of evaluated classifiers.

Library		Model	Type	Domain
aeon	(Spinnato et al., 2024)	BORF RIFC	dictionary-based transform + LGBM classifier interval-based transform + LGBM classifier	regular, ragged partially observed
diffrax	(Kidger et al., 2020)	NCDE	neural controlled differential equations	unevenly sampled
pypots	(Cao et al., 2018) (Che et al., 2018) (Zhang et al., 2022) (Du et al., 2023) (Wu et al., 2022)	BRITS GRU-D RAINDROP SAITS TIMESNET	bidirectional recurrent imputation network gated recurrent unit with decay graph neural network self-attention-based imputation transformer temporal 2d-variation transformer.	partially observed partially observed partially observed partially observed partially observed
sktime	(Ke et al., 2017) (Dempster et al., 2021) (Bagheri et al., 2016)	LGBM ROCKET SVM	gradient boosted tree kernel-based transform + LGBM classifier support vector machine with distance kernel	tabular regular regular, ragged
tslearn	(Sakoe & Chiba, 1978)	KNN	distance-based with dynamic time warping	regular, ragged

**Datasets.** Following established repositories such as UEA and UCR, we compile a diverse collection of datasets that vary in size (small to large), length (short to long), and dimensionality (univariate to multivariate), ensuring broad representativeness. We solely focus on naturally irregular datasets, without artificially inducing irregularity (Tables 1 and 5). First, our collection contains widely used ITS classification datasets: PhysioNet 2012 (P12) (Silva et al., 2012), PhysioNet 2019 (P19) (Reyna et al., 2020), and the MIMIC-III (MI3) clinical database (Johnson et al., 2016) from the medical domain, as well as *Pamap2* (PAM) (Reiss & Stricker, 2012) for physical activity monitoring. Additionally, we include the 11 variable-length univariate time series classification problems (Guna et al., 2014; Caputo et al., 2018; Mezari & Maglogiannis, 2018; Gao et al., 2014) from (Bagnall et al., 2020), the 4 partially observed datasets (Ihler et al., 2006; City of Melbourne, 2020) from (Middlehurst et al., 2024b), and the 7 variable-length multivariate time series classification problems (de Souza, 2018; Williams et al., 2006; Chen et al., 2014; Kudo et al., 1999; Hammami & Bedda, 2010) from (Ruiz et al., 2021). We also provide datasets that, to the best of our knowledge, were never used in these kinds of benchmarks. These include data for trajectory classification of entities such as mammals (AN)(Ferrero et al., 2018), birds (SE) (Browning et al., 2018), and vehicles like buses and trucks (VE), taxis (Moreira-Matias et al., 2013) (TA) and combinations of the previous (Zheng et al., 2010) (GS). Further, we include a small dataset about the productivity prediction for garment employees (Imran et al., 2021) (PGE), and a human activity recognition dataset (Vidulin et al., 2010) (LPA). Finally, inspired by the classical Cylinder-Bell-Funnel benchmark (Saito, 1994) for regular time series classification, we introduce an irregular version called Alembics-Bowls-Flasks (ABF), in which the class depends on the skewness of the time sampling. Where available, we use the default train/test split for training and inference, else we set them based on each dataset description and original paper.

**Models.** The objective of these experiments is to benchmark methods capable of naturally handling *ITS* without introducing bias through imputation. For this reason, and to keep the benchmarks to a reasonable amount, we limit our evaluation to classifiers that inherently support irregular inputs and are available in the aforementioned libraries (Table 2 and Appendix C). As classical baselines, we use K-Nearest Neighbors (KNN) with Dynamic Time Warping (Sakoe & Chiba, 1978), a time series Support Vector Machine (SVM) with a Longest Common Subsequence (LCSS) kernel (Bagheri et al.,



325

326

327

328

331

332

333

334

335 336 337

338

339

340

341

342

343

344 345

346

347

348

349

350

351

352

353

354

355

356

357 358

359 360

361

362

363

364

366

367

368

369

370

371

372

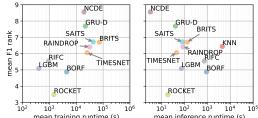
373

374

375

376

377



terms of F1. Best models to the right. Connected models are statistically tied.

Figure 5: CD plot for the benchmarked models in Figure 6: Mean F1 rank against training and inference runtimes for the top 11 models across all datasets. The best models are on the bottom left.

2016), and a LightGBM classifier (LGBM) trained directly on raw ITS, ignoring temporal dependencies. For regular time series models, we include the Bag-Of-Receptive-Fields (BORF) (Spinnato et al., 2024) from aeon, ROCKET (Dempster et al., 2020; 2021) via its MINIROCKET version in sktime, and a Random Interval Feature Classifier (RIFC). These models transform the data and rely on downstream classifiers; we use LGBM to handle possible NaNs. For partially observed data, we benchmark GRU-D (Che et al., 2018), BRITS (Cao et al., 2018), RAINDROP (Zhang et al., 2022), two transformer models, SAITS (Du et al., 2023) and TIMESNET (Wu et al., 2023), from pypots, and a Neural Controlled Differential Equation model (NCDE) (Kidger et al., 2020) from diffrax.

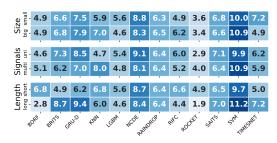
**Experimental Setup.** Following standard practice in similar benchmarking studies (Bagnall et al., 2017; Middlehurst et al., 2024b), all models are trained using the default hyperparameters provided by their respective libraries or those recommended in the original papers. The goal of this benchmark, consistent with prior bake-offs, is to identify the model that best generalizes with a single, reasonable parameter configuration rather than fine-tuning each model for individual datasets. For this reason, the results of these benchmarks do not necessarily highlight the best possible model for a given task, but the model that generalizes best in many. Each model is allocated two weeks ( $\approx 20000$  minutes) for training and inference on each dataset, with access to 32 cores and 512 GB of memory, and to a GPU when the model can use it<sup>2</sup>. Experiments are repeated three times for highly stochastic models, and the average performance is maintained. We use the F1 score with macro averaging as the primary performance metric, as it is robust in the presence of unbalanced data (Japkowicz, 2013), which occurs in some of our datasets. Accuracy results, along with additional metrics and statistical tests, are reported in Appendix D and are consistent with the following findings.

#### 5.1 RESULTS AND DISCUSSION.

We present a comparative analysis of the aggregate results of the benchmark outcomes. We report a critical difference (CD) plot in Figure 5, which ranks models in terms of F1. Models are arranged from right to left, with lower ranks indicating better performance. Models connected by a horizontal bar are statistically tied under a one-sided Holm-corrected Wilcoxon signed-rank test with a significance threshold of 0.05. ROCKET emerged as the clear top-performing model, demonstrating consistent superiority across the datasets. Even if this result aligns with its established reputation as one of the best models for regular time series classification (Middlehurst et al., 2024b), its efficacy on irregular data is somewhat surprising, as ROCKET does not exploit any information about said irregularity. Following ROCKET, a cluster of methods, including BORF, LGBM, RIFC, TIMESNET, exhibits statistically tied performance. Lower ranks are occupied by RAINDROP, KNN, BRITS, followed by GRU-D and NCDE, with SVM distinctly identified as the worst-performing model.

Performance vs. Time. Besides predictive performance, runtime is also a significant factor. In Figure 6, we compare the average F1 rank against training and inference runtimes, discarding SVM for better readability. The better-performing, faster models appear in the bottom-left region of the plot. In terms of training, LGBM is the fastest, followed by RIFC and ROCKET, with ROCKET also being also very fast during inference. For this reason, ROCKET emerges as the best tradeoff between F1 and runtime. Interestingly, despite being designed for tabular data, LGBM performs

<sup>&</sup>lt;sup>2</sup>System: IBM SYSTEM POWER AC922 Compute Nodes with  $2 \times 16$ -core 2.7GHz POWER9 CPUs, 512GB of RAM. NVIDIA Tesla V100 32GB GPU



380

381

382

383

384

385 386

387

388

389 390 391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

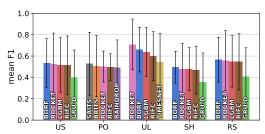
427

428

429

430

431



signals (center), and time series length (bottom).

Figure 7: Mean F1 rank (lower is better) against Figure 8: Mean F1 (higher is better) of the 5 bestdataset size in terms of instances (top), number of performing models for each type of irregularity.

well. This finding aligns with observations in (Tan et al., 2020), where gradient-boosting trees showed strong performance in regular time series regression. LGBM is a compelling choice due to its decent performance and exceptionally fast training time, making it attractive for practitioners needing solid baselines. Neural network-based methods, though designed for ITS, underperform in these bake-off-style benchmarks, except for their competitive inference runtime. Similar patterns appear in regular time series classification (Middlehurst et al., 2024b). We hypothesize that simpler, generalist, models, like ROCKET, excel in bake-off settings due to their low-variance, high-bias inductive bias, making them robust across a wide range of tasks, contrary to specialized models, which exhibit strong performance on specific types of irregularity or dataset characteristics, especially after fine-tuning.

**Performance vs. Dimension.** Figure 7 (top) shows the mean F1 ranks of all benchmarked models (lower is better), stratified by dataset size: small (at most 500 instances) and large (more than 500 instances). KNN and RIFC exhibit a noticeable worsening in rank on larger datasets, indicating limited scalability or reduced robustness as the number of training examples increases. In contrast, LGBM, and especially TIMESNET, improve significantly in rank, suggesting that more complex models, particularly transformer-based ones, benefit from greater data availability to better exploit their capacity. Figure 7 (center) shows the mean F1 ranks for univariate and multivariate time series. While the best-ranked model is again ROCKET, all neural network-based approaches benefit from increased dimensionality, making them particularly suitable for multivariate time series. Figure 7 (bottom) reports the mean F1 ranks stratified by time series length: short (at most 360 observations) and long (more than 360 observations). Here, recurrent models such as GRU-D and BRITS, along with several other neural architectures, tend to struggle on longer sequences. RAINDROP stands out as an exception, likely owing to its graph-based design. Meanwhile, models that rely on localized or interval-based features, such as ROCKET, RIFC, and especially BORF, show improved performance on longer time series, indicating that in this case, simpler is better (more details available in Appendix C).

**Performance vs. Irregularity.** In Figure 8, we report the average F1 score of the top-5 performing models within each irregularity group (higher is better). ROCKET, BORF, and LGBM consistently rank among the top three across unevenly sampled, unequal length, shifted, and ragged sampling time series. GRU-D, while generally ranking lower overall, appears among the top five models in three out of the five groups, showing solid average performance. Partially observed time series exhibit markedly different behavior: here, models designed to handle missing data, such as SAITS and BRITS, outperform ROCKET, BORF, and LGBM. This suggests that explicitly modeling missingness can be highly beneficial, particularly for datasets with structured patterns of missing values.

**Performance after Fine-tuning.** In Table 3, we present the average performance of the top three generalist models, ROCKET, BORF, and LGBM, evaluated in terms of area under the Receiver Operating Characteristic curve (auc) and area under the Precision-Recall curve (aupr) following hyperparameter tuning. These evaluations follow the same 5-fold cross-validation setup and are compared against reference results from (Li et al., 2023; Liu et al., 2024; Zheng et al., 2024) on the two most commonly used irregular medical datasets: P12 (Silva et al., 2012) and P19 (Reyna et al., 2020). This benchmark aims to assess whether *generalist* classifiers can also be effectively fine-tuned for specific tasks, and to compare them with state-of-the-art specialist deep learning models such as CONTIFORMER (Chen et al., 2024), GRU-D (Che et al., 2018), MTSFORMER (Zheng et al., 2024), MUSICNET (Liu et al., 2024), and RAINDROP (Zhang et al., 2022). Results indicate that, when optimally fine-tuned, deep learning-based algorithms outperform simpler regular time series

Table 3: Comparison of best-performing models from the bake-off, against baseline reference results (higher is better). Best values in bold, second best underlined.

		BORF	CONTI FORMER	GRU-D	LGBM	MTS FORMER	MUSIC NET	RAIN DROP	ROCKET
P12								82.8±1.7 44.0±3.0	
P19	auc aupr							$\frac{87.0}{51.8} \pm 2.3$	

classifiers. However, except for ROCKET, which underperforms in this test, this advantage is not always substantial; for instance, LGBM achieves the fourth-best score on P19, outperforming models like CONTIFORMER and GRU-D. Another advantage of models such as ROCKET, BORF, and LGBM is that the performance is very stable, with near-zero standard deviation to a single decimal place. This underscores the value of being able to readily apply standard approaches, as they can offer fast, stable, and non-trivial baselines. However, deep learning offers more flexibility for optimizing on specific tasks, with reasonable inference times when aiming for raw performance for deployment purposes.

**Performance vs. Trustworthiness.** Though not the main focus of this work, we briefly address model trustworthiness, crucial in high-stakes fields like healthcare, where *ITS* are common. The most interpretable models in our benchmark are BORF, which relies on subsequence presence/absence, and RIFC, which uses simple interval-based features, both followed by a tree-based model. Neural models can be interpreted with gradient-based methods, though the reliability of their explanations on *ITS* is unexplored. The top-performing model, ROCKET, offers little interpretability and depends on expensive model-agnostic techniques (Theissler et al., 2022). Robustness to random initialization also matters: models with high variance across seeds hinder reproducibility. Stable methods like LGBM, BORF, and KNN may be preferable in sensitive settings, even at some cost in performance.

#### 6 Conclusion

In this work, we presented pyrregular, a unified framework for addressing the challenges of *ITS*. By introducing a standardized repository for *ITS* classification and structuring the datasets in a common array format, we provided a cohesive way to work with varying forms of irregularity. Our extensive empirical evaluation of 12 state-of-the-art classifiers and baseline methods on 34 datasets emphasizes both the complexity of this domain and the benefits of a shared benchmarking resource. Results indicate that, with appropriate configuration and tuning, specialist models such as neural networks still attain state-of-the-art performance. However, extending their applicability across diverse tasks remains a significant challenge. Interestingly, simple generalist classifiers originally designed for regular time series data, such as ROCKET, perform remarkably well on irregular time series in bake-off-style benchmarks, even without leveraging the irregularity itself. This observation reveals a crucial research gap: the need to develop *generalist* methods capable of explicitly exploiting irregularities, such as missingness and timestamp information.

The construction of this extensive set of benchmarks was greatly facilitated by pyrregular, which abstracts the complexities of *ITS* across diverse libraries. While we aimed to provide a diverse and representative selection of baseline models, our choices were also guided by practical considerations such as library availability and interface compatibility, rather than exhaustive coverage. We acknowledge that several other relevant baselines could further enrich the comparison. Our goal was not to be fully comprehensive, but to establish a robust and extensible starting point for benchmarking within a unified framework. Further, we deliberately limited the scope of the benchmarks to classification, as achieving the same level of detail for other tasks, such as forecasting, anomaly detection, or imputation, would require an effort comparable in scale to what we present here, and is therefore left for future work. Nevertheless, because the proposed array format is task-independent and some curated datasets already include additional target variables, our framework naturally enables exploration of these tasks (see Appendix G for details). Going forward, pyrregular will be extended to such additional tasks, integrated with more datasets, and enriched with methods from a broader selection of time series libraries, increasing its relevance across diverse research domains.

# REFERENCES

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
- Hameer Abbasi. Sparse: A more modern sparse array library. In *SciPy*, pp. 65–68, 2018.
- Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper Schulz, Lorenzo Stella, Ali Caner Türkmen, and Yuyang Wang. GluonTS: Probabilistic and Neural Time Series Modeling in Python. *Journal of Machine Learning Research*, 21(116):1–6, 2020. URL http://jmlr.org/papers/v21/19-820.html.
- Mohammad Ali Bagheri, Qigang Gao, and Sergio Escalera. Support vector machines with time series distance kernels for action classification. In 2016 IEEE Winter Conference on Applications of Computer Vision (WACV), pp. 1–7. IEEE, 2016.
- Anthony Bagnall, Jason Lines, Aaron Bostrom, James Large, and Eamonn Keogh. The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data mining and knowledge discovery*, 31:606–660, 2017.
- Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom, Paul Southam, and Eamonn Keogh. The uea multivariate time series classification archive, 2018. *arXiv* preprint arXiv:1811.00075, 2018.
- Anthony Bagnall, Michael Flynn, James Large, Jason Lines, and Matthew Middlehurst. On the usage and performance of the hierarchical vote collective of transformation-based ensembles version 1.0 (hive-cote v1. 0). In *Advanced Analytics and Learning on Temporal Data: 5th ECML PKDD Workshop, AALTD 2020, Ghent, Belgium, September 18, 2020, Revised Selected Papers 6*, pp. 3–18. Springer, 2020.
- Ella Browning, Mark Bolton, Ellie Owen, Akiko Shoji, Tim Guilford, and Robin Freeman. Predicting animal behaviour using deep learning: Gps data alone accurately predict diving in seabirds. *Methods in Ecology and Evolution*, 9(3):681–692, 2018.
- Wei Cao, Dong Wang, Jian Li, Hao Zhou, Lei Li, and Yitan Li. Brits: Bidirectional recurrent imputation for time series. *Advances in neural information processing systems*, 31, 2018.
- Fabio Marco Caputo, Pietro Prebianca, Alessandro Carcangiu, Lucio Davide Spano, and Andrea Giachetti. Comparing 3d trajectories for simple mid-air gesture recognition. *Comput. Graph.*, 73: 17–25, 2018.
- Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent neural networks for multivariate time series with missing values. *Scientific reports*, 8(1):6085, 2018.
- Yanping Chen, Adena Why, Gustavo Batista, Agenor Mafra-Neto, and Eamonn Keogh. Flying insect classification with inexpensive sensors. *Journal of insect behavior*, 27:657–677, 2014.
- Yuqi Chen, Kan Ren, Yansen Wang, Yuchen Fang, Weiwei Sun, and Dongsheng Li. Contiformer: Continuous-time transformer for irregular time series modeling. *Advances in Neural Information Processing Systems*, 36, 2024.
- ChoroChronos Archive. Trucks dataset dataset and algorithms | chorochronos.org. http://www.chorochronos.org/. Accessed: 2025-01-23.
- Martina Cinquini, Fosca Giannotti, Riccardo Guidotti, and Andrea Mattei. Handling missing values in local post-hoc explainability. In *World Conference on Explainable Artificial Intelligence*, pp. 256–278. Springer, 2023.

- City of Melbourne. Pedestrian counting system. http://www.pedestrian.melbourne.vic.gov.au, 2020. Accessed: 2025-01-23.
- Camila Leite da Silva, Lucas May Petry, and Vania Bogorny. A survey and comparison of trajectory classification methods. In *2019 8th Brazilian conference on intelligent systems (BRACIS)*, pp. 788–793. IEEE, 2019.
  - Hoang Anh Dau, Anthony Bagnall, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, and Eamonn Keogh. The ucr time series archive. *IEEE/CAA Journal of Automatica Sinica*, 6(6):1293–1305, 2019.
  - Vinícius M. A. de Souza. Asphalt pavement classification using smartphone accelerometer and complexity invariant distance. *Eng. Appl. Artif. Intell.*, 74:198–211, 2018.
    - Angus Dempster, François Petitjean, and Geoffrey I Webb. Rocket: exceptionally fast and accurate time series classification using random convolutional kernels. *Data Mining and Knowledge Discovery*, 34(5):1454–1495, 2020.
    - Angus Dempster, Daniel F Schmidt, and Geoffrey I Webb. Minirocket: A very fast (almost) deterministic transform for time series classification. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 248–257, 2021.
    - Wenjie Du. Pypots: A python toolbox for data mining on partially-observed time series. *arXiv* preprint arXiv:2305.18811, 2023.
    - Wenjie Du, David Côté, and Yan Liu. Saits: Self-attention-based imputation for time series. *Expert Systems with Applications*, 219:119619, 2023.
    - Iain S Duff, Albert Maurice Erisman, and John Ker Reid. *Direct methods for sparse matrices*. Oxford University Press, 2017.
    - Carlos Andres Ferrero, Luis Otavio Alvares, Willian Zalewski, and Vania Bogorny. Movelets: Exploring relevant subtrajectories for robust trajectory classification. In *Proceedings of the 33rd Annual ACM symposium on applied computing*, pp. 849–856, 2018.
    - Jingkun Gao, Suman Giri, Emre Can Kara, and Mario Berges. PLAID: a public dataset of high-resoultion electrical appliance measurements for load identification research: demo abstract. In *BuildSys*, pp. 198–199. ACM, 2014.
    - Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.
    - Alessio Gravina, Daniele Zambon, Davide Bacciu, and Cesare Alippi. Temporal graph odes for irregularly-sampled time series. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pp. 4025–4034, 2024.
  - Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *The Tenth International Conference on Learning Representations, ICLR* 2022, *Virtual Event, April* 25-29, 2022. OpenReview.net, 2022. URL https://openreview.net/forum?id=uYLFoz1vlAC.
    - Joze Guna, Grega Jakus, Matevz Pogacnik, Saso Tomazic, and Jaka Sodnik. An analysis of the precision and reliability of the leap motion sensor and its suitability for static and dynamic tracking. *Sensors*, 14(2):3702–3720, 2014.
  - James D Hamilton. *Time series analysis*. Princeton university press, 2020.
  - Nacereddine Hammami and Mouldi Bedda. Improved tree model for arabic speech recognition. In 2010 3rd international conference on computer science and information technology, volume 5, pp. 521–526. IEEE, 2010.

- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2. URL https://doi.org/10.1038/s41586-020-2649-2.
  - Hrayr Harutyunyan, Hrant Khachatrian, David C. Kale, Greg Ver Steeg, and Aram Galstyan. Multitask learning and benchmarking with clinical time series data. *Scientific Data*, 6(1):96, 2019. ISSN 2052-4463. doi: 10.1038/s41597-019-0103-9. URL https://doi.org/10.1038/s41597-019-0103-9.
  - Andrew Harvey, Siem Jan Koopman, and Jeremy Penzer. Messy time series: a unified approach. *Advances in econometrics*, 13:103–144, 1998.
  - Stephan Hoyer and Joe Hamman. xarray: Nd labeled arrays and datasets in python. *Journal of Open Research Software*, 5(1):10–10, 2017.
  - Alexander Ihler, Jon Hutchins, and Padhraic Smyth. Adaptive event detection with time-varying poisson processes. In *KDD*, pp. 207–216. ACM, 2006.
  - Abdullah Al Imran, Md Shamsur Rahim, and Tanvir Ahmed. Mining the productivity data of the garment industry. *Int. J. Bus. Intell. Data Min.*, 19(3):319–342, 2021.
  - Ali Ismail-Fawaz, Angus Dempster, Chang Wei Tan, Matthieu Herrmann, Lynn Miller, Daniel F Schmidt, Stefano Berretti, Jonathan Weber, Maxime Devanne, Germain Forestier, et al. An approach to multiple comparison benchmark evaluations that is stable under manipulation of the comparate set. *arXiv* preprint arXiv:2305.11921, 2023.
  - Nathalie Japkowicz. Assessment metrics for imbalanced learning. *Imbalanced learning: Foundations, algorithms, and applications*, pp. 187–206, 2013.
  - Alistair Johnson, Tom Pollard, and Roger Mark. Mimic-iii clinical database demo (version 1.4). *PhysioNet*, 10:C2HM2Q, 2019.
  - Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
  - Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30, 2017.
  - Patrick Kidger. On Neural Differential Equations. PhD thesis, University of Oxford, 2021.
  - Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural controlled differential equations for irregular time series. *Advances in Neural Information Processing Systems*, 33:6696–6707, 2020.
  - Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3): 455–500, 2009.
  - Mineichi Kudo, Jun Toyama, and Masaru Shimbo. Multidimensional curve classification using passing-through regions. *Pattern Recognit. Lett.*, 20(11-13):1103–1111, 1999.
  - Cristiano Landi, Riccardo Guidotti, Mirco Nanni, and Anna Monreale. The trajectory interval forest classifier for trajectory classification. In *SIGSPATIAL/GIS*, pp. 67:1–67:4. ACM, 2023a.
  - Cristiano Landi, Francesco Spinnato, Riccardo Guidotti, Anna Monreale, and Mirco Nanni. Geolet: An interpretable model for trajectory classification. In *IDA*, volume 13876 of *Lecture Notes in Computer Science*, pp. 236–248. Springer, 2023b.

- Steven Cheng-Xian Li and Benjamin Marlin. Learning from irregularly-sampled time series: A missing data perspective. In *International Conference on Machine Learning*, pp. 5937–5946. PMLR, 2020.
- Zekun Li, Shiyang Li, and Xifeng Yan. Time series as images: Vision transformer for irregularly sampled time series. *Advances in Neural Information Processing Systems*, 36:49187–49204, 2023.
- Jiexi Liu, Meng Cao, and Songcan Chen. Musicnet: A gradual coarse-to-fine framework for irregularly sampled multivariate time series analysis. *arXiv* preprint arXiv:2412.01063, 2024.
- Jiexi Liu, Meng Cao, and Songcan Chen. Timecheat: A channel harmony strategy for irregularly sampled multivariate time series analysis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 18861–18869, 2025.
- Markus Löning, Anthony Bagnall, Sajaysurya Ganesh, Viktor Kazakov, Jason Lines, and Franz J Király. sktime: A unified interface for machine learning with time series. arXiv preprint arXiv:1909.07872, 2019.
- Yicheng Luo, Zhen Liu, Linghao Wang, Binquan Wu, Junhao Zheng, and Qianli Ma. Knowledge-empowered dynamic graph network for irregularly sampled medical time series. *Advances in Neural Information Processing Systems*, 37:67172–67199, 2024.
- Yonghong Luo, Xiangrui Cai, Ying Zhang, Jun Xu, et al. Multivariate time series imputation with generative adversarial networks. *Advances in neural information processing systems*, 31, 2018.
- Giangiacomo Mercatali, Andre Freitas, and Jie Chen. Graph neural flows for unveiling systemic interactions among irregularly sampled time series. *Advances in Neural Information Processing Systems*, 37:57183–57206, 2024.
- Antigoni Mezari and Ilias Maglogiannis. An easily customized gesture recognizer for assisted living using commodity mobile devices. *Journal of Healthcare Engineering*, 2018(1):3180652, 2018.
- Matthew Middlehurst, Ali Ismail-Fawaz, Antoine Guillaume, Christopher Holder, David Guijo-Rubio, Guzal Bulatova, Leonidas Tsaprounis, Lukasz Mentel, Martin Walter, Patrick Schäfer, et al. aeon: a python toolkit for learning from time series. *Journal of Machine Learning Research*, 25(289): 1–10, 2024a.
- Matthew Middlehurst, Patrick Schäfer, and Anthony Bagnall. Bake off redux: a review and experimental evaluation of recent time series classification algorithms. *Data Mining and Knowledge Discovery*, pp. 1–74, 2024b.
- Robin Mitra, Sarah F McGough, Tapabrata Chakraborti, Chris Holmes, Ryan Copping, Niels Hagenbuch, Stefanie Biedermann, Jack Noonan, Brieuc Lehmann, Aditi Shenvi, et al. Learning from data with structured missingness. *Nature Machine Intelligence*, 5(1):13–23, 2023.
- Luis Moreira-Matias, Michel Ferreira, Joao Mendes-Moreira, L. L., and J. J. Taxi Service Trajectory Prediction Challenge, ECML PKDD 2015. UCI Machine Learning Repository, 2013. DOI: https://doi.org/10.24432/C55W25.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Jim Pivarski, Peter Elmer, and David Lange. Awkward arrays in python, c++, and numba. In *EPJ Web of Conferences*, volume 245, pp. 05023. EDP Sciences, 2020.
- Attila Reiss and Didier Stricker. Introducing a new benchmarked dataset for activity monitoring. In 2012 16th international symposium on wearable computers, pp. 108–109. IEEE, 2012.
- Matthew A Reyna, Christopher S Josef, Russell Jeter, Supreeth P Shashikumar, M Brandon Westover, Shamim Nemati, Gari D Clifford, and Ashish Sharma. Early prediction of sepsis from clinical data: the physionet/computing in cardiology challenge 2019. *Critical care medicine*, 48(2):210–217, 2020.

- Yulia Rubanova, Ricky TQ Chen, and David K Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. *Advances in neural information processing systems*, 32, 2019.
  - Donald B Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
    - Alejandro Pasos Ruiz, Michael Flynn, James Large, Matthew Middlehurst, and Anthony Bagnall. The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, 35(2):401–449, 2021.
    - Naoki Saito. Local feature extraction and its applications using a library of bases. Yale University, 1994.
    - Hiroaki Sakoe and Seibi Chiba. Dynamic programming algorithm optimization for spoken word recognition. *IEEE transactions on acoustics, speech, and signal processing*, 26(1):43–49, 1978.
    - Md Abu Hanif Shaikh and KM Azharul Hasan. Efficient storage scheme for n-dimensional sparse array: Gcrs/gccs. In 2015 International Conference on High Performance Computing & Simulation (HPCS), pp. 137–142. IEEE, 2015.
    - Satya Narayan Shukla and Benjamin M. Marlin. Multi-time attention networks for irregularly sampled time series. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?id=4c0J6lwQ4\_.
    - Robert H Shumway, David S Stoffer, and David S Stoffer. *Time series analysis and its applications*, volume 3. Springer, 2000.
    - Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. Predicting in-hospital mortality of icu patients: The physionet/computing in cardiology challenge 2012. In *2012 computing in cardiology*, pp. 245–248. IEEE, 2012.
    - Francesco Spinnato, Riccardo Guidotti, Anna Monreale, and Mirco Nanni. Fast, interpretable and deterministic time series classification with a bag-of-receptive-fields. *IEEE Access*, 2024.
    - Chang Wei Tan, Christoph Bergmeir, Francois Petitjean, and Geoffrey I Webb. Monash university, uea, ucr time series extrinsic regression archive. *arXiv* preprint arXiv:2006.10996, 2020.
    - Romain Tavenard, Johann Faouzi, Gilles Vandewiele, Felix Divo, Guillaume Androz, Chester Holtz, Marie Payne, Roman Yurchak, Marc Rußwurm, Kushal Kolar, and Eli Woods. Tslearn, a machine learning toolkit for time series data. *Journal of Machine Learning Research*, 21(118):1–6, 2020. URL http://jmlr.org/papers/v21/20-091.html.
    - Andreas Theissler, Francesco Spinnato, Udo Schlegel, and Riccardo Guidotti. Explainable ai for time series classification: a review, taxonomy and research directions. *Ieee Access*, 10:100700–100724, 2022.
    - V Vidulin, M Lustrek, B Kaluza, R Piltaver, and J Krivec. Localization data for person activity. *UCI Machine Learning Repository*, 2010.
    - Jun Wang, Wenjie Du, Wei Cao, Keli Zhang, Wenjia Wang, Yuxuan Liang, and Qingsong Wen. Deep learning for multivariate time series imputation: A survey. *arXiv preprint arXiv:2402.04059*, 2024.
    - Philip B Weerakody, Kok Wai Wong, Guanjin Wang, and Wendell Ela. A review of irregular time series data handling with gated recurrent neural networks. *Neurocomputing*, 441:161–178, 2021.
    - Ben H. Williams, Marc Toussaint, and Amos J. Storkey. Extracting motion primitives from natural handwriting data. In *ICANN* (2), volume 4132 of *Lecture Notes in Computer Science*, pp. 634–643. Springer, 2006.
    - Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. *arXiv preprint arXiv:2210.02186*, 2022.

- Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL https://openreview.net/forum?id=ju\_Uqw3840q.
  - Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the (in) fidelity and sensitivity of explanations. *Advances in neural information processing systems*, 32, 2019.
  - Weijia Zhang, Chenlong Yin, Hao Liu, Xiaofang Zhou, and Hui Xiong. Irregular multivariate time series forecasting: A transformable patching graph neural networks approach. In *Forty-first International Conference on Machine Learning*, 2024.
  - Xiang Zhang, Marko Zeman, Theodoros Tsiligkaridis, and Marinka Zitnik. Graph-guided network for irregularly sampled multivariate time series. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. URL https://openreview.net/forum?id=Kwm8I7dU-15.
  - Liangwei Nathan Zheng, Zhengyang Li, Chang George Dong, Wei Emma Zhang, Lin Yue, Miao Xu, Olaf Maennel, and Weitong Chen. Irregularity-informed time series analysis: Adaptive modelling of spatial and temporal dynamics. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pp. 3405–3414, 2024.
  - Yu Zheng, Quannan Li, Yukun Chen, Xing Xie, and Wei-Ying Ma. Understanding mobility based on GPS data. In *UbiComp*, volume 344 of *ACM International Conference Proceeding Series*, pp. 312–321. ACM, 2008.
  - Yu Zheng, Lizhu Zhang, Xing Xie, and Wei-Ying Ma. Mining interesting locations and travel sequences from gps trajectories. In *Proceedings of the 18th international conference on World wide web*, pp. 791–800, 2009.
  - Yu Zheng, Xing Xie, and Wei-Ying Ma. Geolife: A collaborative social networking service among user, location and trajectory. *IEEE Data Eng. Bull.*, 33(2):32–39, 2010.

# A SUMMARY OF NOTATION

We have adopted a tensor-like notation inspired by (Kolda & Bader, 2009). The time series dataset is structured along three dimensions: the instance dimension, which consists of n instances (e.g.,  $X_i$  denotes the i-th time series in the dataset  $\mathbf{X}$ ); the signal dimension, which includes d channels (e.g.,  $\mathbf{x}_{i,j}$  represents the j-th signal in time series  $\mathbf{X}_i$ ); and the time dimension, spanning  $\mathcal{T}$  points (e.g.,  $x_{i,j,t_k}$  represents the  $t_k$  observation of j-th signal in time series  $\mathbf{X}_i$ ). We use tildes to specify the index being referenced (e.g.,  $t_k \in \mathbf{t}$  corresponds to the k-th timestamp at the dataset's level, while  $t_k \in \tilde{\mathbf{t}}$  corresponds to the k-th timestamp at the time series's level). For improved readability, indices are omitted when they are not relevant.

Table 4: Summary of notation.

Notation	
$\mathbf{X}, \mathbf{X}, \mathbf{x}, x$	time series dataset, instance, signal, entry
$\mathbf{t}, \mathbf{ ilde{t}}, \mathbf{ ilde{t}}, t$	timestamps for a time series dataset, instance, signal, entry
k	timestamp index
n	number of instances in a dataset
d	number of signals in a time series
$\mathcal{T}, T,  au$	number timestamps in a time series dataset, instance, signal
i,j,k	indexes for instances, signals, timestamps

### B TAXONOMY OF TIME SERIES IRREGULARITIES

In addition to the well-known missingness taxonomy introduced in (Rubin, 1976) (MCAR, MAR, and MNAR), Mitra et al. (2023) proposed an additional category: structural missingness (SM). While Rubin's framework is typically formulated in terms of univariate patterns, SM highlights situations where missingness is organized across multiple variables and exhibits systematic structure. Our primary aim, distinct from previous works, is to preserve such structural patterns of missingness.

Consider, for instance, daily heart rate signals collected by wearables over three months. Data may be missing completely at random (MCAR) when some days are absent because the device randomly fails to sync, in which case missingness is unrelated to any variable. It may be missing at random (MAR) when data are more frequently absent on weekends, particularly for users with low recorded activity. It may be missing not at random (MNAR) when users remove the device precisely when feeling unwell, so missingness coincides with unrecorded spikes in heart rate. Finally, it may exhibit structural missingness (SM) when devices differ in recording frequency, such as once per second versus once per millisecond, or when a firmware update produces week-long gaps.

In this last case, missingness follows clear temporal patterns tied to device characteristics or design flaws, rather than to a single variable. Addressing such missingness (or raggedness) should therefore be an intentional modeling choice by the practitioner, not the result of routine preprocessing. We provide here formal definitions for each type of time series irregularity and use minimal counterexamples to show that none of these irregularities implies the others.

**Definition B.1** (Uneven Sampling). A signal  $\mathbf{x} = [x_{t_1}, \dots, x_{t_{\tau}}] \in \mathbb{R}^{\tau}$  is said to be *unevenly sampled* if there exists at least one index  $k \in \{1, \dots, \tau-1\}$  such that the time interval between successive observations is not constant, i.e.,  $t_{k+1} - t_k \neq \Delta t$  for some fixed  $\Delta t \in \mathbb{R}$ .

The same definition applies to time series instances and datasets, using their respective indices  $\dot{t}$ ,  $\dot{t}$ .

**Definition B.2** (Partial Observation). A signal  $\mathbf{x} = [x_{t_1}, \dots, x_{t_{\tau}}] \in \mathbb{R}^{\tau}$  is said to be *partially observed* if at least one value  $x_{t_k}$  is missing and represented by a special symbol NaN, indicating the absence of an observation at a timestamp where one was expected, i.e.,  $x_{t_k} = \text{NaN}$  for some  $k \in \{1, \dots, \tau\}$ .

Again, the same definition applies to time series instances and datasets.

**Definition B.3** (Raggedness). Raggedness is a structural irregularity that arises in a multivariate time series  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_d\} \in \mathbb{R}^{d \times T}$  when the component signals do not share a common timestamp

index. Formally, raggedness is present when there exist at least two signals  $\mathbf{x}_a$  and  $\mathbf{x}_b$  such that  $\tilde{\mathbf{t}}_a \neq \tilde{\mathbf{t}}_b$ . It manifests in three independent forms:

- (a) Ragged Length:  $\tau_a \neq \tau_b$ .
- (b) Shift:  $(t_{a,1} < t_{b,1}) \wedge (t_{a,\tau_a} < t_{b,\tau_b})$ .
- (c) Ragged Sampling:  $\Delta t_{a,k} \neq \Delta t_{b,k}$  for some k, where  $\Delta t_{j,k} = t_{j,k+1} t_{j,k}$ . The index k ranges from 1 to  $\min(\tau_a, \tau_b) 1$ , so only intervals that exist in both signals are compared.

The same definition applies to time series datasets.

We now show that the five forms of time series irregularity are mutually independent: none implies any of the others. This is shown through minimal examples of time series that satisfy one irregularity condition while exhibiting none of the others.

#### B.1 UNEVEN SAMPLING

Let  $X = \{\mathbf{x}_a, \mathbf{x}_b\}$  be a time series where both signals share the same timestamp index,  $\tilde{\mathbf{t}} = \tilde{\mathbf{t}}_a = \tilde{\mathbf{t}}_b = [t_1, t_2, t_3]$ , and assume that the sampling intervals are not constant, i.e.,  $t_2 - t_1 \neq t_3 - t_2$ . Then X is unevenly sampled.

UNEVEN SAMPLING  $\Rightarrow$  PARTIAL OBSERVATION. Suppose that all values in both  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are observed (i.e., none are NaN). Then X is unevenly sampled, but not partially observed.

UNEVEN SAMPLING  $\Rightarrow$  RAGGEDNESS. Since  $\tilde{\mathbf{t}}_a = \tilde{\mathbf{t}}_b$ , both signals are aligned on the same timestamps. Therefore, X is not ragged.

#### **B.2** Partial Observation

Let  $X = \{\mathbf{x}_a, \mathbf{x}_b\}$  be a time series where both signals share the same timestamp index,  $\tilde{\mathbf{t}} = \tilde{\mathbf{t}}_a = \tilde{\mathbf{t}}_b = [t_1, t_2, t_3]$ . Suppose that one observation is missing, e.g.,  $x_{a,t_2} = NaN$ . Then X is partially observed.

PARTIAL OBSERVATION  $\Rightarrow$  UNEVEN SAMPLING. Let the timestamps be equally spaced, i.e.,  $t_2 - t_1 = t_3 - t_2 = \Delta t$ . Then  $\boldsymbol{X}$  is partially observed but evenly sampled.

PARTIAL OBSERVATION  $\Rightarrow$  RAGGEDNESS. Since both signals are defined over the same set of timestamps,  $\tilde{\mathbf{t}}_a = \tilde{\mathbf{t}}_b$ ,  $\boldsymbol{X}$  is not ragged.

#### B.3 RAGGED LENGTH

Let  $X = \{\mathbf{x}_a, \mathbf{x}_b\}$  be a time series exhibiting ragged length, with  $\tilde{\mathbf{t}}_a = [t_1, t_2]$  and  $\tilde{\mathbf{t}}_b = [t_1, t_2, t_3]$ . Then the unified timestamp index is  $\tilde{\mathbf{t}} = [t_1, t_2, t_3]$ , and X satisfies  $\tau_a = 2 \neq 3 = \tau_b$ .

RAGGED LENGTH  $\Rightarrow$  UNEVEN SAMPLING. Let the timestamps be evenly spaced, i.e.,  $t_2 - t_1 = t_3 - t_2 = \Delta t$ . Then X exhibits ragged length, but is evenly sampled.

RAGGED LENGTH  $\Rightarrow$  PARTIAL OBSERVATION. Suppose that all values in both  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are observed (i.e., no *NaNs*). Then X exhibits ragged length, but is not partially observed.

RAGGED LENGTH  $\Rightarrow$  SHIFT. Although the signals have different lengths, they both start at the same time,  $t_1$ . Hence, X is not shifted.

RAGGED LENGTH  $\Rightarrow$  RAGGED SAMPLING. The sampling intervals are identical across both signals, i.e.,  $\Delta \tilde{t}_{a,1} = \Delta \tilde{t}_{b,1} = t_2 - t_1$ . Therefore,  $\boldsymbol{X}$  is not raggedly sampled.

#### B.4 SHIFT

Let  $X = \{\mathbf{x}_a, \mathbf{x}_b\}$  be a time series exhibiting shift, with  $\tilde{\mathbf{t}}_a = [t_1, t_2]$  and  $\tilde{\mathbf{t}}_b = [t_2, t_3]$ . Then the unified timestamp index is  $\tilde{\mathbf{t}} = [t_1, t_2, t_3]$ , and X is shifted, as  $\mathbf{x}_a$  starts and ends before  $\mathbf{x}_b$ .

SHIFT  $\Rightarrow$  UNEVEN SAMPLING. Let the timestamps be evenly spaced, i.e.,  $t_2 - t_1 = t_3 - t_2 = \Delta t$ . Then  $\boldsymbol{X}$  exhibits shift, but is evenly sampled.

SHIFT  $\Rightarrow$  PARTIAL OBSERVATION. Suppose that all values in both  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are observed (i.e., no *NaNs*). Then X exhibits shift, but is not partially observed.

SHIFT  $\Rightarrow$  RAGGED LENGTH. Both signals have the same number of observations, i.e.,  $\tau_a = \tau_b = 2$ . Hence, X exhibits shift but not ragged length.

SHIFT  $\Rightarrow$  RAGGED SAMPLING. The sampling intervals within each signal are equal, i.e.,  $\Delta \tilde{t}_{a,1} = t_2 - t_1 = \Delta \tilde{t}_{b,1} = t_3 - t_2$ . Therefore, X is not raggedly sampled.

### B.5 RAGGED SAMPLING

 Let  $X = \{\mathbf{x}_a, \mathbf{x}_b\}$  be a time series exhibiting ragged sampling, with  $\tilde{\mathbf{t}}_a = [t_1, t_2]$  and  $\tilde{\mathbf{t}}_b = [t_1, t_3]$ . Then the unified timestamp index is  $\tilde{\mathbf{t}} = [t_1, t_2, t_3]$ , and the sampling intervals differ across signals:  $\Delta \tilde{t}_{a,1} = t_2 - t_1 \neq t_3 - t_1 = \Delta \tilde{t}_{b,1}$ .

RAGGED SAMPLING  $\Rightarrow$  UNEVEN SAMPLING. Let the global timestamps satisfy  $t_2 - t_1 = t_3 - t_2 = \Delta t$ . Then  $\boldsymbol{X}$  is raggedly sampled but not unevenly sampled.

RAGGED SAMPLING  $\Rightarrow$  PARTIAL OBSERVATION. Suppose that all values in both  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are observed (i.e., no NaNs). Then X exhibits ragged sampling, but is not partially observed.

RAGGED SAMPLING  $\Rightarrow$  RAGGED LENGTH. Both signals contain the same number of observations,  $\tau_a = \tau_b = 2$ . Thus, X is not ragged in length.

RAGGED SAMPLING  $\Rightarrow$  SHIFT. Both signals start at the same time,  $t_1$ , and have the same length. Therefore, X is not shifted.

These examples are minimal and can be easily extended to longer signals and time series. They suffice to establish that all forms of irregularity discussed, both in the main and raggedness subtypes, are pairwise independent. None of them implies any other, as illustrated also in Figure 2. To the best of our knowledge, this taxonomy accounts for all known forms of structural time series irregularity relevant to data modeling and representation.

# C EXPERIMENTAL DETAILS.

In this section, we summarize experimental details regarding the models and datasets.

#### C.1 Models

The objective of these experiments is to benchmark methods capable of naturally handling irregular time series without introducing bias through imputation techniques. To achieve this, we limit our evaluation to classifiers that inherently support missing data in their input and are available in major time series libraries. Below, we describe the implementation details and hyperparameters for each method. Parameters that are not mentioned are left to their default in their library implementation.

**Bag-of-Receptive-Fields** (BORF) The Bag of Receptive Fields (BORF) algorithm (Spinnato et al., 2024) from the aeon library extracts discretized subsequences and counts their appearance in the time series, allowing the presence of missing data. A downstream LightGBM classifier with default parameters is used to handle transformed features. For the fine-tuned benchmark, the hyperparameter was on performed on the *min\_window\_to\_signal\_std\_ratio* in the interval [0, 0.2] with 0.05 increments.

**Bidirectional Recurrent Imputation for Time Series** (BRITS) The BR*ITS* algorithm (Cao et al., 2018), also from the pypots library, employs a bidirectional recurrent network for imputing and classifying incomplete time series. It uses a hidden layer size of 256 and a batch size of 32. Training runs for up to 1000 epochs, with early stopping after 50 epochs of no improvement.

**Gated Recurrent Unit with Decay** (GRU-D) The GRU-D model (Che et al., 2018), available in the pypots library, extends the Gated Recurrent Unit architecture by introducing decay mechanisms that account for missing data patterns. The recurrent hidden layer size is set to 256, with a batch size of 32. Training uses a maximum of 1000 epochs, with early stopping triggered after 50 epochs of no improvement.

**K-Nearest Neighbors with DTW** (KNN) This baseline model employs the tslearn K-Nearest Neighbors algorithm, configured to use the Dynamic Time Warping (DTW) distance measure. DTW incorporates temporal alignment to handle time series of varying lengths effectively. The distance computation uses a Sakoe-Chiba band (Sakoe & Chiba, 1978) of 10 points, which limits the warping window to a fixed radius.

**LightGBM** (LGBM) LightGBM (Ke et al., 2017) is a gradient-boosting framework optimized for speed and efficiency, and can naturally handle missing values. In this baseline, it is trained directly with default parameters on raw time series data transformed into a tabular format using the sktime Tabularizer. For the fine-tuned benchmark, hyperparameter optimization was conducted over a predefined search space that included the number of leaves  $(num\_leaves) \in \{31, 63, 127\}$ , maximum tree depth  $(max\_depth) \in \{-1, 7, 10\}$ ,  $(learning\_rate) \in \{0.05, 0.1\}$ , and the minimum number of samples per leaf  $(min\_data\_in\_leaf) \in \{20, 100\}$ .

Neural Controlled Differential Equation (NCDE) The Neural CDE model (Kidger et al., 2020), implemented via the diffrax library, learns continuous-time representations of time series data using differential equations. It employs an Euler solver with a maximum of 100 steps, with step size equal to the minimum time difference between any two adjacent observations, a hidden layer size of 8, and a width size of 32. Training uses a maximum of 1000 iterations, using Adam as optimizer, with a starting learning rate of 0.01, patience of 200 for early stopping, and a learning rate reduction factor of 0.5 after 50 stagnant iterations.

**Raindrop** (RAINDROP) The Raindrop model (Zhang et al., 2022), a graph-based neural network from pypots, handles irregular time series by sending messages over graphs that are optimized for capturing time-varying dependencies among sensors. This model uses 2 layers, a feed-forward network size of 256, 2 attention heads, and a dropout rate of 0.3. Training employs a batch size of 32, with early stopping after 50 epochs of no improvement.

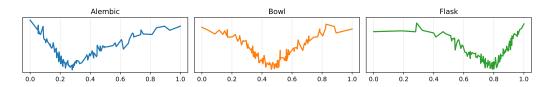


Figure 9: Three examples of instances from the (ABF) dataset, from left to right, Alembic, Bowl, and Flask.

Random Interval Feature Classifier (RIFC) The Random Interval Feature Classifier (RIFC) leverages the RandomIntervalFeatureExtractor from the sktime library to generate simple statistical summaries (mean, standard deviation, minimum, maximum, median, skewness, and kurtosis) from randomly selected intervals within the time series, with the number of intervals being the logarithm of the time series length. These features are subsequently used by a downstream LightGBM classifier to perform classification.

Minimally Random Convolutional Kernel Transform (ROCKET) Rocket, in its Minirocket implementation (Dempster et al., 2021) from the sktime library, employs 10000 fixed convolutional kernels to extract features from time series data. This implementation includes MiniRocketMultivariateVariable, which handles multivariate time series while tolerating missing data. The transformation could include missing data; therefore, instead of the most common ridge classifier, LightGBM with default parameters is used. For the fine-tuned benchmark, hyperparameter optimization was conducted over the number of kernels,  $num\_kernels \in \{100, 500, 1000, 5000, 10000, 50000\}$ .

Self-Attention Imputation for Time Series (SAITS) The SAITS model (Du et al., 2023), implemented in the pypots library, employs a transformer-based architecture specifically tailored for time series imputation. It utilizes a dual self-attention mechanism across temporal dimensions, enabling it to capture both global and local patterns despite missing values. In this configuration, SAITS is trained with 2 attention layers, a model dimension of 256, 4 attention heads, and hidden dimensions  $d_k = 64$ ,  $d_v = 64$ , and  $d_{\rm ffn} = 128$ . A dropout rate of 0.1 is used for both the transformer blocks and attention layers. The model is optimized over a maximum of 1000 epochs, with early stopping triggered after 50 stagnant epochs. Training is performed with a batch size of 32.

**Support Vector Machine with LCSS Kernel** (SVM) This method uses the sktime implementation of a Support Vector Machine, enhanced with the Longest Common Subsequence (LCSS) distance kernel (Bagheri et al., 2016). LCSS is robust to missing values and temporal distortions, as it matches time series subsequences with allowable gaps. The kernel uses a Sakoe-Chiba constraint with a radius of 10. Each time series is standardized using z-score normalization. The model is trained for a maximum of 1000 iterations.

**TimesNet** (TIMESNET) TimesNet (Wu et al., 2023) is a modern transformer-based architecture designed for multivariate time series modeling, emphasizing temporal receptive fields through learnable convolutional kernels. Its implementation here leverages 2 layers and 3 convolutional kernels with dynamic top-k temporal selection. The model dimension is set to 64, with a feed-forward network size of 128. Training is conducted using a batch size of 32 over 1000 epochs, with early stopping after 50 epochs without validation improvement.

# C.2 DATASETS

The repository includes 34 datasets, each briefly described below, along with the data preparation steps applied. <sup>3</sup> For datasets without a predefined train-test split, we created a stratified, instance-based 70-30% train-test split.

<sup>&</sup>lt;sup>3</sup>Data is hosted at **link redacted for double-blind review**.

**Alembics Bowls Flasks.** (ABF) This dataset is inspired by the classical Cylinder-Bell-Funnel (CBF) benchmark (Saito, 1994) for regular time series classification. Similarly to CBF, there are three classes, which are Alembics, Bowls, and Flasks. The classes differ by how much the temporal axis is skewed, i.e., if it has positive (Alembic), negative (Flask), or no skewness (Bowl). For each time series, 128 values are sampled from a circumference and then standardized. There are 10 instances for each class in the training set and 300 for each in the test set. An example is presented in Figure 9.

Animals (AN) This dataset, generated during the Starkey project (Ferrero et al., 2018), consists of trajectories from three animal species—elk, deer, and cattle. The classification task commonly used in the literature (Ferrero et al., 2018; Landi et al., 2023b;a) involves inferring the species based on movement patterns. The target classes in the dataset are balanced, with 38 trajectories for the elk, 30 for the deer, and 34 for the cattle.

Geolife (GS) This dataset was collected during the GeoLife Project (Microsoft Research Asia) from April 2007 to August 2012 (Zheng et al., 2009; 2008; 2010). It contains the trajectories of 182 users and has been preprocessed as detailed in the public *User Guide-1.3*. One of the most common supervised machine-learning tasks using this dataset is to identify (a subset) of the 11 means of transportation. We defined three target variables with a decreasing number of classes. The first target variable includes all the means of transportation in the dataset: airplane, bike, boat, bus, car, motorcycle, run, subway, taxi, train, and walk. The second target variable, used in (Ferrero et al., 2018), groups the transportation modes into six classes: bike, bus/taxi, car, subway, train, and walk. The third target variable, used in (Landi et al., 2023b), simplifies the classification into two categories: private (bike, boat, car, motorcycle, run, walk) and public (the remaining modes of transportation). In Section 5, we benchmark the models against the first target variable. In this setting, each class accounts for approximately 9.1% of the total instances, but the standard deviation is 12.7%, i.e., the target variable is highly imbalanced.

**GPS Data of Seabirds** (SE) This dataset, introduced in (Browning et al., 2018), consists of GPS data collected from 108 seabirds spanning three species: European shag (15), common guillemot (31), and razorbill (62). Similar to the *Animals* dataset, the species has been used to evaluate model performance in inferring species. The target variable is imbalanced, with the majority class (razorbill) comprising 62 individuals, while the minority class (European shag) includes only 15.

**Localization Data for Person Activity** (LPA) Introduced in (Vidulin et al., 2010), this dataset contains data from 5 individuals performing 11 different actions: falling, lying, lying down, on all fours, sitting, sitting down, sitting on the ground, standing up from lying, standing up from sitting, standing up from sitting on the ground, walking. Each action was recorded by tracking the positions of the body's right and left ankles, chest, and belt in a 3-dimensional space, resulting in 12 distinct signals per time series.

MIMIC-III Clinical Database Demo (MI3) Introduced by (Johnson et al., 2016; 2019) on the Physionet platform (Goldberger et al., 2000), the dataset contains health-related data associated with 40,000 patients in critical care at the Beth Israel Deaconess Medical Center from 2001 to 2012. Since the full version is available to credentialed users under strict requirements, we use the publicly available demo version in our work. We preprocess the data in accordance with (Harutyunyan et al., 2019). The classification target involves predicting in-hospital mortality.

PAMAP2 Physical Activity Monitoring (PA2) This dataset, introduced in (Reiss & Stricker, 2012), contains data from 9 subjects (1 female, 8 male) performing 19 different physical activities: ascending stairs, car driving, computer work, cycling, descending stairs, folding laundry, house cleaning, ironing, lying, nordic walking, playing soccer, rope jumping, running, sitting, standing, transient, vacuum cleaning, walking, watching TV. The data includes measurements from 3 inertial measurement units (IMUs) positioned on the dominant arm, chest, and dominant side's ankle. Specifically, from each IMU sensor, the dataset contains information about the temperature, the 3-dimensional acceleration, gyroscope and magnetometer data, and the sensor orientation. Additionally, heart rate observations are included. The two types of sensors record data at different sampling rates: 100 Hz for the IMUs and 9 Hz for the heart rate monitor. We preprocess the data according to the authors' guidelines when downloading the dataset. Data from the "transient" activity, i.e., movements between the end of one

activity and the start of another, was excluded. The remaining 18 activities serve as classification target classes.

**PhysioNet 2012** (P12) Published as data for the "Predicting Mortality of ICU Patients: The PhysioNet/Computing in Cardiology" challenge in 2012 (Silva et al., 2012), the data contains information about the patient, like age, gender, height, and weight, and 37 different types of time series. Similar to the MIMIC-III dataset, the classification target is about predicting in-hospital death.

**PhysioNet 2019** (P19) Published as data for the "Early Prediction of Sepsis from Clinical Data: The PhysioNet/Computing in Cardiology" challenge in 2019 (Reyna et al., 2020), the dataset contains demographic information about the patients, such as age, gender, height, and weight, alongside 34 other time-series variables for vital signs and laboratory test values. The classification task involves predicting whether a patient has sepsis or not.

**Productivity Prediction of Garment Employees** (PGE) Introduced in (Imran et al., 2021), this dataset contains information about garment manufacturing processing on a per-team level. Additionally, this dataset contains a team productivity performance index, which ranges between 0 and 1. As suggested by the authors, we use this index as a classification target. Specifically, we defined a team *efficient* if the productivity performance index is strictly greater than 0.75.

**Taxi** (TA) This dataset, introduced as part of the "ECML/PKDD 15: Taxi Trip Time Prediction (II) Competition" (Moreira-Matias et al., 2013) consists of 121,312 trajectories of Taxis in Porto (Portugal). The classification task is to predict the type of call that generated the run. The types of calls could be: A if this trip was dispatched from the central, B if this trip was demanded directly to a taxi driver on a specific stand C otherwise. The classes are balanced.

**Vehicles** (VE) GPS trajectories about two different types of vehicles -buses and trucks- moving in Athens. This dataset is available from download from the Chorochronos Archive (ChoroChronos Archive).

**UEA and UCR Irregular Datasets.** The other 22 irregular time-series datasets were downloaded from the UEA and UCR dataset repository. In particular, we included the following datasets:

• 11 variable-length univariate time series classification problems from (Bagnall et al., 2020): AllGestureWiimoteX, AllGestureWiimoteY and AllGestureWiimoteZ (GX, GY, GZ) from (Guna et al., 2014); GestureMidAirD1, GestureMidAirD2, and GestureMidAirD3 (GM1, GM2, GM3) from (Caputo et al., 2018); GesturePebbleZ1 and GesturePebbleZ2 (GP1, GP2) from (Mezari & Maglogiannis, 2018); PickupGestureWiimoteZ and ShakeGestureWiimoteZ (PGZ, SGZ) from (Guna et al., 2014); PLAID (PL) from (Gao et al., 2014);

 4 fixed length univariate time series with missing values from (Middlehurst et al., 2024b): DodgerLoopDay, DodgerLoopGame, and DodgerLoopWeekend (DD, DG, DW) from (Ihler et al., 2006); MelbournePedestrian (MP) (City of Melbourne, 2020) extracted from the City of Melbourne website;

7 variable-length multivariate time series from (Ruiz et al., 2021): AsphaltObstaclesCoordinates, AsphaltPavementTypeCoordinates, and AsphaltRegularityCoordinates (AOC, APT, ARC) from (de Souza, 2018); CharacterTrajectories (CT) from (Williams et al., 2006); InsectWingbeat (IW) from (Chen et al., 2014); JapaneseVowels (JV) from (Kudo et al., 1999); SpokenArabicDigits (SAD) from (Hammami & Bedda, 2010);

Table 5 contains the full list of curated datasets at the moment of publication on our repository. The list additionally contains some information about the datasets: the number of instances, #Inst, number of signals, #Sign, and number of observations, #Obs  $(\max_i^n(T_i))$ , the number of target classes #TC and the standard deviation between the number of instances per class (CU). Additionally, the dataset contains information about the time series, like the percentage of missing values (MV)-computed as the ratio between the NaN observations divided by the total number of observations- and the sampling coefficient of variation (SCV), alongside information on the different kind of irregularity in the dataset.

Table 5: Summary of dataset characteristics: the number of instances (#Inst), signals (#Sign), and observations (#Obs); target classes (#TC) and class imbalance (CU); as well as time-series-specific metrics like missing values (MV) and sampling coefficient of variation (SCV), and each type of irregularity, i.e., unevenly sampled (US), partially observed (PO), unequal length (UL), shift (SH), ragged sampling (RS).

Cat	Name	Source	#Inst	#Sign	#Obs	#TC	CU (σ)	MV (%)	SVC	US	РО	UL	SH	RS
ų	MI3	(Johnson et al., 2016)	57	17	145	2	0.20	0.83	0.60	/	/	/	/	/
health	P12	(Silva et al., 2012)	7990	37	203	2	0.36	0.94	0.59	/	/	/	/	/
he	P19	(Reyna et al., 2020)	40334	34	334	2	0.43	0.98	0.18	✓	/	/	/	1
	CT	(Williams et al., 2006)	2858	3	182	20	0.01	0.34	0.00	Х	Х	/	Х	Х
2	GM1	(Caputo et al., 2018)	338	1	360	26	0.00	0.54	0.00	X	X	/	X	X
<i>io</i> ;	GM2	(Caputo et al., 2018)	338	1	360	26	0.00	0.54	0.00	X	X	/	X	X
human activity recognition	GM3	(Caputo et al., 2018)	338	1	360	26	0.00	0.54	0.00	X	X	1	X	X
308	GP1	(Mezari & Maglogiannis, 2018)	304	1	455	6	0.01	0.52	0.00	X	X	/	X	X
rec	GP2	(Mezari & Maglogiannis, 2018)	304	1	455	6	0.01	0.52	0.00	X	X	/	X	X
ity	GX	(Guna et al., 2014)	1000	1	385	10	0.00	0.68	0.00	X	X	/	X	X
άż	GY	(Guna et al., 2014)	1000	1	385	10	0.00	0.68	0.00	X	X	/	X	X
ıα	GZ	(Guna et al., 2014)	1000	1	385	10	0.00	0.68	0.00	X	X	1	X	X
ıaı	LPA	(Vidulin et al., 2010)	273	12	2870	11	0.00	0.95	9.04	/	X	/	/	/
m	PAM	(Reiss & Stricker, 2012)	124	52	110883	16	0.03	0.82	0.01	/	/	1	/	/
_	PGZ	(Guna et al., 2014)	100	1	361	10	0.00	0.60	0.00	X	X	/	X	X
	SGZ	(Guna et al., 2014)	100	1	385	10	0.00	0.57	0.00	X	X	/	X	X
	AN	(Ferrero et al., 2018)	102	2	291	3	0.03	0.50	1.21	/	Х	/	Х	1
	AOC	(de Souza, 2018)	781	3	736	4	0.03	0.59	0.00	X	X	/	X	X
	APT	(de Souza, 2018)	2111	3	2371	3	0.06	0.83	0.00	X	X	1	X	X
ίtλ	ARC	(de Souza, 2018)	1502	3	4201	2	0.01	0.91	0.00	X	X	/	X	X
nobility	GS	(Zheng et al., 2010)	5977	2	96282	11	0.13	0.99	10.27	/	X	/	/	/
шс	MP	(City of Melbourne, 2020)	3633	1	24	10	0.00	0.00	0.01	X	X	/	X	X
	SE	(Browning et al., 2018)	108	4	6048	3	0.18	0.60	0.00	/	X	/	/	/
	TA	(Moreira-Matias et al., 2013)	121312	2	119	3	0.13	0.61	0.00	/	X	/	/	/
	VE	(ChoroChronos Archive)	381	2	1095	2	0.22	0.57	5.29	/	X	✓	X	1
-rc	DD	(Ihler et al., 2006)	158	1	288	7	0.01	0.01	0.00	Х	/	Х	Х	Х
sensor	DG	(Ihler et al., 2006)	158	1	288	2	0.02	0.01	0.00	X	/	X	X	X
se	DW	(Ihler et al., 2006)	158	1	288	2	0.21	0.01	0.00	X	/	X	X	X
	IW	(Chen et al., 2014)	50000	200	22	10	0.00	0.70	0.00	Х	Х	/	Х	Х
-	J۷	(Kudo et al., 1999)	640	12	29	9	0.03	0.46	0.00	X	X	1	X	X
other	PGE	(Imran et al., 2021)	24	9	59	2	0.13	0.19	0.68	1	X	1	/	1
0	PL	(Gao et al., 2014)	1074	1	1344	11	0.05	0.76	0.00	X	X	1	X	X
	SAD	(Hammami & Bedda, 2010)	8798	13	93	10	0.00	0.57	0.00	X	X	✓	X	X
synth	ABF	new!	930	1	128	3	0.00	0.00	1.95	/	Х	Х	Х	Х

Given  $y_h$  as the labels vector containing only the h-th class, CU is defined as follows:

$$CU = \sqrt{\frac{\sum_{h=0}^{c} (y_h - \mu)}{c}} \tag{1}$$

where  $\mu$  is the average number of observations. Given  $\Delta \tilde{\mathbf{t}}$  as the vector of differences between consecutive timestamps of a signal, the SCV is computed as the coefficient of variation (the ratio of the standard deviation to the mean) for each signal, averaged first across each time series and then over the entire dataset.

We divided the dataset into 6 categories based on the type of phenomena captured: *healthcare*, *human activity recognition*, *mobility* (or more generically, geo-temporal motion), *sensors*, *synthetic* data, and *others* for datasets that don't fall in any of the previous categories (like the UCR audio and speech categories).

# D ADDITIONAL RESULTS AND STATISTICAL TESTS

The full result table in terms of F1 is available in Table 7. Further, we provide several other statistical tests, using a diverse range of metrics, and with respect to different dataset subgroups.

**Critical Difference Plots.** Figure 10 shows the CD-plots for common performance metrics and runtimes. *F1*, *accuracy*, *roc-auc*, *precision*, and *recall* yield consistent rankings for the top four models, ROCKET, BORF, LGBM, and RIFC, as well as for the three lowest-performing ones: GRU-D, NCDE, and SVM. In the mid-range, rankings vary slightly across metrics: for instance, KNN performs notably worse in terms of F1 compared to accuracy, whereas TIMESNET shows the opposite trend. As for training time, KNN, being a lazy learner, is the fastest, followed by RIFC and ROCKET. Although LGBM ranks fourth, the previous results in median runtime (Figure 6) suggest that it may be slightly slower on smaller datasets but highly efficient on larger ones, which contributes to its overall strong performance. Neural network-based models generally exhibit longer training times but benefit from faster inference; nevertheless, ROCKET and LGBM maintain a performance edge across both phases.

F1 CD-plots computed for subsets of datasets with specific characteristics, are shown in Figure 11. These plots provide additional and complementary insights to those in Figures 7 and 8. Notably, they reinforce the observation that models explicitly designed for partially observed data tend to outperform more general-purpose approaches, even though the top rankings remain closely contested among SAITS, RIFC, LGBM, BRITS, and ROCKET. BRITS and TIMESNET, in particular, show strong performance on shorter datasets, ranking second and third, respectively, and closely trailing ROCKET. The remaining plots are similar to those discussed in Section 5.

**Multiple Comparison Matrices.** While the widely used CD-plot is effective, it has been criticized in (Ismail-Fawaz et al., 2023) for its susceptibility to manipulation, as the average rank of a model can be influenced by the performance of other comparators. For this reason, we also propose MCM matrix for several metrics in Figures 12 to 14. However, in our case, results are consistent with the CD-plots presented in the previous paragraph, and in the main text, and are presented here in the appendix only due to space limitations. Again, the top four models are always ROCKET, BORF, LGBM, and RIFC, and the lowest-performing are GRU-D, NCDE, and SVM, with mid-range models rankings changing slightly from metric to metric.

**Additional Performance vs Runtime Plots.** We report in Figures 15 to 19 the performance rankings across multiple metrics, dataset subsets, and with respect to both training and inference times. In addition to the insights discussed in the main text, these figures reveal that neural network-based models tend to cluster together in terms of both runtime and performance, regardless of the dataset subset or evaluation metric. This suggests that, although their relative rankings may vary, their overall behavior remains consistent.

Rank Correlation. We report in Figure 20 the F1 rank correlation among models. Models are hierarchically clustered using average linkage applied to the rank correlation matrix. Positive correlations indicate that models tend to perform similarly across datasets, reflecting comparable strengths or weaknesses, while negative correlations suggest divergent performance, highlighting complementary behaviors or differing inductive biases. Reinforcing the categorization proposed in the main text, the plot reveals a strong cluster of generalist methods, LGBM, ROCKET, RIFC, and BORF, which group together at the top hierarchical level. The second major cluster includes the remaining models, with specialist approaches like BRITS and GRU-D showing high correlation, which is expected given their shared RNN architecture. Similarly, TIMESNET and SAITS also form a coherent transformers subgroup. Notable exceptions to the generalist/specialist categorization are SVM, likely due to its overall poor performance across datasets, and KNN, which we hypothesize behaves differently due to its lazy learning paradigm based on distances, which could be more prone to sensitivity to dataset-specific characteristics.

**Model Failures and Limitations.** From these experiments, several model weaknesses become apparent, particularly in relation to specific data characteristics. For example, Figure 7 highlights how RNN-based methods fail to handle long time series effectively, while Table 3 shows that ROCKET underperformed relative to its baseline results after fine-tuning.

 Additional insights arise from the CD plots in Figure 11. Comparing the general rankings in Figure 11a with those on specific subsets reveals which models are most sensitive to dataset properties. For instance, Figure 11g shows that the transformer-based TIMESNET performs worse on smaller datasets, a point also observerd in Section 5. BORF, despite its strong overall performance, ranks third-to-last on partially observed data and declines significantly on short time series (Figures 11c and 11i). KNN also struggles under shift and ragged sampling conditions (Figures 11e and 11f). Notably, KNN was the weakest model in terms of memory consumption, which explodes with longer series (Table 9).

To provide a more fine-grained view, we report in Table 6 each model's worst performance in terms ratio between that worst-case rank and its average rank across all datasets. Higher ratios indicate greater variability, a phenomenon most pronounced among models that otherwise perform strongly on average, such as ROCKET, BORF, and LGBM. Several notable cases emerge. ROCKET, for instance, performs poorly on ABF, a dataset with highly uneven sampling. Similarly, BORF ranks 2.4 times worse than its average on the Mimic3 dataset, which is also highly irregular. Interestingly, LGBM performs unexpectedly poorly on the Garment dataset, whose small size would normally favor tree-based models.

These findings highlight that strong average performance does not necessarily imply robustness across all dataset types. In particular, models often fail on datasets with structural irregularities or atypical sampling patterns.

Table 6: Worst-case dataset performance for each model, along with the ratio between its worst rank and average rank across all datasets. Higher ratios indicate greater variability compared to average performance.

model	worst dataset performance	worst-to-average rank ratio
BORF	Mimic3	2.4
BRITS	AllGestureWiimoteX	1.8
GRU-D	CharacterTrajectories	1.6
KNN	Physionet2012	1.9
LGBM	Garment	2.3
NCDE	ShakeGestureWiimoteZ	1.4
RAINDROP	InsectWingbeat	1.8
RIFC	GeolifeSupervised	2.2
ROCKET	Abf	3.0
SAITS	Animals	1.5
SVM	AllGestureWiimoteY	1.2
TIMESNET	DodgerLoopDay	2.0

Impact of irregularity on explanations. As discussed in Section 5, XAI for irregular time series remains largely unexplored. pyrregular allows researchers to work directly with data while preserving its irregularities, avoiding the bias introduced by imputation choices, which is fundamental since explanations are known to be highly sensitive to input variations (Yeh et al., 2019). This, however, is only a first step. Even when the data retains its irregularity (as in our approach), and even when models can handle irregular inputs, the explainers themselves typically cannot. In line with the observations of Cinquini et al. (2023), we argue that this is primarily an implementation gap on the explainer side. Addressing this limitation would enable our taxonomy of irregularities to be applied to more fine-grained interpretability. For example, it could help distinguish whether a model assigns importance to a missing value because of partial observation or because of raggedness, offering deeper insights into the model's behavior under irregular conditions.

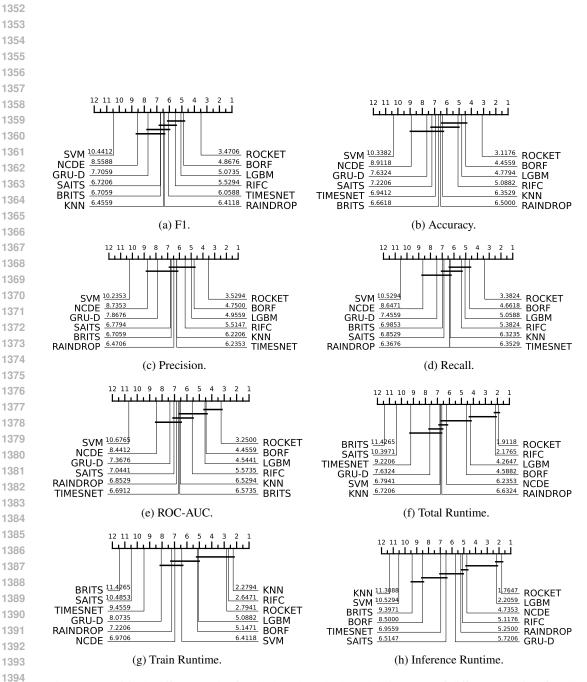


Figure 10: Critical Difference plot for the benchmarked models in terms of different metrics, for all datasets. Best models to the right. The performance of models connected by the bar is statistically tied, using a one-sided Holm-corrected Wilcoxon sign rank test with a critical value of 0.05.

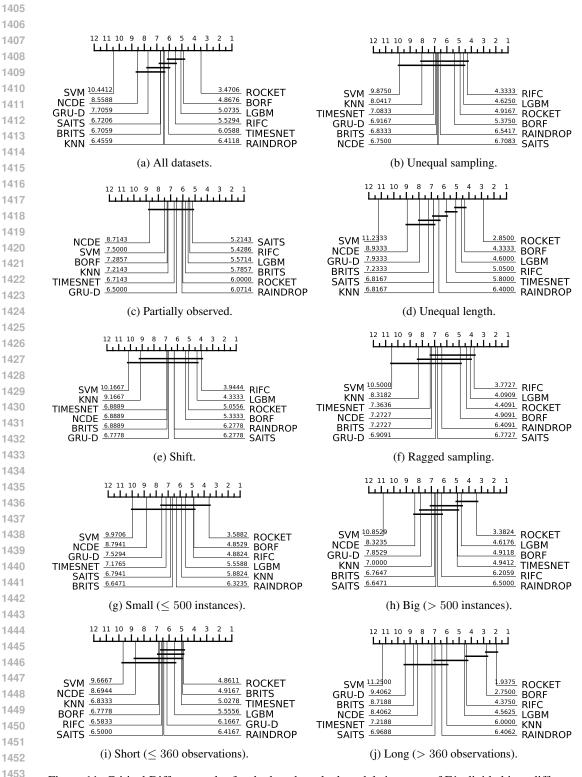


Figure 11: Critical Difference plot for the benchmarked models in terms of F1, divided into different groups. Best models to the right. The performance of models connected by the bar is statistically tied, using a one-sided Holm-corrected Wilcoxon sign rank test with a critical value of 0.05.

1460													
1461													
1462													
1463													
1464		ROCKET	BORF	LGBM	RIFC	TIMESNET	RAINDROP	SAITS	KNN	BRITS	GRU-D	NCDE	SVM
1465	Mean-f.	0.669	0.625	0.606	0.571	0.525	0.520	0.518	0.507	0.482	0.435	0.396	0.245
1466	ROCKET _	r>c / r=c / r <c Wilcoxon p-value</c 	0.044 24 / 3 / 7 0.001	0.063 24 / 3 / 7 0.003	0.098 26 / 4 / 4 ≤ 1e-03	0.144 24 / 0 / 10 ≤ 1e-03	0.149 26 / 0 / 8 ≤ 1e-03	0.151 25 / 0 / 9 0.001	0.162 28/0/6 ≤ 1e-03	0.187 25 / 0 / 9 0.001	0.234 25 / 0 / 9 ≤ 1e-03	0.273 27 / 0 / 7 ≤ 1e-03	0.424 31 / 0 / 3 ≤ 1e-03
1467 1468	BORF _ 0.625	-0.044 7 / 3 / 24 0.999		0.019 19 / 4 / 11 0.084	0.054 22 / 3 / 9 0.015	0.100 20 / 0 / 14 0.008	0.105 25 / 0 / 9 0.002	0.107 23 / 0 / 11 0.009	0.118 20 / 0 / 14 0.024	0.143 21 / 0 / 13 0.006	0.190 23 / 0 / 11 0.001	0.229 28 / 0 / 6 ≤ 1e-03	0.380 29 / 1 / 4 ≤ 1e-03
1469	LGBM _ 0.606	-0.063 7 / 3 / 24 0.997	-0.019 11 / 4 / 19 0.916	-	0.035 18 / 3 / 13 0.114	0.081 20 / 0 / 14 0.007	0.086 23 / 0 / 11 0.002	0.088 23 / 0 / 11 0.008	0.099 21 / 0 / 13 0.054	0.124 23 / 0 / 11 0.003	0.170 25 / 0 / 9 ≤ 1e-03	0.210 29 / 0 / 5 ≤ 1e-03	0.361 30 / 1 / 3 ≤ 1e-03
1470	RIFC - 0.571	-0.098 4 / 4 / 26 1.000	-0.054 9 / 3 / 22 0.985	-0.035 13 / 3 / 18 0.886		0.046 20 / 0 / 14 0.114	0.051 21 / 0 / 13 0.035	0.053 22 / 0 / 12 0.039	0.064 21 / 0 / 13 0.098	0.089 23 / 0 / 11 0.019	0.136 25 / 0 / 9 0.002	0.175 26 / 0 / 8 ≤ 1e-03	0.326 31 / 0 / 3 ≤ 1e-03
1471	TIMESNET _ 0.525	-0.144 10 / 0 / 24 1.000	-0.100 14 / 0 / 20 0.992	-0.081 14 / 0 / 20 0.994	-0.046 14 / 0 / 20 0.889		0.005 16 / 2 / 16 0.634	0.007 19/2/13 0.316	0.018 15 / 2 / 17 0.575	0.043 20 / 3 / 11 0.053	0.090 23 / 2 / 9 0.004	0.129 23 / 1 / 10 0.001	0.280 27 / 2 / 5 ≤ 1e-03
1472	RAINDROP _ 0.520	-0.149 8/0/26 1.000	-0.105 9 / 0 / 25 0.998	-0.086 11 / 0 / 23 0.998	-0.051 13 / 0 / 21 0.967	-0.005 16 / 2 / 16 0.366	-	0.002 20 / 2 / 12 0.269	0.013 16 / 2 / 16 0.595	0.038 16 / 2 / 16 0.152	0.085 20 / 2 / 12 0.014	0.124 26 / 1 / 7 ≤ 1e-03	0.275 28 / 3 / 3 ≤ 1e-03
1473	SAITS _ 0.518	-0.151 9 / 0 / 25 0.999	-0.107 11/0/23 0.992	-0.088 11 / 0 / 23 0.992	-0.053 12 / 0 / 22 0.962	-0.007 13 / 2 / 19 0.684	-0.002 12 / 2 / 20 0.731	0.209	0.011 12 / 2 / 20 0.714	0.036 17 / 2 / 15 0.196	0.082 20 / 2 / 12 0.022	0.122 27 / 1 / 6 ≤ 1e-03	0.273 29 / 2 / 3 ≤ 1e-03
1474 1475	KNN -	-0.162 6 / 0 / 28 1.000	-0.118 14/0/20 0.977	-0.099 13 / 0 / 21 0.948	-0.064 13 / 0 / 21 0.905	-0.018 17 / 2 / 15 0.425	-0.013 16/2/16 0.405	-0.011 20 / 2 / 12 0.286	0.714	0.025 14/3/17 0.422	0.022 0.071 18/3/13 0.101	0.111 23 / 2 / 9 0.016	0.262 25 / 5 / 4 ≤ 1e-03
1476	BRITS	1.000 -0.187 9 / 0 / 25 0.999	-0.143 13/0/21 0.994	0.948 -0.124 11 / 0 / 23 0.997	-0.089 11 / 0 / 23 0.982	-0.043 11 / 3 / 20 0.947	-0.038 16 / 2 / 16 0.848	-0.036 15 / 2 / 17 0.804	-0.025 17 / 3 / 14 0.578	0.422	0.101 0.047 19 / 4 / 11 0.032	0.016 0.086 22 / 1 / 11 0.004	≤ 1e-03 0.237 27/3/4 ≤ 1e-03
1477	0.482 GRU-D 0.435	-0.234 9 / 0 / 25 1.000	-0.190	-0.170	-0.136	-0.090 9 / 2 / 23	-0.085	-0.082	-0.071	-0.047	0.032	0.004 0.039 16 / 1 / 17	0.190
1478			11/0/23 1.000 -0.229	9 / 0 / 25 1.000 -0.210	9 / 0 / 25 0.998 -0.175	0.996	12 / 2 / 20 0.986 -0.124	12 / 2 / 20 0.978 -0.122	13/3/18 0.899 -0.111	11 / 4 / 19 0.968 -0.086	-0.039	0.182	27 / 2 / 5 ≤ 1e-03 0.151
1479	NCDE _ 0.396	-0.273 7 / 0 / 27 1.000	-0.229 6/0/28 1.000	-0.210 5 / 0 / 29 1.000	-0.175 8 / 0 / 26 1.000	-0.129 10 / 1 / 23 0.999	-0.124 7/1/26 1.000	-0.122 6/1/27 1.000	-0.111 9/2/23 0.984	11 / 1 / 22 0.996	17 / 1 / 16 0.818	0.151	0.151 27 / 1 / 6 ≤ 1e-03
1480	SVM - 0.245 -	-0.424 3 / 0 / 31 1.000	-0.380 4/1/29 1.000	-0.361 3 / 1 / 30 1.000	-0.326 3 / 0 / 31 1.000	-0.280 5 / 2 / 27 1.000	-0.275 3 / 3 / 28 1.000	-0.273 3 / 2 / 29 1.000	-0.262 4 / 5 / 25 1.000	-0.237 4/3/27 1.000	-0.190 5/2/27 1.000	-0.151 6/1/27 1.000	If in bold, then p-value < 0.05
1481							(a) F1	score.					
1482		ROCKET 0.742	BORF 0.700	LGBM 0.674	RIFC 0.644	RAINDROP 0.581	TIMESNET 0.576	SAITS 0.568	BRITS 0.554	KNN 0.542	GRU-D 0.508	NCDE 0.456	SVM 0.300
1483 1484	ROCKET _	Mean-Difference	0.042 24 / 3 / 7	0.068 25 / 3 / 6	0.098 25 / 4 / 5 ≤ 1e-03	0.161 29 / 0 / 5 ≤ 1e-03	0.166 29 / 0 / 5 ≤ 1e-03	0.174 28 / 0 / 6 ≤ 1e-03	0.188 25 / 0 / 9 ≤ 1e-03	0.200 27 / 1 / 6 ≤ 1e-03	0.234	0.286	0.442 29 / 2 / 3 ≤ 1e-03
1485	0.742 BORF	Wilcoxon p-value	0.001	0.001	0.056	0.119	0.124	0.132	0.146	0.158	≤ 1e-03	≤ 1e-03	0.400
1486		7 / 3 / 24 0.999 -0.068	-0.026	18 / 7 / 9 0.052	21 / 3 / 10 0.012 0.030	25 / 0 / 9 ≤ 1e-03 0.094	26 / 0 / 8 0.001 0.098	25 / 0 / 9 0.001 0.106	21 / 0 / 13 0.008 0.120	20 / 1 / 13 0.014 0.132	25 / 0 / 9 ≤ 1e-03 0.166	31 / 0 / 3 ≤ 1e-03 0.218	30 / 1 / 3 ≤ 1e-03
1487	LGBM _ 0.674	6/3/25 0.999 -0.098	-0.026 9 / 7 / 18 0.948 -0.056	-0.030	0.030 19/3/12 0.098	25 / 0 / 9 ≤ 1e-03	25 / 0 / 9 ≤ 1e-03	0.106 24 / 0 / 10 0.001	0.120 24 / 0 / 10 0.001	0.132 20 / 1 / 13 0.041 0.103	0.166 26 / 0 / 8 ≤ 1e-03	0.218 30 / 0 / 4 ≤ 1e-03	0.375 30 / 1 / 3 ≤ 1e-03
1488	RIFC 0.644	-0.098 5 / 4 / 25 1.000	-0.056 10 / 3 / 21 0.988	-0.030 12 / 3 / 19 0.902		0.064 24 / 0 / 10 0.012	0.069 25 / 0 / 9 0.007	0.076 24 / 1 / 9 0.005	0.091 24 / 0 / 10 0.011	0.103 21 / 0 / 13 0.045	0.137 26 / 0 / 8 ≤ 1e-03	0.188 27 / 0 / 7 ≤ 1e-03	0.345 31 / 1 / 2 ≤ 1e-03
1489	RAINDROP _ 0.581	-0.161 5 / 0 / 29 1.000	-0.119 9/0/25 1.000	-0.094 9 / 0 / 25 1.000	-0.064 10 / 0 / 24 0.988	-	0.005 18 / 2 / 14 0.252	0.012 21/2/11 0.114	0.027 18 / 2 / 14 0.226	0.039 15 / 3 / 16 0.537	0.073 20 / 2 / 12 0.021	0.124 28 / 1 / 5 ≤ 1e-03	0.281 27 / 2 / 5 ≤ 1e-03
1490	TIMESNET _ 0.576	-0.166 5 / 0 / 29 1.000	-0.124 8 / 0 / 26 0.999	-0.098 9 / 0 / 25 1.000	-0.069 9 / 0 / 25 0.994	-0.005 14 / 2 / 18 0.748		0.007 18 / 4 / 12 0.255	0.022 17 / 4 / 13 0.198	0.034 14 / 2 / 18 0.614	0.068 22 / 2 / 10 0.016	0.120 22 / 2 / 10 0.002	0.276 25 / 2 / 7 ≤ 1e-03
1491 1492	SAITS _ 0.568	-0.174 6 / 0 / 28 1.000	-0.132 9 / 0 / 25 0.999	-0.106 10 / 0 / 24 0.999	-0.076 9 / 1 / 24 0.995	-0.012 11 / 2 / 21 0.886	-0.007 12 / 4 / 18 0.745	-	0.015 14 / 3 / 17 0.591	0.027 12 / 2 / 20 0.684	0.061 18/2/14 0.071	0.112 25 / 1 / 8 0.001	0.269 28 / 2 / 4 ≤ 1e-03
1492	BRITS _ 0.554	-0.188 9 / 0 / 25 1.000	-0.146 13 / 0 / 21 0.992	-0.120 10 / 0 / 24 0.999	-0.091 10 / 0 / 24 0.989	-0.027 14 / 2 / 18 0.774	-0.022 13 / 4 / 17 0.802	-0.015 17 / 3 / 14 0.409	-	0.012 16/3/15 0.510	0.046 19 / 4 / 11 0.023	0.098 22 / 2 / 10 0.004	0.254 28 / 3 / 3 ≤ 1e-03
1494	KNN _ 0.542	-0.200 6 / 1 / 27 1.000	-0.158 13 / 1 / 20 0.986	-0.132 13 / 1 / 20 0.959	-0.103 13 / 0 / 21 0.956	-0.039 16 / 3 / 15 0.463	-0.034 18 / 2 / 14 0.386	-0.027 20 / 2 / 12 0.316	-0.012 15 / 3 / 16 0.490		0.034 17 / 4 / 13 0.158	0.086 23 / 2 / 9 0.028	0.242 26 / 5 / 3 ≤ 1e-03
1495	GRU-D 0.508	-0.234 8/1/25 1.000	-0.192 9 / 0 / 25 1.000	-0.166 8/0/26 1.000	-0.137 8/0/26 1.000	-0.073 12 / 2 / 20 0.979	-0.068 10 / 2 / 22 0.984	-0.061 14 / 2 / 18 0.929	-0.046 11 / 4 / 19 0.977	-0.034 13 / 4 / 17 0.842	0.158	0.052 19 / 1 / 14	0.208 27 / 3 / 4 ≤ 1e-03
1496	NCDE - 0.456	-0.286 5 / 0 / 29 1.000	-0.244 3 / 0 / 31 1.000	-0.218 4 / 0 / 30 1.000	-0.188 7/0/27 1.000	-0.124 5/1/28 1.000	-0.120 10 / 2 / 22 0.998	-0.112 8/1/25 0.999	-0.098 10 / 2 / 22 0.996	-0.086 9 / 2 / 23 0.972	-0.052 14 / 1 / 19 0.888	0.112	0.157 25 / 1 / 8 5 1e-03
1497	SVM	1.000 -0.442 3/2/29 1.000	-0.400 3 / 1 / 30 1.000	1.000 -0.375 3 / 1 / 30 1.000	1.000 -0.345 2/1/31 1.000	-0.281 5 / 2 / 27 1.000	0.998 -0.276 7 / 2 / 25 1.000	0.999 -0.269 4 / 2 / 28 1.000	0.996 -0.254 3/3/28 1.000	0.972 -0.242 3 / 5 / 26 1.000	0.209	-0.157 8/1/25 1.000	≤ 1e-03  If in bold, then p-value < 0.05
1498	0.300	1.000	1.000	1,000	1,000	1.000	1.000	1.000	1.000	1.000	4/3/27 1.000	1.000	p-value < 0.05
1499													

(b) Accuracy.

Figure 12: Summary performance statistics for the 12 classifiers on 34 datasets, generated using the multiple comparison matrix (MCM). The MCM shows pairwise comparisons. Each cell shows the mean difference in performance, wins/draws/losses, and Wilcoxon p-value for two comparates. The best models on the top left are sorted based on the average performance. The more intense the color, the higher the mean accuracy difference w.r.t. the comparate, positive (red) or negative (blue).

1513													
1514													
1515													
1516													
1517													
1518													
1519	Mean-precis	ROCKET 0.677	BORF 0.650	LGBM 0.619	RIFC 0.600	TIMESNET 0.550	SAITS 0.548	KNN 0.548	RAINDROP 0.543	BRITS 0.507	GRU-D 0.464	NCDE 0.415	SVM 0.269
1520		Mean-Difference r>c/r=c/r <c Wilcoxon p-value</c 	0.027 24 / 3 / 7 0.002	0.058 24 / 3 / 7 0.006	0.077 26 / 4 / 4 0.001	0.127 24 / 0 / 10 0.002	0.128 26 / 0 / 8 0.005	0.128 27 / 1 / 6 0.001	0.133 26 / 0 / 8 0.001	0.169 25 / 0 / 9 0.002	0.213 24 / 1 / 9 0.001	0.261 26 / 0 / 8 ≤ 1e-03	0.408 29 / 2 / 3 ≤ 1e-03
1521	BORF _	-0.027 7 / 3 / 24 0.998	0.002	0.031 19 / 4 / 11 0.074	0.050 20 / 3 / 11 0.024	0.099 22 / 0 / 12 0.016	0.101 24 / 0 / 10 0.014	0.101 21 / 0 / 13 0.047	0.106 24 / 0 / 10 0.005	0.142 22 / 0 / 12 0.008	0.186 25 / 0 / 9 0.001	0.234 28 / 0 / 6 ≤ 1e-03	0.381 29 / 1 / 4 ≤ 1e-03
1522	0.650 LGBM		-0.031	0.074							0.155		
1523	0.619	-0.058 7/3/24 0.994	-0.031 11 / 4 / 19 0.926	0.010	0.019 19/3/12 0.104	0.069 21 / 0 / 13 0.005	0.071 23 / 0 / 11 0.014	0.071 21 / 0 / 13 0.081	0.076 23 / 0 / 11 0.005	0.112 24 / 0 / 10 0.005	26 / 0 / 8 0.001	0.204 29 / 0 / 5 ≤ 1e-03	0.350 30 / 1 / 3 ≤ 1e-03
1524	RIFC _ 0.600	-0.077 4 / 4 / 26 0.999	-0.050 11 / 3 / 20 0.976	-0.019 12 / 3 / 19 0.896	-	0.050 19 / 0 / 15 0.107	0.051 22 / 0 / 12 0.045	0.051 21 / 0 / 13 0.176	0.056 23 / 0 / 11 0.032	0.092 22 / 0 / 12 0.024	0.136 25 / 0 / 9 0.002	0.184 27 / 0 / 7 0.001	0.331 29 / 1 / 4 ≤ 1e-03
1525	TIMESNET _ 0.550	-0.127 10 / 0 / 24 0.998	-0.099 12 / 0 / 22 0.985	-0.069 13 / 0 / 21 0.995	-0.050 15 / 0 / 19 0.896		0.002 19 / 2 / 13 0.316	0.002 14 / 2 / 18 0.684	0.007 13 / 2 / 19 0.684	0.043 19 / 3 / 12 0.095	0.087 22 / 2 / 10 0.015	0.135 23 / 1 / 10 0.001	0.281 29 / 2 / 3 ≤ 1e-03
1526	SAITS _ 0.548	-0.128 8 / 0 / 26 0.995	-0.101 10 / 0 / 24 0.986	-0.071 11 / 0 / 23 0.987	-0.051 12 / 0 / 22 0.956	-0.002 13 / 2 / 19 0.684		0.000 12 / 2 / 20 0.748	0.005 13 / 2 / 19 0.608	0.041 16/2/16 0.201	0.085 19 / 2 / 13 0.029	0.133 28 / 1 / 5 ≤ 1e-03	0.279 29 / 2 / 3 ≤ 1e-03
1527	KNN _ 0.548	-0.128 6 / 1 / 27 0.999	-0.101 13 / 0 / 21 0.955	-0.071 13 / 0 / 21 0.921	-0.051 13 / 0 / 21 0.829	-0.002 18 / 2 / 14 0.316	-0.000 20 / 2 / 12 0.252	-	0.005 18 / 2 / 14 0.236	0.041 16/3/15 0.278	0.085 20 / 4 / 10 0.037	0.133 23 / 2 / 9 0.008	0.279 26 / 5 / 3 ≤ 1e-03
1528 1529	RAINDROP _ 0.543	-0.133 8 / 0 / 26	-0.106 10 / 0 / 24	-0.076 11 / 0 / 23	-0.056 11 / 0 / 23	-0.007 19 / 2 / 13	-0.252 -0.005 19 / 2 / 13	-0.005 14 / 2 / 18	0.236	0.278 0.036 15 / 2 / 17	0.037 0.080 20 / 2 / 12	0.128 26 / 1 / 7 \$ 1e-03	0.274 28 / 3 / 3 ≤ 1e-03
1530	0.543 BRITS	0.999	0.995	0.995	0.969	0.316	-0.392	0.764	-0.036	0.211	0.018	0.092	0.238
1531	0.507	9 / 0 / 25 0.998	-0.142 12 / 0 / 22 0.992	10 / 0 / 24 0.995	12 / 0 / 22 0.977	12 / 3 / 19 0.905	16 / 2 / 16 0.799	15 / 3 / 16 0.722	17 / 2 / 15 0.789 -0.080	-0.044	0.044 20 / 4 / 10 0.059	23 / 1 / 10 0.004	25 / 3 / 6 ≤ 1e-03 0.195
1532	GRU-D _ 0.464	-0.213 9 / 1 / 24 0.999	-0.186 9 / 0 / 25 0.999	-0.155 8 / 0 / 26 0.999	-0.136 9 / 0 / 25 0.998	-0.087 10 / 2 / 22 0.985	-0.085 13 / 2 / 19 0.971	-0.085 10 / 4 / 20 0.963	12 / 2 / 20 0.982	-0.044 10 / 4 / 20 0.941		0.048 17 / 1 / 16 0.137	24/3/7 ≤ 1e-03
1533	NCDE _ 0.415 _	-0.261 8 / 0 / 26 1.000	-0.234 6 / 0 / 28 1.000	-0.204 5 / 0 / 29 1.000	-0.184 7 / 0 / 27 0.999	-0.135 10 / 1 / 23 0.999	-0.133 5 / 1 / 28 1.000	-0.133 9 / 2 / 23 0.992	-0.128 7 / 1 / 26 1.000	-0.092 10 / 1 / 23 0.996	-0.048 16 / 1 / 17 0.863		0.146 24 / 1 / 9 0.001
1534	SVM _ 0.269	-0.408 3 / 2 / 29 1.000	-0.381 4 / 1 / 29 1.000	-0.350 3 / 1 / 30 1.000	-0.331 4 / 1 / 29 1.000	-0.281 3 / 2 / 29 1.000	-0.279 3 / 2 / 29 1.000	-0.279 3 / 5 / 26 1.000	-0.274 3 / 3 / 28 1.000	-0.238 6 / 3 / 25 1.000	-0.195 7 / 3 / 24 1.000	-0.146 9 / 1 / 24 0.999	If in bold, then p-value < 0.05
1535							( ) D						
1536							(a) Prec						
1537	Mean-rec		BORF 0.645	LGBM 0.625	RIFC 0.591	TIMESNET 0.540	RAINDROP 0.538	KNN 0.535	SAITS 0.533	BRITS 0.506	GRU-D 0.466	NCDE 0.423	SVM 0.283
1538	ROCKET	Mean-Difference r>c / r=c / r <c Wilcoxon p-value</c 	0.047 22 / 4 / 8 0.001	0.067 24 / 4 / 6 0.001	0.100 26 / 4 / 4 ≤ 1e-03	0.151 25 / 1 / 8 ≤ 1e-03	0.153 27 / 1 / 6 ≤ 1e-03	0.156 27 / 1 / 6 ≤ 1e-03	0.158 26 / 0 / 8 ≤ 1e-03	0.186 26 / 0 / 8 ≤ 1e-03	0.226 24 / 1 / 9 ≤ 1e-03	0.268 27 / 0 / 7 ≤ 1e-03	0.409 30 / 2 / 2 ≤ 1e-03
1539	BORF _ 0.645	-0.047 8 / 4 / 22 0.999		0.020 19 / 4 / 11 0.074	0.053 22 / 4 / 8 0.013	0.104 22 / 1 / 11 0.004	0.106 24 / 1 / 9 0.001	0.110 20 / 0 / 14 0.026	0.111 24 / 0 / 10 0.003	0.139 22 / 0 / 12 0.006	0.179 23 / 0 / 11 ≤ 1e-03	0.221 28 / 0 / 6 ≤ 1e-03	0.362 30 / 1 / 3 ≤ 1e-03
1540	LGBM _ 0.625	-0.067 6 / 4 / 24 0.999	-0.020 11 / 4 / 19 0.926	-	0.033 19 / 4 / 11 0.072	0.085 21 / 1 / 12 0.002	0.087 23 / 1 / 10 0.001	0.090 19 / 1 / 14 0.071	0.092 22 / 0 / 12 0.004	0.119 24 / 0 / 10 0.001	0.159 24 / 0 / 10 ≤ 1e-03	0.202 29 / 0 / 5 ≤ 1e-03	0.342 30 / 1 / 3 ≤ 1e-03
1541	RIFC	-0.100 4/4/26	-0.053 8 / 4 / 22	-0.033 11 / 4 / 19 0.928	-	0.051 22 / 1 / 11 0.044	0.053 22 / 1 / 11 0.026	0.056 21 / 0 / 13 0.098	0.058 23 / 1 / 10 0.020	0.086 24 / 0 / 10 0.012	0.125 25 / 0 / 9 0.001	0.168 26 / 0 / 8 ≤ 1e-03	0.309 31 / 1 / 2 ≤ 1e-03
1542	0.591 TIMESNET	1.000 -0.151	0.987	0.928 -0.085 12 / 1 / 21 0.998	-0.051 11 / 1 / 22 0.956	0.044	0.026 0.002 15 / 3 / 16 0.585	0.098 0.005 14 / 2 / 18 0.678	0.007	0.034	0.001 0.074 22 / 2 / 10 0.012	0.117	≤ 1e-03 0.257 26 / 3 / 5 ≤ 1e-03
1543 1544	0.540 RAINDROP	8 / 1 / 25 1.000 -0.153	11 / 1 / 22 0.996 -0.106	0.998 -0.087 10 / 1 / 23			0.585		19 / 2 / 13 0.252 0.005	21 / 3 / 10 0.065 0.032		23 / 1 / 10 0.001 0.115	
1544	0.538	-0.153 6/1/27 1.000 -0.156	-0.106 9 / 1 / 24 0.999 -0.110	0.999	-0.053 11 / 1 / 22 0.974	-0.002 16 / 3 / 15 0.415	-	0.003 16 / 2 / 16 0.640	0.005 21 / 2 / 11 0.231 0.002	0.032 19 / 2 / 13 0.148 0.029	0.072 20 / 2 / 12 0.016 0.069	0.115 26 / 1 / 7 ≤ 1e-03 0.112	0.255 28/3/3 ≤ 1e-03 0.252
1546	KNN _ 0.535	6/1/27 1.000	14 / 0 / 20 0.975	-0.090 14 / 1 / 19 0.929	-0.056 13 / 0 / 21 0.905	-0.005 18 / 2 / 14 0.322	-0.003 16 / 2 / 16 0.360	-	19 / 2 / 13 0.269	15 / 4 / 15 0.356	17 / 4 / 13 0.095	23 / 2 / 9 0.011	26 / 6 / 2 ≤ 1e-03
1547	SAITS _ 0.533	-0.158 8 / 0 / 26 1.000	-0.111 10 / 0 / 24 0.997	-0.092 12 / 0 / 22 0.996	-0.058 10 / 1 / 23 0.980	-0.007 13 / 2 / 19 0.748	-0.005 11 / 2 / 21 0.769	-0.002 13 / 2 / 19 0.731	-	0.027 16 / 2 / 16 0.280	0.067 19 / 2 / 13 0.039	0.110 27 / 1 / 6 ≤ 1e-03	0.250 29 / 2 / 3 ≤ 1e-03
1548	BRITS _ 0.506	-0.186 8 / 0 / 26 1.000	-0.139 12 / 0 / 22 0.994	-0.119 10 / 0 / 24 0.999	-0.086 10 / 0 / 24 0.988	-0.034 10 / 3 / 21 0.935	-0.032 13 / 2 / 19 0.852	-0.029 15 / 4 / 15 0.644	-0.027 16 / 2 / 16 0.720		0.040 17 / 4 / 13 0.082	0.083 23 / 1 / 10 0.005	0.223 27 / 3 / 4 ≤ 1e-03
1549	GRU-D _ 0.466	-0.226 9 / 1 / 24	-0.179 11 / 0 / 23	-0.159 10 / 0 / 24	-0.125 9 / 0 / 25	-0.074 10 / 2 / 22	-0.072 12 / 2 / 20	-0.069 13 / 4 / 17	-0.067 13 / 2 / 19	-0.040 13 / 4 / 17		0.043 18 / 1 / 15	0.183 27 / 3 / 4 ≤ 1e-03
1550	NCDE	1.000 -0.268 7/0/27 1.000	1.000 -0.221 6/0/28 1.000	1.000 -0.202 5 / 0 / 29 1.000	0.999 -0.168 8 / 0 / 26 1.000	0.988 -0.117 10 / 1 / 23 0.999	0.984 -0.115 7/1/26 1.000	0.905 -0.112 9/2/23 0.989	0.961 -0.110 6/1/27 1.000	0.918 -0.083 10 / 1 / 23 0.995	-0.043 15 / 1 / 18 0.863	0.137	0.140 27 / 1 / 6 ≤ 1e-03
1551	0.423 SVM	-0.409			-0.309				-0.250		0.863 -0.183 4/3/27	-0.140 6/1/27	If in bold, then
1552	0.283	2 / 2 / 30 1.000	-0.362 3 / 1 / 30 1.000	-0.342 3 / 1 / 30 1.000	2 / 1 / 31 1.000	-0.257 5/3/26 1.000	-0.255 3/3/28 1.000	-0.252 2 / 6 / 26 1.000	3 / 2 / 29 1.000	-0.223 4/3/27 1.000	1.000	1.000	p-value < 0.05
4550													

(b) Recall.

Figure 13: Summary performance statistics for the 12 classifiers on 34 datasets, generated using the multiple comparison matrix (MCM). The MCM shows pairwise comparisons. Each cell shows the mean difference in performance, wins/draws/losses, and Wilcoxon p-value for two comparates. The best models on the top left are sorted based on the average performance. The more intense the color, the higher the mean accuracy difference w.r.t. the comparate, positive (red) or negative (blue).

Mean-total		ROCKET 1485.083	RIFC 1656.918	BORF 5298.659	NCDE 6877.340	GRU-D 9351.653	TIMESNET 12925.685	KNN 16523.164	RAINDROP 18921.227	SVM 21772.918 I	SAITS 30426.722	BRITS 60025.019
LGBM _ 472.791	Mean-Difference - r>c / r=c / r <c Wilcoxon p-value</c 	-1012.291 29 / 0 / 5 ≤ 1e-03	-1184.127 22 / 0 / 12 0.058	-4825.868 14 / 0 / 20 0.845	-6404.549 13 / 0 / 21 0.902	-8878.861 3 / 0 / 31 1.000	-12452.894 0 / 0 / 34 1.000	-16050.373 8 / 0 / 26 1.000	-18448.436 12 / 0 / 22 0.988	-21300.127 10 / 0 / 24 0.999	-29953.930 0 / 0 / 34 1.000	-59552.228 0 / 0 / 34 1.000
ROCKET _ 1485.083	1012.291 5 / 0 / 29 1.000		-171.836 18 / 0 / 16 0.407	-3813.576 5 / 0 / 29 1.000	-5392.258 1 / 0 / 33 1.000	-7866.570 1 / 0 / 33 1.000	-11440.602 0 / 0 / 34 1.000	-15038.081 1 / 0 / 33 1.000	-17436.145 0 / 0 / 34 1.000	-20287.836 0 / 0 / 34 1.000	-28941.639 0 / 0 / 34 1.000	-58539.936 0 / 0 / 34 1.000
RIFC - 1656.918	1184.127 12 / 0 / 22 0.944	171.836 16 / 0 / 18 0.600	-	-3641.741 2 / 0 / 32 1.000	-5220.422 2 / 0 / 32 1.000	-7694.734 1 / 0 / 33 1.000	-11268.767 1/0/33 1.000	-14866.246 2 / 0 / 32 1.000	-17264.309 1/0/33 1.000	-20116.000 1/0/33 1.000	-28769.803 1 / 0 / 33 1.000	-58368.101 1 / 0 / 33 1.000
BORF - 5298.659	4825.868 - 20 / 0 / 14 0.159	3813.576 29 / 0 / 5 ≤ 1e-03	3641.741 32 / 0 / 2 ≤ 1e-03		-1578.682 13 / 0 / 21 0.837	-4052.994 3 / 0 / 31 1.000	-7627.026 2 / 0 / 32 1.000	-11224.505 8 / 0 / 26 1.000	-13622.568 5 / 0 / 29 1.000	-16474.260 7 / 0 / 27 1.000	-25128.063 2 / 0 / 32 1.000	-54726.360 1 / 0 / 33 1.000
NCDE _ 6877.340	6404.549 - 21 / 0 / 13 0.101	5392.258 33 / 0 / 1 ≤ 1e-03	5220.422 32 / 0 / 2 ≤ 1e-03	1578.682 21 / 0 / 13 0.167		-2474.312 10 / 1 / 23 0.998	-6048.345 5 / 1 / 28 1.000	-9645.823 15 / 2 / 17 0.984	-12043.887 17 / 1 / 16 0.822	-14895.578 14 / 1 / 19 0.993	-23549.381 4 / 1 / 29 1.000	-53147.678 2 / 1 / 31 1.000
GRU-D _ 9351.653	8878.861 31 / 0 / 3 ≤ 1e-03	7866.570 33 / 0 / 1 ≤ 1e-03	7694.734 33 / 0 / 1 ≤ 1e-03	4052.994 31 / 0 / 3 ≤ 1e-03	2474.312 23 / 1 / 10 0.002		-3574.032 4 / 2 / 28 1.000	-7171.511 19 / 2 / 13 0.672	-9569.575 21 / 2 / 11 0.335	-12421.266 23 / 2 / 9 0.196	-21075.069 1 / 2 / 31 1.000	-50673.366 0 / 2 / 32 1.000
TIMESNET _ 12925.685	12452.894 34 / 0 / 0 ≤ 1e-03	11440.602 34 / 0 / 0 ≤ 1e-03	11268.767 33 / 0 / 1 ≤ 1e-03	7627.026 32 / 0 / 2 ≤ 1e-03	6048.345 28 / 1 / 5 ≤ 1e-03	3574.032 28 / 2 / 4 ≤ 1e-03		-3597.479 24 / 2 / 8 0.097	-5995.542 27 / 2 / 5 0.001	-8847.233 27 / 2 / 5 0.002	-17501.037 4 / 2 / 28 1.000	-47099.334 2 / 2 / 30 1.000
KNN - 16523.164	16050.373 26 / 0 / 8 ≤ 1e-03	15038.081 33 / 0 / 1 ≤ 1e-03	14866.246 32 / 0 / 2 ≤ 1e-03	11224.505 26 / 0 / 8 ≤ 1e-03	9645.823 17 / 2 / 15 0.016	7171.511 13 / 2 / 19 0.328	3597.479 8 / 2 / 24 0.903		-2398.063 16 / 2 / 16 0.221	-5249.755 13 / 3 / 18 0.551	-13903.558 3 / 2 / 29 1.000	-43501.855 1 / 2 / 31 1.000
RAINDROP _ 18921.227	18448.436 - 22 / 0 / 12 0.012	17436.145 34 / 0 / 0 ≤ 1e-03	17264.309 33 / 0 / 1 ≤ 1e-03	13622.568 29 / 0 / 5 ≤ 1e-03	12043.887 16 / 1 / 17 0.178	9569.575 11 / 2 / 21 0.665	5995.542 5 / 2 / 27 0.999	2398.063 16 / 2 / 16 0.779	-	-2851.691 16 / 2 / 16 0.684	-11505.494 2 / 2 / 30 1.000	-41103.792 1/2/31 1.000
SVM - 21772.918	21300.127 - 24 / 0 / 10 0.001	20287.836 34 / 0 / 0 ≤ 1e-03	20116.000 33 / 0 / 1 ≤ 1e-03	16474.260 27 / 0 / 7 ≤ 1e-03	14895.578 19 / 1 / 14 0.007	12421.266 9 / 2 / 23 0.804	8847.233 5 / 2 / 27 0.998	5249.755 18 / 3 / 13 0.449	2851.691 16 / 2 / 16 0.316	-	-8653.803 3 / 2 / 29 1.000	-38252.100 1/2/31 1.000
SAITS = 30426.722	29953.930 34 / 0 / 0 ≤ 1e-03	28941.639 34 / 0 / 0 ≤ 1e-03	28769.803 33 / 0 / 1 ≤ 1e-03	25128.063 32 / 0 / 2 ≤ 1e-03	23549.381 29 / 1 / 4 ≤ 1e-03	21075.069 31 / 2 / 1 ≤ 1e-03	17501.037 28 / 2 / 4 ≤ 1e-03	13903.558 29 / 2 / 3 ≤ 1e-03	11505.494 30 / 2 / 2 ≤ 1e-03	8653.803 29 / 2 / 3 ≤ 1e-03	-	-29598.297 4/2/28 1.000
BRITS _ 60025.019	59552.228 34 / 0 / 0 ≤ 1e-03	58539.936 34 / 0 / 0 ≤ 1e-03	58368.101 33 / 0 / 1 ≤ 1e-03	54726.360 33 / 0 / 1 ≤ 1e-03	53147.678 31 / 1 / 2 ≤ 1e-03	50673.366 32 / 2 / 0 ≤ 1e-03	47099.334 30 / 2 / 2 ≤ 1e-03	43501.855 31 / 2 / 1 ≤ 1e-03	41103.792 31 / 2 / 1 ≤ 1e-03	38252.100 31 / 2 / 1 ≤ 1e-03	29598.297 28 / 2 / 4 ≤ 1e-03	If in bold, then p-value < 0.05
					(a	) Total F	Runtime					
	ROCKET 0.851	BORF 0.844	LGBM 0.832	RIFC 0.815	SAITS 0.748	RAINDROP 0.744	TIMESNET 0.741	BRITS 0.736	GRU-D 0.718	KNN 0.701	NCDE 0.693	SVM 0.493
Mean-roc_ ROCKET 0.851	auc   Mean-Difference r>c/r=c/r <c Wilcoxon p-value</c 	0.007 23 / 4 / 7 0.007	0.019 21 / 4 / 9 0.021	0.036 27 / 4 / 3 ≤ 1e-03	0.104 27 / 0 / 7 0.002	0.108 27 / 1 / 6 0.001	0.111 27 / 0 / 7 ≤ 1e-03	0.115 25 / 1 / 8 0.002	0.133 25 / 0 / 9 0.002	0.151 30 / 0 / 4 ≤ 1e-03	0.158 27 / 0 / 7 ≤ 1e-03	0.358 31/1/2 ≤ 1e-03
BORF -	-0.007 - 7 / 4 / 23	0.007	0.012 18 / 4 / 12	0.029 22 / 4 / 8 0.007	0.002 0.097 27 / 0 / 7 0.001	0.101 25/1/8 0.001	0.104 24 / 0 / 10 ≤ 1e-03	0.108 22 / 1 / 11 0.004	0.126 22 / 0 / 12	0.144 23 / 0 / 11 0.005	0.151 29 / 0 / 5 ≤ 1e-03	0.351 30 / 1 / 3 ≤ 1e-03
LGBM _ 0.832 _	0.993 -0.019 - 9/4/21 0.979	-0.012 12 / 4 / 18 0.872	0.128	0.007 0.017 20 / 4 / 10 0.057	0.001 0.085 24 / 0 / 10 0.004	0.001 0.089 27/1/6 ≤ 1e-03	≤ 1e-03 0.092 25 / 0 / 9 ≤ 1e-03	0.004 0.096 23 / 1 / 10 0.002	0.001 0.114 26 / 0 / 8 0.001	0.132 21 / 0 / 13 0.025	0.139 29 / 0 / 5 ≤ 1e-03	0.339 30 / 1 / 3 ≤ 1e-03
RIFC _ 0.815	-0.036 - 3 / 4 / 27 1.000	-0.029 8 / 4 / 22 0.993	-0.017 10 / 4 / 20 0.943	-	0.068 24 / 0 / 10 0.026	0.072 21 / 1 / 12 0.045	0.075 21 / 0 / 13 0.030	0.002 0.080 21 / 1 / 12 0.016	0.001 0.098 24 / 0 / 10 0.003	0.025 0.115 22 / 0 / 12 0.023	0.122 26 / 0 / 8 0.002	0.322 31/1/2 ≤ 1e-03
SAITS _ 0.748	-0.104 7/0/27	-0.097 7 / 0 / 27	-0.085 10 / 0 / 24	-0.068 10/0/24	-	0.004 14 / 2 / 18 0.678	0.007 13 / 2 / 19 0.702	0.012 14 / 2 / 18 0.425	0.030 20 / 2 / 12 0.091	0.047 13 / 2 / 19 0.541	0.054 25 / 1 / 8 0.010	0.254 29 / 2 / 3 ≤ 1e-03
RAINDROP _ 0.744	-0.108 - 6/1/27 0.999	0.999 -0.101 8 / 1 / 25 0.999	0.996 -0.089 6/1/27 1.000	0.975 -0.072 12 / 1 / 21 0.955	-0.004 18 / 2 / 14 0.322	0.678	0.702 0.003 17 / 2 / 15 0.231	0.425 0.008 16/3/15 0.307	0.026 19 / 2 / 13 0.088	0.043 13 / 2 / 19 0.527	0.010 0.050 23 / 1 / 10 0.010	0.250 27 / 4 / 3 ≤ 1e-03
TIMESNET =	-0.111 - 7/0/27 1.000	-0.104 10/0/24 1.000	-0.092 9 / 0 / 25 1.000	-0.075 13 / 0 / 21 0.971	-0.007 19 / 2 / 13 0.298	-0.003 15 / 2 / 17 0.769	0.231	0.005 18 / 2 / 14 0.298	0.023 21 / 2 / 11 0.058	0.527 0.040 13/2/19 0.601	0.047 23/1/10 0.052	0.247 26 / 2 / 6 ≤ 1e-03
	-0.115 - 8/1/25	-0.108 11 / 1 / 22	-0.096 10 / 1 / 23	-0.080 12 / 1 / 21	-0.012 18 / 2 / 14 0.575	-0.008 15 / 3 / 16	-0.005 14 / 2 / 18 0.702	0.298	0.018 20 / 2 / 12	0.601 0.035 15 / 2 / 17 0.527	0.052 0.043 23 / 1 / 10 0.014	0.243 29/3/2 ≤ 1e-03
BRITS _ 0.736		0.996	0.998	0.984		0.693		0.019	0.122	0.017		
0.736 GRU-D	-0.133 - 9/0/25	-0.126 12 / 0 / 22	-0.114 8/0/26	-0.098 10 / 0 / 24	-0.030 12 / 2 / 20	13 / 2 / 19	11/2/21	12 / 2 / 20	-	15 / 2 / 17	0.025 20 / 1 / 13	0.225 29 / 2 / 3
0.736 GRU-D _ 0.718 _ KNN	-0.133 - 9/0/25 0.998 -0.151 - 4/0/30	-0.126 12 / 0 / 22 0.999	-0.132	-0.098 10 / 0 / 24 0.997 -0.115 12 / 0 / 22	-0.030 12/2/20 0.909 -0.047 19/2/13	-0.026 13 / 2 / 19 0.912 -0.043 19 / 2 / 13	-0.023 11/2/21 0.942 -0.040 19/2/13	-0.018 12 / 2 / 20 0.878 -0.035 17 / 2 / 15	-0.017 17/2/15	15 / 2 / 17 0.708	20 / 1 / 13 0.161	29 / 2 / 3 ≤ 1e-03
0.736  GRU-D = 0.718 = KNN = 0.701 = NCDE = 0.736	-0.133 - 9/0/25 -0.998 -0.151 -4/0/30 1.000	-0.126 12 / 0 / 22 0.999 -0.144 11 / 0 / 23 0.995	0.999 -0.132 13/0/21 0.976	-0.115 12 / 0 / 22 0.978	-0.047 19 / 2 / 13 0.459	-0.043 19 / 2 / 13 0.473	-0.040 19 / 2 / 13 0.399	-0.035 17 / 2 / 15 0.473	17 / 2 / 15 0.292	15/2/17 0.708	0.025 20/1/13 0.161 0.007 20/2/12 0.152	29 / 2 / 3 ≤ 1e-03 0.207 26 / 6 / 2 ≤ 1e-03 0.200
0.736 GRU-D = 0.718 = KNN = 0.701 =	-0.133 - 9/0/25 0.998 -0.151 - 4/0/30 1.000	-0.126 12 / 0 / 22 0.999 -0.144 11 / 0 / 23 0.995	0.999 -0.132 13 / 0 / 21 0.976	-0.115 12 / 0 / 22 0.978	-0.047 19 / 2 / 13 0.459	-0.043 19 / 2 / 13 0.473	-0.040 19 / 2 / 13 0.399	-0.035 17 / 2 / 15 0.473	17 / 2 / 15 0.292	15/2/17 0.708	20 / 1 / 13 0.161	29 / 2 / 3 ≤ 1e-03 0.207 26 / 6 / 2 ≤ 1e-03

(b) ROC-AUC.

Figure 14: Summary performance statistics for the 12 classifiers on 34 datasets, generated using the multiple comparison matrix (MCM). The MCM shows pairwise comparisons. Each cell shows the mean difference in performance, wins/draws/losses, and Wilcoxon p-value for two comparates. The best models on the top left are sorted based on the average performance. The more intense the color, the higher the mean accuracy difference w.r.t. the comparate, positive (red) or negative (blue).

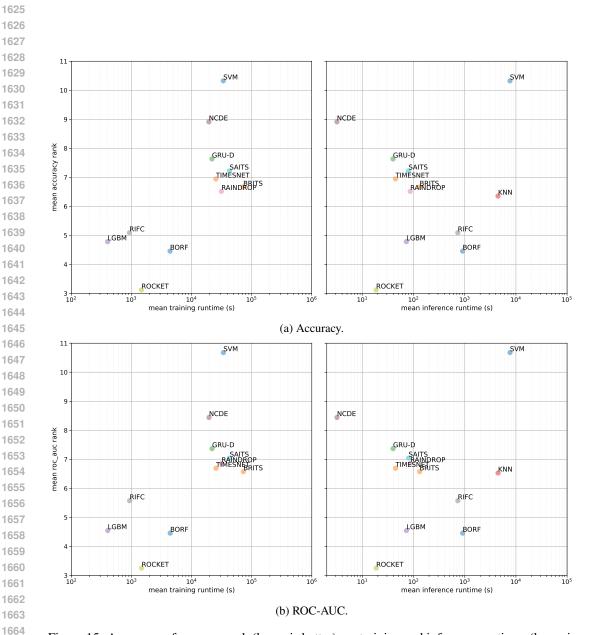


Figure 15: Average performance rank (lower is better) vs. training and inference runtimes (lower is better). Best values are on the bottom-left of each plot.

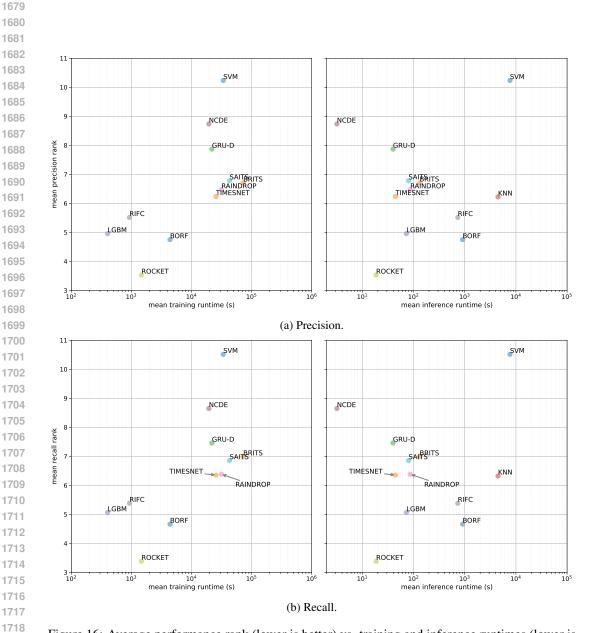


Figure 16: Average performance rank (lower is better) vs. training and inference runtimes (lower is better). Best values are on the bottom-left of each plot.

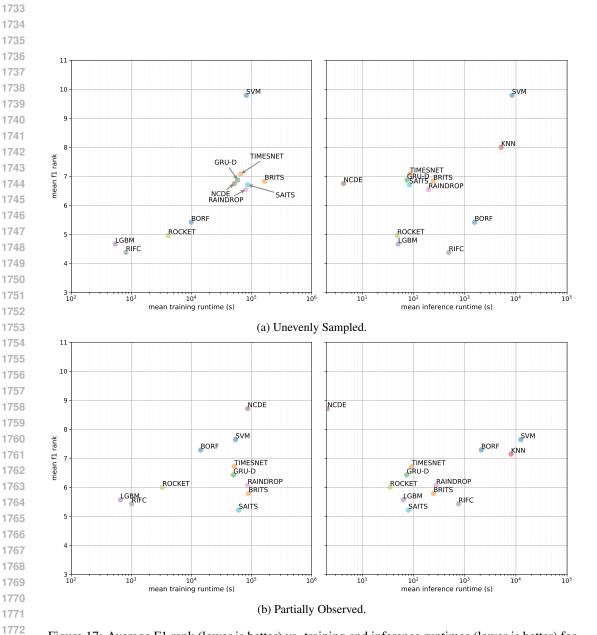


Figure 17: Average F1 rank (lower is better) vs. training and inference runtimes (lower is better) for subsets of datasets. Best values are on the bottom-left of each plot.

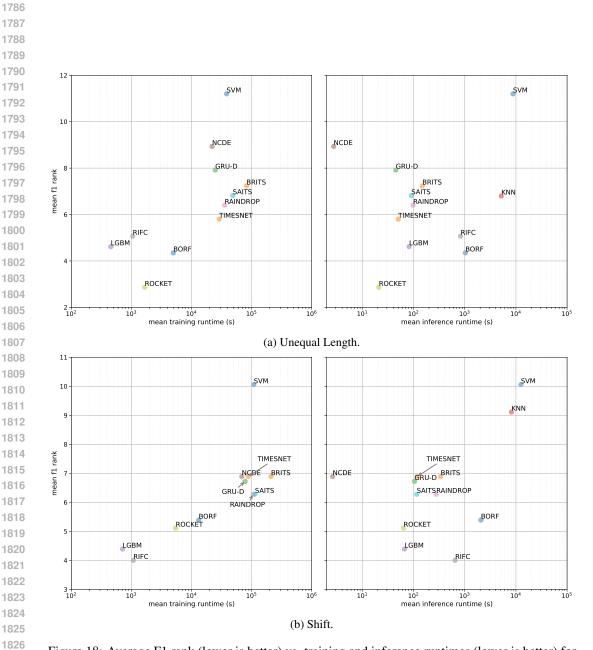


Figure 18: Average F1 rank (lower is better) vs. training and inference runtimes (lower is better) for subsets of datasets. Best values are on the bottom-left of each plot.

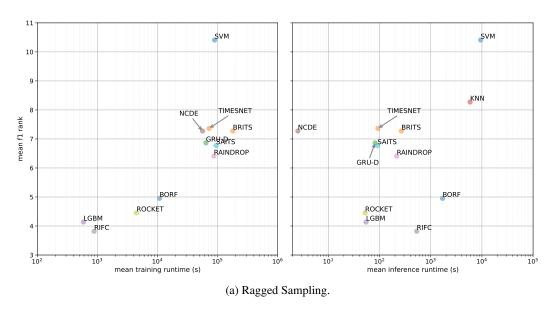


Figure 19: Average F1 rank (lower is better) vs. training and inference runtimes (lower is better) for subsets of datasets. Best values are on the bottom-left of each plot.

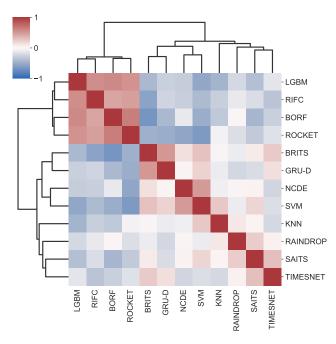


Figure 20: F1 rank correlation between models. Models are hierarchically clustered using average linkage applied to the rank correlation matrix. Positive correlations indicate that models tend to perform similarly across datasets, reflecting comparable strengths or weaknesses. Negative correlations suggest that models excel on different datasets, revealing complementary behaviors or distinct inductive biases.

Table 7: F1 score on the test set for each dataset and each classifier. The average of 3 runs is taken for methods that highly depend upon initialization, i.e., all approaches besides BORF, KNN, LGBM, and SVM. Missing values are due to exceeding memory or maximum runtime. The best values for each dataset are in bold. 

F BRITS	N LGBM NCDE F		ROCKET	SAITS	SVM	TIMESNET
$0.33\pm0.01$ $0.28\pm0.09$ $0.31$	<b>0.41</b> ±0.02		$0.17\pm0.00$	$0.30\pm0.03$	0.32	$0.31\pm0.01$
$0.80  0.65\pm0.00  0.68\pm0.03  0.80  0.80 \ 0.80  0.81$	$0.43\pm0.11  0.64\pm0.05$	$0.88\pm0.02$	0.90±0.04	$0.59\pm0.03$	0.07	$0.61\pm0.01$
0.78+0.02 0.49+0.31 0.34	0.69+0.00		0.9±0.01 0.96±0.00	$0.73\pm0.02$	0.05	0.86+0.02
$0.99\pm0.00$ $0.63\pm0.10$ $0.36$	$0.81 \pm 0.00$		<b>0.99</b> ±0.01	$0.99\pm0.00$	0.10	0.99±0.00
$0.96\pm0.05$ $0.64\pm0.30$ <b>0.98</b>		$0.94\pm0.02$	$0.98\pm0.00$	$0.94\pm0.05$	89.0	$0.95\pm 0.02$
$0.52\pm0.02$ $0.46\pm0.07$ $0.44$	$0.29\pm0.12$		$0.54\pm0.03$	$0.43\pm0.06$	0.23	$0.20\pm0.02$
$0.72\pm0.07$ $0.71\pm0.07$ <b>0.90</b>	$0.51 \pm 0.04$	_	$0.34\pm0.00$	$0.57\pm0.10$	0.85	$0.42\pm0.03$
$0.93\pm0.02$ $0.63\pm0.27$ <b>0.97</b>		_	$0.42\pm0.00$	$0.96\pm0.01$	96.0	$0.63\pm 0.02$
$0.47\pm0.04$ $0.24\pm0.08$ $0.33$	$0.23\pm 0.02$	_	$0.66\pm0.02$	$0.41\pm 0.02$	0.04	$0.50\pm0.04$
$0.32\pm0.05$ $0.40\pm0.08$ $0.26$	_		$0.57\pm0.05$	$0.24\pm0.25$	0.13	$0.45\pm0.03$
$0.22\pm0.02$ $0.06\pm0.02$ $0.14$	•		$0.48 \pm 0.03$	$0.27\pm0.11$	0.01	$0.32\pm0.03$
$0.27\pm0.06$ $0.23\pm0.14$ $0.75$	_		$0.89\pm0.02$	$0.75\pm0.02$	0.16	$0.73\pm0.06$
$0.31\pm0.04$ $0.56\pm0.24$ $0.74$	Ī	10 $0.76\pm0.05$	$0.85 \pm 0.05$	$0.54\pm0.10$	0.43	$0.58\pm 0.04$
1	<b>0.52</b> ±0.29	$0.07 \pm 0.02$	$0.31\pm0.15$	1	,	1
$0.09\pm0.09$ $0.11\pm0.04$ $0.65$	$0.13\pm0.05$ $0.44\pm0.08$		$0.70\pm0.01$	$0.33\pm0.09$	0.11	$0.44\pm0.00$
$0.18\pm0.07$ $0.13\pm0.07$ $0.65$	$0.26\pm0.01$ $0.46\pm0.04$		$0.70\pm0.02$	$0.43\pm0.10$	0.13	$0.44\pm0.02$
$0.17\pm0.09$ $0.10\pm0.04$ $0.62$	_	_	$0.69\pm0.01$	$0.19\pm0.12$	90.0	$0.33\pm0.02$
$0.65\pm0.01$ $0.61\pm0.00$ -	_	_	$0.53\pm0.00$	$0.13\pm0.07$	0.02	$0.60\pm0.00$
$0.96\pm0.00$ $0.96\pm0.01$	$0.57\pm0.02$ $0.94\pm0.02$		$0.94\pm0.01$	$0.96\pm0.01$	0.47	$0.97 \pm 0.01$
$0.28\pm0.20$ $0.21\pm0.14$ $0.02$	_		$0.02\pm0.01$	$0.26\pm0.06$	0.02	$0.27\pm0.03$
$0.42\pm0.11$ $0.35\pm0.00$ $0.35$			$0.35\pm0.00$	$0.46\pm0.10$	0.35	$0.42\pm0.11$
$0.92\pm0.00$ $0.92\pm0.01$ $0.88$	_	_	$0.94\pm0.00$	$0.67\pm0.34$	0.44	$0.93\pm0.01$
$0.46\pm0.00$ $0.61\pm0.02$ $0.12$	$0.49\pm0.01$ $0.56\pm0.02$	_	$0.47\pm0.01$	$0.55\pm 0.01$	0.46	$0.56\pm0.01$
$0.49\pm0.00$ $0.55\pm0.02$ -	- 0.69±0.01	_	$0.71\pm0.01$	$0.72\pm0.00$	0.05	$0.71\pm0.01$
1	1	$0.37 \pm 0.32$	$0.66 \pm 0.10$	1	ı	ı
<b>0.78</b> ±0.00 <b>0.78</b> ±0.00 <b>0.78</b>	$0.57\pm0.29$ $0.48\pm0.26$	$0.40\pm0.00$	$0.40\pm0.00$	$0.65\pm 0.22$	0.40	$0.43\pm0.05$
$0.34\pm0.03$ $0.26\pm0.14$ $0.61$	$0.30\pm0.02$ $0.46\pm0.34$	34 0.60±0.12	$0.73\pm0.00$	$0.65\pm 0.07$	0.16	$0.57\pm0.06$
$0.36\pm0.04$ $0.20\pm0.10$ $0.64$	$0.28\pm0.01$ $0.53\pm0.03$	$0.47\pm0.02$	$0.85\pm0.01$	$0.39\pm0.06$	0.20	$0.45\pm 0.02$
76.0 76.0 0.99±0.00 0.99±0.00 86	$0.74\pm0.01$ $0.98\pm0.00$	$0.65\pm0.42$	$0.98\pm0.00$	$0.95\pm 0.01$	0.62	$0.99\pm0.00$
$17  0.48 \pm 0.15  0.49 \pm 0.17  0.38  0.42$	$0.33\pm0.08$ $0.40\pm0.15$	15 <b>0.82</b> ±0.04	$0.80\pm0.08$	$0.27\pm0.02$	0.26	$0.24\pm0.06$
$0.34\pm0.05$ $0.38\pm0.08$ $0.77$	$0.12\pm0.05$ $0.46\pm0.04$	$0.75\pm0.08$	$0.88 \pm 0.02$	$0.57 \pm 0.11$	0.13	$0.49\pm0.09$
$0.23\pm0.00$ $0.23\pm0.00$ -	$0.33\pm0.01$ $0.25\pm0.02$	_	$0.57\pm0.01$	$0.25\pm 0.01$	ı	$0.25\pm0.01$
<b>77</b> 0.50±0.08 0.60±0.01 0.88 0.94	$0.63\pm0.08$ $0.65\pm0.02$	0.90±0.05	$0.94\pm0.02$	$0.62 \pm 0.03$	0.40	$0.59\pm0.01$

# E COMPLEXITY AND PIPELINE COSTS

 In this section, we analyze the efficiency of our proposed library and the benchmarked classifiers in terms of both time and space complexity. The end-to-end cost of the classification pipeline consists of three components: (i) loading the dataset from disk into memory, (ii) converting it into a dense representation, and (iii) running the models. The last component is an *external* cost, since pyrregular wraps existing state-of-the-art classifiers from other libraries, and is reported separately in Table 9.

In Table 8 we report instead the *internal* costs for datasets with a size greater than 10MB. The first two columns of report the empirical times needed for dataset loading and conversion. Theoretically, the dominant cost arises when converting the sparse COO representation into dense form, which requires ranking the timestamps (Section 4). This amounts to sorting within each time series, leading to a complexity that scales linearly with the number of time series and log-linearly with the number of non-null observations per series. Thus, in practice, runtimes are efficient: for example, P19 takes less than 3 seconds end-to-end, while the largest dataset, PA2, is converted in under one minute.

The third and fourth columns of Table 8 compare disk usage of our proposed array format with that of the raw data. In most cases, the proposed format either matches or substantially reduces disk requirements. For instance, GS decreases from 0.24GB in raw form to 0.09GB with our approach, while the reduction is even more pronounced for TA, which shrinks from 1.81GB to only 0.08GB. These reductions are especially valuable for large-scale datasets where disk I/O is a bottleneck.

The last three columns of Table 8 detail the memory footprint of different representations. The sparse COO representation incurs a cost of four times the number of non-null observations, accounting for the storage of coordinates and values. Conversion into a minimally ragged dense format leads to a worst-case memory complexity of  $O(n \times d \times T)$ , where  $T = \max_i^n(T_i)$  is the longest series length. If the dataset is instead expanded into a fully ragged dense array, the worst-case complexity becomes  $O(n \times d \times T)$ , which grows quickly with irregularity. Empirical results illustrate these trends. For example, on PA2, the sparse representation required only 3.93GB, compared to 5.33GB for a minimally ragged dense format. The largest savings are seen in highly irregular datasets: for TA, the sparse format used 0.34GB, while the fully ragged dense array would require over 4TB of memory, an impractical cost.

Table 8: Loading and conversion times (in seconds) for datasets using the proposed array format, along with disk size consumption (GB) compared to the raw data. Memory usage (GB) of the sparse representation is also reported relative to dense alternatives.

	tin	ne (s)	disk siz	ze (GB)		memory (GB)						
	Loading	Conversion	ours	raw	ours	dense w/o raggedness	dense with raggedness					
ABF	0.03	0.06	$\sim$ 0.00	0.01	$\sim$ 0.00	~0.00	0.81					
AOC	0.09	0.12	0.01	$\sim$ 0.00	0.02	0.01	0.01					
APT	0.30	0.42	0.02	0.02	0.08	0.11	0.11					
ARC	0.21	0.28	0.01	0.01	0.05	0.14	0.14					
CT	0.12	0.16	0.01	0.01	0.03	0.01	0.01					
GS	1.25	3.62	0.09	0.24	0.29	8.58	377.15					
IW	6.93	13.01	0.36	0.31	2.00	1.64	1.64					
LPA	0.07	0.10	$\sim\!\!0.00$	0.02	0.01	0.07	4.02					
MI3	0.13	0.01	$\sim\!\!0.00$	0.04	$\sim\!$ 0.00	$\sim$ 0.00	0.03					
P12	0.35	0.65	0.01	0.08	0.1	0.45	6.35					
P19	0.96	2.08	0.03	0.24	0.31	3.41	3.43					
PA2	13.46	21.35	0.83	1.61	3.93	5.33	21.47					
PL	0.06	0.07	$\sim 0.00$	$\sim 0.00$	0.01	0.01	0.01					
SAD	0.49	0.65	0.02	0.02	0.14	0.08	0.08					
SE	0.15	0.17	0.01	0.06	0.03	0.02	0.84					
TA	1.49	2.34	0.08	1.81	0.34	0.22	4135.02					
VE	0.05	0.08	$\sim$ 0.00	0.01	0.01	0.01	0.17					

Table 9: Total runtime (seconds) for each dataset and each classifier. The average of 3 runs is taken for methods that highly depend upon initialization, i.e., all approaches besides BORF, KNN, LGBM, and SVM. Missing values are due to exceeding memory or maximum runtime. The best values for each dataset are in bold. 

TIMESNET	65±12	$203\pm29$	$2601 \pm 438$	$21283\pm1363$	$18775\pm 1290$	$1989 \pm 292$	$232\pm 30$	$64\pm 12$	$310\pm 378$	$586 \pm 104$	$1078\pm 266$	$634\pm 26$	$575\pm 165$	$797 \pm 149$	1	$1202\pm373$	$1151\pm 198$	$1194\pm225$	$5825 \pm 1060$	$337 \pm 42$	$35634\pm26815$	$81 \pm 12$	$489 \pm 90$	$7817\pm1145$	$51343\pm6196$	1	$21\pm3$	$197\pm 15$	$7031\pm1167$	$6991 \pm 936$	$3122\pm594$	$239 \pm 34$	92470±29583	$2339 \pm 196$
SVM	20	11	1196	14178	8021	2898	22	7	7	102	103	103	26	105	,	307	311	306	279477	41	66	4	950	5764	144163	,	7	∞	2666	17106	2019	10	ı [	926
SAITS																											$39 \pm 25$	$365 \pm 63$	$26156 \pm 2239$	$6956 \pm 536$	$59800 \pm 14909$	$325{\pm}20$	$197788 \pm 30354$	$10987 \pm 6396$
ROCKET	$1\pm 0$	$1\pm 0$	$23\pm1$	$38\pm0$	$52{\pm}1$	$191 \pm 3$	7±0	$0\pm0$	$0\pm0$	$25\pm1$	$26\pm1$	$28\pm1$	$5\pm0$	$5\pm0$	$26010\pm136$	$17\pm1$	$17\pm0$	$17\pm1$	$118 {\pm} 2$	$0 \mp 0$	$27{\pm0}$	$0\pm0$	$142\pm 1$	$22{\pm}1$	$206 {\pm} 2$	$22963\pm92$	$0\pm0$	$1\pm 0$	$40\pm1$	$103\pm1$	$39\pm1$	$2\pm 0$	<b>350</b> ±9	0∓ <b>9</b>
RIFC	$5\pm0$	$2\pm 0$	$21{\pm}1$	78∓6	<b>68</b> ±5	$72{\pm}4$	$2_{\pm 0}$	$1\pm0$	$1\pm0$	<b>4</b> ±0	<b>4</b> ±0	<b>4</b> ±0	$3 \pm 0$	$3\pm 0$	$1826\pm 22$	<b>7</b> ±0	$7\pm 0$	$7\pm0$	$39775\pm 851$	$35\pm 2$	$46\pm1$	$0\mp9$	$18_{\pm 1}$	$1907 \pm 34$	$9243\pm919$	$1179 \pm 97$	$2\pm 0$	$1\pm 0$	$11\pm 0$	$675\pm53$	<b>0</b> ∓6	$1\pm 0$	$1305\pm 19$	8∓0
RAINDROP	$16\pm 2$	$23\pm4$	94±2	$13014\pm21270$	$88376\pm7613$	$472\pm 84$	$39\pm9$	$19\pm 9$	$40\pm 37$	$76\pm 19$	$156\pm 37$	$80\pm18$	$6 \pm 95$	$50\pm 9$	ı	85±5	<b>78</b> ±2	<b>97</b> ±3	$5937 \pm 1596$	89±7	$10318\pm17149$	$14\pm 2$	$550 \pm 207$	$7582 \pm 12025$	$294468 \pm 72194$	ı	$9\pm1$	$30\pm 15$	$5162 \pm 8392$	$1966 \pm 398$	$20918\pm2925$	25±5	$20576 \pm 7435$	$108 \pm 27$
NCDE	$1290\pm2051$	$93\pm 24$	$92\pm 28$	$180 \pm 5$	$175\pm60$	$243\pm 85$	$111\pm 29$	$111\pm 68$	$243\pm 141$	$158\pm 42$	$188 \pm 100$	$136 \pm 39$	$91\pm 27$	$90\pm 25$	$4815\pm6530$	$52\pm 25$	$59\pm11$	$87\pm 24$	$38767 \pm 8867$	$92\pm27$	$186 \pm 10$	$76\pm 16$	$205 \pm 85$	$1204\pm205$	1	ı	$197\pm 123$	$110\pm 28$	$106\pm 2$	$930\pm 143$	$123\pm 31$	$115\pm 69$	$10592\pm6509$	$115\pm 49$
LGBM	2	12	139	302	14 4	613	32	1	1	346	364	440	61	71	1358	373	374	230	4533	133	278	3	306	239	751	4004	_	-	452	87	17	28	352	32
KNN	10	∞	1318	32482	26775	2563	4	$\epsilon$	$\epsilon$	95	94	93	82	79		942	939	879	,	48	447	3	816	40226	,		7	11	4017	17309	164	12	- (	362
GRU-D	$33\pm6$	$121\pm 24$	$1274\pm503$	$9416 \pm 3823$	$16236\pm5490$	$1542\pm719$	$90 \pm 24$	$41\pm 6$	$67 \pm 39$	$477 \pm 163$	$466 \pm 86$	$467 \pm 222$	$338\pm 14$	$998\pm53$	1	$430\pm133$	$612\pm77$	$533\pm67$	$1460 \pm 62$	$51\pm0$	$3854 \pm 303$	$41\pm 10$	$122\pm13$	$4182 \pm 948$	$45315\pm6110$	1	$50 \pm 69$	$245 \pm 142$	$4615\pm 3922$	$1829 \pm 133$	$1893\pm 120$	$267 \pm 29$	$46781\pm20524$	$1311\pm589$
BRITS	$230\pm 8$	$5952\pm510$	$5572\pm906$	$47501 \pm 10144$	$246725 \pm 140403$	$18061\pm4779$	$1989 \pm 478$	$1063\pm215$	$1526 \pm 904$	$10400\pm 2849$	$15746 \pm 1939$	$7442\pm1118$	$6802\pm5500$	$6579\pm1563$	1	$4846\pm 1652$	5314±447	$6439\pm 3059$	$7067 \pm 138$	$289 \pm 79$	$31534 \pm 11204$	$1421 \pm 335$	$6230\pm7$	$72508\pm 32359$	244447±112157		$784 \pm 165$	$2182\pm670$	S		$80101 \pm 34115$	$1821 \pm 119$	$861475 \pm 86400$	$76064\pm5836$
BORF	16	18	295	1000	629	300	20	12	12	18	20	22	25	27	5718	62	54	63	31587	23	411	12	56	5455	30858	77647	9	∞	108	8487	204	6	16825	127
	ABF	AN	AOC	APT	ARC	CI	DD	DG	DW	GM1	GM2	GM3	GP1	GP2	GS	ďΧ	ĞΥ	ЗS	ΜI	λſ	$\Gamma$ PA	MI3	MP	P12	P19	PA2	PGE	PGZ	PL	SAD	SE	SGZ	TA	ΛE

# F ARRAY STRUCTURES

We report a summary of the main formats used to represent regular and irregular time series data in the literature in Table 10.

Table 10: Overview of the main formats used to represent regular and irregular time series data in the literature, categorized by tensor type. The table details the underlying data structures (classes), the software libraries that implement them, their usage across the time series libraries considered in this study, and their support for timestamps and tensor operations.

Type	Format	Library	Class	Usage	Timestamps	Tensor Ops.
		numpy	Array	aeon	X	/
		numpy	Array	sktime	X	✓
se		numpy	Array	tslearn	X	✓
Dense	3D Tensor	numpy	MaskedArray	-	X	✓
Ω		jax	Array	diffrax	<b>/</b> *	✓
		tensorflow	Array	-	X	✓
		torch	Tensor	pypots	X	✓
		awkward	AwkwardArray	_	X	✓
Ragged		tensorflow	RaggedTensor	-	X	✓
88	3D Tensor	torch	NestedTensor	-	X	✓
$R_{\hat{s}}$		zarr	RaggedArray	-	X	✓
		pyarrow	ListArray	-	X	✓
		sparse	GCXS	_	Х	✓
Sparse	3D Tensor	sparse	DOK	-	X	✓
S		sparse	COO	-	×	✓
	Nested List	python	List[Array]	aeon	Х	Х
Other	3D tensor**	xarray	Dataset	-	✓	✓
Ö	Long	pandas	DataFrame	sktime	✓	X
	MultiIndex	pandas	DataFrame	sktime	✓	X

<sup>\*</sup> only as a separate channel

# G EXTENDING PYRREGULAR TO OTHER TASKS

As noted in Section 6, our framework is designed to extend naturally to several additional tasks beyond classification. In particular, we highlight regression, forecasting, and anomaly detection, which are already supported at the representation level and require only minor adjustments to dataset metadata or the inclusion of auxiliary variables.

**Regression.** This task involves predicting continuous outcomes and is directly supported by our framework. Examples include SAPS-I (Simplified Acute Physiology Score) in PhysioNet 2012 or raw productivity in the Garment dataset.

**Forecasting.** Here the objective is to predict future values of a time series given its history. We plan to introduce a static variable with a cutoff point to indicate the train/test split, and to extend the accessor method to provide users with a straightforward mechanism for performing this split.

**Anomaly detection.** This task aims to identify unusual or irregular patterns in the data. Since anomalies may have the same shape as the underlying dataset, they cannot be indicated via static variables. Instead, leveraging the support for additional data arrays in xarray, we will represent anomalies using sparse binary masks that flag anomalous regions in the time series.

**Model support.** We also plan to support a set of representative models for such tasks. A non-comprehensive list includes recent work introducing dynamic graph networks for medical data (Luo et al., 2024), image-based transformers for irregular series (Li et al., 2023), channel harmony strategies (Liu et al., 2025), graph neural flows (Mercatali et al., 2024), temporal graph ODEs (Gravina et al., 2024), state space models (Gu et al., 2022), and patching graph neural networks for forecasting (Zhang et al., 2024), to name a few.

<sup>\*\*</sup> with additional tensors for static variables

# H QUICK GUIDE

2106

21072108

2109

211021112112

2113 2114

21152116

211721182119

2120

212121222123

2124

2125

212621272128

2129

213021312132

21332134

2135

21512152

2153

Extensive documentation and examples were removed from the Supplementary Materials due to double-blind constraints. Thus, we provide a quick start guide and simple workflow notebooks below.

```
pip install pyrregular[models]
```

#### H.1 LIST DATASETS

If you want to see all the datasets available, you can use the list\_datasets function:

```
from pyrregular import list_datasets
df = list_datasets()
```

#### H.2 LOAD A DATASET

To load a dataset, you can use the load\_dataset function. For example, to load the "Garment" dataset, you can do:

```
from pyrregular import load_dataset
df = load_dataset("Garment.h5")
```

### H.3 CLASSIFICATION

To use the dataset for classification, you can just "densify" it:

```
2136
       from pyrregular import load_dataset
2137
2138
       df = load_dataset("Garment.h5")
2139
       X, _ = df.irr.to_dense()
2140
       y, split = df.irr.get_task_target_and_split()
2141
2142
      X_train, X_test = X[split != "test"], X[split == "test"]
2143
      y_train, y_test = y[split != "test"], y[split == "test"]
2144
2145
       # We have ready-to-go models from various libraries:
       from pyrregular.models.rocket import rocket_pipeline
2146
2147
      model = rocket_pipeline
2148
      model.fit(X_train, y_train)
2149
      model.score(X_test, y_test)
2150
```

The dataset can be also easily used in pytorch

```
2160
2161
2162
2163
2164
           Notebook: Basic Workflow
2165
2166
      [1]: import xarray as xr
2167
2168
           List available datasets
2169
2170
           To view available datasets, you can use the list_datasets function.
2171
      [2]: from pyrregular import list_datasets
2172
2173
      [3]: print(list_datasets())
2174
2175
           ['Abf.h5', 'AllGestureWiimoteX.h5', 'AllGestureWiimoteY.h5',
           'AllGestureWiimoteZ.h5', 'Animals.h5', 'AsphaltObstaclesCoordinates.h5',
2176
           \verb|'AsphaltPavementTypeCoordinates.h5', |'AsphaltRegularityCoordinates.h5'|,
2177
           'CharacterTrajectories.h5', 'DodgerLoopDay.h5',
2178
           'DodgerLoopGame.h5', 'DodgerLoopWeekend.h5', 'Garment.h5',
2179
           'GeolifeSupervised.h5', 'GestureMidAirD1.h5', 'GestureMidAirD2.h5',
2180
           'GestureMidAirD3.h5', 'GesturePebbleZ1.h5', 'GesturePebbleZ2.h5',
2181
           'JapaneseVowels.h5', 'Ldfpa.h5', 'MelbournePedestrian.h5', 'Mimic3.h5',
2182
           'PLAID.h5', 'Pamap2.h5', 'Physionet2012.h5', 'Physionet2019.h5',
2183
           'PickupGestureWiimoteZ.h5', 'Seabirds.h5', 'ShakeGestureWiimoteZ.h5',
2184
           'SpokenArabicDigits.h5', 'Taxi.h5', 'Vehicles.h5']
2185
2186
           Loading the dataset from the online repository
2187
           Loading a dataset is as from the online repo is as simple as calling the load_dataset function with
2188
           the dataset name.
2189
      [4]: from pyrregular import load_dataset
2190
2191
           ds = load_dataset("Garment.h5")
2192
2193
           The dataset is loaded as an xarray dataset. The dataset is saved in the default os cache directory,
2194
           which can be found with:
2195
           import pooch
2196
           print(pooch.os_cache("pyrregular"))
2197
           You can also use xarray to directly load a local file. In this case, you have to specify our backend
2198
           as pyrregular in the engine argument.
2199
2200
           import xarray as xr
           ds = xr.load_dataset("path/to/file.h5", engine="pyrregular")
2201
2202
           You can view the underlying DataArray by calling the data variable.
2203
2204 [65]: da = ds.data
2205
     [66]: da
2206
2207
2208
2209
2210
2211
2212
```

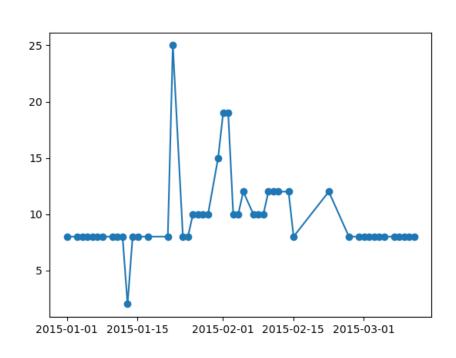
```
2214
2215
2216
2217
2218
2219 [66]: <xarray.DataArray 'data' (ts_id: 24, signal_id: 9, time_id: 59)> Size: 329kB
           <COO: shape=(24, 9, 59), dtype=float64, nnz=10267, fill_value=nan>
2220
           Coordinates:
2221
               day
                                        (time_id) <U9 2kB 'Thursday' ... 'Wednesday'
2222
               department
                                        (ts_id) <U9 864B 'finishing' ... 'sweing'
2223
               productivity_binary
                                        (ts_id) int32 96B 1 0 1 1 1 1 1 1 ... 1 1 0 0 0 0 1
2224
               productivity_class
                                        (ts_id) <U4 384B 'high' 'low' ... 'low' 'high'
2225
                                        (ts_id) float32 96B 0.8126 0.6283 ... 0.7005 0.7503
               productivity_numerical
2226
                                        (time_id) <U8 2kB 'Quarter1' ... 'Quarter2'</pre>
               quarter
2227
             * signal_id
                                        (signal_id) <U21 756B 'idle_men' ... 'wip'
               split
                                        (ts_id) <U5 480B 'train' 'train' ... 'train' 'train'
2228
                                        (ts_id) int32 96B 1 10 11 12 2 3 4 ... 3 4 5 6 7 8 9
               t.eam
2229
             * time_id
                                        (time_id) datetime64[ns] 472B 2015-01-01T01:00:00...
2230
             * ts_id
                                        (ts_id) <U12 1kB 'finishing_1' ... 'sweing_9'
2231
           Attributes:
2232
                           2024-12-04T21:50:44.408790-12:00
               _fixed_at:
2233
               _is_fixed:
2234
                           [Abdullah Al Imran, Md Shamsur Rahim, Tanvir Ahmed]
               author:
2235
                           {'default': {'task': 'classification', 'split': 'split', 'tar...
               configs:
2236
               license:
                           CC BY 4.0
2237
               source:
                           https://archive.ics.uci.edu/dataset/597/productivity+predicti...
               title:
                           Productivity Prediction of Garment Employees
2238
2239
    [67]: # the shape is (n_time_series, n_channels, n_timestamps)
2240
           da.shape
2241
2242
    [67]: (24, 9, 59)
2243
2244 [68]: # the array is stored as a sparse array
2245
           da.data
2246
2247 [68]: <COO: shape=(24, 9, 59), dtype=float64, nnz=10267, fill_value=nan>
2248
     [69]: # dimensions contain the time series ids, signal ids and timestamps
2249
           da.dims
2250
2251 [69]: ('ts_id', 'signal_id', 'time_id')
2252
2253 [70]: # e.g., these are the time series ids
           da["ts_id"].data
2254
2255
     [70]: array(['finishing_1', 'finishing_10', 'finishing_11', 'finishing_12',
2256
                  'finishing_2', 'finishing_3', 'finishing_4', 'finishing_5',
2257
                  'finishing_6', 'finishing_7', 'finishing_8', 'finishing_9',
2258
                  'sweing_1', 'sweing_10', 'sweing_11', 'sweing_12', 'sweing_2',
2259
                  'sweing_3', 'sweing_4', 'sweing_5', 'sweing_6', 'sweing_7',
2260
                  'sweing_8', 'sweing_9'], dtype='<U12')
2261
2262
2263
2264
2265
2266
2267
```

```
2268
2269
2270
2271
2272
2273 [72]: # there are also static variables, such as the class
           da["productivity_binary"].data
2274
     [72]: array([1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0,
2276
                  0, 1], dtype=int32)
2277
2278 [74]: # the train/test split
2279
           da["split"].data
2280
2281 [74]: array(['train', 'train', 'test', 'train', 'train', 'test', 'train',
                  'train', 'train', 'test', 'train', 'train', 'test', 'train',
2282
                  'train', 'test', 'train', 'train', 'train', 'train', 'test',
2283
                  'train', 'train', 'train'], dtype='<U5')
2284
2285 [75]: # all the coordinates can be accessed via the `coords` variable
2286
           da.coords
2287
2288 [75]: Coordinates:
               day
                                        (time_id) <U9 2kB 'Thursday' ... 'Wednesday'
2289
               department
                                        (ts_id) <U9 864B 'finishing' ... 'sweing'
2290
               productivity_binary
                                        (ts_id) int32 96B 1 0 1 1 1 1 1 1 ... 1 1 0 0 0 0 1
2291
               productivity_class
                                        (ts_id) <U4 384B 'high' 'low' ... 'low' 'high'
               productivity_numerical (ts_id) float32 96B 0.8126 0.6283 ... 0.7005 0.7503
2293
                                        (time_id) <U8 2kB 'Quarter1' ... 'Quarter2'
               quarter
2294
                                        (signal_id) <U21 756B 'idle_men' ... 'wip'
             * signal_id
2295
                                        (ts_id) <U5 480B 'train' 'train' ... 'train' 'train'
               split
2296
                                        (ts id) int32 96B 1 10 11 12 2 3 4 ... 3 4 5 6 7 8 9
               team
2297
                                        (time_id) datetime64[ns] 472B 2015-01-01T01:00:00...
             * time_id
                                        (ts_id) <U12 1kB 'finishing_1' ... 'sweing_9'
             * ts_id
2298
2299
     [76]: # metadata contains informations about the datasets and tasks
2300
           da.attrs
2301
2302 [76]: {'_fixed_at': '2024-12-04T21:50:44.408790-12:00',
2303
            '_is_fixed': True,
2304
            'author': [Abdullah Al Imran, Md Shamsur Rahim, Tanvir Ahmed],
            'configs': {'default': {'task': 'classification',
2305
              'split': 'split',
2306
              'target': 'productivity_binary'},
2307
             'regression': {'task': 'regression',
2308
              'split': 'split',
2309
              'target': 'productivity_numerical'}},
2310
            'license': 'CC BY 4.0',
2311
            'source': 'https://archive.ics.uci.edu/dataset/597/productivity+prediction+of+g
2312
           arment+employees',
2313
            'title': 'Productivity Prediction of Garment Employees'}
2314
2315
2316
2317
2318
2319
2320
2321
```

```
2322
2323
2324
2325
2326
          Data Handling and Plotting
2327
2328
          Data can be accessed with standard xarray methods.
2329
2330 [77]: import matplotlib.pyplot as plt
           import numpy as np
2331
2332
     [78]: # the first time series
2333
           da[0]
2334
2335 [78]: <xarray.DataArray 'data' (signal_id: 9, time_id: 59)> Size: 9kB
2336
           <COO: shape=(9, 59), dtype=float64, nnz=392, fill_value=nan>
           Coordinates:
2337
                                        (time_id) <U9 2kB 'Thursday' ... 'Wednesday'
               day
2338
               department
                                        <U9 36B 'finishing'
2339
               productivity_binary
                                        int32 4B 1
2340
               productivity_class
                                        <U4 16B 'high'
2341
               productivity_numerical
                                        float32 4B 0.8126
2342
                                        (time_id) <U8 2kB 'Quarter1' ... 'Quarter2'
               quarter
2343
                                        (signal_id) <U21 756B 'idle_men' ... 'wip'
             * signal_id
2344
               split
                                        <U5 20B 'train'
2345
                                        int32 4B 1
               team
                                        (time_id) datetime64[ns] 472B 2015-01-01T01:00:00...
             * time_id
2346
                                        <U12 48B 'finishing_1'
               ts_id
2347
           Attributes:
2348
               _fixed_at: 2024-12-04T21:50:44.408790-12:00
2349
               _is_fixed: True
2350
                            [Abdullah Al Imran, Md Shamsur Rahim, Tanvir Ahmed]
               author:
2351
               configs:
                            {'default': {'task': 'classification', 'split': 'split', 'tar...
2352
               license:
                           CC BY 4.0
2353
                           https://archive.ics.uci.edu/dataset/597/productivity+predicti...
               source:
2354
               title:
                           Productivity Prediction of Garment Employees
2355
2356 [79]: # the first channel of the first time series
           da[0, 0]
2357
2358
     [79]: <xarray.DataArray 'data' (time_id: 59)> Size: 784B
2359
           <COO: shape=(59,), dtype=float64, nnz=49, fill_value=nan>
2360
           Coordinates:
2361
                                        (time_id) <U9 2kB 'Thursday' ... 'Wednesday'
               dav
2362
               department
                                        <U9 36B 'finishing'
2363
               productivity_binary
                                        int32 4B 1
               productivity_class
                                        <U4 16B 'high'
2364
               productivity_numerical float32 4B 0.8126
2365
                                        (time_id) <U8 2kB 'Quarter1' ... 'Quarter2'</pre>
               quarter
2366
               signal_id
                                        <U21 84B 'idle men'
2367
               split
                                        <U5 20B 'train'
2368
                                        int32 4B 1
               team
2369
2370
2371
2372
2373
2374
```

```
2376
2377
2378
2379
2380
                                     (time_id) datetime64[ns] 472B 2015-01-01T01:00:00...
2381
            * time_id
              ts_id
                                     <U12 48B 'finishing_1'
2382
          Attributes:
2383
              _fixed_at:
                         2024-12-04T21:50:44.408790-12:00
2384
              _is_fixed: True
2385
              author:
                          [Abdullah Al Imran, Md Shamsur Rahim, Tanvir Ahmed]
2386
                          {'default': {'task': 'classification', 'split': 'split', 'tar...
              configs:
2387
              license:
                         CC BY 4.0
2388
                         https://archive.ics.uci.edu/dataset/597/productivity+predicti...
              source:
2389
                         Productivity Prediction of Garment Employees
              title:
2390
2391 [80]: # to access the underlying sparse vector
          da[0, 0].data
2392
    [80]: <COO: shape=(59,), dtype=float64, nnz=49, fill_value=nan>
2394
2395 [87]: # to access the underlying dense vector
2396
          da[0, 4].data.todense()
8., nan, nan, nan, 8., 25., 8., 8., 10., 10., 10., 10., 15.,
2399
                 19., 19., 10., 10., 12., 10., 10., 10., 12., 12., 12., 12., 8.,
2400
                 nan, nan, nan, nan, 12., nan, nan, nan, 8., 8., 8., 8.,
2401
                 8., 8., 8., 8., 8., 8.]
2402
2403 [89]: # this vector contains a lot of nans, which are the padding necessary to have
2404
           \hookrightarrowshared timestamps w.r.t. the whole dataset
2405
          np.isnan(da[0, 4].data.todense()).sum()
2406
2407 [89]: 10
2408
    [90]: plt.plot(da[0, 4]["time_id"], da[0, 4], marker="o")
2409
2410 [90]: [<matplotlib.lines.Line2D at 0x14eb06990>]
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
```

```
2431
2432
2433
2434
2435
                   25
2436
2437
2438
2439
                   20
2440
2441
2442
                   15
2443
2444
2445
2446
                   10
2447
2448
2449
                    5
2450
2451
2452
2453
                    2015-01-01 2015-01-15
                                               2015-02-01 2015-02-15 2015-03-01
2454
2455
2456
2457
2458 [92]: # using the custom ".irr" accessor, we can filter out the nans to the minimum_
            \hookrightarrowamount possible due to raggedness
2459
           np.isnan(da.irr[0, 4].data.todense()).sum()
2460
2461
     [92]: 0
2462
2463 [93]: plt.plot(da.irr[0, 4]["time_id"], da.irr[0, 4], marker="o")
2465 [93]: [<matplotlib.lines.Line2D at 0x14eb6b230>]
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
```



2511
2512 [94]: # the fourth channel first 10 time series of the dataset, as a heatmap
2513 da.irr[:10, 4].plot()

2514 [94]: <matplotlib.collections.QuadMesh at 0x14dcf3680>

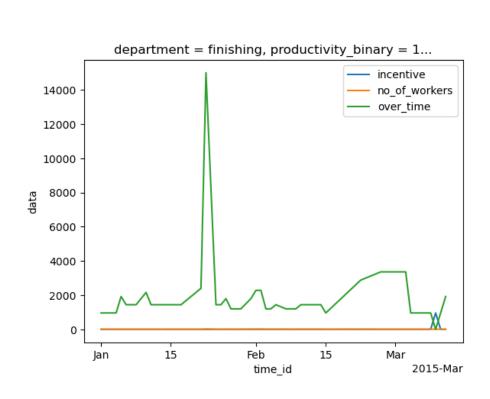
```
signal id = no of workers
    finishing_7
    finishing_6
                                                                                 25
    finishing_5
                                                                                 20
    finishing_4
    finishing_3
ts id
                                                                               - 15 g
    finishing_2
   finishing_12
                                                                                - 10
   finishing_11
                                                                                 5
   finishing_10
    finishing_1
                     08
                          15
                                22
                                        Feb
                                              08
                                                          22
                                                                    08
                                                    15
                                                               Mar
                                         time_id
                                                               2015-Mar
```

```
2566
2567

# plotting some channels
da.irr[0, 2].plot(label=da.coords["signal_id"][2].item())
da.irr[0, 4].plot(label=da.coords["signal_id"][4].item())
da.irr[0, 5].plot(label=da.coords["signal_id"][5].item())
plt.legend()

2571
```

2572[103]: <matplotlib.legend.Legend at 0x16ea32870>



#### Downstream Tasks

The xarray is nice, but not supported by basically any downstream library. Thus, we can convert it into a numpy array.

```
2625
2626[104]: %%time
           # time series data, timestamps
2627
           X, T = da.irr.to_dense(
2628
               normalize_time=True,
                                      # normalize the time index to [0, 1]
2629
           )
2630
2631
          CPU times: user 2.23 s, sys: 79 ms, total: 2.31 s
2632
          Wall time: 2.34 s
2633
    106]: # the shape is (n_time_series, n_channels, n_timestamps), timestamps are
2634
            →returned as a separate channel, for downstream methods that are able to use_
2635
            \hookrightarrow them
2636
           X.shape, T.shape
2637
2638[106]: ((24, 9, 59), (24, 1, 59))
2639
2640
2641
2642
2643
2644
```

```
2646
2647
2648
2649
2650
2651[107]: # static variables
           Z = da.coords.to_dataset()[["split", "productivity_binary"]].to_pandas()
2652
           Z.head()
2653
2654[107]:
                          split productivity_binary department productivity_class \
2655
           ts_id
2656
           finishing_1
                                                    1 finishing
                                                                                 high
                          train
2657
           finishing_10 train
                                                    0 finishing
                                                                                 low
2658
           finishing_11
                         test
                                                    1 finishing
                                                                                 high
2659
           finishing_12 train
                                                    1 finishing
                                                                                 high
2660
           finishing_2 train
                                                    1 finishing
                                                                                 high
2661
                          productivity_numerical team
2662
           ts id
2663
           {\tt finishing\_1}
                                        0.812625
                                                      1
2664
           finishing_10
                                         0.628333
                                                     10
2665
           finishing_11
                                         0.874028
                                                     11
2666
           finishing_12
                                         0.922840
                                                     12
2667
           finishing_2
                                         0.819271
2668
2669[108]: # target and split
           y, split = da.irr.get_task_target_and_split()
2670
2671
2672
           Train-test split
2673
2674[111]: X_train, X_test = X[split != "test"], X[split == "test"]
           y_train, y_test = y[split != "test"], y[split == "test"]
2675
           X_train.shape, y_train.shape, X_test.shape, y_test.shape
2676
2677[111]: ((18, 9, 59), (18,), (6, 9, 59), (6,))
2679
           Classification
2680
           We have several ready-to-use classifiers in the pyrregular package. Be sure to install the required
2681
           dependencies.
2682
2683[118]: from pyrregular.models.rocket import rocket_pipeline
2685[119]: | %%time
           model = rocket_pipeline
2686
           model.fit(X_train, y_train)
2687
           model.score(X_test, y_test)
2688
<sup>2689</sup>[119]: 0.6666666666666666
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
```

## Notebook: Dataset Conversion

#### The "Long Format"

2705 2706

2707 2708

2709 2710

2711

2712

2713 2714

2715

The basic format to convert any dataset to our representation is the long format. The long format is simply a tuple:

```
(time_series_id, channel_id, timestamp, value, static_var_1, static_var_2, ...).
```

If your dataset contains rows that are in this format, you are almost good to go. Else, there will be a little bit of preprocessing to do.

## Case 1. (easy) Your dataset is already in the long format

```
Let's assume for now your dataset is already in this form. Here is a minimal working example.
2716
2717 [28]: import pandas as pd
2718
           import numpy as np
2719
2720 [29]: df = pd.DataFrame(
               {
2721
                   "time_series_id": np.random.choice(["A", "B", "C"], size=100),
2722
                   "channel_id": np.random.choice(["X", "Y", "Z"], size=100),
2723
                   "timestamp": pd.date_range("2023-01-01", periods=100, freq="H"),
2724
                   "value": np.random.randn(100),
2725
               }
2726
           )
2727
           df["labels"] = df["time_series_id"].map(
2728
               {"A": 0, "B": 1, "C": 1}
2729
             # let's say we have labels
           df.head()
2730
2731
            time_series_id channel_id
                                                                value labels
    [29]:
                                                  timestamp
2732
                                     Y 2023-01-01 00:00:00 0.105162
           0
                          В
                                                                             1
2733
                          В
                                     Z 2023-01-01 01:00:00 -0.573337
           1
                                                                             1
2734
           2
                                     X 2023-01-01 02:00:00 -1.973967
                          В
                                                                             1
2735
                          С
           3
                                     Y 2023-01-01 03:00:00 0.656065
                                                                             1
2736
           4
                          Α
                                     Y 2023-01-01 04:00:00 -0.500246
                                                                             0
2737
2738 [30]: # Let's save this dataframe to a CSV file
           df.to_csv("your_original_dataset.csv", index=False)
2739
2740
     [31]: # the csv file can be converted to our format using our interface
2741
2742
           from pyrregular.io_utils import read_csv
2743
           from pyrregular.reader_interface import ReaderInterface
2744
           from pyrregular.accessor import IrregularAccessor
2745
2746
           class YourDataset(ReaderInterface):
2747
2748
```

```
2755
2756
2757
2758
               @staticmethod
2759
               def read_original_version(verbose=False):
2760
                   return read_csv(
2761
                        filenames="your_original_dataset.csv",
2762
                        ts_id="time_series_id",
2763
                        time_id="timestamp",
2764
                        signal_id="channel_id",
2765
                        value_id="value",
2766
                        dims={
2767
                            "ts_id": [
2768
                                "labels"
2769
                            ], # static variable that depends on the time series id
                            "signal_id": [],
2770
                            "time_id": [],
2771
                        },
2772
                        time_index_as_datetime=False,
2773
                        verbose=verbose,
2774
                   )
2775
2776 [32]: da = YourDataset.read_original_version(True)
2777
           da
2778
           Getting dataset metadata: Oit [00:00, ?it/s]
2779
2780
           Reading dataset:
                              0%1
                                            | 0/100 [00:00<?, ?it/s]
2781
     [32]: <xarray.DataArray (ts_id: 3, signal_id: 3, time_id: 100)> Size: 3kB
2782
           <COO: shape=(3, 3, 100), dtype=float64, nnz=100, fill_value=nan>
2783
           Coordinates:
2784
             * time_id
                           (time_id) <U19 8kB '2023-01-01 00:00:00' ... '2023-01-05 03:00...
2785
               labels
                           (ts_id) int64 24B 0 1 1
2786
                           (ts_id) <U1 12B 'A' 'B' 'C'
             * ts_id
             * signal_id (signal_id) <U1 12B 'X' 'Y' 'Z'
2787
2788
          If you don't know if a variable is static, or to which dimension it depends from, you can check it.
2789
2790 [33]: from pyrregular.data_utils import infer_static_columns
2791
           infer_static_columns(df, "time_series_id")
2792
2793
     [33]: ['labels']
2794
2795
           The dataset can be saved with our custom accessor
2796
2797 [34]: da.irr.to_hdf5("your_dataset.h5")
2798
           And then loaded directly with xarray
2799
2800 [35]: import xarray as xr
2801
2802
2803
2804
2805
2806
```

```
2808
2809
2810
2811
2812
2813 [36]: da2 = xr.load_dataset("your_dataset.h5", engine="pyrregular")
           da2
2814
2815
     [36]: <xarray.Dataset> Size: 11kB
2816
           Dimensions:
                          (ts_id: 3, signal_id: 3, time_id: 100)
2817
           Coordinates:
2818
                           (ts_id) int32 12B 0 1 1
               labels
2819
             * signal_id (signal_id) <U1 12B 'X' 'Y' 'Z'
             * time_id
                           (time_id) <U19 8kB '2023-01-01 00:00:00' ... '2023-01-05 03:00...
             * ts_id
2821
                           (ts_id) <U1 12B 'A' 'B' 'C'
2822
           Data variables:
               data
                           (ts_id, signal_id, time_id) float64 3kB <COO: nnz=100,
2823
           fill_value=nan>
2824
2825
          Case 2. Your dataset is not in the long format
2826
2827
```

Let's say you have a 3d numpy array, containing the time series, and a numpy array containing only the labels.

```
2829
    [37]: import numpy as np
2830
2831
           shape = (10, 2, 100) # 10 time series, 2 channels, 100 timestamps
           data = np.full(shape, np.nan)
2833
           mask = np.random.rand(*shape) < 0.35
           data[mask] = np.random.randn(mask.sum())
2835
           labels = np.random.randint(0, 2, shape[0])
2836
2837
           np.save("your_more_complex_dataset.npy", data)
           np.save("your_more_complex_dataset_labels.npy", labels)
2838
2839
           data.shape, labels.shape
2840
```

2841 [37]: ((10, 2, 100), (10,)) 2842

2828

2843

2844

2856 2857

2859 2860 2861 You need only a function that takes the data and the labels, and returns a dataframe in the long format, yielding it row by row.

```
2845
     [38]: def read_your_dataset(filenames):
2846
               data = np.load(filenames["data"])
2847
               labels = np.load(filenames["labels"])
2848
               ts_ids, signal_ids, timestamps = np.indices(shape)
2849
               ts_ids, signal_ids, timestamps = ts_ids.ravel(), signal_ids.ravel(),
            →timestamps.ravel()
2850
2851
               for ts_id, signal_id, timestamp in zip(ts_ids, signal_ids, timestamps):
2852
                   value = data[ts_id, signal_id, timestamp]
2853
                   if np.isnan(value):
2854
                       continue
2855
```

```
2863
2864
2865
2866
                   label = labels[ts_id]
2867
                   yield dict(
2868
                        time_series_id=ts_id,
2869
                        channel_id=signal_id,
2870
                        timestamp=timestamp,
2871
                        value=value,
2872
                        labels=label,
2873
                   )
2874
2875 [39]: from pyrregular.io_utils import read_csv
2876
           from pyrregular.reader_interface import ReaderInterface
           from pyrregular.accessor import IrregularAccessor
2877
2878
           class YourDataset(ReaderInterface):
2879
               Ostaticmethod
2880
               def read_original_version(verbose=False):
2881
                   return read_csv(
2882
                       filenames={
2883
                            "data": "your_more_complex_dataset.npy",
2884
                            "labels": "your_more_complex_dataset_labels.npy",
2885
                       },
2886
                       ts_id="time_series_id",
                       time_id="timestamp",
2887
                       signal_id="channel_id",
2888
                       value_id="value",
2889
                       dims={
2890
                            "ts_id": [
2891
                                "labels"
2892
                            ], # static variable that depends on the time series id
2893
                            "signal_id": [],
2894
                            "time_id": [],
2895
2896
                        reader_fun=read_your_dataset,
2897
                        time_index_as_datetime=False,
                        verbose=verbose,
2898
                        attrs={
2899
                            "authors": "Bond, James Bond", # you can add any attribute you_
2900
            \rightarrow want
2901
                       }
2902
                   )
2903
2904 [40]: da = YourDataset.read_original_version(True)
2905
2906
          Getting dataset metadata: Oit [00:00, ?it/s]
2907
2908
                                            | 0/720 [00:00<?, ?it/s]
          Reading dataset:
                              0%1
2909
2910
2911
2912
2913
2914
2915
```

```
2916
2917
2918
2919
2920
2921 [40]: <xarray.DataArray (ts_id: 10, signal_id: 2, time_id: 100)> Size: 23kB
           <COO: shape=(10, 2, 100), dtype=float64, nnz=720, fill_value=nan>
2922
           Coordinates:
2923
                          (time_id) int64 800B 0 1 2 3 4 5 6 7 \dots 92 93 94 95 96 97 98 99
             * time_id
2924
                          (ts_id) int64 80B 0 0 0 1 1 1 0 1 1 0
               labels
2925
                          (ts_id) <U21 840B '0' '1' '2' '3' '4' '5' '6' '7' '8' '9'
             * ts_id
2926
             * signal_id (signal_id) <U21 168B '0' '1'
2927
           Attributes:
2928
               authors: Bond, James Bond
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
```