

Missing Value Imputation: With Application to Handwriting Data

Zhen Xu and Sargur N. Srihari

State University of New York at Buffalo, Buffalo, USA

ABSTRACT

Missing values make pattern analysis difficult, particularly with limited available data. In longitudinal research, missing values accumulate, thereby aggravating the problem. Here we consider how to deal with temporal data with missing values in handwriting analysis. In the task of studying development of individuality of handwriting, we encountered the fact that feature values are missing for several individuals at several time instances. Six algorithms, i.e., random imputation, mean imputation, most likely independent value imputation, and three methods based on Bayesian network (static Bayesian network, parameter EM, and structural EM), are compared with children's handwriting data. We evaluate the accuracy and robustness of the algorithms under different ratios of missing data and missing values, and useful conclusions are given. Specifically, static Bayesian network is used for our data which contain around 5% missing data to provide adequate accuracy and low computational cost.

Keywords: Missing Value Imputation, Bayesian Network, Parameter EM, Structural EM

1. INTRODUCTION

Missing values are pervasive and important in pattern recognition and machine learning. It has been described as "one of the most important statistical and design problems in research".¹ In real life learning applications, the available data always contain missing values. In some cases, data are missing by human mistake or by accident. For example, people may forget to record the data at certain time. In addition, temporary power outage and sensor failure, or even communication channel block may generate missing values. In other cases, data are missing systematically or constantly. For instance, in medical diagnosis, patients usually do not take all possible tests. Moreover, some dust and scratches on the lens can produce missing values until they are being taken care of. Although missing values are common, many existing multivariate analysis methods, such as principle component analysis² and singular value decomposition,³ require complete data. There are two simple ways to solve this problem. Firstly, extra human efforts are required for missing values imputation. For example, the doctor can ask patients to take extra tests. This is the most accurate yet expensive way, and it only suits for important applications with relatively static nature and small size. Specifically, we cannot observe the original value at different times for time sensitive applications. Secondly, data with missing values are simply discarded.⁴ This is a very wasteful manner dealing with valuable data. One missing pixel within one million makes the whole image useless. In addition, the wastefulness is aggravated in case of longitudinal research,⁵ which involves multiple observations on the same object over a long time period. In conclusion, we need a good solution for missing value imputation with adequate accuracy and moderate computational cost.

In our longitudinal research, the children's handwriting data, including samples for the same student at different grades, are used for studying the development of handwriting individuality along the time.⁶ It studies children's handwriting traits' developmental trend from the time they begin to learn writing to when the traits are established. The handwriting data are sampled from a large number of students beginning from 2nd grade to 4th grade. Within this time period, the students either begin to learn or just have learned how to write. The collection occurs at every spring, and records how students' handwriting habits develop as handwriting skills continue to mature. Writing samples collected for the same student from different years are examined and studied. The word *and* is chosen because it is one of the few words that children write frequently and repeatedly. In the process of collecting samples, all students received the same instructions. In addition, all writing samples are gathered under the same forms to control variations coming from document size and format. After collection,

Further author information: (Send correspondence to Sargur N. Srihari)
Sargur N. Srihari: E-mail: srihari@cedar.buffalo.edu

X_1	# strokes in "a"	one		two		three		uppercase		X_1	Initial stroke of "a"	staff right		staff left		staff center			
X_2	formation of "a" staff	tented		retraced		looped		no staff		X_2	Formation of "a" staff	tented		retraced		looped		no staff	
X_3	# strokes in "n"	one		two		three		uppercase		X_3	Number of "n" arches	one		two					
X_4	formation of "n" staff	tented		retraced		looped		no staff		X_4	Shape of "n" arches	pointed		rounded		retraced		combination	
X_5	shape of arch of "n"	pointed		rounded						X_5	Location of "n" mid	above base		below base		at base			
X_6	# strokes in "d"	one		two		three		uppercase		X_6	Formation of "d" staff	tented		retraced		looped			
X_7	formation of "d" staff	tented		retraced		looped		no staff		X_7	Formation of "d" initial	overhand		underhand		straight across			
X_8	initial stroke of "d"	staff top		bulb						X_8	Formation of "d" terminal	curved up		straight		curved down		no obvious end stroke	
X_9	Unusual formations	formation		symbol						X_9	Symbol	unusual				symbol			
X_{10}	a-n relationship	a taller		a equal		a smaller				X_{10}	a-n relationship	a taller		a equal		a smaller			
X_{11}	a-d relationship	a taller		a equal		a smaller				X_{11}	a-d relationship	a taller		a equal		a smaller			
X_{12}	n-d relationship	n taller		n equal		n smaller				X_{12}	n-d relationship	n taller		n equal		n smaller			

Figure 1: Twelve features of handwriting data. (a) Printed *and*. (b) Cursive *and*.

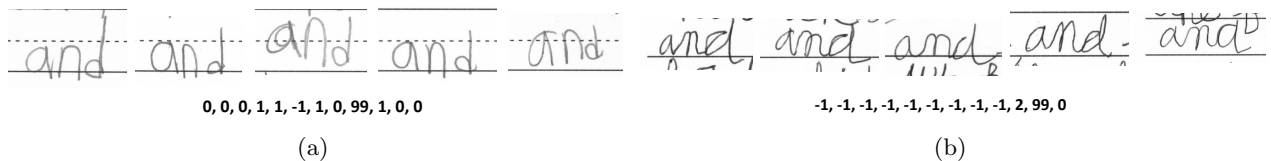


Figure 2: Examples of two sets of images corresponding to two feature values, which contain incommensurate (99) and missing (-1) values. (a) Printed *and*. (b) Cursive *and*.

the questioned document examiners (QDEs) manually assign different features to handwriting. Specifically, the word *and* is represented by 12 features $\mathbf{X} = \{X_i\}$, $i = 1, \dots, 12$, as shown in Figure 1. Ordinal values are assigned to each feature as $\{0, 1, \dots, k_i - 1\}$ with two exceptions: the value -1 indicates the missing value and the value 99 shows the inconsistency within collected samples. Each feature value is extracted from five samples of the same person for the reason of stability. In collected data, we encountered the problem that some feature values are missing for several individuals at several time instances. In addition, there are some incommensurate values within the data. Some image examples are shown in Figure 2. Because our longitudinal research is intended to continue for 11 years, and follows students all the way until they graduate from high school, missing value problem becomes very important. Each year, even when there is a very small amount of missing values, 11 years' accumulation will contribute to a large portion of the data. The missing and incommensurate values deteriorate the qualities of data and the results of analysis algorithms, especially for distance related ones, as shown in Figure 3. So our goal is to preprocess these data to make them complete and commensurate for future use.

In this paper, we present six missing value imputation algorithms and compare their performance in our children's handwriting data. The proposed Bayesian network framework unifies missing value imputation with other inference tasks. In addition, we conduct extensive experiments and give useful conclusions on different ratios of missing data and missing value. These different scenarios suit for longitudinal research with different lengths. It is a fairly complete comparative study of missing value imputation problem. Specifically, the missing data are defined as data with missing values, while the complete data are data without missing values. The whole data contain complete data and missing data. The missing data ratio measures the percentage of missing data within the whole data while the missing value ratio measures the percentage of missing values within the missing data.

2. BASIC ALGORITHMS

In this section, we introduce some basic algorithms for data commensuration and missing value imputation. Firstly, we solve the data commensuration problem. It is a simple yet meaningful move. Secondly, we introduce three highly efficient imputation algorithms. These algorithms assume all features of the data are independent from each other. They can be used when there are large data sets yet with very limited computational power.

	Year	Grade	Sample ID	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	X ₁₁	X ₁₂
1	2012	1	MaL2023Kr	0	0	0	1	-1	0	1	0	99	0	0	0
2	2012	1	OIC2023Ka	-1	0	0	1	1	0	3	1	99	0	99	99
3	2012	4	EtM2020Ga	0	1	1	5	-1	0	3	0	99	2	0	0
4	2013	2	BrB2023Bu	0	0	0	0	-1	0	1	1	99	0	0	0
5	2013	2	CaG2023Ma	0	-1	0	1	1	0	1	0	99	99	0	0
6	2013	2	DaM2023Sa	0	-1	0	99	1	0	1	1	99	0	0	0
7	2013	2	LiD2023Ru	0	-1	0	1	1	0	99	1	99	2	2	0
8	2012	2	ViG2022Kl	0	0	0	1	1	1	3	-1	0	0	0	0
9	2013	3	BrE2022Ba	0	0	0	1	1	-1	1	0	99	1	0	0
10	2013	3	LuP2022La	2	4	0	1	1	2	3	-1	99	0	0	0
11	2013	3	MoA2022Ch	0	0	0	1	0	1	3	-1	99	99	0	0
12	2013	3	JoM2022Le	0	1	0	1	1	0	-1	1	2	0	0	2
13	2012	4	AlE2020BR	5	99	0	1	1	0	3	-1	0	0	0	0
14	2012	4	LaS2020Gu	0	1	0	1	1	-1	3	1	99	0	0	0
15	2013	5	An_2020Wo	0	99	0	1	1	0	-1	1	99	2	0	0
16	2013	2	IsS2023We	0	0	0	0	0	0	3	0	99	-1	-1	-1
17	2013	2	ZoA2023Wa	0	-1	0	-1	-1	0	1	1	99	0	0	0
18	2012	3	EIC2021Sh	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	2
19	2012	3	IsR2021Br	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	2
20	2012	3	Jar2021Be	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
21	2012	3	KaE2021Re	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
22	2012	3	KeM2021Fl	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
23	2012	3	LeR2021Ha	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	1	2
24	2012	3	OwH2021Br	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
25	2012	3	Sam2021Ri	-1	-1	-1	-1	-1	-1	-1	-1	1	-1	-1	-1

(a)

	Year	Grade	Sample ID	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	X ₁₁	X ₁₂
1	2013	3	AlM2022Ca	1	-1	1	0	0	99	99	99	99	0	0	0
2	2013	3	AlK2022Ko	1	99	1	1	0	-1	0	0	99	2	99	99
3	2013	3	JoC2022Kl	-1	1	1	0	1	1	0	0	99	0	0	2
4	2013	3	JoO2022Fi	2	1	1	1	0	1	-1	2	99	0	0	0
5	2013	3	Ma_2022Ka	1	-1	0	0	0	1	0	0	99	2	0	0
6	2013	3	MaG2022Po	-1	1	1	1	1	1	0	0	99	0	0	0
7	2013	3	NaM2022We	2	1	1	1	1	2	-1	0	0	99	99	99
8	2013	3	SaJ2022Ad	-1	1	1	99	2	1	0	0	99	99	99	99
9	2013	4	SoA2021He	1	1	1	0	0	-1	1	3	2	2	0	0
10	2013	3	CaR2022Ka	-1	-1	1	3	99	99	0	0	99	99	99	99
11	2013	4	Cyl2021Va	2	1	0	0	99	1	0	0	99	-1	-1	-1
12	2013	4	DyB2021Lo	1	1	1	1	0	1	1	0	99	-1	-1	-1
13	2013	4	EmM2021Qu	2	1	1	1	0	1	0	2	99	-1	-1	-1
14	2012	3	AsM2021Gr	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
15	2012	3	MaR2012Sh	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
16	2012	3	MiG2021Pa	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	2
17	2012	3	NaL2021Ge	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	2	2
18	2012	4	AlG2020Fi	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
19	2012	4	AnG2020Ma	-1	-1	-1	-1	-1	-1	-1	-1	-1	2	99	0
20	2012	4	EMM2020VO	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
21	2012	4	HaK2020Sm	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
22	2012	4	IrO2020Yo	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	2	2
23	2012	4	JAC2020EG	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	0
24	2012	4	Mal2020Ca	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	0	2
25	2013	4	ClK2021Ko	-1	-1	-1	-1	-1	-1	-1	-1	-1	1	-1	-1
26	2013	5	JoW2020Li	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
27	2012	4	KaR2020Re	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
28	2012	4	LaE2020Jo	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
29	2012	4	MAE2020Mi	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
30	2012	4	OIE2020Fr	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
31	2012	4	SaL2020Ha	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
32	2012	4	StW2020Ba	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
33	2012	4	TyP2020Hu	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
34	2012	4	TyR2020CH	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1

(b)

Figure 3: Examples of missing and incommensurate feature values for (a) printed, and (b) cursive data, where 99 means inconsistent value and -1 means missing value.

2.1 Data Commensuration

As we have mentioned in the last section, there are some incommensurate 99 values in our data, which deteriorate the distance related data analysis algorithms. In order to solve this problem, we treat 99 value as a new commensurate feature value. For example, for feature X_i with k_i values, we reassign the 99 value as k_i . It is reasonable because the 99 value does not fit in any defined category.

2.2 Random Imputation (Rand)

In random imputation, we randomly choose discrete values within the feasible scope to fill in the missing values. Specifically, $Rand_{X_i} = \forall x_i \in Val(X_i)$, where $Val(X_i) = \{0, 1, \dots, k_i - 1\}$. If feature X_i contains incommensurate 99 value, the scope of X_i contains one more value k_i .

2.3 Mean Imputation (Mean)

Even though random imputation is very efficient and robust, we want to give more educated guess. In mean imputation,⁷ we use the mean values of the features to fill in the missing values. It can be computed as $Mean_{X_i} = \frac{\sum_{m=1}^N X_{im}}{N}$. If the computed mean is not a discrete value, we use nearest discrete value as the mean.

2.4 Most Likely Independent Value Imputation (MLIV)

Intuitively, the mean imputation should be better than random imputation. However, there are still many cases that mean value does not work very well. For example, when the underlying probability distribution is bimodal with equal probability mass, the mean value will lie in the middle and it will give a bad imputation. Motivated by this example, our third algorithm, most likely independent value imputation, uses the most probable value of each feature to fill in the missing values. It can be computed as $MLIV_{x_i} = argmax_{x_i} M[x_i]$, where $M[x_i]$ is the number of times $X_i = x_i$ in data \mathcal{D} . This is our best solution when all features of the data are independent.

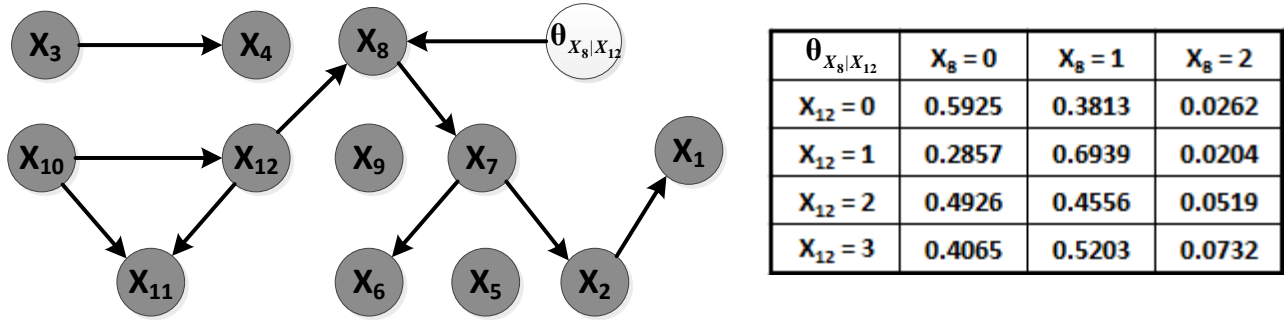


Figure 4: An example of Bayesian network that is composed of structure and parameters. Only one parameter $\theta_{X_8|X_{12}}$ is shown for the interest of space.

2.5 Correlation & Imputation

In real life, it is a rare case that all features are independent. Instead, correlation widely exists within data. It is a very powerful and helpful tool for predictions. For example, we can even use the frequency of storks to predict the human birth rate.⁸ Naturally we can use the observed values and the correlation between observed and missing features to perform a better imputation. For example, if we know one student is smart and works hard, then there is a very high probability that he gets an A in an easy course even when we miss his grade.

In recent years, some machine learning approaches, such as K-nearest neighbors and single value decomposition based method, have been used for missing value imputation.^{9,10} Most of them implicitly use the correlation for imputation. However, these methods highly depend on tuning parameters, which is a very hard task without prior knowledge. The quality of parameter tuning significantly affects the performance. More importantly, because the imputed data are used for probabilistic inference in the future, i.e., answer conditional probability queries, we want a coherent framework for both missing value imputation and probabilistic inference. Intuitively, by this way, the imputed value will give better inference performance. Here we introduce a Bayesian network (BN) based framework for missing value imputation. It explicitly explores the correlations to improve imputation accuracy. BN is probabilistic graphical model that compactly represents a joint distribution and corresponding conditional independence assumptions for a set of variables.^{3,11} The BN structure is a directed acyclic graph (DAG) where nodes represent variables and directed edges represent statistical dependent relationships. The BN structure encodes the conditional independency relationships and BN parameters specify conditional probability distributions (CPDs). Both graph structure and parameters together determine the joint distribution. The learned distribution can be used for missing value imputation by maximum likelihood principle. The BN is chosen for several reasons. Firstly, BN is an intuitive way and a general framework of representing and visualizing the statistical relationships between different variables. Secondly, BN is a compact way of encoding conditional independent relationships and it scales well with high dimensional data. Thirdly, BN allows efficient inference algorithms, which can be further used to perform temporal analysis for children’s handwriting in different grades.

3. BAYESIAN NETWORK BASED MISSING VALUE IMPUTATION ALGORITHMS

3.1 Bayesian Network

3.1.1 Representation

BN can compactly represent a joint distribution and encode the corresponding conditional independencies.^{3,11} An example of BN is shown in Figure 4. Edge $X_3 \rightarrow X_4$ means X_3 is a parent of X_4 , and variable X_4 has statistical dependence relationship with variable X_3 . Accordingly, X_4 is a child of X_3 . It is an intuitive way to visualize the statistical relationships and the DAG structure compactly encodes the conditional independencies. For example, we can find X_4 depends on X_3 , and X_8 is independent of X_{11} when conditioning on X_{12} . These dependencies and independencies are hard to find and represent in the joint distribution or data itself. Because of conditional independencies, the joint distribution can be decomposed as the product of CPDs, which can be

represented as parameters of BN.¹¹ One example CPD $\theta_{X_8|X_{12}} = P(X_8|X_{12})$ is given in Figure 4, which is a tabular representation of CPDs. Both graph structure \mathcal{G} and parameters θ of BN together determine the joint distribution.

3.1.2 Score Function

One widely used score function is negative log-likelihood, which is also called as log-loss.¹¹ It is one way of measuring how good the specific BN fits the data. Given a data set \mathcal{D} with N i.i.d samples, the log-loss for a BN with structure \mathcal{G} and parameters θ is defined as:

$$l(\mathcal{D}|\mathcal{G}, \theta) = - \sum_{m=1}^N \log P(\mathbf{X}_m) = - \sum_{i=1}^n \sum_{m=1}^N \log P(X_{im}|\mathbf{pa}_{im}) = - \sum_{i=1}^n \sum_{m=1}^N \log \theta_{X_{im}|\mathbf{pa}_{im}}, \quad (1)$$

where \mathbf{X}_m is a single data point that contains the values of all variables, X_{im} is the value of i_{th} variable in the m_{th} sample, \mathbf{pa}_{im} are values of parents of X_i in the m_{th} sample, n is the number of variables, and N is the number of samples. $\theta_{X_{im}|\mathbf{pa}_{im}} = P(X_{im}|\mathbf{pa}_{im})$, which is the BN parameter and represents a single element of CPD for variable X_i . The log-loss measures the loss of using the BN to represent the true probability distribution. Lower log-loss score signifies a better BN. We need to learn both the structure and parameters to construct a good BN.

3.1.3 Parameter Learning

Given a BN structure \mathcal{G} , we need to learn the parameters to fit the data well. For notation simplicity, we define $x_i \in Val(X_i)$ and $\mathbf{p}_i \in Val(\mathbf{pa}_{X_i})$, which are possible values of X_i and possible values of parents of X_i , the log-loss in Equation 1 can be redefined as:

$$l(\mathcal{D}|\mathcal{G}, \theta) = - \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{p}_i} \log \theta_{x_i|\mathbf{p}_i}^{M[x_i, \mathbf{p}_i]}, \quad (2)$$

where $M[x_i, \mathbf{p}_i]$ is the number of times $X_i = x_i$ and $\mathbf{pa}_{X_i} = \mathbf{p}_i$ in \mathcal{D} and $\theta_{x_i|\mathbf{p}_i}$ is the corresponding conditional probability. The score transformation is got by grouping all the occurrences of $\theta_{x_i|\mathbf{p}_i}$ over the whole data. For each \mathbf{p}_i , $\sum_{x_i} \theta_{x_i|\mathbf{p}_i} = 1$. Considering these constraints using Lagrange multiplier, the maximum likelihood estimation (MLE) of parameters with respect to the log-loss score can be computed as $\hat{\theta}_{x_i|\mathbf{p}_i} = \frac{M[x_i, \mathbf{p}_i]}{M[\mathbf{p}_i]}$, where $M[\mathbf{p}_i] = \sum_{x_i} M[x_i, \mathbf{p}_i]$. So $M[x_i, \mathbf{p}_i]$ are sufficient statistics to compute the MLE parameters.

3.1.4 Structure Learning

Given data, learning the BN structure is a harder problem than parameter learning. Note that after learning the structure, BN parameters are also determined given complete data as we saw in the previous section. Automatic BN structure learning is a very active research area, and there are a large number of papers working on this topic. There are mainly three approaches. Firstly, the constraint-based approach^{12, 13} starts from a complete graph and adopts the independent tests to find a set of conditional independencies which can be used to prune edges. Since independence information is only partial property of the data, the constraint-based approach usually cannot learn the full structure. Secondly, the score-based approach^{11, 14} begins from an empty graph and searches possible graph space to find the BN structure with the highest score. Because the possible searching space contains a super-exponential number of structures, finding the best structure is NP-hard¹⁵ and many heuristic algorithms have been designed. Thirdly, the Bayesian model averaging approach learns and averages a number of possible BN structures.^{16, 17} It outputs an ensemble of structures instead of a single one. Different possible structures and edges have their own confidence level.

3.2 Static BN (SBN)

As shown in section 3.1, BN can be used to represent joint distribution and hence correlation structures among the features. For the static BN algorithm, we first learn BN using complete data, then use the learned BN to do inference on missing data. In BN inference, given observed values of a set of variables, we can evaluate probability distribution over missing values of another set of variables. After learning probability distribution of the missing values, we use maximum likelihood principle to choose the most probable value. Note that, in the

learned BN, when the node has no connections, i.e., there is no correlation between this node and the others, static BN algorithm can reduce to the most likely independent value imputation. This algorithm has static nature in the sense that the learned BN does not change as the missing data increase. However, from another point of view, the observed values in missing data are still valid observations, which can be used to learn a better BN. In general, more learning data generate better learning results. This idea motivates our next two algorithms.

3.3 Parameter EM (ParaEM)

In parameter EM, we assume the size of missing data is moderate compared with the size of complete data; the amount of observed values within missing data can change the parameters of BN, and thus the probability distribution of missing values. Considering the observed values within the missing data significantly complicates the parameter learning, because we need to marginalize out the missing values. For a single data point \mathbf{X}_m , $m = 1, \dots, N$, we divide it into two parts, i.e., the observed part \mathbf{O}_m and the missing part \mathbf{H}_m . The probability of a single point should equal to the probability of the observed part, so the log-loss can be computed as:

$$l(\mathcal{D}|\mathcal{G}, \theta) = - \sum_{m=1}^N \log P(\mathbf{O}_m) = - \sum_{m=1}^N \log \sum_{\mathbf{H}_m} P(\mathbf{O}_m, \mathbf{H}_m) = - \sum_{m=1}^N \log \sum_{\mathbf{H}_m} P(\mathbf{X}_m). \quad (3)$$

Compared with Equation 1, there are some complications involved. Firstly, the sum within the log prevents the computation of log-loss score decomposes into variable level. Secondly, for complete data, the log-loss is unimodal with unique global maximum in parameter space, as shown in section 3.1.3. However, for log-loss in Equation 3, the sum of unimodal distributions becomes multimodal. In order to optimize the multimodal log-loss in high dimensional space, EM guarantees to find a local minimum using iterative methods.¹⁸ In the expectation step, we compute the expected sufficient statistics using BN inference on current parameters. Specifically, $M_t[x_i, \mathbf{p}_i] = \sum_{m=1}^N P(x_i, \mathbf{p}_i | \mathbf{O}_m, \theta_t)$, and $M_t[\mathbf{p}_i] = \sum_{x_i} M_t[x_i, \mathbf{p}_i]$. In the maximization step, we view the expected sufficient statistics as real observed numbers, and compute new parameters using maximum likelihood estimation. Specifically, $\theta_{x_i|\mathbf{p}_i}^{t+1} = \frac{M_t[x_i, \mathbf{p}_i]}{M_t[\mathbf{p}_i]}$. Every iteration of EM will decrease the log-loss until it converges. Compare with static BN, parameter EM considers the influences of observed values within the missing data and potentially perform better because of increased size of learning data. In addition, the static BN can be viewed as one iteration of parameter EM including collecting sufficient statistics and optimizing the parameters. Furthermore, in static BN, when we observe $X_i = x_i$ and $\mathbf{pa}_{\mathbf{X}_i} = \mathbf{p}_i$, we add 1 to $M[x_i, \mathbf{p}_i]$. However, in parameter EM, we add $P(x_i, \mathbf{p}_i | \mathbf{O}_m, \theta_t)$, which is a probabilistic and soft assignment. The finer change in sufficient statistics may result in lower log-loss. Finally, EM algorithm may converge to a very poor local maximum. In practice, we can use multiple random restart to find a better result.³

3.4 Structural EM (StructEM)

In structural EM, we assume there are lots of missing data and the valid observations within the missing data change both the structure and parameters of BN. When considering missing data, in addition to the harder parameter learning shown in section 3.3, the structure search in structure learning also becomes considerably more complex. In searching procedure, we apply different search operations including edge addition, edge deletion, and edge reversal. In order to measure the goodness of these operations, we need to define a delta score δ to measure the score difference before and after the search operation, and prefer the operations minimizing the log-loss score. When dealing with the complete data, as shown in Equation 1, after applying one edge operation, only one conditional probability changes because of the new parent set. All other terms keep the same. So the local change keeps in local and the delta score can be cached for further similar change. With missing data, as shown in Equation 3, the local change propagates to all parameters and the computed score can not be cached for further use. It will significantly increase the computation costs and make the structure learning intractable. Structural EM makes a good trade-off between computation complexity and accuracy.¹⁹ In the expectation step, we use the current learned BN, including structure and parameters, to produce the complete data and compute the expected sufficient statistics. In the maximization step, we use the complete data and expected sufficient statistics to learn a better BN structure and parameters. Every iteration of structural EM will decrease the log-loss score until it converges. Note that in structural EM, BIC score is used to guarantee the convergence.¹¹ We use log-loss here only for illustration purpose. In addition, in maximization step, any score-based Bayesian structure learning algorithm can be applied.

4. COMPARISON OF ALGORITHMS

4.1 Experimental Setting and Performance Metric

Our experiments focus on the children’s handwriting data and six algorithms are tested under a variety of conditions. Given the statistics of our data, best algorithm is selected in terms of imputation accuracy and computational cost, and it will be used later in section 5. In order to accurately measure the performance, we use complete data as ground truth and randomly make some of them as missing values. For the comparison completeness, we vary the missing data ratio for 10% to 90% and the missing value ratio from 1% to 50%. In our experiments, we choose the K2 algorithm as BN structure learning algorithm²⁰ and implement the missing value imputation algorithms based on Bayes net toolbox²¹ and BNT structure learning package.²²

Since the feature values are categorical, we use the imputation error rate as performance metric in experiments instead of using distance based metrics. It can be defined as the ratio of the number of wrong imputations N_{Wrong} and the total number of missing values N_{total} , $Error = \frac{N_{Wrong}}{N_{total}}$.

4.2 Data Properties

Here we show some properties of our children’s handwriting data. Although we use more complex techniques to solve the imputation problem, these properties are simple and good indicators of how good the imputations would be. In real applications, we do not have the ground truth of missing values. However, these data properties can give us some confidence levels of imputed values. Here we only show the results of the printed data for the interest of space, the cursive data analysis follows the same principle. The Figure 5(a) shows the entropy of every feature. It would roughly signify how good our imputation algorithms are when assuming each feature is independent, especially for most likely independent value imputation. The error rate of most likely independent value imputation for each feature is shown in Figure 5(b). Note that the relationship between entropy and error rate follows the same trend except feature 4 and feature 8. We observe feature 4 has more possible categorical values, which naturally increase its entropy. So we want a normalized version of entropy. The normalized entropy can be defined as $H_N(X_i) = -\frac{\sum_{x_i} p(x_i) \ln p(x_i)}{k_i}$, where x_i is the possible discrete values of X_i , and k_i is its number. The normalized entropy is shown in Figure 5(c). We can see the trend of normalized entropy and error rate of most likely independent value imputation are quite similar. A slight deviation lies in feature 8. This is because the distribution of feature 8 is bimodal with probability mass dispersing between two feature values, so its error rate is higher. The deviation of unimodality can be further measured using dip test.²³ In conclusion, the normalized entropy and modal information decide the trend of error rate for single feature imputation.

In order to measure the performance improvement of the imputation considering correlation, Figure 6 shows the correlations between features. The error rate will decrease as the correlation becomes stronger. It is a single indicator of how much the performance improves considering correlation. From Figure 6, we can see the correlations between features in our data are not very strong. For both cursive and printed data, within 66 possible pairwise correlations, there are only 4 pairs that the absolute values of correlation coefficients are larger than 0.2. However, even in this situation, using the correlation information can still result in moderate improvement. For printed data, the four strongest correlations around the diagonal in Figure 6(a) result in feature-wise performance improvement from 5.32% to 17.53%, corresponding to four edges $X_1 \rightarrow X_2$, $X_3 \rightarrow X_4$, $X_6 \rightarrow X_7$, and $X_{12} \rightarrow X_{11}$ in Figure 7(a). For cursive data, the four strongest correlations in 6(b) result in feature-wise performance improvement from 8.17% to 21.53%, corresponding to four edges $X_3 \rightarrow X_7$, $X_3 \rightarrow X_8$, $X_{10} \rightarrow X_{11}$, and $X_{11} \rightarrow X_{12}$ in Figure 7(b).

4.3 Performance Comparison of Missing Value Imputation Algorithms

In this section, we measure and compare the performance of imputation algorithms under a variety of conditions and make some useful observations and conclusions. Specifically, we set the missing data ratio from 90% to 10% at 10% interval, and missing value ratio as 9 different values within 1%, 2%, 3%, 5%, 7%, 10%, 20%, 30%, 50%. The missing data ratio is set in 10% interval because it mainly affects the performance of parameter EM and structural EM. When the increased number of valid observations within the missing data is too small, the performance will keep the same and the results will be redundant. The missing value ratio prefers small values with unbalanced distribution because in reality, the missing value ratio is usually small. Furthermore, when the

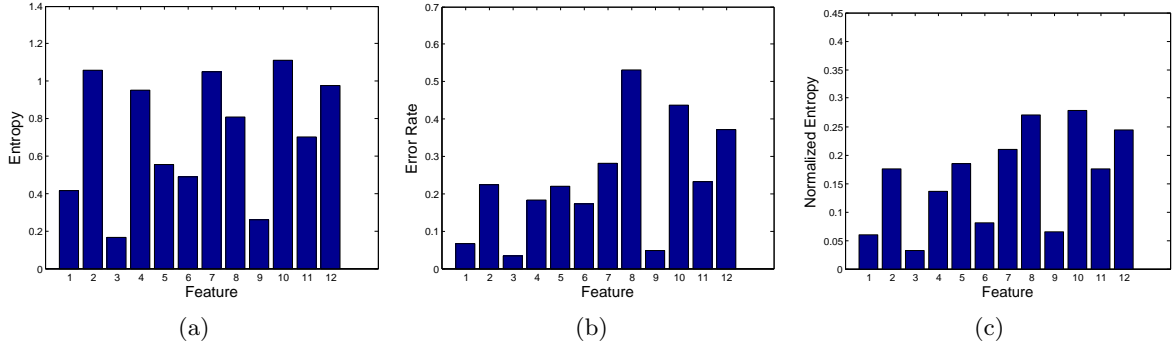


Figure 5: Feature-wise entropy and error rate of printed data. (a) Entropy for every feature. (b) Error rate for every feature using most likely independent value imputation. (c) Normalized entropy for every feature.

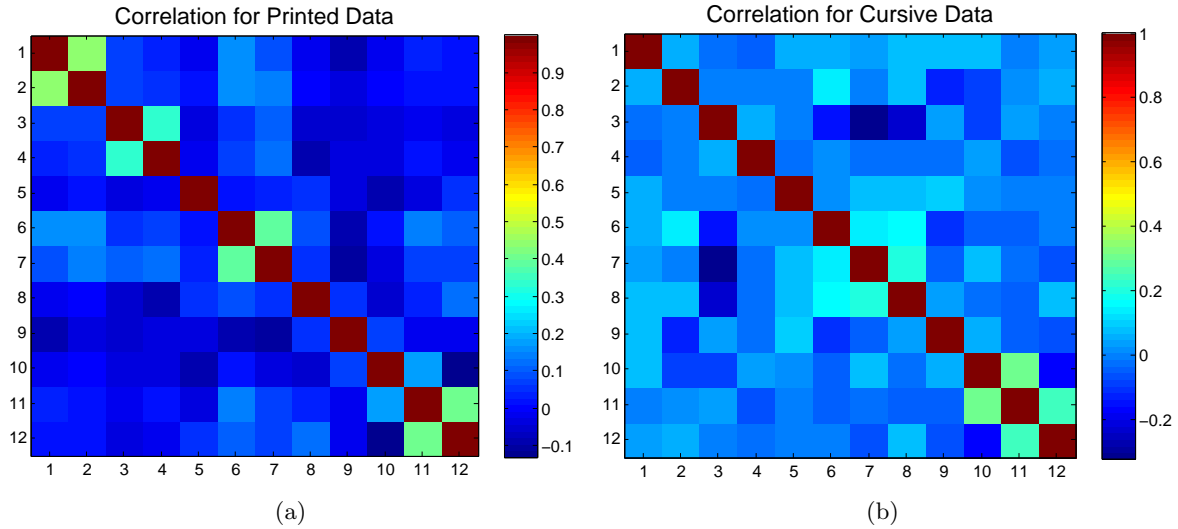


Figure 6: Correlation matrices of complete data. (a) For printed data, the four strongest correlations are between features 1 and 2, 3 and 4, 6 and 7, 11 and 12, along the diagonal. (b) For cursive data, the four strongest correlations are between features 3 and 7, 3 and 8, 10 and 11, 11 and 12.

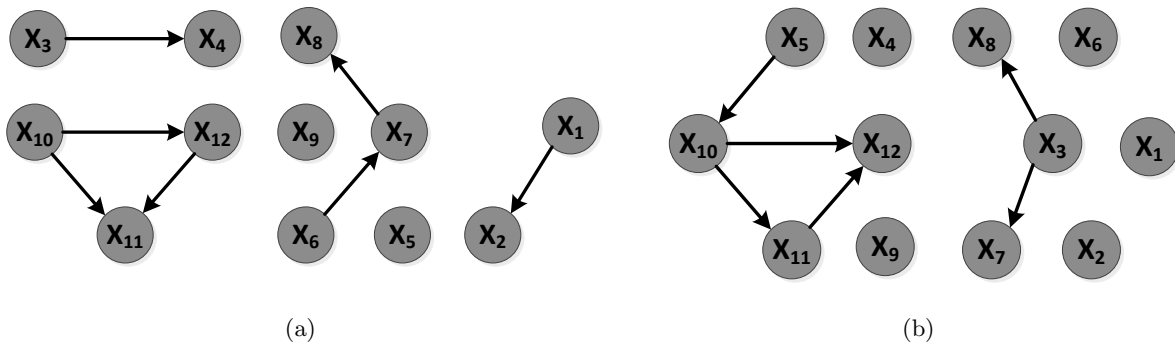


Figure 7: The learned BNs from data. All strongest correlations in Figure 6 have corresponding edges. (a) BN for printed data. (b) BN for cursive data.

missing value ratio is too high, the missing data contain very few valid observations and the imputation accuracy usually decreases a lot. In this case, people recommend to simply discard the data.¹ We conduct extensive

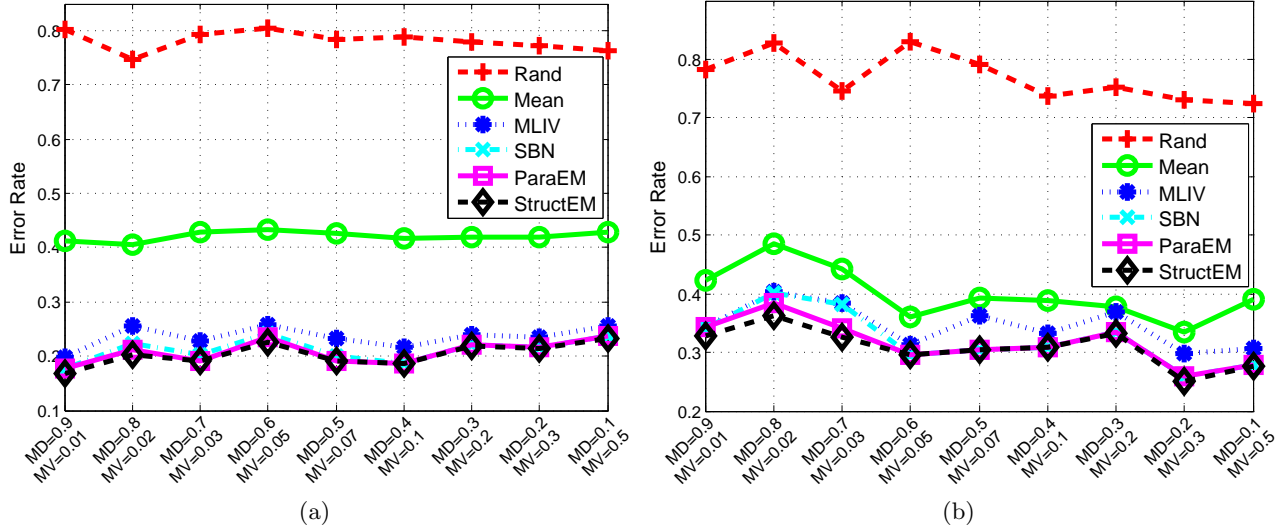


Figure 8: Error rate of six imputation algorithms. Rand means random imputation. Mean denotes mean imputation. MLIV denotes most likely independent value imputation. SBN means static BN algorithm. ParaEM mean parameter EM algorithm. StructEM denotes structural EM algorithm. MD denotes missing data ratio. MS denotes missing value ratio. (a) Printed data. (b) Cursive data.

experiments on all possible combinations of missing data and missing value ratios. In order to emphasize the key points and save space, we only show some representative results. At first, we show the performance using the diagonal view of ratio matrix, within which rows are different missing data ratios and columns are different missing value ratios. For example, the first case is 90% missing data and 1% missing value, the second case is 80% missing data and 2% missing value, and so on. The general conclusions will be given under these cases. Secondly, we show the performance of row and column views of ratio matrix, i.e., fixing the missing data ratio and varying the missing value ratios, and vice versa. Some finer observations are given. Figure 8 is the performance of six imputation algorithms, from which we can see:

- Most likely independent value imputation is the best algorithm when assuming each feature is independent. For printed data, the average accuracy improvement of most likely independent value imputation is 18.58% compared to mean imputation, and 54.67% for random imputation. For cursive data, the improvement is 5.36% and 40.54%. The improvement is smaller for cursive data because there are much more diversities, so they have higher entropy and the probability mass is more decentralized on different feature values.
- Using correlations between features can indeed improve the performance of imputation. The static BN is better or at least equal to the most likely independent value imputation. For printed data, the average accuracy improvement is 2.32%, while for cursive data, the improvement is 2.27%. The highest improvements happen when there are enough complete data to learn the correct correlation structure and the missing value ratio is low, so missing values do not corrupt the whole correlation structure. Note that when the correlation between features are stronger, the accuracy improvement can be higher.
- The lowest imputation error is 16.75% for printed data and 25.13% for cursive data. The cursive data has lower accuracy because they have more diversities. For the same reason, their results have less robustness.

In order to highlight the differences between the last four algorithms, Figure 9 shows the detailed performance of four imputation algorithms, from which we can observe:

- When the missing data ratio is high, and the missing value ratio is low, observations within missing data become very helpful for us to learn a better BN and improve imputation performance. Parameter EM improves the accuracy compared with static BN under this condition. For printed data, the highest accuracy improvement is 1.21%, while for cursive data, the highest improvement is 4.14%. Note that the

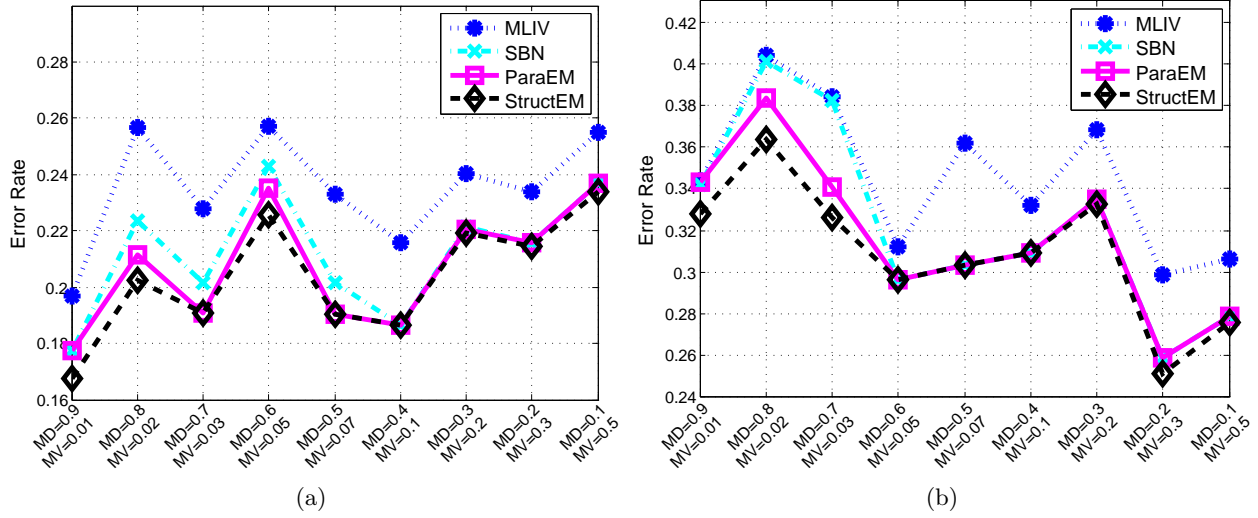


Figure 9: Highlighted error rate of four imputation algorithms. Here MLIV, SBN, ParaEM, StructEM, MD, and MV share the same meanings as in Figure 8. (a) Printed data. (b) Cursive data.

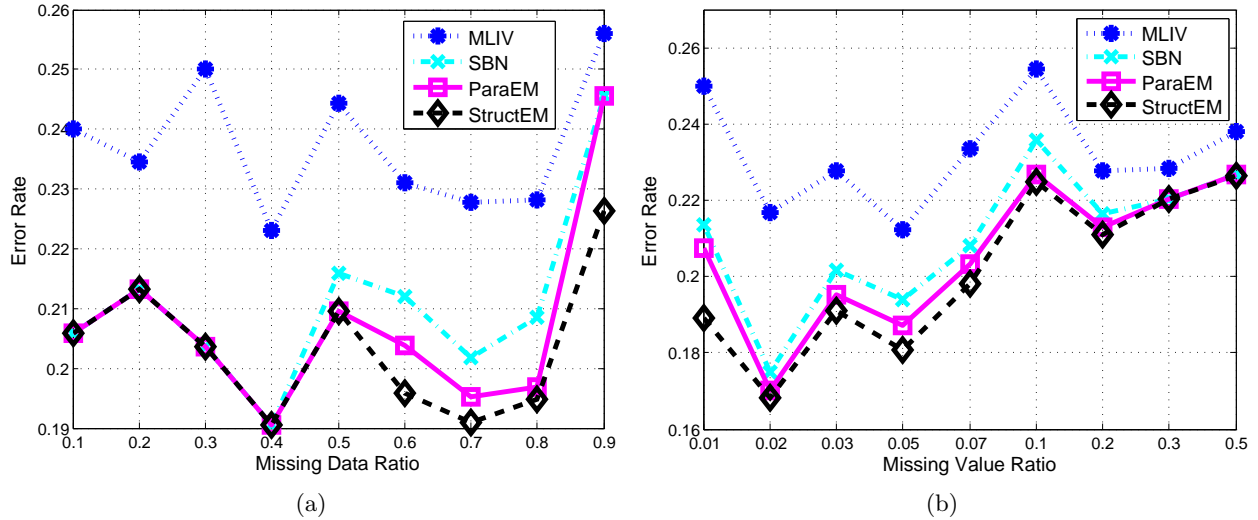


Figure 10: Highlighted error rate of four imputation algorithms. Here MLIV, SBN, ParaEM, and StructEM share the same meanings as in Figure 8. (a) Fixed missing value ratio, varying missing data ratios. (b) Fixed missing data ratio, varying missing value ratios.

accuracy improvement of cursive data is better because there are more diversities, and we need more data to reveal the whole distribution. So the increased amount of observations really helps a lot.

- When the missing data ratio is high while the missing value ratio is low, structural EM outperforms static BN at most 2.12% for printed data and 5.6% for cursive data, and correspondingly 1.08% and 2.02% for parameter EM. The higher accuracy improvements come from more flexibilities, i.e., both allowed changes in structures and parameters. Similarly, more diverse cursive data enjoy more improvements.

In order to get more insights of imputation performance, we show the row and column views of the results of ratio matrix. Specifically, we fix the missing value ratio to 3% and vary missing data ratios, and fix the missing data ratio to 30% and vary missing value ratios. Only the performance of printed data is given in Figure 10 for the interest of space. From the figure 10 we can conclude:

	Year	Grade	Sample ID	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	X ₁₁	X ₁₂
1	2012	1	MaL2023Kr	0	0	0	1	1	0	1	0	3	0	0	0
2	2012	1	OIC2023Ka	0	0	0	1	1	0	3	1	3	0	3	3
3	2012	4	ETM2020Ga	0	0	0	0	0	0	1	1	3	0	0	0
4	2013	2	BrB2023Bu	0	1	0	1	1	0	1	0	3	3	0	0
5	2013	2	CaG2023Ma	0	1	0	6	1	0	1	1	3	0	0	0
6	2013	2	DaM2023Sa	0	1	0	1	1	0	4	1	3	2	2	0
7	2013	2	LiD2023Ru	0	0	0	1	1	1	3	0	0	0	0	0
8	2012	2	ViG2022Kl	0	0	0	1	1	0	1	0	3	1	0	0
9	2013	3	BrE2022Ba	2	4	0	1	1	2	3	1	3	0	0	0
10	2013	3	LuP2022La	0	0	0	1	0	1	3	0	3	3	0	0
11	2013	3	MoA2022Ch	0	1	0	1	1	0	1	1	2	0	0	2
12	2013	3	JoM2022Le	5	5	0	1	1	0	3	0	0	0	0	0
13	2012	4	AlE2020BR	0	1	0	1	1	0	3	1	3	0	0	0
14	2012	4	LaS2020Gu	0	5	0	1	1	0	1	1	3	2	0	0
15	2013	5	An_2020Wo	0	0	0	0	0	0	3	0	3	0	0	0
16	2013	2	IsS2023We	0	1	0	1	1	0	1	1	3	0	0	0

	Year	Grade	Sample ID	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	X ₁₁	X ₁₂
1	2013	3	AIM2022Ca	1	1	1	0	0	3	3	4	3	0	0	0
2	2013	3	AlK2022Ko	1	3	1	1	0	1	0	0	3	2	3	3
3	2013	3	JoC2022Kl	2	1	1	0	1	1	0	0	3	0	0	2
4	2013	3	JoO2022Fi	2	1	1	1	0	1	0	2	3	0	0	0
5	2013	3	Ma_2022Ka	1	1	0	0	0	1	0	0	3	2	0	0
6	2013	3	MaG2022Po	2	1	1	1	1	1	0	0	3	0	0	0
7	2013	3	NaM2022We	2	1	1	1	2	1	0	0	3	3	3	2
8	2013	3	Saj2022Ad	2	1	1	4	2	1	0	0	3	3	3	3
9	2013	4	SoA2021He	1	1	1	0	0	1	1	3	2	2	0	0
10	2013	3	CaR2022Ka	2	1	1	3	3	3	0	0	3	3	3	3
11	2013	4	Cyl2021Va	2	1	0	0	3	1	0	0	3	0	0	0
12	2013	4	DyB2021Lo	1	1	1	1	0	1	1	0	3	0	0	0
13	2013	4	EmM2021Qu	2	1	1	1	0	1	0	2	3	0	0	0

(a)

(b)

Figure 11: Feature values after data commensuration and missing value imputation for (a) printed, and (b) cursive data, where the sample IDs correspond to Figure 3.

- As the missing value ratio is fixed, when the missing data ratio is low, parameter EM and structural EM have the same performance as static BN. As the missing data ratio increases higher than 40%, parameter EM and structural EM perform better. Because parameter EM only allows for parameter change according to new observations, it cannot change BN structure to accommodate large amount of them. So parameter EM works best when there are moderate amount of new observations. When the amount of observations within the missing data is too small or too large, its performance improvement becomes small. Instead, structural EM outperforms parameter EM in most cases.
- When the missing data ratio is fixed, as the missing value ratio increases, the performance of static BN, parameter EM, and structural EM become the same. This is because the number of valid observations within missing data becomes less. In addition, as more missing values corrupt the correlation structure, the improvement of static BN compared with most likely independent value becomes smaller.

5. APPLICATION TO MISSING HANDWRITING DATA

While the whole longitudinal handwriting research is intended for 11 years, currently we only collect samples from 2012 to 2013. For the purpose of temporal analysis, we only keep the individual samples having valid values for 2 years. In total, there are 888 individual printed handwriting samples. There are 96.4% incommensurate data, within which the incommensurate value ratio is 11.32%. The missing data ratio is 1.35% for 2012, 1.46% for 2013 and 2.81% for 2 years. Similarly, there are 282 individual cursive handwriting samples. There are 95.92% incommensurate data, within which the incommensurate value ratio is 13%. The missing data ratio is 5.32% for 2012, 1.06% for 2013 and 6.38% for 2 years. Note that as we observed, the missing mechanism is missing at random,⁴ so the two years' missing data accumulation is the sum of individual year's. We predict, even in the future collections, each year's missing data will accumulate with rare overlapping. As we have discussed in section 4.3 and shown in Figure 3, when the missing value ratio is too high, the number of valid observations within missing data are too small, which makes the sample less useful and decreases imputation accuracy. Under this situation, the imputed results are similar to samples generated by sampling techniques. So we delete the samples with more than 50% missing values. Note that in case that simulated samples are needed for inference purpose, BN is a generative model, which can generate enough number of samples. After the adjustment, the missing value ratio is 10.42% for printed data and 25% for cursive data. In our real children's handwriting data, because of low missing data ratio, we use static BN to perform the imputation for the adequate accuracy and low computational cost. Partial results are shown in Figure 11. From these examples, we can see our algorithm indeed improves the quality of data by data commensuration and missing value imputation.

6. CONCLUSION

In this paper, we solve the problem of missing values that we have encountered in our children's handwriting data, which are used to study the development of handwriting individuality along the time. This problem

becomes more serious in our longitudinal research and temporal analysis as more missing values accumulated with the time. We present and implement several missing value imputation algorithms on the handwriting data, and compare their performance under different ratios of missing data and missing value. We can conclude, for independent features, most likely independent value imputation performs best. The correlation usually improves the imputation accuracy no matter how weak it is. When missing data ratio is low, static BN works best because of the adequate accuracy and low computational cost. When the missing data ratio is high and the missing value ratio is low, structural EM works best because of the inclusion of more observations and higher model flexibilities. Given the statistics of our data, we use the static BN to perform missing value imputation.

REFERENCES

- [1] A. N. Baraldi and C. K. Enders, “An introduction to modern missing data analyses,” *Journal of School Psychology* **48**(1), pp. 5–37, 2010.
- [2] I. Jolliffe, *Principal component analysis*, Wiley Online Library, first ed., 2005.
- [3] C. M. Bishop *et al.*, *Pattern recognition and machine learning*, Springer, first ed., 2006.
- [4] J. W. Graham, “Missing data analysis: Making it work in the real world,” *Annual Review of Psychology* **60**, pp. 549–576, 2009.
- [5] S. Menard, *Longitudinal research*, Sage, first ed., 2002.
- [6] M. Puri, S. N. Srihari, and L. Hanson, “Probabilistic modeling of children’s handwriting,” in *Proc. of Document Recognition and Retrieval XXI*, 2014.
- [7] A. R. T. Donders, G. J. van der Heijden, T. Stijnen, and K. G. Moons, “Review: a gentle introduction to imputation of missing values,” *Journal of Clinical Epidemiology* **59**(10), pp. 1087–1091, 2006.
- [8] D. Janzing, J. Peters, E. Sgouritsa, K. Zhang, J. M. Mooij, and B. Schölkopf, “On causal and anticausal learning,” in *Proc. of ICML*, 2012.
- [9] O. Troyanskaya, M. Cantor, G. Sherlock, P. Brown, T. Hastie, R. Tibshirani, D. Botstein, and R. B. Altman, “Missing value estimation methods for dna microarrays,” *Bioinformatics* **17**(6), pp. 520–525, 2001.
- [10] S. Oba, M.-a. Sato, I. Takemasa, M. Monden, K.-i. Matsubara, and S. Ishii, “A bayesian missing value estimation method for gene expression profile data,” *Bioinformatics* **19**(16), pp. 2088–2096, 2003.
- [11] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*, MIT press, 2009.
- [12] J. Pearl, *Causality: models, reasoning and inference*, Cambridge Univ Press, 2000.
- [13] C. Meek, “Causal inference and causal explanation with background knowledge,” in *Proc. of UAI*, 1995.
- [14] N. Friedman, I. Nachman, and D. Peér, “Learning bayesian network structure from massive datasets: the sparse candidate algorithm,” in *Proce. of UAI*, 1999.
- [15] D. M. Chickering, D. Heckerman, and C. Meek, “Large-sample learning of bayesian networks is np-hard,” *The Journal of Machine Learning Research* **5**, pp. 1287–1330, 2004.
- [16] J. A. Hoeting, D. Madigan, A. E. Raftery, and C. T. Volinsky, “Bayesian model averaging: a tutorial,” *Statistical Science* , pp. 382–401, 1999.
- [17] A. E. Raftery, D. Madigan, and J. A. Hoeting, “Bayesian model averaging for linear regression models,” *Journal of the American Statistical Association* **92**(437), pp. 179–191, 1997.
- [18] T. K. Moon, “The expectation-maximization algorithm,” *Signal Processing Magazine* **13**(6), pp. 47–60, 1996.
- [19] N. Friedman, “The bayesian structural em algorithm,” in *Proc. of UAI*, 1998.
- [20] G. F. Cooper and E. Herskovits, “A bayesian method for the induction of probabilistic networks from data,” *Machine Learning* **9**(4), pp. 309–347, 1992.
- [21] K. Murphy *et al.*, “The bayes net toolbox for matlab,” *Computing Science and Statistics* **33**(2), pp. 1024–1034, 2001.
- [22] P. Leray and O. Francois, “Bnt structure learning package: Documentation and experiments,” *PSI Technical Reports* , 2004.
- [23] J. A. Hartigan and P. Hartigan, “The dip test of unimodality,” *The Annals of Statistics* , pp. 70–84, 1985.