

General Table Completion using a Bayesian Nonparametric Model

Authored by:

Zoubin Ghahramani
Isabel Valera

Abstract

Even though heterogeneous databases can be found in a broad variety of applications, there exists a lack of tools for estimating missing data in such databases. In this paper, we provide an efficient and robust table completion tool, based on a Bayesian nonparametric latent feature model. In particular, we propose a general observation model for the Indian buffet process (IBP) adapted to mixed continuous (real-valued and positive real-valued) and discrete (categorical, ordinal and count) observations. Then, we propose an inference algorithm that scales linearly with the number of observations. Finally, our experiments over five real databases show that the proposed approach provides more robust and accurate estimates than the standard IBP and the Bayesian probabilistic matrix factorization with Gaussian observations.

1 Paper Body

A full 90% of all the data in the world has been generated over the last two years and this expansion rate will not diminish in the years to come [17]. This extreme availability of data explains the great investment that both the industry and the research community are expending in data science. Data is usually organized and stored in databases, which are often large, noisy, and contain missing values. Missing data may occur in diverse applications due to different reasons. For example, a sensor in a remote sensor network may be damaged and transmit corrupted data or even cease to transmit; participants in a clinical study may drop out during the course of the study; or users of a recommendation system rate only a small fraction of the available books, movies, or songs. The presence of missing values can be challenging when the data is used for reporting, information sharing and decision support, and as a consequence, missing data treatment has captured the attention in diverse areas of data science such as machine learning, data mining, and data warehousing and management. Several studies have shown that probabilistic modeling can help to estimate missing values, detect errors in databases, or provide probabilistic responses to queries

[19]. In this paper, we exclusively focus on the use of probabilistic modeling for missing data estimation, and assume that the data are missing completely at random (MCAR). There is extensive literature in probabilistic missing data estimation and imputation in homogeneous databases, where all the attributes that describe each object in the database present the same (continuous or discrete) nature. Most of the work assumes that databases contain only either continuous data, usually modeled as Gaussian variables [21], or discrete, that can be either modeled by discrete likelihoods [9] or simply treated as Gaussian variables [15, 21]. However, there still exists a lack of work dealing with heterogeneous databases, which in fact are common in real applications and where the standard approach is to treat all the attributes, either continuous or discrete, as Gaussian variables. As a motivating example, consider a database that contains the answers to a survey, including diverse information about the participants such as age (count data), gender (categorical data), salary (continuous non negative data), etc. 1

In this paper, we provide a general Bayesian approach for estimating and replacing the missing data in heterogeneous databases (being the data MCAR), where the attributes describing each object can be either discrete, continuous or mixed variables. Specifically, we account for real-valued, positive real-valued, categorical, ordinal and count data. To this end, we assume that the information in the database can be stored in a matrix (or table), where each row corresponds to an object and the columns are the attributes that describe the different objects. We propose a novel Bayesian nonparametric approach for general table completion based on feature modeling, in which each object is represented by a set of latent variables and the observations are generated from a distribution determined by those latent features. Since the number of latent variables needed to explain the data depends on the specific database, we use the Indian buffet process (IBP) [8], which places a prior distribution over binary matrices where the number of columns (latent variables) is unbounded. The standard IBP assumes real-valued observations combined with conjugate likelihood models that allow for fast inference algorithms [4]. Here, we aim at dealing with heterogeneous databases, which may contain mixed continuous and discrete observations. We propose a general observation model for the IBP that accounts for mixed continuous and discrete data, while keeping the properties of conjugate models. This allows us to propose an inference algorithm that scales linearly with the number of observations. The proposed algorithm does not only infer the latent variables for each object in the table, but it also provides accurate estimates for its missing values. Our experiments over five real databases show that our approach for table completion outperforms, in terms of accuracy, the Bayesian probabilistic matrix factorization (BPMF) [15] and the standard IBP which assume Gaussian observations. We also observe that the approach based on treating mixed continuous and discrete data as Gaussian fails in estimating some attributes, while the proposed approach provides robust estimates for all the missing values regardless of their discrete or continuous nature. The main contributions in this paper are: i) A general observation model (for mixed continuous and discrete data) for the IBP that allows us to derive an inference

algorithm that scales linearly with the number of objects, and its application to build ii) a general and scalable tool to estimate missing values in heterogeneous databases. An efficient C-code implementation for Matlab of the proposed table completion tool is also released on the authors website.

2

Related Work

In recent years, probabilistic modeling has become an attractive option for building database management systems since it allows estimating missing values, detecting errors, visualizing the data, and providing probabilistic answers to queries [19]. BayesDB, for instance, is a database management system that resorts to Crosscat [18], which originally appeared as a Bayesian approach to model human categorization of objects. BayesDB provides missing data estimates and probabilistic answer to queries, but it only considers Gaussian and multinomial likelihood functions. In the literature, probabilistic low-rank matrix factorization approaches have been broadly applied to table completion (see, e.g., [14, 15, 21]). In these approaches, the table database X is approximated by a low-rank matrix representation $X \approx ZB$, where Z and B are usually assumed to be Gaussian distributed. Most of the works in this area have focused on building automatic recommendation systems, which appears as the most popular application of missing data estimation [14, 15, 21]. More specific models to build recommendation systems can be found in [7, 22], where the authors assume that the rates each user assign to items are generated by a probabilistic generative model which, based on the available data, accounts for similarities among users and among items to provide good estimates of the missing rates. Probabilistic matrix factorization can also be viewed as latent feature modeling, where each object is represented by a vector of continuous latent variables. In contrast, the IBP and other latent feature models (see, e.g., [16]) assume binary latent features to represent each object. Latent feature models usually assume homogeneous databases with either real [14, 15, 21] or categorical data [9, 12, 13], and only a few works consider heterogeneous data, such as mixed real and categorical data [16]. However, up to our knowledge, there are no general latent feature models (nor table completion tools) to directly deal with heterogeneous databases. To fill this gap, in this paper we provide a general table completion approach for heterogeneous databases, based on a generalized IBP, that allows for efficient inference. 1

<http://probcomp.csail.mit.edu/bayesdb/>

2

3

Model Description

Let us assume a table with N objects, where each object is defined by D attributes. We can store the data in an $N \times D$ observation matrix X , in which each D -dimensional row vector is denoted by $\mathbf{x}_n = [x_{1n}, \dots, x_{Dn}]$ and each entry is denoted by x_{dn} . We consider that column vectors \mathbf{x} (i.e., each dimension in the observation matrix X) may contain the following types of data:

- Continuous variables: 1. Real-valued, i.e., $x_{dn} \in \mathbb{R}$; 2. Positive real-valued, i.e., $x_{dn} \in \mathbb{R}^+$.
- Discrete variables: 1. Categorical data, i.e., x_{dn} takes values

in a finite unordered set, e.g., $x_{dn} \in \{\text{blue}, \text{red}, \text{black}\}$. 2. Ordinal data, i.e., x_{dn} takes values in a finite ordered set, e.g., $x_{dn} \in \{\text{never}, \text{sometimes}, \text{often}, \text{usually}, \text{always}\}$. 3. Count data, i.e., $x_{dn} \in \{0, \dots, ?\}$,

We assume that each observation x_{dn} can be explained by a K -length vector of latent variables associated to the n -th data point $z_n = [z_{n1}, \dots, z_{nK}]$ and a weighting vector $B_d = [b_{d1}, \dots, b_{dK}]$ (being K the number of latent variables), whose elements b_{dk} weight the contribution of k -th the latent feature to the d -th dimension of X . We gather the latent binary feature vectors z_n in a $N \times K$ matrix Z , which follows an IBP with concentration parameter α , i.e., $Z \sim \text{IBP}(\alpha)$ [8]. We place a 2 Gaussian distribution with zero mean and covariance matrix Σ_B over the weighting vectors B_d . For convenience, z_n is a K -length row vector, while B is a K -length column vector. To accommodate for all kinds of observed random variables described above, we introduce an auxiliary Gaussian variable y_{nd} , such that when conditioned on the auxiliary variables, the latent variable model behaves as a standard IBP with Gaussian observations. In particular, we assume y_{nd} is Gaussian distributed with mean $z_n B_d$ and variance σ^2 , i.e., $p(y_{nd} | z_n, B_d) = \mathcal{N}(y_{nd} | z_n B_d, \sigma^2)$, and assume that there exists a transformation function over the variables y_{nd} to obtain the observations x_{dn} , mapping the real line \mathbb{R} into the observation space. The resulting generative model is shown in Figure 1, where Z is the IBP latent matrix, and Y_d and B_d contain, respectively, the auxiliary Gaussian variables y_{nd} and the weighting factors b_{dk} for the d -dimension of the data. Additionally, \mathcal{D} denotes the set of auxiliary random variables needed to obtain the observation vector x_d given Y_d , and H_d contains the hyper-parameters associated to the random variables in \mathcal{D} . This model assumes that the observations x_{dn} are independent given the latent matrix Z , the weighting matrices B_d and the auxiliary variables \mathcal{D} . Therefore, the likelihood can be factorized as

$$\begin{aligned} & p(X | Z, \{B_d, \mathcal{D}_d\}_{d=1}^D) \\ &= \prod_{d=1}^D \prod_{n=1}^N p(x_{dn} | z_n, B_d, \mathcal{D}_d) \end{aligned}$$

Note that, if we assume Gaussian observations and set $Y_d = x_d$, this model resembles the standard IBP with Gaussian observations [8]. In addition, conditioned on the variables Y_d , we can infer the latent matrix Z as in the standard IBP. We also remark that auxiliary Gaussian variables to link a latent model with the observations have been previously used in Gaussian processes for multi-class classification [6] and for ordinal regression [2]. However, up to our knowledge, this simple approach has not been used to account for mixed continuous and discrete data, and the existent approaches for the IBP with discrete observations propose non-conjugate likelihood models and approximate inference algorithms [12, 13].

Likelihood Functions

Now, we define the set of transformations that map from the Gaussian variables y_{nd} to the corresponding observations x_{dn} . We consider that each dimension in the table X may contain any of the discrete or continuous variables detailed above, provide a likelihood function for each kind of data and, in turn, also a likelihood function for mixed data. 2

For convenience, we capitalized here the notation for the weighting vectors B_d .

3

Real-valued Data. In this case, we assume that $x_d = Y_d$ in the model in Figure 1 and consider the standard approach when dealing with real-valued observations, which consist of assuming a Gaussian likelihood function. In particular, as in the standard linear-Gaussian IBP [8], we assume that each observation x_{dn} is distributed as $p(x_{dn} | z_n, B_d) = N(x_{dn} | z_n B_d, \sigma_y^2)$. **Positive Real-valued Data.** In order to obtain positive real-valued observations, i.e., $x_{dn} \in \mathbb{R}^+$, we apply a transformation over y_{nd} that maps from the real numbers to the positive real numbers, i.e., $x_{dn} = f(y_{nd} + u_{dn})$, where u_{dn} is a Gaussian noise variable with variance σ_u^2 , and $f: \mathbb{R} \rightarrow \mathbb{R}^+$ is a monotonic differentiable function. By change of variables, we obtain the likelihood function for positive real-valued observations as

$$p(x_{dn} | z_n, B_d) = q \cdot f' \left(f^{-1}(x_{dn}) \right) \cdot N(y_{nd} | z_n B_d, \sigma_y^2 + \sigma_u^2), \quad (1)$$

where $f^{-1}: \mathbb{R}^+ \rightarrow \mathbb{R}$ is the inverse function of the transformation $f(\cdot)$, i.e., $f^{-1}(f(v)) = v$. Note that in this case we resort to the Gaussian variable u_{dn} in order to obtain x_{dn} from y_{nd} , and therefore, $\sigma_d^2 = \sigma_u^2$ and $H_d = \sigma_u^2$. **Categorical Data.** Now we account for categorical observations, i.e., each observation x_{dn} can take values in the unordered index set $\{1, \dots, R_d\}$. Hence, assuming a multinomial probit model, we can write $d_{dn} = \arg \max_{r \in \{1, \dots, R_d\}} (z_n B_d)_r + u_{dnr}$, where u_{dnr} is a Gaussian noise variable with variance σ_y^2 , and therefore, we can obtain the probability of each element x_{dn} taking value $r \in \{1, \dots, R_d\}$ as

$$p(x_{dn} = r | z_n, B_d) = \frac{\exp(z_n B_d)_r + u_{dnr}}{\sum_{r'=1}^{R_d} \exp(z_n B_d)_{r'} + u_{dnr'}} \quad (2)$$

denotes the K -length weighting vector, in which each b_{dkr} where k being weights the influence of the k -th feature for the observation x_{dn} taking value r . Note that, under this d likelihood model, since we have a Gaussian auxiliary variable y_{nr} and a weighting factor b_{dkr} for each possible value of the observation $r \in \{1, \dots, R_d\}$, we need to gather all the weighting factors d in the $N \times R_d$ matrix Y_d . B_d in a $K \times R_d$ matrix B_d , and all the Gaussian auxiliary variables y_{nr} .

$d = z_n B_d + u_{dn}$, where u_{dn} is a Gaussian noise variable with variance σ_y^2 , and therefore, we can obtain the probability of each element x_{dn} taking value $r \in \{1, \dots, R_d\}$ as

$$p(x_{dn} = r | z_n, B_d) = E_p(u) \cdot \frac{\exp(z_n B_d)_r + u_{dnr}}{\sum_{r'=1}^{R_d} \exp(z_n B_d)_{r'} + u_{dnr'}} \quad (3)$$

where subscript r in b_{dkr} states for the column in B_d ($r \in \{1, \dots, R_d\}$), $\Phi(\cdot)$ denotes the cumulative density function of the standard normal distribution and $E_p(u)$ denotes expectation with respect to the distribution $p(u) = N(0, \sigma_y^2)$.

). Ordinal Data. Consider ordinal data, in which each element x_{dn} takes values in the ordered index set $\{1, \dots, R_d\}$. Then, assuming an ordered probit model, we can write $\mathbb{P}(x_{dn} = r | \mathbf{y}_{nd}, \mathbf{B}_d) = \Phi(\mathbf{z}_{nd}^T \mathbf{B}_d - \tau_{r-1}) - \Phi(\mathbf{z}_{nd}^T \mathbf{B}_d - \tau_r)$ where again \mathbf{y}_{nd} is Gaussian distributed with mean $\mathbf{z}_n^T \mathbf{B}_d$ and variance σ^2 , and τ_r for $r \in \{1, \dots, R_d - 1\}$ are the thresholds that divide the real line into R_d regions. We assume the thresholds τ_r are sequentially generated from the truncated Gaussian distribution $\tau_r \sim N(\tau_{r-1} - \sigma^2, \sigma^2) I(\tau_r \in [\tau_{r-1}, \tau_r])$, where $\tau_0 = -\infty$ and $\tau_{R_d} = +\infty$. As opposed to the categorical case, now we have a unique \mathbf{B}_d

weighting vector \mathbf{B}_d and a unique Gaussian variable \mathbf{y}_{nd} for each observation x_{dn} . Hence, the value of x_{dn} is determined by the region in which \mathbf{y}_{nd} falls. Under the ordered probit model [2], the probability of each element x_{dn} taking value $r \in \{1, \dots, R_d\}$ can be written as $\mathbb{P}(x_{dn} = r | \mathbf{y}_{nd}, \mathbf{B}_d) = \Phi(\mathbf{z}_{nd}^T \mathbf{B}_d - \tau_{r-1}) - \Phi(\mathbf{z}_{nd}^T \mathbf{B}_d - \tau_r)$. (5) Let us remark that, if the d -dimension of the observation matrix contains ordinal data, the set of d auxiliary variables reduces to the Gaussian thresholds $\tau_d = \{\tau_{1d}, \dots, \tau_{R_d}\}$ and $\mathbf{H}_d = \tau_d$. \mathbf{d} Count Data. In count data each observation x_{dn} takes non-negative integer values, i.e., $x_{dn} \in \{0, 1, 2, \dots\}$. Then, we assume $x_{dn} = \text{bvc}(\mathbf{y}_{nd})$, (6) where bvc returns the floor of v , that is the largest integer that does not exceed v , and $f: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is a monotonic differentiable function that maps from the real numbers to the positive real numbers. We can therefore write the likelihood function as $\mathbb{P}(x_{dn} = r | \mathbf{y}_{nd}, \mathbf{B}_d) = \frac{1}{r!} \frac{d^r}{d\mathbf{y}_{nd}^r} f(\mathbf{y}_{nd}) \exp(-f(\mathbf{y}_{nd}))$ (7) where $f: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is the inverse function of the transformation f .

?

Z

\mathbf{Y}_d

$2 \times B$

\mathbf{B}_d

\mathbf{X}_d

$\mathbf{H}_d, d = 1, \dots, D$

Figure 1: Generalized IBP for mixed continuous and discrete observations.

4

Inference Algorithm

In this section we describe our algorithm for inferring the latent variables given the observation matrix. Under our model, detailed in Section 3, the probability distribution over the observation matrix is fully characterized by the latent matrices \mathbf{Z} and $\{\mathbf{B}_d\}_{d=1}^D$ (as well as the auxiliary variables τ_d). Hence, if we assume the latent vector \mathbf{z}_n for the n -th datapoint and the weighting factors \mathbf{B}_d (and the auxiliary variables τ_d) to be known, we have a probability distribution over missing observations x_{dn} from which we can obtain estimates for x_{dn} by sampling from this distribution,³ or by simply taking either its mean, mode or median value. However, this procedure requires the latent matrix \mathbf{Z} and the latent weighting factors \mathbf{B}_d (and τ_d) to be known. We use Markov Chain Monte Carlo (MCMC) methods, which have been broadly applied to infer the IBP matrix (see, e.g., in [8, 23, 20]). The proposed inference algorithm is summarized

in Algorithm 1. This algorithm exploits the information in the available data to learn the similarities among the objects (captured in our model by the latent feature matrix Z), and how these latent features show up in the attributes that describe the objects (captured in our model by B_d). In Algorithm 1, we first need to update the latent matrix Z . Note that conditioned on $\{Y_d\}_{d=1}^D$, both the latent are independent of the observation matrix X . Admatrix Z and the weighting matrices $\{B_d\}_{d=1}^D$ dditionally, since $\{B_d\}_{d=1}^D$ and $\{Y_d\}_{d=1}^D$ are Gaussian distributed, we can analytically marginalize $\{B_d\}_{d=1}^D$ out the weighting matrices $\{B_d\}_{d=1}^D$ to obtain $p(\{Y_d\}_{d=1}^D | Z)$. Therefore, to infer the matrix Z , we d=1 d=1 can apply the collapsed Gibbs sampler which presents better mixing properties than the uncollapsed 3 Note that sampling from this distribution might be computationally expensive. In this case, we can easily obtain samples of x_{dn} by exploiting the structure of our model. In particular, we can simply sample the auxiliary Gaussian variables y_{nd} given z_n and B_d , and then obtain an estimate for x_{dn} by applying the corresponding transformation, detailed in Section 3.1.

5

Algorithm 1 Inference Algorithm. Input: X Initialize: initialize Z and $\{Y_d\}_{d=1}^D$ 1: for each iteration do 2: Update Z given $\{Y_d\}_{d=1}^D$. 3: for $d = 1, \dots, D$ do 4: Sample B_d given Z and Y_d according to (8). 5: Sample Y_d given X , Z and B_d (as shown in the Supplementary Material). 6: Sample $?_d$ if needed (as shown in the Supplementary Material). 7: end for 8: end for d D Output: Z , $\{B_d\}_{d=1}^D$ and $\{?_d\}_{d=1}^D$

Gibbs sampler and, in consequence, is the standard method of choice in the context of the standard linear-Gaussian IBP [8]. However, this algorithm suffers from a high computational cost (being complexity per iteration cubic with the number of data points N), which is prohibitive when dealing with large databases. In order to solve this limitation, we resort to the accelerated Gibbs sampler [4] instead. This algorithm presents linear complexity with the number of datapoints and is detailed in the Supplementary Material. Second, we need to sample the weighting factors in B_d , which is a $K \times R_d$ matrix in the case of categorical attributes, and a K -length column vector otherwise. We denote each column vector in B_d by b_{dr} . The posterior over the weighting vectors are given by $p(b_{dr} | y_{rd}, Z) = N(b_{dr} | P^{-1}_{?d} y_{rd}, P^{-1}_{?d})$,

(8)

2 where $P = Z_i^T Z + 1/2 B_i^T B_i$ and $?_{dr} = Z_i^T y_{rd}$. Note that the covariance matrix $P^{-1}_{?d}$ depend neither on the dimension d nor on r , so we only need to invert the $K \times K$ matrix P once at each iteration. We describe in the Supplementary Material how to efficiently compute P after changes in the Z matrix by rank one updates, without the need of computing the matrix product $Z_i^T Z$.

Once we have updated Z and B_d , we sample each element in Y_d from the distribution $d \sim x_{dn}, z_n, b_d$ spec $z_n b_d, ?_d$ if the observation x_{dn} is missing, and from the posterior $p(y_{nr} | N(y_{nr} | \text{ifid in the Supplementary Material, otherwise. Finally, we sample the auxiliary variables in } ?_d \text{ from their posterior distribution (detailed in the Supplementary Material) if necessary. This two latter steps involve, in the worst case, sampling from a doubly$

truncated univariate normal distribution (see the Supplementary Material for further details), for which we make use of the algorithm in [11].

5

Experimental evaluation

We now validate the proposed algorithm for table completion on five real databases, which are summarized in Table 1. The datasets contain different numbers of instances and attributes, which cover all the discrete and continuous variables described in Section 3. We compare, in terms of predictive log-likelihood, the following methods for table completion: ? The proposed general table completion approach denoted by GIBP (detailed in Section 3). ? The standard linear-Gaussian IBP [8] denoted by SIBP, treating all the attributes as Gaussian. ? The Bayesian probabilistic matrix factorization approach [15] denoted by BPMF, that also treats all the attributes in X as Gaussian distributed. For the GIBP, we consider for the real positive and the count data the following transformation, that maps from the real numbers to the real positive numbers, $f(x) = \log(\exp(wx) + 1)$, where w is a user hyper-parameter. Before running the SIBP and the BPMF methods we normalize each column in matrix X to have zero-mean and unit-variance. Then, in order to provide estimates for the missing data, we denormalize the inferred Gaussian variable. Additionally, since both the SIBP and the BPMF assume continuous observations, when dealing with discrete data, we estimate each missing value as the closest integer value to the (denormalized) Gaussian variable. 6

Dataset Statlog German credit dataset [5] QSAR biodegradation dataset [10] Internet usage survey dataset [1] Wine quality Dataset [3]

N 1,000

6,497

D 20 (10 C + 4 O + 6 N) 41 (2 R + 17 P + 4 C + 18 N) 32 (23 C + 8 O + 1 N) 12 (11 P + 1 N)

NESARC dataset [13]

43,000

55 C

1,055 1,006

Description Collects information about the credit risks of the applicants. Contains molecular descriptors of biodegradable and non-biodegradable chemicals. Contains the responses of the participants to a survey related to the usage of internet. Contains the results of physicochemical tests realized to different wines. Contains the responses of the participants to a survey related to personality disorders.

0

?2

?2

?3 ?4

GIBP SIBP BPMF

?5 ?6 10

20

30 40 % of missing data

50

(a) Statlog.

?1 Log?likelihood

?1 Log?likelihood

Log?likelihood

Table 1: Description of datasets. ?R? states for real-valued variables, ?P? for positive real-valued variables, ?C? for categorical variables, ?O? for ordinal variables and ?N? for count variables

?4 GIBP SIBP BPMF

?6 ?8 ?10 10

30 40 % of missing data

50

(b) QSAR biodegradation.

10

20

30

40 50 60 70 % of missing data

80

90

(c) Internet usage survey.

?0.5 Log?likelihood

Log?likelihood

GIBP SIBP BPMF

?2

?2.5 20

0 GIBP SIBP BPMF

?5

?10

?0.6 ?0.7 GIBP SIBP

?0.8 10

?1.5

20

30

40 50 60 70 % of missing data

80

90

10

(d) Wine quality.

20

30

40 50 60 70 % of missing data

80

90

(e) Nesarc database

Figure 2: Average test log-likelihood per missing datum. The ?whiskers? show a standard deviations from the average test log-likelihood. In Figure 2,

we plot the average predictive log-likelihood per missing value as a function of the percentage of missing data. Each value in Figure 2 has been obtained by averaging the results in 20 independent sets where the missing values have been randomly chosen. In Figures 2a and 2b, we cut the plot in 50% because, in these two databases, the discrete attributes present a mode value that is present for more than 80% of the instances. As a consequence, the SIBP and the BPMF algorithms assign probability close to one to the mode, which results in an artificial increase in the average test log-likelihood for larger percentages of missing data. For the BPMF model, we have used different numbers of latent features (in particular, 10, 20 and 50), although we only show the best results for each database, specifically, $K = 10$ for the NESARC and the wine databases, and $K = 50$ for the remainder. Both the GIBP and the SIBP have not inferred a number of (binary) latent features above 25 in any case. Note that in Figure 2e, we only plot the test log-likelihood for the GIBP and the SIBP because the BPMF provides much lower values. As expected, we observe in Figure 2 that the average test log-likelihood decreases for the three models when the number of missing values increases (flat shape of the curves are due to the y-axis scale). In this figure, we also observe that the proposed general IBP model outperforms the SIBP and the BPMF for four of the the databases, being the SIBP slightly better for the Internet database. The BPMF model presents the lowest test-log-likelihood in all the databases. Now, we analyze the performance of the three models for each kind of discrete and continuous variables. Figure 3 shows average predictive likelihood per missing value for each attribute in the table, i.e., for each dimension in X . In this figure we have grouped the dimensions according to the kind of data that they contain, showing in the x-axis the number of considered categories for the case of categorical and ordinal data. In this figure, we observe that the GIBP presents similar performance

for all the attributes in the five databases, while for the SIBP and the BPMF models, the test-loglikelihood falls drastically for some of the attributes, being this effect worse in the case of the BPMF (it explains the low log-likelihood in Figure 2). This effect is even more evident in Figures 2b and 2d. We also observe, in Figures 2 and 3, that both IBP based approaches (the GIBP and the SIBP) outperform the BPMF, with the proposed GIBP being the one that best performs across all the databases. We can conclude that, unlike to the BPMF and the GIBP, the SIBP provides accurate estimates for the missing data regardless of their discrete or continuous nature.

6

Conclusions

In this paper, we have proposed a table completion approach for heterogeneous databases, based on an IBP with a generalized likelihood that allows for mixed discrete and continuous data. We have then derived an inference algorithm that scales linearly with the number of observations. Finally, our experimental results over five real databases have shown that the proposed approach outperforms, in terms of robustness and accuracy, approaches that treat all the attributes as Gaussian variables. Log-likelihood

0 ?10 GIBP SIBP BPMF

?20 ?30
 C5
 C10
 C5
 C3
 C4
 C3
 C3
 C4
 C2
 C2 O4 Attribute
 O5
 O5
 O2
 N
 N
 N
 N
 N
 N
 (a) Statlog. Log?likelihood
 10 0 ?10 ?20 ?30
 GIBP SIBP BPMF R R P P P P P P P P P P P P P P C2C2C4C2
 N Attribute
 (b) QSAR biodegradation. Log?likelihood
 0 ?2 ?4 GIBP SIBP BPMF
 ?6 ?8
 C3 C3 C3 C3 C3 C3 C4 C4 C4 C5 C5 C6 C6 C6 C6 C6 C6 C5 C5 C3 C2 C2
 C2 C9 O6 O7 O7 O7 O7 O7 O8 O6 N Attribute
 (c) Internet usage survey. Log?likelihood
 10 0 ?10 GIBP SIBP BPMF
 ?20 ?30
 P
 P
 P
 P
 P
 P P Attribute
 P
 P
 P
 P
 N
 (d) Wine quality. Log?likelihood
 0 ?10 GIBP SIBP BPMF
 ?20 ?30

CC
Attribute

(e) Nesarc database

Figure 3: Average test log-likelihood per missing datum in each dimension of the data with 50% of missing data. In the x-axis ?R? states for real-valued variables, ?P? for positive real-valued variables, ?C? for categorical variables, ?O? for ordinal variables and ?N? for count variables. The number that accompanies the ?C? or ?O? corresponds to the number of categories. Acknowledgments Isabel Valera acknowledge the support of Plan Regional-Programas I+D of Comunidad de Madrid (AGES-CM S2010/BMD-2422), Ministerio de Ciencia e Innovaci?on of Spain (project DEIPRO TEC2009-14504-C02-00 and program Consolider-Ingenio 2010 CSD2008-00010 COMONSENS). Zoubin Ghahramani is supported by the EPSRC grant EP/I036575/1 and a Google Focused Research Award. 8

2 References

- [1] Pew Research Centre. 25th anniversary of the web. Available on: <http://www.pewinternet.org/datasets/2014-25th-anniversary-of-the-web-omnibus/>.
- [2] W. Chu and Z. Ghahramani. Gaussian processes for ordinal regression. *J. Mach. Learn. Res.*, 6:1019? 1041, December 2005.
- [3] P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems*. Dataset available on: <http://archive.ics.uci.edu/ml/datasets.html>, 47(4):547?553, 2009.
- [4] F. Doshi-Velez and Z. Ghahramani. Accelerated sampling for the indian buffet process. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML ?09*, pages 273?280, New York, NY, USA, 2009. ACM.
- [5] J. Eggermont, J. N. Kok, and W. A. Kusters. Genetic programming for data classification: Partitioning the search space. In *Proceedings of the 2004 Symposium on applied computing (ACM SAC04)*. Dataset available on: <http://archive.ics.uci.edu/ml/datasets.html>, pages 1001?1005. ACM, 2004.
- [6] M. Girolami and S. Rogers. Variational Bayesian multinomial probit regression with Gaussian process priors. *Neural Computation*, 18:2006, 2005.
- [7] P. Gopalan, F. J. R. Ruiz, R. Ranganath, and D. M. Blei. Bayesian Non-parametric Poisson Factorization for Recommendation Systems. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2014.
- [8] T. L. Griffiths and Z. Ghahramani. The Indian buffet process: an introduction and review. *Journal of Machine Learning Research*, 12:1185?1224, 2011.
- [9] X.-B. Li. A Bayesian approach for estimating and replacing missing categorical data. *J. Data and Information Quality*, 1(1):3:1?3:11, June 2009.
- [10] K. Mansouri, T. Ringsted, D. Ballabio, R. Todeschini, and V. Consonni. Quantitative structureactivity relationship models for ready biodegradability of chemicals. *Journal of Chemical Information and Modeling*. Dataset available on: <http://archive.ics.uci.edu/ml/datasets.html>.
- [11] C. P. Robert. Simulation of truncated normal variables. *Statistics and computing*, 5(2):121?125, 1995.
- [12] F. J. R. Ruiz, I. Valera, C. Blanco, and F. Perez-Cruz. Bayesian nonpara-

metric modeling of suicide attempts. *Advances in Neural Information Processing Systems*, 25:1862–1870, 2012. [13] F. J. R. Ruiz, I. Valera, C. Blanco, and F. Perez-Cruz. Bayesian nonparametric comorbidity analysis of psychiatric disorders. *Journal of Machine Learning Research* (To appear). Available on <http://arxiv.org/pdf/1401.7620v1.pdf>, 2013. [14] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. In *Advances in Neural Information Processing Systems*, 2007. [15] R. Salakhutdinov and A. Mnih. Bayesian probabilistic matrix factorization using Markov Chain Monte Carlo. In *Proceedings of the 25th International Conference on Machine Learning, ICML ’08*, pages 880–887, New York, NY, USA, 2008. ACM. [16] E. Salazar, M. Cain, E. Darling, S. Mitroff, and L. Carin. Inferring latent structure from mixed real and categorical relational data. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12), ICML ’12*, pages 1039–1046, New York, NY, USA, July 2012. Omnipress. [17] ScienceDaily. Big data, for better or worse: 90% of world’s data generated over last two years. [18] P. Shafto, C. Kemp, Mansinghka V., and Tenenbaum J. B. A probabilistic model of cross-categorization. *Cognition*, 120(1):1 – 25, 2011. [19] S. Singh and T. Graepel. Automated probabilistic modelling for relational data. In *Proceedings of the ACM of Information and Knowledge Management, CIKM ’13*, New York, NY, USA, 2013. ACM. [20] M. Titsias. The infinite gamma-Poisson feature model. *Advances in Neural Information Processing Systems*, 19, 2007. [21] A. Todeschini, F. Caron, and M. Chavent. Probabilistic low-rank matrix completion with adaptive spectral regularization algorithms. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 845–853. Curran Associates, Inc., Dec. 2013. [22] C. Wang and D. M. Blei. Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’11*, pages 448–456, New York, NY, USA, 2011. ACM. [23] S. Williamson, C. Wang, K. Heller, and D. Blei. The IBP compound Dirichlet process and its application to focused topic modeling. *Proceedings of the 27th Annual International Conference on Machine Learning*, 2010.