

GRU-M: A Joint Impute and Learn Approach for Sequential Prediction under Missing Data

3 Soumen Pachal*

S.PACHAL@TCS.COM

4 Avinash Achar*

ACHAR.AVINASH@TCS.COM

5 Nancy Bhutani†

NANCYBHUTANI96@GMAIL.COM

6 Akash Das†

DAS.AKASH3@TCS.COM

7 *Tata Consultancy Services Research, IITM Research Park, Chennai, India*

8 **Editors:** Vu Nguyen and Hsuan-Tien Lin

9 Abstract

10 Sequential Prediction in presence of missing data is an old research problem. Classically,
 11 researchers have tackled this by imputing data first and then building predictive models.
 12 This 2-stage process is typically prone to errors and to circumvent this, researchers have
 13 provided a variety of techniques which employ a joint impute and learn approach before
 14 prediction. Among these, Recurrent Neural Networks (RNNs) have been very popular given
 15 their natural ability to tackle sequential data efficiently. Existing state-of-art approaches
 16 either (i) do not impute (ii) do not completely factor the available information around a gap,
 17 (iii) ignore position information within a gap and so on. Our approach intelligently addresses
 18 these gaps by proposing a novel architecture which jointly imputes and learns by taking
 19 into account (i) information from either end of the gap (ii) proximity to the left/right-end of
 20 a gap (iii) the length of the gap. In context of this work, prediction means either sequence
 21 classification or forecasting. In this paper, we demonstrate the utility of the proposed
 22 architecture on forecasting tasks. We benchmark against a range of state-of-art baselines
 23 and in scenarios where data is either (a) naturally missing or (b) synthetically masked.

24 **Keywords:** Sequential Prediction; Missing Data; RNN; GRU; Time-series;

25 1. Introduction

26 Sequence based prediction (classification [Xing et al. \(2010\)](#) and forecasting [Lim and Zohren](#)
 27 [\(2021\)](#)), is a classic research problem with numerous applications. Given constant method-
 28 ological enhancements and newer emerging applications, it continues to be very relevant.
 29 Recurrent Neural Networks (RNNs) have been a popular state-of-the-art approach for the
 30 same. Deep RNNs have been highly successful in diverse domains like computer vision,
 31 NLP, time series classification/forecasting, speech and audio processing and so on. RNNs
 32 have exhibited state-of-the-art performance (and much more) in tasks like machine transla-
 33 tion [Sutskever et al. \(2014\)](#), handwriting recognition [Graves et al. \(2009\)](#) etc.

34 When data is partly missing, predictive modelling becomes further challenging. This is
 35 the problem tackled in this paper. Data could be missing due to many factors like sensor
 36 failure, noise (under high noise levels, its better to ignore the data), cost of sensing, mainte-
 37 nance and so on. Sequence prediction under missing data also has a vast literature. RNNs
 38 in particular have also been explored for this. This paper addresses sequence prediction

*. Both authors contributed equally.

†. Both authors contributed equally.

39 under missing data using RNNs in a novel way. Our proposed ideas can be employed for
 40 both classification and forecasting. In the classification context, data can appear as variable
 41 length sequences, where each sequence is mapped to a class. In forecasting, data is typically
 42 in the form of one long sequence. We consider a general multi-step forecasting scenario with
 43 possibly additional exogenous variables to be handled.

44 1.1. Problem & Motivation

45 We consider real-valued time series with missing values. For sequence classification appli-
 46 cations, time-stamps can be either discrete or continuous (real-valued). Here data comes
 47 as multiple sequences, where each sequence is real-valued (possibly vector valued). Further,
 48 each of these sequences can have different lengths. The classification task is to categorize
 49 each sequence into one of many predefined classes. For example, each vector sequence could
 50 indicate essential parameter measurements like blood pressure, heart rate, sugar level etc.
 51 from a patient, while the classification task may mean detecting a certain disease.

52 If the task is forecasting, then the time-stamps are assumed discrete (integers or natural
 53 numbers). The forecasting task considered here involves, predicting one or more endoge-
 54 nous variables over a multi-step forecast horizon, in the presence of possible exogenous
 55 inputs, which influence the evolution of the endogenous variables. Such forecasting tasks
 56 arise in diverse applications. A retailer may be interested in forecasting one or more prod-
 57 ucts sales (modelled as endogenous variables) in the presence of exogenous price variation.
 58 In electricity markets, forecasting power demand (endogenous) by factoring temperature
 59 influences/fluctuations (exogenous) is another example.

60 1.2. Contributions

61 Early approaches for prediction under missing data have restricted themselves to a 2-step
 62 approach where data is first imputed and subsequently a predictive model is built. However
 63 this 2-step approach is very sensitive to the imputation quality and can suffer from excess
 64 errors due to poor imputation. To address this, diverse approaches which avoid an explicit
 65 imputation step have been proposed.

66 RNN approaches form a dominant class of techniques for sequential prediction under
 67 missing data. A small subclass of these techniques build a predictive model by avoiding any
 68 form of imputation during learning by capturing the pattern of missingness directly [Lipton
 69 et al. \(2016\)](#); [Pachal and Achar \(2022\)](#), where the missingness pattern to be captured could
 70 be very complex. On the other hand, most other techniques employ a joint impute and
 71 train strategy. Here, missing gaps in data are imputed during predictive model building, by
 72 utilizing the neighboring data that is available. Our proposed approach falls under this class
 73 of approaches. Some of these approaches ignore information from the right-end of a data gap
 74 [Che et al. \(2016\)](#), while others which consider information from either end, do not explicitly
 75 consider time distances from either side of the data gap [Cao et al. \(2018\)](#). Our approach
 76 proposes a novel joint and learn strategy while addressing the above issues. Specifically,

- 77 • We propose a novel RNN based joint impute and learn strategy which factors both
 78 the (closest) left-end and right-end values of a gap, for imputation.

- 79 • Our novel strategy takes into account the position of the missing point and its distance
80 from both the left and right end of a gap. In the regions closer to the data gap end, it
81 can impute by maximally factoring the influence from the respective end values. For
82 large gaps, it can ignore the right and left end values while imputing in the mid-region.
- 83 • We demonstrate effectiveness of the proposed architecture on real time-series where
84 data is either (i)naturally missing or (ii)artificially masked.

85 2. Related Work

86 Prediction under missing data has a long literature. A natural approach to tackle this is
87 to impute first before doing predictive modeling. Data imputation itself has an extensive
88 literature. The early and simpler methods for data imputation include smoothing, inter-
89 polation, splines [Johnson \(2013\)](#) etc. which capture straightforward patterns in the data.
90 The advanced imputation approaches include spectral analysis [Mondal and Percival \(2010\)](#),
91 matrix factorization [Koren et al. \(2009\)](#), kernel methods [Rehfeld et al. \(2011\)](#), EM algo-
92 rithm [García-Laencina et al. \(2010\)](#), GANs (Generative Adversarial Networks) [Luo et al.](#)
93 [\(2018\)](#), Transformers [Nie et al. \(2024\)](#) etc. Hence, the imputation step can be computationally
94 expensive. Further, imputation methods typically make strong, restrictive assumptions
95 like missing at random, low missing rates and so on which may not hold in practice. The
96 two step sequential process of imputation followed by prediction can suffer from imputation
97 errors which can in-turn render poor predictive models.

98 To circumvent this, literature has seen techniques where data imputation is not sep-
99 arately carried out. When performed, it is coupled with the predictive model building
100 process. There have been a wide range of techniques proposed along these lines. These
101 include RNNs, GANs, Gaussian processes (GP) etc. Gaussian processes essentially do a
102 Bayesian estimation where the predictive distribution given the observed values is inferred.
103 They are a natural data-driven model to consider when data is irregularly sampled with
104 many missing values. There has been some recent work along this direction [Futoma \(2018\)](#);
105 [Li \(2020\)](#) mostly applied in the health care domain (clinical time series). In the rest of the
106 section, we elaborate on RNN techniques only, as our work is RNN based.

107 RNNs have been explored to tackle missing data scenarios even before the deep learning
108 surge [Bengio and Gingras \(1995\)](#); [Tresp and Briegel \(1997\)](#); [Parveen and D. Green \(2001\)](#).
109 [Bengio and Gingras \(1995\)](#) is perhaps the earliest approach which unfolds the RNN to allow
110 the missing values to settle down (converge) while the other inputs are fixed at the observed
111 values. [Parveen and D. Green \(2001\)](#) employs a slightly modified architecture of [Bengio](#)
112 [and Gingras \(1995\)](#). While both use non-gated RNNs, [Bengio and Gingras \(1995\)](#) uses a
113 feedback structure based on a Jordan network while [Parveen and D. Green \(2001\)](#) uses an
114 Elmann network.

115 Among the RNN approaches, a couple of them impute neither before nor during model
116 building. [Lipton et al. \(2016\)](#) encodes the missingness pattern as a simple binary vector
117 indicating the presence or absence of a data point and then trains. A more involved ap-
118 proach of encoding the missingness pattern using two RNN layers in a lossless fashion was
119 proposed in [Pachal and Achar \(2022\)](#). Both these approaches do not bother to impute even
120 during model building. To capture all the complex dependencies based on the patterns of
121 missingness can be hard on these RNN models.

122 There is a recent attention-based RNN approach which first imputes using some padding
 123 (from left) OR some form of interpolation [Cinar et al. \(2018\)](#). During subsequent model
 124 building it doesn't use the imputed values as they are, but assumes some form of weighted
 125 influence from these imputed values on the prediction. These weights are suitably parame-
 126 terized and the associated parameters are learnt parallelly with the predictive model.

127 Among the recent joint impute and learn strategies, GRU-D [Che et al. \(2016\)](#) imputes
 128 a missing point as weighted combination of the overall (time-series) mean and the closest
 129 available point at the left end of a data gap. The weight that controls the influence from
 130 the left end is assumed to decay exponentially with the time distance from the left end. The
 131 decay/weight factor is learnt in conjunction with the predictive model.

132 While GRU-D totally ignores the right-end of a data gap, BRITS is a more recent
 133 approach incorporating influence of both right and left end of a gap towards imputation,
 134 using a bidirectional layer. In the state evolution from left, there is a state decay (inspired
 135 from GRU-D) incorporated which is a function of how far the current event is to the closest
 136 available point to its left. The state at each time-step is non-linearly transformed via a
 137 feed-forward network to capture the input at the next time-step (to the right). Further,
 138 the state of the backward layer is transformed to capture the input at the next step to the
 139 left. If the input at the current time-step is missing, then its predicted to be the average of
 140 the output of the previous step of both the forward and backward layer. In BRITS, equal
 141 importance is given to both forward and backward imputed values (computed by passing the
 142 available information from either end of the gap through bidirectional layer) irrespective of
 143 its position in the gap. In our GRU-M model, we consider the left-end and right-end available
 144 values directly. Further, we have considered a convex combination of left-end, right-end and
 145 mean values where the coefficients are learned along with the other model parameters during
 146 training. BRITS has been recently extended to tackle multivariate time-series with selective
 147 dependencies, using graph neural networks in GRIN [Cini et al. \(2021\)](#).

148 Recently, Bi-GAN, a GAN approach [Gupta et al. \(2021\)](#) was considered using a joint
 149 impute and learn strategy. It also factors both the left and right end of a gap for imputation
 150 similar to BRITS. But unlike BRITS, it adopts a generator-discriminator framework where
 151 the generator is a bidirectional RNN as in BRITS, while training is done in an adversarial
 152 setting jointly with a discriminator which predicts the presence or absence of data. In
 153 addition, Bi-GAN differs from BRITS in the last step where a missing input is modelled as
 154 a weighted average (instead of standard average) of the outputs of the previous step from
 155 the forward and backward layers. These weights are assumed to be a function of the distance
 156 from the left-end/right-end gap.

157 2.1. In Perspective of Other RNN Approaches

158 Unlike GRU-D, our approach (denoted as GRU-M) factors information from both the left
 159 and right end of a data gap in a novel and intelligent fashion. As explained in detail in the
 160 next section, our approach can be viewed as a non-trivial extension of GRU-D. The attention-
 161 based approach of [Cinar et al. \(2018\)](#) while intelligent in learning a weighted influence of
 162 the imputed values during model building, can still be prone to imputation errors due to its
 163 separate imputation step.

164 Our approach uses the closest left and right available inputs directly to impute. BRITS
 165 and Bi-GAN use the closest left and right values in an indirect fashion by (a)passing the
 166 closest left available input through the forward layer (b)passing the closest right available
 167 input through the backward layer and (c) by taking an average of these two arrived values
 168 from either direction. Unlike our method, both are prone to error accumulation as imputa-
 169 tion is carried out sequentially from either direction. Further, unlike BRITS or Bi-GAN, our
 170 approach additionally allows for using the mean as a possible impute option, which makes
 171 sense when imputing in the middle region of large data gaps.

172 3. Proposed Methodology

173 3.1. GRU

174 We start by describing the GRU unit gating equations in detail as our proposed approach
 175 builds on it. Also, GRU unit as the building block for RNNs is currently very popular
 176 across sequence prediction applications [Ravanelli et al. \(2018\)](#); [Che et al. \(2016\)](#); [Gruber
 177 and Jockisch \(2020\)](#). GRU based cell computes its hidden state (for one layer) at the i^{th}
 178 time-step as follows.

$$z_i = \sigma(W^z u_i + U^z h_{i-1} + b_z) \quad (1)$$

$$r_i = \sigma(W^r u_i + U^r h_{i-1} + b_r) \quad (2)$$

$$\tilde{h}_i = \tanh(U(r_i \odot h_{i-1}) + W u_i + b) \quad (3)$$

$$h_i = (1 - z_i) \odot h_{i-1} + z_i \odot \tilde{h}_i \quad (4)$$

179 where h_{i-1} is state at the previous time-instant, (u_i) is current input, $\sigma(\cdot)$ denotes
 180 sigmoid function, z_i is update gate vector and r_i is the reset gate vector. \tilde{h}_i is the additional
 181 memory (summary of all inputs so far) which is a function of u_i and h_{i-1} , the previous hidden
 182 state. The reset signal controls the influence of the previous state on the new memory. The
 183 final current hidden state is a convex combination (controlled by z_i) of the new memory and
 184 the memory at the previous step, h_{i-1} . All associated weight matrices W^z , W^r , W , U^z , U^r ,
 185 U and vectors b_z , b_r and b are trained using back-propagation through time.

186 3.2. Proposed GRU-M Unit

187 As explained earlier in related work, GRU-D in addition to using the masking binary vector
 188 also adopts a novel joint impute-learn strategy. We first describe the essentials of GRU-D,
 189 then point out its drawbacks and finally describe our method, GRU-M which can be viewed
 190 as a non-trivial extension of GRU-D.

191 A multivariate time-series X with D variables of length N is denoted as

$$\langle (x_1, t_1), (x_2, t_2), \dots (x_i, t_i) \dots (x_N, t_N) \rangle \in \mathbb{R}^{D \times T} \quad (5)$$

192 Each $x_i \in \mathbb{R}^D$ represents the i^{th} observation, while x_i^d represents its d^{th} component. As
 193 evident, $t_i \in \mathbb{R}$ denote the time-stamp of x_i , the time instant at which the measurement
 194 happens.

195 To capture which variables are missing at which time-instant, we define formally the
 196 binary masking variable, m_i^d , as

$$m_i^d = \begin{cases} 1, & \text{if } x_i^d \text{ is observed} \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

197 GRU-D also maintains a time-interval δ_i^{dL} , denoting the distance from the closest avail-
 198 able data-point on the left for the d^{th} variable. More formally,

$$\delta_i^{dL} = \begin{cases} t_i - t_{i-1} + \delta_{i-1}^{dL}, & \text{if } i > 1, m_{i-1} = 0 \\ t_i - t_{i-1}, & \text{if } i > 1, m_{i-1} = 1 \\ 0, & \text{if } i = 1. \end{cases} \quad (7)$$

199 GRU-D essentially models a *decay* mechanism from the left end of a gap and works with
 200 δ_i^L ($D \times 1$ vector of δ_i^{dL} s) only. When data is missing at time-step i at the d^{th} input variable,
 201 the input is hypothesized to be a weighted combination of the closest available input to the
 202 left and the time average of the d^{th} component (across the sequence). The decay/weight
 203 factor is modelled using a monotonically decreasing function of δ_i^L and ranges between 0
 204 and 1. The vector of decay/weight factors ($D \times 1$ vector) at the i^{th} time-step would be

$$\gamma_i^L = \exp\{-\max(0, W_\gamma \delta_i^L + b_\gamma)\} \quad (8)$$

205 where b_γ is $D \times 1$ vector and W_γ is a $D \times D$ matrix assumed to be diagonal. This assumes
 206 a component-wise independence of decay rates which is pretty realistic. The modified input
 207 that is input to the GRU unit is now

$$\hat{x}_i^d = m_i x_i^d + (1 - m_i)(\gamma_i^{dL} x_i^{dL} + (1 - \gamma_i^{dL}) \tilde{x}^d) \quad (9)$$

208 where x_i^{dL} is the last available data point to the left of the d^{th} component of i^{th} observation,
 209 \tilde{x}^d is the empirical mean across all available data points among the N observations of the
 210 d^{th} variable. Replacing u_i by \hat{x}_i^d in eqns.(1)-(4) is essentially the GRU-D architecture by
 211 incorporating input decay.

212 **Drawback of GRU-D:** GRU-D's hypothesis is that closer a missing observation is
 213 to the left end of the gap ($\delta_i^{dL} \rightarrow 0$), closer will \hat{x}_i^d be to x_i^{dL} . The farther away the
 214 missing observation is from gap's left end, closer will \hat{x}_i^d be to \tilde{x}^d . In this process, GRU-D
 215 totally ignores information from right end of the gap. This means even though the missing
 216 observation is far away from the left end, it may actually be very close to the right end. In
 217 which case, \hat{x}_i^d must be close to x_i^{dR} , (the closest available observation to the right of i^{th}
 218 observation at d^{th} component) instead of \tilde{x}^d , especially when the gap length is large. In a
 219 bid to address this lacunae in GRU-D and beyond, we propose our novel approach termed
 220 GRU-M.

221 GRU-M unlike GRU-D takes into account the right-end of a gap as well while computing
 222 \hat{x}_i^d , the modified input at each time-step. Towards this, we additionally maintain another
 223 time-interval vector δ_i^{dR} denoting the distance from the closest available data-point on the
 224 right of any observation. This can be efficiently computed recursively as follows.

$$\delta_i^{dR} = \begin{cases} t_{i+1} - t_i + \delta_{i+1}^{dR}, & \text{if } i < N, m_{i+1} = 0 \\ t_{i+1} - t_i, & \text{if } i < N, m_{i+1} = 1 \\ 0, & \text{if } i = N. \end{cases} \quad (10)$$

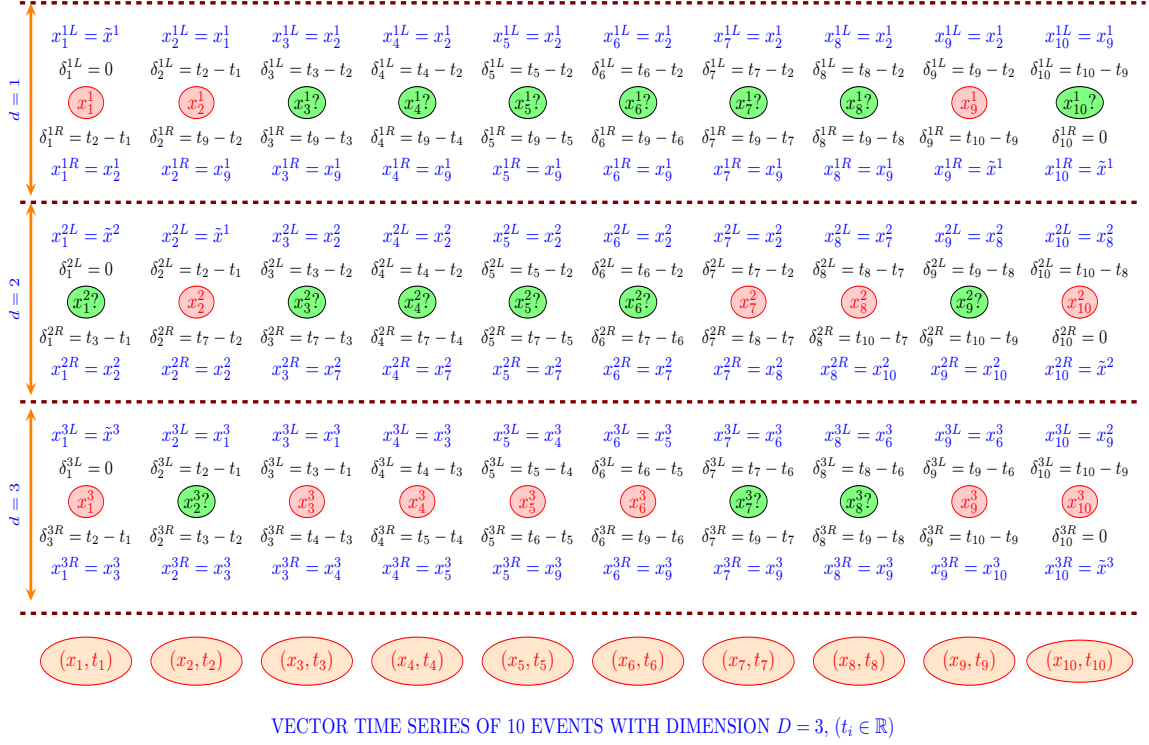


Figure 1: Illustration of various variables (δ_i^{dL} , δ_i^{dR} , x_i^{dL} , x_i^{dR}) that determine the modified input into GRU-unit (in GRU-M) at each time-step. Green - missing, Red - available.

225 At the right end of the time-series, the closest available data point must be at zero distance,
 226 hence $\delta_N^{dR} = 0$. For any other observation i , if the right immediate observation is available
 227 ($m_{i+1} = 0$), then the time distance $\delta_i^{dR} = t_{i+1} - t_i$. If the right immediate observation
 228 is not available ($m_{i+1} = 0$), then one can recursively compute this as the sum of (i)time
 229 distance to the immediate right observation ($t_{i+1} - t_i$), (ii) distance from the immediate
 230 right observation to *its* immediately available data point to the right, δ_{i+1}^{dR} .

231 Please refer to Fig. 1 for a detailed illustration of δ_i^{dR} and δ_i^{dL} on an example data
 232 sequence. For instance, consider the first component ($d = 1$) of the input at time-step 4
 233 (x_4^1), in Fig. 1. The associated entry is marked in green indicating a missing entry. The
 234 closest available data point to the left in the first component is at time-step 2. Note that
 235 available data points are marked in red. Hence δ_4^{dL} , the time distance to the closest available
 236 data point to the left is $(t_4 - t_2)$. The closest available data point to the right in the first
 237 component is at time-step 9. Hence δ_4^{dR} , the time distance to the closest available data point
 238 to the right is $(t_9 - t_4)$.

239 Instead of learning one decay or weight factor between 0 and 1, we learn three factors
 240 as follows. Let $\delta_i^d = [\delta_i^{dL}, \delta_i^{dR}]$.

$$\Gamma_i^d = F(\delta_i^d), \text{ where } F(\cdot) \text{ denotes a feed-forward map, } \Gamma_i^d \text{ is } (3 \times 1) \text{ vector.} \quad (11)$$

241

$$\gamma_i^{dL} = \frac{\exp(\Gamma_i^d(1))}{\sum_{j=1}^3 \exp(\Gamma_i^d(j))}, \quad \gamma_i^{dR} = \frac{\exp(\Gamma_i^d(2))}{\sum_{j=1}^3 \exp(\Gamma_i^d(j))}, \quad \gamma_i^{dm} = \frac{\exp(\Gamma_i^d(3))}{\sum_{j=1}^3 \exp(\Gamma_i^d(j))} \quad (12)$$

242 where

$$\gamma_i^{dL} + \gamma_i^{dR} + \gamma_i^{dm} = 1 \quad (13)$$

243 $F(\cdot)$ denotes a general feed-forward map with a linear activation at the output layer. Hence,
 244 $F(\cdot)$ in its simplest form with no hidden nodes would be a linear transformation, namely

$$F(\delta_i^d) = W_M^d \delta_i^d + b_M^d. \quad (\text{where } W_M^d \text{ is } (3 \times 2), b_M^d \text{ is } (3 \times 1)) \quad (14)$$

245 Essentially, the above equations model a (potentially) general feed-forward structure (FFN
 246 blocks in Fig. 2) with two inputs $[\delta_i^L, \delta_i^R]$, and three outputs $[\Gamma_i^d(1), \Gamma_i^d(2), \Gamma_i^d(3)]$ followed by
 247 a *softmax* transformation. The modified input at time i to the GRU unit is now a convex
 248 combination of three quantities, x_i^{dL} , x_i^{dR} (the closest observation to the right of i) and \tilde{x}^d ,
 249 where the co-efficients come from the above described 3-output softmax layer.

$$\hat{x}_i^d = m_i^d x_i^d + (1 - m_i^d)(\gamma_i^{dL} x_i^{dL} + \gamma_i^{dR} x_i^{dR} + \gamma_i^{dm} \tilde{x}^d) \quad (15)$$

250 Entities x_i^{dR} , x_i^{dL} can be efficiently computed recursively as follows.

$$x_i^{dL} = \begin{cases} x_{i-1}^d & \text{if } i > 1, m_{i-1} = 1 \\ x_{i-1}^{dL} & \text{if } i > 1, m_{i-1} = 0 \\ \tilde{x}^d, & i = 1 \end{cases} \quad (16)$$

251

$$x_i^{dR} = \begin{cases} x_{i+1}^d & \text{if } i < N, m_{i+1} = 1 \\ x_{i+1}^{dR} & \text{if } i < N, m_{i+1} = 0 \\ \tilde{x}^d, & i = N \end{cases} \quad (17)$$

252 At the right end, x_i^{dR} has strictly no closed observation to the right. It would not be
 253 unreasonable to assume x_N^{dR} to be \tilde{x}^d , mean of the d^{th} variable and we adopt this convention
 254 even on the left (in particular on x_1^{dL}). This can be verified in Fig. 1 where x_1^{dL} and x_{10}^{dR}
 255 ($N = 10$ here) are chosen to be \tilde{x}^d . If the immediate observation to the right is available,
 256 then $x_i^{dR} = x_{i+1}^d$. If not, then x_i^{dR} must be same as the immediate observation to the right
 257 of the $(i + 1)^{\text{th}}$ observation, and hence is updated recursively to x_{i+1}^{dR} . An analogous update
 258 is performed from the left for x_i^{dL} .

259 As an illustration look at time-step 4 in the first component ($d = 1$) of the input in
 260 Fig. 1. $x_4^{1L} = x_2^1$ as the closest available point to the left is at t_2 . Similarly $x_4^{1R} = x_9^1$ as
 261 the closest available point to the right is at t_9 . Fig. 1 illustrates the various variables on an
 262 example time-series with 10 observations. The input dimension, D is chosen to be 3. The
 263 red circles represent available observations while the green ones indicated missing values.
 264 The various variables δ_i^{jL} , δ_i^{jR} (eqn. (7) & eqn. (10)) and x_i^{dL} , x_i^{dR} (eqn. (16) & eqn. (17))
 265 are clearly explained in this figure via the explicit values it takes on in all 3 components.

266 Fig. 2 explains the overall scheme of the modified input at the i^{th} step of GRU-M. It
 267 considers a 2-dimensional ($D = 2$) input. If $m_i^d = 1$, i.e. data is available, then \hat{x}_i^d is just
 268 x_i^d . If $m_i^d = 0$, i.e. data is missing, then the entire architecture comes into play. Lets stick

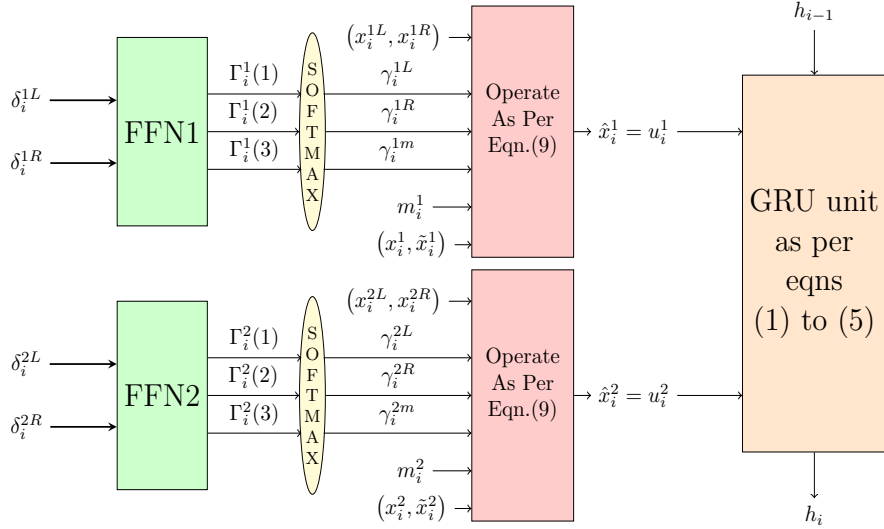


Figure 2: Illustration of computation of the modified input (\hat{x}_i^d) which is fed at the input (u_i^d) position of the GRU unit at i^{th} time-step of unfolded RNN, Input Dimension, $D = 2$ here, h_{i-1} comes from previous ($(i-1)^{th}$) time-step of unfolded RNN.

269 to $d = 1$ for illustration. Both δ_i^{1L} and δ_i^{1R} are processed by a feed forward network (FFN1)
 270 which implements eqn. (11) with a linear activation at the output. The 3 outputs of the
 271 FFN1 are passed through a softmax layer to obtain the 3 factors γ_i^{dL} , γ_i^{dR} and γ_i^{dm} . These
 272 3 factors in conjunction with x_i^{1L} (closest available value to the left), x_i^{1R} (closest available
 273 value to the right), \tilde{x}_i^1 (mean of the 1st component) and m_i^1 (masking variable) are fed to
 274 the next block which implements eqn. (15). The output of this block is \tilde{x}_i^1 which is fed at
 275 u_i^1 , the 1st input position of the GRU block. A similar story holds for the 2nd input position
 276 of the GRU-block at the i^{th} time-step.

277 **Advantages of GRU-M:** The idea of the above parametrization is as follows. For large
 278 gaps, when the missing observation (where $m_i = 0$) is closer to the left end of the gap, \hat{x}_i^d
 279 must be close to x_i^{dL} . Similarly when missing observation closer to the right end, \hat{x}_i^d should
 280 be close to x_i^{dR} . When missing observation is around the center of the gap, the true value
 281 may have no bearing with the left-end or right-end value. Hence in the mid-gap region, it
 282 may be best to impute with \tilde{x}^d , which our GRU-M scheme can capture. On the other hand,
 283 when the gap length ($\delta_i^L + \delta_i^R$) is small or medium, the above parametrization should ideally
 284 ignore \tilde{x}^d completely, as the time-series true value may actually not come anywhere close to
 285 the mean \tilde{x}^d . For instance, when both x_i^{dL} and x_i^{dR} are close in value, while \tilde{x}^d is relatively
 286 much farther from either of them, \tilde{x}^d must be ignored while evaluating \hat{x}_i^d . *Our proposed*
 287 *parametrization can capture all such important features.*

288 **Model Parameters:** Compared to GRU-D, while GRU-M captures many more in-
 289 fluences during imputation as explained above, this comes at a cost of some additional
 290 parameters. While the standard GRU unit parameters are same in both GRU-D and GRU-
 291 M, the difference lies in how \hat{x}_i is evaluated. GRU-D, uses two $2D$ extra parameters for
 292 this, where D parameters come from the bias vector b_γ and remaining D parameters from

293 the diagonal of W_γ matrix (refer eqn. (8)). In contrast, in GRU-M, for each component
 294 (or the d^{th} component), we need to additionally store at least a matrix W_M^d and a vector
 295 b_M^d . Recall from eqn. (14), that W_M^d is (3×2) while b_M^d is (3×1) , hence contributing to
 296 9 parameters. Hence in total GRU-M uses $9D$ extra parameters at least. In comparison
 297 to GRU-D, GRU-M needs $7D$ extra parameters to capture the additional influences that
 298 GRU-D completely misses. This is not much, while training with these extra parameters
 299 gives us superior predictions as demonstrated in our experiments.

300 GRU-D in addition to input decay captures missingness patterns by incorporating a
 301 hidden state decay from the left. This aspect of GRU-D can also be extended. However in
 302 this paper, for purposes of testing GRU-M and GRU-D, we stick to input-decay only.

303 3.3. GRU-M unit for sequence prediction (classification/forecasting)

304 We have till now discussed how sequential inputs with missing entries from a (vector) time-
 305 series are processed by an RNN layer with the proposed GRU-M unit. Suppose one wants
 306 to perform classification of a sequence (with missing entries). This sequence needs to be
 307 processed by an RNN layer with GRU-M unit. One needs to feed the state output from the
 308 last time-step to a soft-max layer and train using a standard cross-entropy loss, for instance.

309 **Forecasting:** As mentioned earlier we use a multi-step forecasting setting. *For accurate*
 310 *multi-step forecasting, we use an encoder-decoder Sutskever et al. (2014) model which consists*
 311 *of two RNN layers of the proposed GRU-M units. For forecasting Wen et al. (2017), a*
 312 *historical window (of some fixed size) into the past starting from the forecast horizon is used*
 313 *as input for forecasting. This window also referred to as input window is fed sequentially*
 314 *into the encoder or the first RNN layer. The last state of the encoder RNN layer is typically*
 315 *fed as the initial state into the second RNN layer, which is the decoder. The decoder is*
 316 *unfolded into as many steps as the width of the forecast horizon. Exogenous inputs (if*
 317 *any) of the input window are placed at the encoder inputs, while the exogenous inputs in the*
 318 *forecast horizon are placed as decoder inputs in a sequential fashion.*

319 4. Results

320 In this paper, we restrict our experimental validation of the GRU-M unit to multi-step
 321 forecasting task using an ED architecture as described above. *The data sets we use are*
 322 *all publicly available. We are unfortunately not in a position to share our codes due to*
 323 *proprietary reasons, but we believe we have shared full details of our proposed architecture*
 324 *using which any reader can implement the same.*

325 4.1. Data Sets, Error Metrics and Hyper-parameters

326 **D1:** Air Quality¹ data-set was recorded for a period of 1 year (Mar 2004 to Feb 2005) on an
 327 hourly basis. The data is *naturally missing* here. It contains 9358 points of average response
 328 from 5 metal oxide sensors. We use CO(GT) data, from one of these sensors (denoted as
 329 D1 from now on) to measure the performance of our model.

1. github.com/mehak25/BiGAN/tree/master/data/air/initial

330 **D2:** Electricity data-set² is a publicly available data-set which measured electricity con-
 331 sumption in kWh of 321 clients for every hour. We divided the whole data-set into 3 groups
 332 namely low, medium and high based on average consumption. From each group, we choose
 333 one sequence randomly and denote this collection of 3 sequences as D2.

334 **D3:** M5³ is a publicly available data-set distributed across 12 aggregation levels which con-
 335 tains daily sales of different products for 5.4 years. We pick sequences from level 10 which
 336 is the aggregated sales for each product across all stores and states. While this level of
 337 data contains 3049 sequences, we pick 4 sequences (referred to as *D3* from now on) based
 338 on sufficient total variation. The idea is higher the total variation of the sequence, more
 339 challenging would be the forecasting. We employ synthetic masking on these 4 sequences.
 340 The total variation we employ is defined below:

$$TV(f) = \frac{1}{L-h+1} \sum_{i=1}^{L-h+1} \frac{1}{h} \sum_{j=i}^{i+h-1} |f(j+1) - f(j)|,$$

341 where L is the total length of the sequence and h is the size of the moving window ($h \ll L$).
 342 TV can be defined without a moving window with just the inner summation, but a high TV
 343 based on this definition will not necessarily capture sufficient variation across the sequence.
 344 We use a small/local moving window to check for sufficient variation across the length of
 345 the sequence.

346 We consider following well-known error metrics to measure performance: (1)**MSE** (Mean
 347 Square Error) (2)**MAPE** (Mean Absolute Percentage Error) (3)**MASE** (Mean Absolute
 348 Scale Error Hyndman and Koehler (2006)). While MSE is a scale dependent metric, the
 349 other two are scale independent, while being complementary to each other.

350 The APE is relative error (RE) expressed in percentage. If \hat{X} is predicted value, while
 351 X is the true value, $RE = (\hat{X} - X)/X$. In the multi-step setting, APE is computed for each
 352 step and is averaged over all steps to obtain the MAPE for one window of the prediction
 353 horizon. APE while has the advantage of being a scale independent metric, can assume
 354 abnormally high values and can be misleading when the true value is very low. An alternative
 355 complementary error metric which is scale-free could be MASE.

356 The MASE is computed with reference to a baseline metric. The choice of baseline is
 357 typically the *copy previous* predictor, which just replicates the previous observed value as
 358 the prediction for the next step. For a given window of one prediction horizon of K steps
 359 ahead, let us denote the i^{th} step error by $|\hat{X}_i - X_i|$. The i^{th} scaled error is defined as

$$e_s^i = \frac{|\hat{X}_i - X_i|}{\frac{1}{n-K} \sum_{j=K+1}^n |X_j - X_{j-K}|} \quad (18)$$

360 where n is no. of data points in the training set. The normalizing factor is the average i^{th}
 361 step-ahead error of the copy-previous baseline on the training set. Hence the MASE on a
 362 multi-step prediction window w of size K will be

$$MASE(w, K) = \frac{1}{K} \sum_{j=1}^K e_s^j \quad (19)$$

2. <https://github.com/laignuokun/multivariate-time-series-data>

3. <https://www.kaggle.com/competitions/m5-forecasting-accuracy/data>

Table 1 represents the broad choice of hyper-parameters used for training.

Table 1: Hyper parameters during training.

Parameters	Description (D1/D2/D3)
Batch size	256/512/256
Learning rate	0.01/0.05/0.001
Hidden state vector size	16/16/20
Optimizer	Adam/RMSProp/Adam

363

364 4.2. Baselines and Masking

365 We consider a wide range of baselines to benchmark our method, GRU-M against: (1)
 366 EDC - An Encoder-Decoder (Seq2Seq) predictive model which first imputes using cubic
 367 spline interpolation (using a piece-wise cubic polynomial that is twice differentiable) before
 368 model building. (2) BRITS Cao et al. (2018) - A bi-directional RNN approach for joint
 369 imputation and prediction exploiting information from both ends of a data gap. (3) RITS
 370 Cao et al. (2018) - unidirectional version of BRITS. (4) Bi-GAN Gupta et al. (2021) - A
 371 GAN based jointly impute and learn technique incorporating information from either ends
 372 of a gap. (5) GRU-D Che et al. (2016) - An ED model with a GRU-D unit (adopting a joint
 373 impute and learn strategy) in both encoder and decoder.

374 The first approach above is based on imputing first (using a strong technique) followed
 375 by predictive model building using an ED model on complete imputed data. The rest of
 376 the baselines adopt a joint impute and learn approach using some distinct form of an RNN
 377 architecture and training approach. Hence our bench-marking enables a diverse comparison
 378 with state-of-art baselines.

379 4.2.1. Assessing significance of mean error differences statistically

380 We have conducted a Welch t-test (unequal variance) based significance assessment (across
 381 all relevant experiments) under all the mean metrics (MSE, MASE, MAPE) differences (Pro-
 382 posed vs Baseline) with a significance level of 0.05 for null hypothesis rejection. The best
 383 performing method’s error is highlighted in bold if its MASE/MAPE improvement over ev-
 384 ery other method is statistically significant. We allow for highlighting the second/third best
 385 errors in situations when the mean error differences between the best and second/third best
 386 errors are statistically insignificant.

387 4.2.2. Synthetic Masking

388 We traverse the data sequentially and at each step we ask whether to mask at the current
 389 step or not. To perform this, we toss a coin at every step with an adaptive/varying q
 390 (probability of head). On seeing a tail, we do nothing and move ahead. If heads, we need
 391 to mask, in particular decide the number of consecutive steps ahead (τ_w) that need to be
 392 masked. For this, we consider a bag of window lengths T_w , from which we uniformly sample
 393 to obtain τ_w . Depending on the length τ_w sampled, the heads probability q is updated as
 394 $q = \frac{c}{\tau_w}$, where c is a control parameter. If $c = 1$, the strategy masks about 50% of the points,

395 because for every τ_w points masked, the next τ_w points (in expectation) are retained. In
 396 both D2 and D3, we have chosen $c = 0.2$, which means we mask about 16% of the points.

397 **Synthetic masking advantage:** Since underlying actual data is known, one can test
 398 the performance of trained models on *all* points of a separate test set in the forecast horizon.

399 4.3. Results on D1 (data naturally missing)

400 For D1, input window size was set as 20 while the prediction horizon was varied over 8, 12
 401 and 16 steps. We separately kept aside 10% data for testing our model. This means we
 402 measure the performance of our method on 909, 905 and 901 test examples for prediction
 403 horizon 8, 12 and 16 respectively. Table 2 gives a detailed account of the various errors,
 404 vindicating GRU-M’s superior performance.

405 Specifically, in MASE terms, GRU-M shows a minimum (statistically significant) im-
 406 provement of 0.15 over all baselines in the best case scenario (Prediction horizon = 12).
 407 Please note in other two scenarios also, the minimum improvement over all other baselines
 408 is a very healthy 0.14. In MAPE terms, GRU-M shows a minimum improvement of 8%
 409 across all baselines in the best case scenario (Prediction horizon = 8). In other two scenar-
 410 ios too, the minimum improvement over all baselines is a statistically significant 7% and 5%
 411 respectively. In MSE terms, GRU-M outperforms all baselines across all prediction horizons.

Table 2: Results on D1 (naturally missing) for different forecast horizons.

Method	Prediction horizon = 8			Prediction horizon = 12			Prediction horizon = 16		
	MASE	MAPE	MSE	MASE	MAPE	MSE	MASE	MAPE	MSE
EDC	1.07	39	1.78	0.91	35	1.72	0.86	37	1.63
Bi-GAN	1.14	32	2.48	1.22	38	3.33	1.39	48	4.39
RITS	1.05	32	1.97	0.94	32	1.95	0.88	32	1.95
BRITS	1.05	32	1.97	0.94	32	1.97	0.88	32	1.95
GRU-D	0.79	42	1.76	0.78	42	1.76	0.74	42	1.76
GRU-M	0.65	24	0.77	0.63	25	0.87	0.62	27	0.94

412

413 4.4. Results on D2 (data synthetically masked)

414 We perform masking on D2 and D3 to generate missingness as discussed above. For D2,
 415 we set width of input window as 20 and prediction horizon at 12. This implies we are
 416 predicting for the next 12 hours ahead. A test size of 10% was separately set aside for
 417 each of the 3 series to measure performance of our model. Table 3 demonstrates GRU-Ms
 418 superior performance compared to the baselines. *In both the current and next experiment,*
 419 *performance of Bi-GAN was very poor (in particular MAPE was > 100%). Hence we have*
 420 *not reported Bi-GAN results in Tab. 3 and Tab. 4.*

421 In particular, in terms of MASE/MAPE, GRU-M shows a minimum improvement of
 422 0.05 (Medium consumption) and 13% (Medium) respectively over all baselines in the best
 423 case scenario for the medium category. Even in high category, GRU-M outperforms all
 424 baselines based on all metrics in a statistically significant fashion. In the low category too,
 425 except based on MASE metric where the performance is comparable to GRU-D, GRU-M
 426 outperforms all baselines in a statistically significant fashion.

Table 3: Results on D2: Masking Window Lengths $T_w = \{5, 6, \dots, 12\}$, Parameter $c = 0.2$

Method	Low			Medium			High		
	MASE	MAPE	MSE	MASE	MAPE	MSE	MASE	MAPE	MSE
EDC	0.55	23	1135	0.68	25	14941	0.37	14	110529
RITS	0.89	31	3222	1.27	39	45221	1.36	45	977696
BRITS	0.88	33	2804	1.03	33	35518	1.04	36	637439
GRU-D	0.43	23	1167	0.36	24	14468	0.32	22	198071
GRU-M	0.51	22	1046	0.31	11	4739	0.28	11	74424

4.5. Results on D3 (data synthetically masked)

Here, the input window size is set to 20, while forecast horizon was chosen to be 5 days ahead. A separate 97 days out of 1941 was kept aside for testing which implies we use 73 examples to measure the performance of our method. Table 4 represents the errors of all relevant methods, which indicates the superior or comparable performance of GRU-M in terms of all three error metrics, compared to all baselines.

In particular, under MSE metric, GRU-M though is one of the best performing methods, it has comparable performance with EDC (also indicated in bold) on sequences 1, 2. In terms of MASE/MAPE, GRU-M shows a minimum improvement of 0.04 and 7% respectively over all baselines, in the best case scenario (Seq 3). Our results based on MASE/MAPE indicate that GRU-D and EDC have comparable performance with GRU-M on some sequences.

Table 4: Results on D3: Masking Window Lengths $T_w = \{5, 6, \dots, 12\}$, Parameter $c = 0.2$

Method	Seq 1			Seq 2			Seq 3			Seq 4		
	MASE	MAPE	MSE	MASE	MAPE	MSE	MASE	MAPE	MSE	MASE	MAPE	MSE
EDC	0.65	19	3293	0.43	14	1732	0.62	22	9366	0.63	41	1759
RITS	0.72	24	3740	1.45	50	13219	1.08	36	33177	0.71	65	2143
BRITS	0.85	30	5325	1.17	40	8750	0.98	31	33945	0.54	49	1273
GRU-D	0.59	23	4098	0.38	16	1874	0.40	19	7198	0.42	39	1106
GRU-M	0.56	17	2589	0.34	10	1706	0.38	12	4647	0.48	39	1064

437

5. Conclusions

We proposed a novel architecture (GRU-M) for sequential learning under missing data. Its a joint impute and learn approach where the parameterized functions which impute and perform predictive modelling are simultaneously learnt. It factors in information not only from both ends of a data gap, but also the distance from the left and right-end of the gap in a novel fashion, distinct from existing approaches. It can also incorporate state-decay from the right if necessary whenever there are bidirectional layers, which can arise during forecasting. Overall, it can be viewed as a non-trivial generalization of GRU-D [Che et al. \(2016\)](#), which is a state-of-art technique for the same problem. Our proposed approach can be employed for both sequence classification and sequence forecasting. Based on 3 diverse metrics, we bench-marked our approach on data sets where data was either naturally missing or synthetically masked. We compared against a range of diverse, but closely related state-of-art RNN approaches for sequence forecasting under missing data. Our results clearly vindicate the viability of our proposed approach. As future work, we would like to test our proposed architecture on sequential classification.

452

453 **References**

- 454 Yoshua Bengio and Francois Gingras. Recurrent neural networks for missing or asynchronous
 455 data. In *Proceedings of the 8th International Conference on Neural Information Processing*
 456 *Systems*, page 395–401, Cambridge, MA, USA, 1995. MIT Press.
- 457 Wei Cao, Dong Wang, Jian Li, Hao Zhou, Lei Li, and Yitan Li. Brits: Bidirectional recurrent
 458 imputation for time series. *Advances in Neural Information Processing Systems 31*, pages
 459 6775–6785, 2018.
- 460 Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recur-
 461 rent neural networks for multivariate time series with missing values. *Scientific Reports*,
 462 8, Jun 2016. doi: 10.1038/s41598-018-24271-9.
- 463 Yagmur Gizem Cinar, Hamid Mirisaee, Parantapa Goswami, Eric Gaussier, and Ali Ait-
 464 Bachir. Period-aware content attention rnns for time series forecasting with missing values.
 465 *Neurocomputing*, 312:177–186, 2018. ISSN 0925-2312.
- 466 Andrea Cini, Ivan Marisca, and Cesare Alippi. Filling the g_ap_s: Multivariate time
 467 series imputation by graph neural networks. In *International Conference on Learning*
 468 *Representations*, 2021.
- 469 Joseph David Futoma. *Gaussian Process-Based Models for Clinical Time Series in Health-*
 470 *care*. PhD thesis, Statistical Science, Duke University, 2018.
- 471 Pedro J. García-Laencina, José-Luis Sancho-Gómez, and Aníbal R. Figueiras-Vidal. Pattern
 472 classification with missing data: A review. *Neural Comput. Appl.*, 19(2):263–282, March
 473 2010. ISSN 0941-0643. doi: 10.1007/s00521-009-0295-6.
- 474 Alex Graves, Marcus Liwicki, Santiago Fernández, Roman Bertolami, Horst Bunke, and
 475 Jürgen Schmidhuber. A novel connectionist system for unconstrained handwriting recog-
 476 nition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(5):855–868, May 2009.
- 477 Nicole Gruber and Alfred Jockisch. Are gru cells more specific and lstm cells more sensitive
 478 in motive classification of text? *Front. Artif. Intell.*, 2020. doi: 10.3389/frai.2020.00040.
- 479 Mehak Gupta, Thao-Ly T. Phan, H. Timothy Bunnell, and Rahmatollah Beheshti. Con-
 480 current imputation and prediction on ehr data using bi-directional gans: Bi-gans for ehr
 481 imputation and prediction. In *Proceedings of the 12th ACM Conference on Bioinformatics,*
 482 *Computational Biology, and Health Informatics (BCB '21)*, pages 1–9, 2021.
- 483 Rob J. Hyndman and Anne B. Koehler. Another look at measures of forecast accuracy.
 484 *International Journal of Forecasting*, 22(4):679–688, 2006.
- 485 Sharon A. Johnson. *Splines*, pages 1443–46. Springer US, 2013. ISBN 978-1-4419-1153-7.
- 486 Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recom-
 487 mender systems. *Computer*, 42(8):30–37, August 2009.
- 488 Steven Cheng-Xian Li. *Learning from Irregularly-Sampled Time Series*. PhD thesis, Steven
 489 Cheng-Xian Li, University of Massachusetts Amherst, 2020.

- 490 Bryan Lim and Stefan Zohren. Time-series forecasting with deep learning: a survey. *Phil.*
491 *Trans. R. Soc. A.*, 379(2194):40–48, 2021.
- 492 Zachary C Lipton, David Kale, and Randall Wetzel. Directly modeling missing data in
493 sequences with rnns: Improved classification of clinical time series. In *Proc. of the 1st*
494 *Machine Learning for Healthcare Conference*, volume 56, pages 253–270. PMLR, 2016.
- 495 Yonghong Luo, Xiangrui Cai, Ying ZHANG, Jun Xu, and Yuan xiaojie. Multivariate time
496 series imputation with generative adversarial networks. In *Advances in Neural Information*
497 *Processing Systems*, volume 31, pages 1603–1614, 2018.
- 498 Debashis Mondal and Donald Percival. Wavelet variance analysis for gappy time series.
499 *Annals of the Institute of Statistical Mathematics*, 62(5):943–966, October 2010.
- 500 Tong Nie, Guoyang Qin, Wei Ma, Yuewen Mei, and Jian Sun. Imputeformer: Low rankness-
501 induced transformers for generalizable spatiotemporal imputation. In *Proceedings of the*
502 *30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '24,
503 page 2260–2271, New York, NY, USA, 2024.
- 504 Soumen Pachal and Avinash Achar. Sequence prediction under missing data: An rnn ap-
505 proach without imputation. In *Proc. of the 31st ACM Int. Conf. on Information &*
506 *Knowledge Management*, CIKM '22, page 1605–1614, 2022.
- 507 S. Parveen and P. D. Green. Speech recognition with missing data using recurrent neural
508 nets. In *Proceedings of the 14th International Conference on Neural Information Process-*
509 *ing Systems: Natural and Synthetic*, page 1189–1195, Cambridge, MA, USA, 2001.
- 510 Mirco Ravanelli, Philemon Brakel, Maurizio Omologo, and Yoshua Bengio. Light gated
511 recurrent units for speech recognition. *IEEE Transactions on Emerging Topics in Com-*
512 *putational Intelligence*, 2(2):92–102, Apr 2018.
- 513 Kira Rehfeld, Norbert Marwan, Jobst Heitzig, and Jürgen Kurths. Comparison of correlation
514 analysis techniques for irregularly sampled time series. *Nonlinear Processes in Geophysics*,
515 18(3):389–404, July 2011. doi: 10.5194/npg-18-389-2011.
- 516 Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural
517 networks. In *Proceedings of the 27th International Conference on Neural Information*
518 *Processing Systems - Volume 2*, pages 3104–3112, 2014.
- 519 Volker Tresp and Thomas Briegel. A solution for missing data in recurrent neural networks
520 with an application to blood glucose prediction. In *Proceedings of the 10th International*
521 *Conference on Neural Information Processing Systems*, page 971–977, Cambridge, MA,
522 USA, 1997. MIT Press.
- 523 Ruofeng Wen, Kari Torkkola, Balakrishnan Narayanaswamy, and Dhruv Madeka. A multi-
524 horizon quantile recurrent forecaster. In *The 31st Conference on Neural Information*
525 *Processing Systems (NIPS 2017), Time Series Workshop*, 2017.
- 526 Zhengzheng Xing, Jian Pei, and Eamonn Keogh. A brief survey on sequence classification.
527 *SIGKDD Explor. Newsl.*, 12(1):40–48, November 2010. ISSN 1931-0145.