

MAIN: Multihead-Attention Imputation Networks

Spyridon Mouselinos*, Kyriakos Polymenakos*[†], Antonis Nikitakis*, Konstantinos Kyriakopoulos*[‡]

{s.mouselinos, k.polymenakos, a.nikitakis, k.kyriakopoulos}@deepsea.ai

*DeepSea Technologies, Athens, Greece

[†] University of Oxford, Department of Engineering Science

[‡] University of Cambridge, Department of Engineering

Abstract—The problem of missing data, usually absent in curated and competition-standard datasets, is an unfortunate reality for most machine learning models used in industry applications. Recent work has focused on understanding the nature and the negative effects of such phenomena, while devising solutions for optimal imputation of the missing data, using both discriminative and generative approaches. We propose a novel mechanism based on multi-head attention which can be applied effortlessly in any model and achieves better downstream performance without the introduction of the full dataset in any part of the modeling pipeline. Our method inductively models patterns of missingness in the input data in order to increase the performance of the downstream task. Finally, after evaluating our method against baselines for a number of datasets, we found performance gains that tend to be larger in scenarios of high missingness.

Index Terms—Imputation, Attention, Deep Learning, Encoding

I. INTRODUCTION

Corrupted or missing data often raise the following questions: If the missing values are imputed, what value is to be used for "padding"? What is the underlying relation with the normalization scheme of the input values? How should categorical inputs be treated? Popular approaches include replacing the missing values with their mean, in the case of continuous variables, or the median, in the case of discrete variables. Other approaches introduce a binary indicator of missingness that is jointly trained with the rest of the model parameters.

The simplest, although widely used, method is *mean/mode impute*, where missing numerical values are replaced by their mean (0 in the Zeta Normalized Scale) and categorical values are replaced by their mode (the most frequent class of the dataset). One disadvantage of mean/mode impute is that it can introduce ambiguity: variables that take values close to the mean or mode, might not be distinguished from variables whose missing values were imputed with the value of the mean or median. To resolve this ambiguity, a bitwise indicator (also referred to as *mask*) was introduced, that assigns 1 to missing instances and 0 to non-missing instances. For simplicity, we refer to this method as Zero Impute (with) Mask Concat (**ZIMC**).

A straightforward extension to ZIMC is Sparsity Normalization (**SN**), which creates a per-example (horizontal) scaling by a constant factor [1], based on the global missingness factor. This method resolves the Variable Sparsity Problem, but fails

to capture the dependencies in-between the features. In more detail, no adjustment is made based on what feature is missing and the values the existing features have.

Since the problem of missing data naturally defines a context of missing information we are inspired from feature-wise transformations [2] as methods of context-based processing. More specifically the Feature Wise Linear Modulation (**FiLM**), naturally augments the idea of SN [1] (**SN**) since it can be seen as a conditional scaling method. In order to efficiently capture the missing context we combine both the feature vector (i.e values) as well as the bitwise missingness indicator mask into contextual embeddings. In this regard, our method is capable of producing independent scaling coefficients for each individual feature, surpassing the **SN** method on various datasets and thus confirms that vertical (per-feature) conditioned imputation surpasses horizontal (per-instance) methods. Furthermore, we extend the attention mechanism of AimNet [3] into one that uses a richer key-value embedding representation, capable of encapsulating position, missingness and value information. On the contrary, in AimNet, only the missingness context is captured. That enables our model to continuously improve the embedding representations with both existing and missing features.

Our approach, Multihead-Attention Imputation Networks (**MAIN**), builds on **FiLM** as it creates multiple coefficients for each feature (multiple vertical imputation), based on the attention mechanism. In this modeling scenario, we create N sets of N coefficients where N is the number of input features, with each set representing a conditioned modulation of a feature based on the value and existence of all features. This learnable combination mechanism captures both global feature-value and feature existence effects.

Briefly our method extends the SN method of Yi et al. [1] in two ways:

- 1) learning a different weight coefficient for each individual feature - instead of the whole feature vector, thus dismissing the concept of linear covariation as proposed in SN[1].
- 2) conditioning the coefficients not only on the existence vector of the input features but also on their observed values.

. Also our attention mechanism differs from [3] as we use multi-head attention and we build embeddings with both the value and the missing context of all the features.

II. RELATED WORK

Imputation methods are often classified as *discriminative* and *generative*, based on the formulation of the missingness modeling mechanism. In discriminative modeling, missing features are approximated directly through their conditional distributions, given the values of the existing features. On the contrary, in generative modeling, missing features are often imputed in a two-step fashion: initially, the underlying joint distribution of all features is modeled, followed by a second step of conditionally generating the missing features from the existing information in the feature vector.

Discriminative methods depend on certain assumptions about the missingness mechanism. Specifically, the data must conform to the Missing Completely At Random (MCAR) or - at a modest degree - Missing At Random (MAR) assumption [4, 5]. For MCAR, the probability of an entry missing is independent from the values of the other features, while for MAR, that probability depends on the values of the other features. In both cases, the probability of a feature missing is assumed to be independent from the value of that feature.

Generative models focus on the joint probability distribution of the data instead of the conditional [6, 7]. In [6] the authors deal with problem of missing data with the use of deep latent variable models (DLVMs). Their approach is based on an importance-weighted autoencoder (IWAE) which maximises a potentially tight lower bound of the log-likelihood of the observed data. In [7] the authors propose a general framework implementing VAEs to fit incomplete heterogenous data. Generative models in general are more flexible and can effectively capture multimodal distributions, while discriminative models rely on point estimates. The extra flexibility of the generative approach allows these models to deal with strongly MAR scenarios, as well as cases where the probability of a feature missing is dependent on the value of that feature (Missing Not At Random, or MNAR). On the other hand they are difficult to train and specifically, GANs can suffer from slow convergence and mode collapse problems.

It should be pointed that, joint modeling may excel in scenarios where the underlying connections between features are of spatial and/or auto-regressive nature (e.g image completion) [5, 8, 9], but in case of tabular data this can be equally approached by (multiple) conditional modeling, since there is not such a variety of modalities that can only be approached by manifold walks.

Furthermore, when the final objective is to build a predictive task based on specific input (e.g regression, classification) the flexibility that generative modeling provides is unfit for time-critical applications. In a scenario of new incoming information, a generative model will act in a "two-step fashion". Firstly, it will impute the missing data producing a complete feature vector and as a second step will apply an auxiliary or different architecture/model on the downstream task. Discriminative methods on the other hand can operate both in single and two-step fashion.

The work by Yi et al. [1] gave an in-depth insight to the first

two questions about data imputation, and re-introduced us with the phenomenon of Variable Sparsity Problem (VSP), while making significant progress towards the "one-step" approach. The VSP problem stresses an undesired phenomenon where the model's performance drops as its output significantly varies due to the rate of missingness in the given input also identified in [10]. The solution proposed by Yi et al. [1], named Sparsity Normalization (SN), seems to give a significant boost in the reconstruction capabilities of existing networks, while also increasing the robustness of a model's downstream task when dealing with missing inputs.

In this work, we opt for a discriminative modeling schema, supporting the paradigm shift towards "one-step" methods, and study the effect of query-key-value attention performed in a multi-head fashion. We propose a custom key representation encoding, called Positional Encoded Vector (PEV). We show that even in heavily corrupted data our model incorporates the best attributes of discriminative and generative modeling, being an one-step method that is valid for both numeric and categorical inputs. The contribution of this work can be summarized around 4 key ideas:

- 1) Due to its main application in real-time industrial scenarios, this work focuses on the performance gains of the downstream task(s), in contrast with other works that seem to sacrifice performance in favour of reconstruction of the full feature vector.
- 2) We show that our method has the best of both worlds between discriminative and generative methods as far as the downstream task performance is concerned. It can handle both MCAR and MAR assumptions while not suffering from the limitation of training with complete vectors, that some discriminative and generative approaches silently imply. After testing its performance against 7 Datasets it shows consistent performance margin, as well as convergence without any instability issues.
- 3) We achieve a single-step, missingness-agnostic behaviour: our scheme can train on missing data directly, exploiting all the examples (including partially missing ones), and test on a different set never seen by the model, with completely untreated data (including missing features). This is not the case in autoencoder-based (i.e reconstruction) schemes like AimNet or generative methods [11], where the model can train only on complete data where the missingness/reconstruction is artificially induced.
- 4) The PEV-based attention mechanism models both the missingness and the value of each individual feature in relation to the other features on the raw data distribution. AimNet [3] models only the missing context in a missing-induced distribution, while in VSP [1] the feature vector is regularized horizontally, not taking into account each feature's individual variation.

III. PROBLEM FORMULATION

A. Preliminaries

Following the formulation of VSP [1], let $\mathbf{x} \in \mathbb{R}^l$ and $\mathbf{y} \in \mathbb{R}^d$ represent a training input/output pair instance of a

model with D tasks. Without loss of generality we will set $|D| = 1$ for the rest of our analysis.

To formulate missingness, we introduce $\mathbf{m} \in \{0, 1\}^l$, as the binary mask indicating missing values in \mathbf{x} . Hence, we define $\mathbf{x}_{\text{miss}} = \mathbf{x} \odot \mathbf{m}$, where \odot denotes the Hadamard product (element-wise multiplication) of two vectors, as the corrupted input the model observes after the effect of the missingness mechanism. Finally, for the model let's assume a N -layer feed forward network, with convex, non-decreasing non-linearities σ in each layer but the last. Each layer i contains n_i units, and we use $\mathbf{W}^i \in \mathbb{R}^{n_i \times n_{i-1}}$ to denote the weight matrix, $\mathbf{b}^i \in \mathbb{R}^{n_i}$ to denote the bias, and $\mathbf{h}^i \in \mathbb{R}^{n_i}$ for the post-activation output vector. The calculation executed at each layer can then be written as $\mathbf{h}^i = \sigma^i(\mathbf{W}^i \mathbf{h}^{i-1} + \mathbf{b}^i)$.¹

B. Modeling Assumptions

Our method works under the following two assumptions:

- 1) Each element m^k of the binary mask \mathbf{m} is MAR, meaning it may not depend on the other mask elements or the value of the element x^k , but it does depend on the values of the original input vector $x_i, i \neq k$, as described in [4]. We will denote the means of the mask vector as μ_m and the means of the input vector as μ_{missing} .
- 2) The coordinates of \mathbf{b}^i and the elements of \mathbf{W}^i are mutually independent and follow the same distribution with means μ_b^i and μ_w^i respectively, as in [1, 12, 13].

C. Methods

Theorem 1. *The expected value of the output layer of an N -Layer FFN with l -dimensional input and convex, non-decreasing, non-linearities under the MCAR assumption is bounded by: $E[\mathbf{h}^N] \geq f_N \circ f_{N-1} \circ \dots \circ f_1(\mu_m \mu_x)$, where $f_i(x) = \sigma(\mathbf{W}^i \cdot \mathbf{x}^{i-1} + \mathbf{b}^i)$.*

Proof of 1. Proof is result of [1].

In Theorem 2 we extend the result of Theorem, for the more general MAR case.

Theorem 2. *The expected value of the first layer of an N -Layer FFN of Theorem 1, under the MAR assumption, is bounded by: $E[\mathbf{h}^1] \geq \sigma(n_0 \mu_w^1 \mu_m \mu_x + \mu_b^1) + T_1$, where $T_1 = \sigma(n_0 \mu_w^1 \text{Cov}(\mathbf{x}^1, \mathbf{m}^1))$.*

Proof of 2. For the output activation vector of the first layer \mathbf{h}^1 it holds that:

$$\mathbf{h}^1 = \sigma(\mathbf{W}^1 \times \mathbf{x}_{\text{missing}} + \mathbf{b}) = \sigma(\mathbf{W}^1 \times \mathbf{x} \odot \mathbf{m} + \mathbf{b}).$$

Then, by taking the expected value of both sides:

$$E[\mathbf{h}^1] = E[\sigma(\mathbf{W}^1 \times \mathbf{x} \odot \mathbf{m} + \mathbf{b})],$$

And due to non-decreasing convexity of σ :

$$E[\mathbf{h}^1] \geq \sigma(E[\mathbf{W}^1 \times \mathbf{x} \odot \mathbf{m} + \mathbf{b}]).$$

¹Throughout the paper we use uppercase and bold notation for matrices, lower and bold notation for vectors and lower and italic notation for elements. Superscripts without parentheses as \mathbf{h}^i are used to denote a quantity at layer i , and superscripts in parentheses, as $h^{(i)}$ to denote the i th element of \mathbf{h} .

For the second part of the inequality it holds that: $\sigma(E[\mathbf{W}^1 \times \mathbf{x} \odot \mathbf{m} + \mathbf{b}]) = \sigma(E[\mathbf{W}^1] \times E[\mathbf{x} \odot \mathbf{m}] + E[\mathbf{b}])$, based on assumption 2.

Focusing on elements of $E[\mathbf{x} \odot \mathbf{m}]$ we have that:

$$\begin{aligned} E[\mathbf{x} \odot \mathbf{m}] &= [E[x^{(1)} m^{(1)}], \dots, E[x^{(n_0)} m^{(n_0)}]]^T = \\ &= [\mu_m^{(1)} \mu_x^{(1)} + \text{cov}(x^{(1)}, m^{(1)}), \dots \\ &\quad \dots, \mu_m^{(n_0)} \mu_x^{(n_0)} + \text{cov}(x^{(n_0)}, m^{(n_0)})]^T \end{aligned}$$

where $x^{(i)}, m^{(i)}$ refers to the i th component of the vector \mathbf{x} and \mathbf{m} respectively.

For the conditional probabilities $p(x_i | m_i)$, and $p(x_i)$ we have under:

$$\text{MCAR: } p(x_i | m_i, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_l) =$$

$$p(x_i | m_i) = p(x_i)$$

$$\text{MAR: } p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_l) =$$

$$p(x_i | m_i, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_l)$$

but

$$p(x_i | m_i) \neq p(x_i)$$

Recall that under the MCAR assumption $\text{cov}(x^{(i)}, m^{(i)}) = 0$, as \mathbf{x}, \mathbf{m} are independent. Under the MAR assumption we are working with, $x^{(i)}, m^{(i)}$ are conditionally independent, conditioned on the values of the rest of the features $x^{(j)}, i \neq j$, and in this case $\text{cov}(x^{(i)}, m^{(i)})$ is not necessarily 0.

In that fashion the elementwise calculations lead to:

$$\begin{aligned} \sigma(E[\mathbf{W}^1] \times E[\mathbf{x} \odot \mathbf{m}] + E[\mathbf{b}]) &= \sigma(E[\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n_1)}] \\ &\quad [\mu_m^{(1)} \mu_x^{(1)} + \text{cov}(x^{(1)}, m^{(1)}), \dots, \mu_m^{(n_0)} \mu_x^{(n_0)} + \\ &\quad + \text{cov}(x^{(n_0)}, m^{(n_0)})]^T + E[\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(n_0)}]) \end{aligned}$$

Finally for elements i, j with $j = 0, \dots, n_0$ and $\mathbf{w}^{(i)} \in \mathbb{R}^{n_0}$:

$$\begin{aligned} \sigma(E[(\mathbf{w}^{(i)}) (\mu_m^{(j)} \mu_x^{(j)} + \text{cov}(x^{(j)}, m^{(j)})) + E[\mathbf{b}^{(j)}]]) \\ \geq \sigma(n_0 \mu_w^1 \mu_m \mu_x + \mu_b^1) + \sigma(n_0 \mu_w^1 \text{cov}(x^{(j)}, m^{(j)})) \end{aligned}$$

Now let $T_1 = \sigma(n_0 \mu_w^1 \text{Cov}(\mathbf{x}^1, \mathbf{m}^1))$, where $\text{Cov}(a, b)$ denotes the vector of element-wise covariances between items of a and b then:

$$E[\mathbf{h}^1] \geq \sigma(n_0 \mu_w^1 \mu_m \mu_x + \mu_b^1) + T_1.$$

This term is propagated through all the network layers, adding bias to the model's output. (Full proof in the Appendix). A straightforward way to tackle the issue would be altering the model's input, by a subtraction debias:

$$\mathbf{x}_{\text{new}} = \mathbf{x} - [\text{cov}(x^{(1)}, m^{(1)}), \dots, \text{cov}(x^{(n)}, m^{(n)})]^T$$

thus alleviating the MAR-derived effect (T_1 term). Subsequently, one could apply sparsity normalisation, as proposed by Yi et al. [1]. This way, the mean of the network's output is not explicitly dependent on the the missingness rate over all data instances. However, that term is intractable since it requires the full features $x^{(i)}$ to be known. Instead, we propose a simple mechanism that takes into account both the input

feature values and their missingness pattern, and evaluates the similarity between data instances with the same, or similar, structure.

Core Idea: Latent factors conditioned on the value and missingness vector, including the covariance between each input feature and the respective missingness indicator can be approximated through transformer-like similarity scores between non-linear embedding projections under the MAR assumption.

Instead of the transformation :

$$\mathbf{x}_{\text{new}} = (\mathbf{x} - \beta)/\alpha$$

where β is the per-feature covariance term and α the per-instance sparsity normalization scaling, we opt for an embedding mechanism. The mechanism is conditioned on the fused value-existence input vectors $(\mathbf{x} \odot \mathbf{m}, \mathbf{m})$, and creates a common representation embedding space $K \subset \mathbb{R}^{emb}$ for the input instances. Multihead self-attention is employed towards that goal, since, it can use the per-instance available context to retrieve the appropriate non-affine transformations from the embedding space ($\mathbf{trans} \in K$) and apply them to the provided input. In this way both the per-instance context (similar to α), and the global context (similar to β) are incorporated in the final solution:

$$\text{lookup}(\mathbf{x}) = \mathbf{trans}$$

$$\mathbf{x}_{\text{new}} = \text{transform}(\mathbf{x}, \mathbf{trans})$$

where the lookup function, $\text{lookup} : \mathbb{R}^N \rightarrow \mathbb{R}^{emb}$ is implemented with multihead self-attention, and the transform function, $\text{transform} : \mathbb{R}^{(N+emb)} \rightarrow \mathbb{R}^N$ by non-linear feed-forward components.

IV. MAIN ALGORITHM

The MAIN algorithms consists of two complementary steps: The Positional Encoded Vector (PEV) creation, and the multihead attention step with opacity gating. Let $\mathbf{x} \in \mathbb{R}^N$ denote an input vector of N features, and x_1, x_2, \dots, x_N the scalar values representing the value of each specific feature in \mathbf{x} . PEV encoding augments the initial feature vector in order to:

- 1) Differentiate between the case of a "filled value", in case of a missing feature, i.e zero in zeta normalization scenarios, and the existing equivalent of that feature value in a feature vector - zero because the feature in that specific instance is equal to the mean of the distribution that generates it.
- 2) Create an orthogonal basis of feature existence. From this perspective, the original feature vector is treated as a linear combination of inter-changeable vectors whose direction is denoted by the feature position and magnitude by their respective feature's value.

Here, for the specified input length N , the minimum number of bits required to describe all available positions, bw is first calculated. Then, for each training example \mathbf{x} the following components are calculated:

- a) A binary mask, indicating the existence of each feature in \mathbf{x} .
- b) A binary positional encoding mask, marking the position of each feature in \mathbf{x} .

Finally, we concatenate a),b) feature-wise with the original input. That results to a 2D input matrix for each training example $\mathbf{x}_{\text{augmented}} \in \mathbb{R}^{N \times \log_2(N+1)}$. In that way, we provide a simple way for the model to capture the similarity between data instances based on both their values and their missingness patterns. The aforementioned PEV procedure is described in Algorithm 1.

Algorithm 1 PEV creation

```

1: function PEVMASKGENERATOR( $x$ )
2:    $bw \leftarrow \lceil \log_2 n \rceil$ 
3:   for  $i \leftarrow 1$  to  $n$  do
4:      $bin_i \leftarrow bin_{bw}(i)$ 
5:      $m_i \leftarrow 0$  if  $x[i] = \emptyset$  else 1
6:      $pev_i \leftarrow [m_i | bin_i]$ 
7:    $pev \leftarrow [pev_1, pev_2, \dots, pev_n]$ 
8:    $m \leftarrow [m_1, m_2, \dots, m_n]$ 
9:   return  $pev, m$ 

```

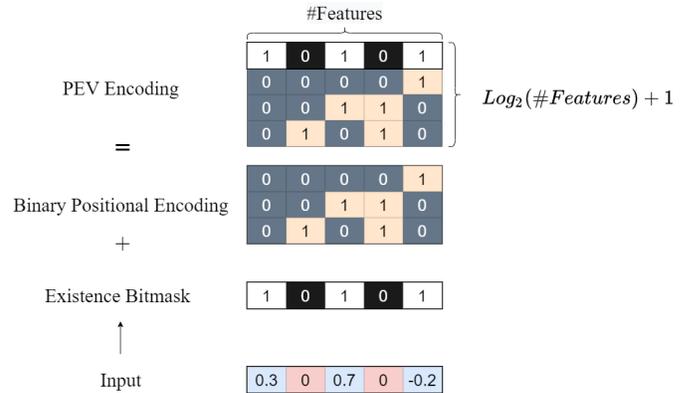


Fig. 1 Creation of PEV encoding Matrix.

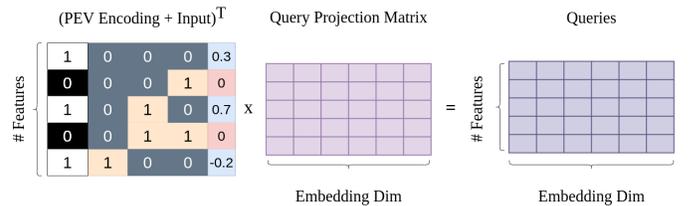


Fig. 2 Creation of Query Projections.

Algorithm 2 Opaque Multi Head Attention

```

1: while convergence is not achieved do
2:   for batch in  $X$  do
3:     for  $(x, y)$  in batch do
4:        $pev, m \leftarrow PevMaskGenerator(x)$ 
5:        $query \leftarrow proj_n(dq([x : pev]))$ 
6:        $key \leftarrow proj_k(dk([x : pev]))$ 
7:        $value \leftarrow proj_n(dv([x : pev]))$ 
8:        $imputed_x \leftarrow mha(query, key, value)$ 
9:        $\hat{x} \leftarrow \gamma \cdot imputed_x + (1 - \gamma) \cdot query$ 
10:      for  $dsm_j$  in Model Tasks do
11:         $\hat{y}_j \leftarrow dsm_j(\hat{x})$ 
12:       $\hat{y} \leftarrow [\hat{y}_0 : \hat{y}_1 : \dots : \hat{y}_l]$ 
13:       $loss \leftarrow e(y, \hat{y})$ 
14:       $batch\ loss \leftarrow \sum loss$ 
15:       $backprop(batch\ loss)$ 
16:   return

```

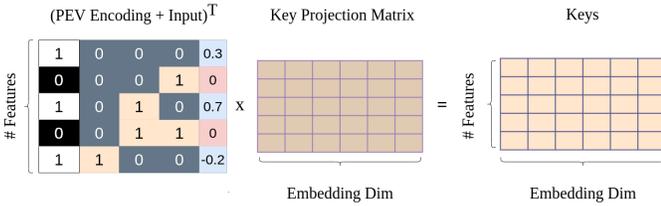


Fig. 3 Creation of Key Projections.

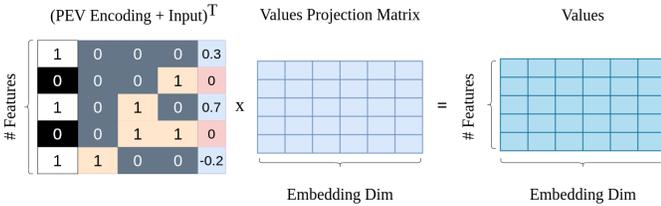


Fig. 4 Creation of Value Projections.

The next step is a Luong-style [14] attention mechanism of each individual augmented feature over the total augmented feature vector, effectively serving as self-attention. Note that the attention function is considered a mapping procedure between a Query and a set of Key-Value pairs [15]. In that fashion, we project the augmented input (PEV + features) into 3 embedding vector spaces, Query, Key and Value, using trainable projections. Here, the PEV augmentation will act as a mechanism capable to guide the focus of the self-attention heads into retrieving saliency from the existence-position perspective. The example of Query projections creation is depicted in Figure 1, and the Keys and Values are created similarly (Figures 3, 2). That of course, could not be possible without the projection of all features to a common embedding space. This is enabled through the use of PEV embeddings, that transform each feature from a mere scalar to a querable vector. This is the key link with the T_N bound of Theorem 2. The correlation between each

feature value and the fact that it is missing, is conditioned on the values of the other features in the original measurement. That is in the end composed out of the each feature's value, and whether it exists or not in the specific instance.

As a next step, we perform multi-head attention over each augmented feature, noted as *mha* in Algorithm 2, resulting in the imputed feature vector *imputed_x*.

Finally, we employ a trainable gating mechanism called opacity gate, $\gamma(\cdot)$. The gate acts as a "trainable sigmoid knob" between the original input of the model, and the imputed input produced by the MAIN mechanism. The gate is unconditioned on the input features, and initialized at 0.5, giving equal importance to both inputs. During training, the gate is left free to decide through backprop whether our method provides an input more "salient" towards minimizing the downstream loss (Gate \simeq 0) or if the original input is better (Gate \simeq 1). Furthermore, the gate acts in a twofold manner, providing an abstraction buffer towards the next layers. During the first iterations, the next layers will try to adapt towards an unoptimal solution to the task, based on the original inputs while the MAIN mechanism is still in early training stage. This is not permanent however, since the gate will shift towards the MAIN inputs, "fine-tuning" the rest of the model into the optimal solution.

The algorithm finishes with collecting all the opacity weighted feature vectors, \hat{x} and feeding them to an arbitrary number of downstream models, noted as *dsm* in Algorithm 2. They are jointly trained with the MAIN component to perform the required predictive tasks, and during this analysis are all considered to be 1, since no-multitask dataset was used.

V. EXPERIMENTS

Our method was evaluated on 7 Tabular datasets, 6 of which come from the UCI Dataset Repository and 1 maritime dataset related to vessel's performance which is property of DeepSea Technologies (i.e DeepSea V9). For all experiments, in order to be comparable with other methods that simulate missingness, binary sampling was performed per feature using the same technique as in [1] and [11]. Our method though is not bound to this restriction and can work directly on missing data, exploiting the full capacity of the dataset. The downstream task in the UCI datasets was binary classification with imbalanced classes, while in the maritime dataset was regression. It is noteworthy to outline here that out of the 6 datasets, Breast, Credit, Spam and Heart are injected with MCAR missingness at 4 different levels: [20%, 40%, 60%, 80%], while Pima, Mammographic and DeepSea V9 have intrinsic MAR missingness at levels: [12.2%, 3.35%, 7.28%]. In order for all the experiments to be compared at a common basis, the full rows of those datasets were polluted so that the total missing rate matched the [20%, 40%, 60%, 80%] scale, as a mixture of MCAR + MAR missingness.

It is important to outline that the scope of the experiments is towards to a **single-step** end-to-end training scheme. We found that the comparison with classical **two-step**, methods that first-impute then train/test formulations (i.e KNN, MICE,

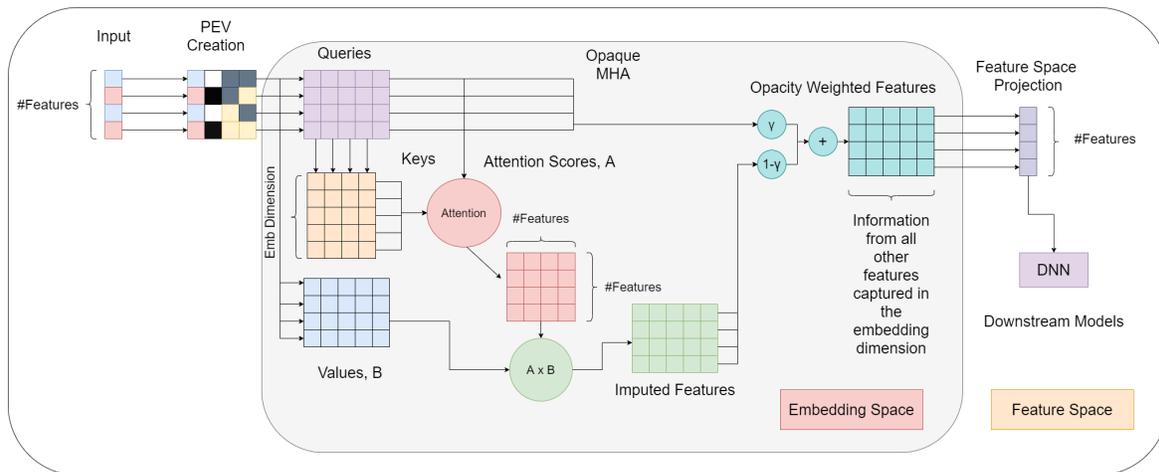


Fig. 5 An overview of MAIN’s architecture. It is comprised of three basic modules: a) PEV creation, b) Opaque Multi-head Attention, c) Downstream Tasks/Models

MCMC) is not only less relevant to the scope of this work but also unfair to them since they don’t optimize for a specific task and give inferior results. For the sake of completeness however, we opt to compare with the **two-step tabular** method, GAIN [11], since its authors shared our common interest of optimizing towards a downstream task. We also compared with SN [1] as our method builds on the same ideas and can be seen as an extension to theirs. ZIMC (i.e zero imputation with mask concat) is serving as a baseline and FiLM [2] although not an imputation method, since it also inspired our work can be seen as an intermediate step between SN and our work MAIN. We choose not to compare with AimNet [3], as there was no official codebase, making it very difficult to replicate results. Furthermore, it was a DataBase-Oriented imputation method for the project HoloClean, not a general purpose Machine Learning Solution like the other compared methods.

The experimental design is similar to [1] and [3], reporting the test AUROC (UCI) and MSE (DeepSea V9) score of 5 runs. In order to keep the comparison fair, the GAIN [11] method was trained and tested on different splits of the dataset unlike the original implementation where training and evaluation was in the same set. The Breast, Spam, Credit and DeepSea V9 datasets were normalized in the range of [0,1], while Heart, Pima were normalized in the range [-1,1] as in [1], [11].

We used no explicit preprocessing for categorical values in the above datasets since our scheme: a) creates automatically trainable embeddings for both scalar and categorical/ordinal variables and b) we don’t reconstruct the missing feature and thus we don’t have to map reconstructed logits to corresponding classes as in [11].

Regarding train/validation/test splits, in most imputation-only approaches the RMSE reconstruction metric is usually reported on a single dataset without any splits. Since no official splits exist in any UCI datasets, instead of the 70-20-10 split used in [1] we opt for a 70-30 split with 10-fold stratified cross-validation.

Models Setup: In order to produce comparable results in

UCI and test the performance gains of our imputation method as a base layer, we tried to keep the total number of layers (imputation + downstream) of our scheme close to what proposed in [16] as the most appropriate for the UCI and also used in [1]; namely 4 Hidden layers @ 256 units and Adam Optimizer. The same principle applied to all the compared schemes where we had: imputation method + 4 layers dedicated to the downstream task.

Due to class imbalance, in our training scheme we pre-calculate the class weights of the target classes and use them into a weighted binary crossentropy loss, penalizing misclassification on the minority class more heavily than misclassifications on the majority class. Finally, the reported metric is the AUROC curve interpolated at 200 points with Riemann summation method.

TABLE I : UCI Breast-Wisconsin

| MR | 20% | 40% | 60% | 80% |
|-------------|-----------------------|-----------------------|-----------------------|-----------------------|
| ZIMC | 0.9501 ± 0.005 | 0.9485 ± 0.007 | 0.9182 ± 0.010 | 0.8470 ± 0.012 |
| GAIN* | 0.9872 ± 0.008 | 0.9475 ± 0.011 | 0.9171 ± 0.035 | 0.8443 ± 0.017 |
| SN | 0.9683 ± 0.007 | 0.9341 ± 0.009 | 0.8593 ± 0.029 | 0.8640 ± 0.027 |
| FiLM | 0.9818 ± 0.003 | 0.9513 ± 0.005 | 0.9263 ± 0.008 | 0.8757 ± 0.015 |
| MAIN | 0.9821 ± 0.003 | 0.9786 ± 0.005 | 0.9693 ± 0.013 | 0.9241 ± 0.014 |

TABLE II : UCI Credit Dataset

| MR | 20% | 40% | 60% | 80% |
|-------------|-----------------------|-----------------------|-----------------------|-----------------------|
| ZIMC | 0.7297 ± 0.008 | 0.7051 ± 0.007 | 0.6833 ± 0.010 | 0.6349 ± 0.013 |
| GAIN* | 0.7412 ± 0.008 | 0.7173 ± 0.013 | 0.6849 ± 0.008 | 0.6019 ± 0.032 |
| SN | 0.7396 ± 0.005 | 0.7145 ± 0.007 | 0.6826 ± 0.003 | 0.6332 ± 0.012 |
| FiLM | 0.7443 ± 0.005 | 0.7187 ± 0.007 | 0.7025 ± 0.005 | 0.6528 ± 0.011 |
| MAIN | 0.7456 ± 0.004 | 0.7209 ± 0.006 | 0.7032 ± 0.004 | 0.6577 ± 0.009 |

As Tables I-IV show, both GAIN and SN provide a far better alternative to the ZIMC baseline in every dataset and missing rate. While FiLM is of subpar performance towards both the generative GAIN and the discriminative SN in most cases, it is

Original Input Vector / Missing Vector Bitmask



Fig. 6 An illustrated example of a training example with missing features. (feat.value/missing)

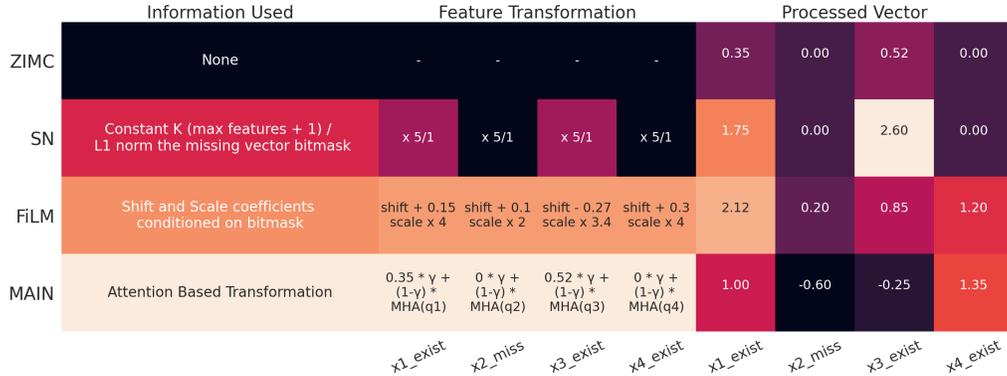


Fig. 7 An illustrated example of the methods under test.

TABLE III : UCI Spam Dataset

| MR | 20% | 40% | 60% | 80% |
|-------------|-----------------------|-----------------------|-----------------------|-----------------------|
| ZIMC | 0.9740 ± 0.005 | 0.9519 ± 0.007 | 0.9243 ± 0.008 | 0.8675 ± 0.011 |
| GAIN* | 0.9645 ± 0.004 | 0.9451 ± 0.008 | 0.9213 ± 0.001 | 0.8623 ± 0.010 |
| SN | 0.9798 ± 0.002 | 0.9568 ± 0.002 | 0.9270 ± 0.002 | 0.8707 ± 0.005 |
| FiLM | 0.9735 ± 0.005 | 0.9571 ± 0.006 | 0.9276 ± 0.007 | 0.8725 ± 0.012 |
| MAIN | 0.9817 ± 0.005 | 0.9589 ± 0.007 | 0.9298 ± 0.008 | 0.8732 ± 0.011 |

TABLE VII : DeepSea V9

| MR | 20% | 40% | 60% | 80% |
|-------------|------------------------|------------------------|------------------------|------------------------|
| ZIMC | 0.0013 ± 0.0005 | 0.0030 ± 0.0006 | 0.0124 ± 0.0010 | 0.0280 ± 0.0013 |
| GAIN* | 0.0009 ± 0.0002 | 0.0027 ± 0.0002 | 0.0086 ± 0.0001 | 0.0267 ± 0.0004 |
| SN | 0.0012 ± 0.0002 | 0.0048 ± 0.0003 | 0.0116 ± 0.0005 | 0.0310 ± 0.0009 |
| FiLM | 0.0010 ± 0.0004 | 0.0035 ± 0.0005 | 0.0091 ± 0.0005 | 0.0270 ± 0.0009 |
| MAIN | 0.0009 ± 0.0003 | 0.0027 ± 0.0002 | 0.0085 ± 0.0002 | 0.0261 ± 0.0008 |

TABLE IV : UCI PIMA DATASET

| MR | 20% | 40% | 60% | 80% |
|-------------|-----------------------|-----------------------|-----------------------|-----------------------|
| ZIMC | 0.8131 ± 0.016 | 0.7864 ± 0.014 | 0.7585 ± 0.015 | 0.7001 ± 0.015 |
| GAIN* | 0.8074 ± 0.017 | 0.7861 ± 0.020 | 0.7503 ± 0.023 | 0.6987 ± 0.022 |
| SN | 0.8121 ± 0.013 | 0.7851 ± 0.021 | 0.7589 ± 0.017 | 0.7006 ± 0.015 |
| FiLM | 0.8017 ± 0.017 | 0.7725 ± 0.015 | 0.7475 ± 0.016 | 0.6913 ± 0.010 |
| MAIN | 0.8442 ± 0.009 | 0.7981 ± 0.007 | 0.7628 ± 0.011 | 0.7183 ± 0.012 |

TABLE V : UCI Heart Dataset

| MR | 20% | 40% | 60% | 80% |
|-------------|-----------------------|-----------------------|-----------------------|-----------------------|
| ZIMC | 0.8256 ± 0.009 | 0.8002 ± 0.012 | 0.6752 ± 0.028 | 0.6702 ± 0.033 |
| GAIN* | 0.8419 ± 0.039 | 0.7586 ± 0.065 | 0.6399 ± 0.102 | 0.6357 ± 0.094 |
| SN | 0.8533 ± 0.007 | 0.7974 ± 0.012 | 0.7570 ± 0.007 | 0.6619 ± 0.014 |
| FiLM | 0.8530 ± 0.007 | 0.7915 ± 0.014 | 0.7516 ± 0.009 | 0.6422 ± 0.017 |
| MAIN | 0.8702 ± 0.008 | 0.8371 ± 0.007 | 0.7851 ± 0.008 | 0.6802 ± 0.011 |

TABLE VI : UCI Mammographic Dataset

| MR | 20% | 40% | 60% | 80% |
|-------------|-----------------------|-----------------------|-----------------------|-----------------------|
| ZIMC | 0.7851 ± 0.014 | 0.7433 ± 0.027 | 0.6712 ± 0.026 | 0.6049 ± 0.026 |
| GAIN* | 0.8010 ± 0.024 | 0.7501 ± 0.036 | 0.7126 ± 0.028 | 0.6448 ± 0.057 |
| SN | 0.7794 ± 0.011 | 0.7367 ± 0.021 | 0.6753 ± 0.014 | 0.6063 ± 0.019 |
| FiLM | 0.8198 ± 0.012 | 0.7472 ± 0.023 | 0.6860 ± 0.042 | 0.6350 ± 0.032 |
| MAIN | 0.8807 ± 0.005 | 0.8535 ± 0.009 | 0.7998 ± 0.012 | 0.7682 ± 0.014 |

noteworthy, that in the case of the heavy MAR dataset (Pima), FiLM gives worse results to SN, while with none of the current

state-of-the-art methods solving explicitly the MAR case.

On the contrary, the MAIN method seems to be vastly outperforming other methods, especially in the case of high missing rates, where it maintains considerably higher scores than its counterparts. In the case of the DeepSea V9 regression task, GAIN and MAIN have the least test MSE and perform equally in the low missing scenarios, while MAIN slightly outperforms GAIN in the high missing scenarios.

VI. DISCUSSION

In this section we will give an illustrated example of a partially missing feature vector in order to provide more insights on how our method works compared to other well-known methods. The setup is the following: in figure 6, we present a z-normalized 4D feature vector that will be provided as an input to the different methods that were compared against our final solution, MAIN.

The 4D input vector consists of 3 non-missing features: $[x1, x3, x4]$ with corresponding values $[0.35, 0.52, 0.00]$, and one missing feature: $x2$. The respective existence bitmask of the input is depicted with a binary indicator after the "slash" symbol (/) of each feature, where '1' means missing feature, thus the corresponding mask is $m = [0, 1, 0, 0]$.

In figure 7 we break-down how the input information is potentially transformed (i.e Feature Transformation) by 4

different approaches, leading to a the Processed Vector on the right in the same figure. The processed vector can be seen as the result of a learnable scaling factor (for the most algorithms) which forms the input to the following-up machine learning method (e.g downstream model). The actual numbers are figurative just to illustrate the underlying mechanics of each method.

We observe that in the case of ZIMC as is the most naive method features are not explicitly affected since it is a mere zero imputation technique. It must be noted here, that the mask concatenation in ZIMC, potentially provides to the underline model with joint knowledge of missingness but that could happen in a deeper layer interactions in the case of neural nets. The even simpler zero-imputation couldn't discriminate at all, between a missing feature from a feature placed at the mode of the distribution in the z-scaled space.

In SN, the inverse L1 norm of the bitmask scales all the non-zero features, by multiplying by their number + 1 (K) which is an improvement over zero-fill. On the other hand the scaling of existing features at the mode of the distribution (i.e zero-mean like x_4) are affected in the same way as the missing ones (i.e x_2), re-scaling only the rest of them in order to alleviate the VSP problem[1]. In FiLM, since all the features are subject to affine transformation conditioned in the missing context (feature + mask), the missing x_2 and non-missing x_4 zero-valued features are treated differently. Thus we claim that FiLM can be seen as a feature-wise rescaling extension of the SN which learns how to condition missing context to independent feature scalars. Following this logic and going one step further, if we explicitly define the conditioning mechanism (i.e attention) we have MAIN which offers even more expressivity to learn the underline mechanics of the missing distribution.

VII. CONCLUSION

In this work, we proposed MAIN a novel method based on multi-head attention to deal with missing data in continuous or discrete datasets. Our method works in a single step by implicitly imputing missing data and is optimized directly on the downstream task, offering an end-to-end trainable system. We demonstrate that MAIN significantly outperforms state-of-the-art methods in a variety of open dataset and also in a proprietary one.

REFERENCES

- [1] Joonyoung Yi, Juhyuk Lee, Kwang Joon Kim, Sung Ju Hwang, and Eunho Yang. Why not to use zero imputation? correcting sparsity bias in training neural networks. In *7th International Conference on Learning Representations*, 2020.
- [2] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. Feature-wise transformations. *Distill*, 2018.
- [3] Richard Wu, Aoqian Zhang, Ihab Ilyas, and Theodoros Rekatsinas. Attention-based learning for missing data imputation in holoclean. 2020.
- [4] Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*. John Wiley & Sons, 2014.
- [5] Steven Cheng-Xian Li, Bo Jiang, and Benjamin M. Marlin. Misgan: Learning from incomplete data with generative adversarial networks. In *7th International Conference on Learning Representations*, 2019.
- [6] Pierre-Alexandre Mattei and Jes Frellsen. MIWAE: Deep generative modelling and imputation of incomplete data sets. In *Proceedings of the 36th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 4413–4423, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- [7] Alfredo Nazábal, Pablo M. Olmos, Zoubin Ghahramani, and Isabel Valera. Handling incomplete heterogeneous data using vaes. *CoRR*, 2018.
- [8] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [9] Ashish Bora, E. Price, and A. Dimakis. Ambientgan: Generative models from lossy measurements. In *ICLR*, 2018.
- [10] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. In *Advances in Neural Information Processing Systems 31*, pages 10727–10737. 2018.
- [11] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. Gain: Missing data imputation using generative adversarial nets. In *6th International Conference on Learning Representations*, 2018.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [13] Marek Śmieja, Łukasz Struski, Jacek Tabor, Bartosz Zieliński, and Przemysław Spurek. Processing of missing data by neural networks. In *Advances in Neural Information Processing Systems 31*, pages 2719–2729. 2018.
- [14] Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Lisbon, Portugal, 2015.
- [15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30*. 2017.
- [16] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, Long Beach, CA, USA, 2017*.