

Journal of Computational Biology: <http://mc.manuscriptcentral.com/liebert/jcb>

**A model-based approach to gene clustering with missing observation reconstruction in a Markov Random Field framework**

Journal:	<i>Journal of Computational Biology</i>
Manuscript ID:	JCB-2008-0078.R1
Manuscript Type:	Original Paper
Date Submitted by the Author:	n/a
Complete List of Authors:	Blanchet, Juliette; INRIA Rhone-Alpes, MISTIS; SFL Vignes, Matthieu; RRI - University of Aberdeen, BioSS
Keyword:	computational molecular biology, STATISTICS, GENE NETWORKS, GENE EXPRESSION
Abstract:	<p>The different measurement techniques that interrogate biological systems provide means for monitoring the behaviour of virtually all cell components at different scales and from complementary angles. However data generated in these experiments are difficult to interpret. A first difficulty arises from high-dimensionality and inherent noise of such data. Organizing them into meaningful groups is then highly desirable to improve our knowledge of biological mechanisms. A more accurate picture can be obtained when accounting for dependencies between components (e.g. genes) under study. A second difficulty arises from the fact that biological experiments often produce missing values. When it is not ignored, the latter issue has been solved by imputing the expression matrix prior to applying traditional analysis methods. Although helpful, this practice can lead to unsound results. We propose in this paper a statistical methodology that integrates individual dependencies in a missing data framework. More explicitly, we present a clustering algorithm dealing with incomplete data in a Hidden Markov Random Field context.</p>

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



For Peer Review

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

# A model-based approach to gene clustering with missing observations **reconstruction** in a Markov Random Field framework

29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Juliette Blanchet <sup>a,1</sup>, Matthieu Vignes <sup>b,\*</sup>

<sup>a</sup>*INRIA Rhône-Alpes, 655 av. de l'Europe, 38330 Saint Ismier Cedex, France*

<sup>b</sup>*Biomathematics and Statistics Scotland at the RRI, Bucksburn, Aberdeen, AB21  
9SB, Scotland, UK*

---

## Abstract

The different measurement techniques that interrogate biological systems provide means for monitoring the behaviour of virtually all cell components at different scales and from complementary angles. However data generated in these experiments **are difficult to interpret**. A first difficulty arises from high-dimensionality and inherent noise of such data. Organizing them into meaningful groups is then highly desirable to improve our knowledge of biological mechanisms. A more accurate picture can be obtained when accounting for dependencies between components (*e.g.* genes) under study. A second difficulty arises from the fact that biological experiments often **produce missing values**. When it is not ignored, the latter issue has been solved by imputing the expression matrix prior to applying traditional analysis methods. Although helpful, this practice can lead to unsound results.

We propose in this paper a statistical methodology that integrates individual dependencies in a missing data framework. More explicitly, we present a clustering algorithm **dealing with incomplete data in a Hidden Markov Random Field context**.

*Preprint submitted to Journal of Computational Biology*      *6th November 2008*

1  
2  
3 This tackles the missing value issue in a probabilistic framework and still allows us to  
4 reconstruct missing observations *a posteriori* without imposing any pre-processing of  
5 the data. Experiments on synthetic data validate the gain in using our method and  
6  
7  
8  
9  
10 real biological data analysis present its potential to extract biological knowledge.

11  
12 *Key words:* Biological interaction network; Gene clustering; Markov Random  
13  
14 Field; Mean field-like approximation; Missing data.  
15

---

## 16 17 18 19 20 21 **1 Introduction**

22  
23  
24  
25 A vast continuously increasing amount of functional data is now available  
26 thanks to recent high-throughput techniques: whole-genome sequences, gene  
27 expression or localization, mass-spectrometry analysis, *etc.*. However, at present,  
28 these complex data are difficult to interpret because of such features as their  
29 high-dimensionality, their inherent noise or even bias, the absence of stan-  
30 dardized representation. Organizing data into meaningful structures is highly  
31 desirable as a first step in unsupervised exploration of the large number of  
32 genes. Most biological mechanisms involve groupings of genes, gene products  
33 or proteins (*e.g.* enzymes) that act in a coordinated manner. Many clustering  
34 algorithms have been proposed over the last decade to decipher the message  
35 contained in DNA microarray data (see a list in the introduction of Kim et al.  
36 (2007)). In particular, Yeung et al. (2001) proposed a Gaussian mixture model  
37 to tackle this issue. This latter method and many others have the drawback to  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51

52  
53 \* Corresponding author.

54 *Email addresses:* juliette.blanchet@inrialpes.fr (Juliette Blanchet),  
55  
56 matthieu@bioass.ac.uk (Matthieu Vignes).

57  
58 <sup>1</sup> Present address: SLF, Fluelastrasse, 11, 7260 Davos Dorf, Switzerland  
59  
60

1  
2  
3 consider gene measurements to be independent. Hence we proposed in a pre-  
4 vious publication (Vignes and Forbes (2007)) an extension of this approach  
5 to account for individual features (*e.g.* microarray data) and dependencies be-  
6 tween genes in a united framework based on *Hidden Markov Random Fields*  
7 (HMRF).  
8  
9

10  
11  
12  
13  
14  
15 All clustering methods above use a full matrix of expression data as an input.  
16 An unfortunate feature of microarray experiments and other high-throughput  
17 technologies is that they often produce multiple missing values (McLachlan  
18 et al. (2004)). Most of the time, these missing entries appear because of vari-  
19 ous experimental issues (Troyanskaya et al. (2001), Bo et al. (2004)): dust or  
20 scratches on the slide, corrupted images, difficulties in measuring fluorescence  
21 intensity, systematic error of the robot that drops the probes, problem with  
22 precise gene spotting on the array...  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

33 A common practice –*case deletion*– is to remove genes and/or arrays from the  
34 analysis to end up with a fully-observed matrix on which classical approaches  
35 can be applied. However this approach can lose important information. Up  
36 to 90% of genes (rows) or experimental conditions (columns) can be affected  
37 (Ouyang et al. (2004)). It can also conceal interactions in a network. An al-  
38 ternative approach is to replace the missing values by zeros or by column/row  
39 means. Such a naive filling-in strategy is a particular case of *single impu-*  
40 *tation*. It is known to cause spurious estimation of summary statistics. The  
41 subsequent clustering results can be misleading (Little and Rubin (2002)).  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 Several more sophisticated methods have been proposed since the pioneering  
54 work of Troyanskaya et al. (2001). Most of them propose single imputation  
55 methods to transform the data matrix into a full matrix as needed by subse-  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

quent classical statistical analysis (see Troyanskaya et al. (2001), Oba et al. (2003), Bo et al. (2004), Ouyang et al. (2004)...). More recently, promising approaches make use of multiple imputation (Sehgal et al. (2005)) or iterative alternate blended clustering and missing values estimation (Kim et al. (2007)). Hu et al. (2006) proposed to improve classical estimation procedures by incorporating a large reference microarray dataset to define a general context for each gene. Nevertheless, none of these approaches take into account relationships imposed by the biological system between genes.

We propose to tackle **both issues of clustering and missing data imputation in a statistical framework. To our knowledge, these two issues have never been tackled simultaneously for dependent data. The clustering methodology has already been presented in a previous work (Blanchet and Vignes (2007)). In this paper, its efficiency is highlighted on varied datasets. The methodology is also expanded to the important issue of missing value reconstruction.** Instead of imputing values prior to the analysis, our integrated approach makes the best use of the statistical framework we consider. Estimation of missing values is *made a posteriori*, based on the network, the observed individual data and the clustering pattern. We are hence able both to quantitatively compare the quality of missing data estimations and to assess the biological significance of our results in regards of approaches cited above. Given the huge amount of such algorithms, we tested algorithms reported to work well (Brock et al. (2008)) and for which we were able to retrieve the corresponding algorithms. We emphasize that our model can be useful in a great range of applications for clustering biological entities of interest such as genes, proteins, metabolites in post-genomics studies. It requires individual possibly incomplete measurements taken on these entities related by a relevant interaction network. Hence

our method is neither organism- nor data-specific.

The present paper is organized as follows: the statistical model is presented in Section 2 with the EM-like estimation procedure, the classification framework and the reconstruction of missing observations. It provides *a posteriori* probabilities of entity (*e.g.* gene) classification given the observed data. This can be seen as a confidence measure of assignment. Experiments on synthetic data are reported in Section 3 while results on real yeast cell-cycle data combined with a network of interacting proteins are presented in Section 4.

## 2 Markovian model for clustering and imputation with missing data

### 2.1 Model

In what follows, we assume that values are *Missing At Random* (MAR, Little and Rubin (2002)). The fact that a datum is missing is not related to its actual unobserved value. A particular case of MAR is when the missingness process does not depend at all on the data, as for example when a dust on the slide produces missing values. Data are then said to be *Missing Completely At Random* (MCAR, Little and Rubin (2002)). An advantage of MAR hypothesis is that maximum likelihood can be estimated independently on the missingness process (Little and Rubin (2002)). In real applications, the MAR assumption might not be true as regards the phenomenon generating missing values. Just think of censorship issues due to machine limits of detection. Data are then said to be *Not Missing At Random* (NMAR). Methods based on MAR assumption can however produce satisfactory results if observed values contain enough

1  
2  
3 information to predict missing values with a likelihood approach. Simulations  
4  
5 in Section 3 show that inferences made by our model under MAR assumption  
6  
7 lead to satisfactory results even on NMAR data.  
8  
9

10  
11 We present in this section a statistical model for clustering and imputing  
12 incomplete dependent data. We refer to the entities of interest (pixels in Sec-  
13 tion 3, genes in Section 4) as *sites*, which we assume to be in interaction.  
14  
15 **These interactions can be due to spatial proximity as for the pixel images  
16 of Section 3, or due to biological relationships as for the genes of Section 4.**  
17  
18 **We further assume that experiments conducted on these sites create incom-  
19 plete data.** We denote  $\mathcal{S}$  the set of  $N$  sites and  $\mathbf{x} = \{\mathbf{x}_i \in \mathbb{R}^D\}$  the  $N \times D$   
20 matrix of observations, for which some entries are missing. For each  $i \in \mathcal{S}$ ,  
21 we write  $o_i \subset \llbracket 1, D \rrbracket$  the indices corresponding to the observed values  $x_{id}$   
22 and  $m_i$  the complementary indices for missing values ( $o_i \cup m_i = \llbracket 1, D \rrbracket$ ). We  
23 shall denote  $\mathbf{x}_i^{o_i} = \{x_{id}, d \in o_i\}$  the vector of observed data at site  $i$ ,  
24  $\mathbf{x}_i^{m_i} = \{x_{id}, d \in m_i\}$  the vector of missing data at site  $i$ ,  $\mathbf{x}^o = \{\mathbf{x}_i^{o_i}, i \in \mathcal{S}\}$   
25 the set of observed data and  $\mathbf{x}^m = \{\mathbf{x}_i^{m_i}, i \in \mathcal{S}\}$  the set of missing data.  
26  
27 We address the issue of clustering *i.e.* distinguishing meaningful groups in a  
28 dataset. In other words, each site  $i \in \mathcal{S}$  has to be assigned one of the  $K$  labels  
29  $z_i \in \llbracket 1, K \rrbracket$ . Dependencies between sites are maintained by an interaction net-  
30 work defining a neighbourhood structure. The model we consider is a *Hidden*  
31 *Markov Random Field* (HMRF), meaning that the hidden labels (or clusters)  
32 follow a Markov Random Field distribution. This is the generalization of the  
33 one-dimensional *Hidden Markov Chains* (also referred to as Hidden Markov  
34 Models, HMM) to higher dimensions, needed to deal with a graph of interac-  
35 tions. In this paper, we restrict to the widely-used Potts model for which the  
36 joint Markovian (or Gibbs) distribution of labels  $\mathbf{Z} = \{Z_i, i \in \mathcal{S}\}$  is:  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



$$P_G(\mathbf{z}) = W^{-1} \exp(\beta \sum_{i \sim j} 1_{z_i = z_j}) \quad (1)$$

where  $i \sim j$  denotes neighbours sites in the network (*i.e.* linked by an edge). Note that the distribution  $P_G(\mathbf{z})$  above depends on a single parameter  $\beta$  controlling the "smoothness" of the classification: the higher  $\beta$ , the more likely two neighbouring sites are to be assigned to the same cluster.

We eventually assume that data are independent conditionally on classes, this is to say that  $P(\mathbf{x}|\mathbf{z}) = \prod_{i \in \mathcal{S}} P(\mathbf{x}_i|z_i)$ . In this paper, class-dependent distributions are considered to be Gaussian:  $P(\cdot|Z_i = k) \sim \mathcal{N}(\mu_k, \Sigma_k)$ .

## 2.2 Parameter estimation

For sake of clarity, we denote  $\theta_k = (\mu_k, \Sigma_k)$  the parameters of the  $k$ -th Gaussian distribution. The single parameter of the Potts distribution is, as in equation (1), denoted by  $\beta$ . Without *a priori* knowledge, the full set of parameters  $\Psi = (\theta_1, \dots, \theta_K, \beta)$  is unknown and has to be estimated.

The method we propose is a maximum likelihood-based approach. The principle is to choose the most likely parameters  $\Psi$  for the observed data. Bayesian techniques would offer an alternative way to draw inference from the likelihood function. Such methods are not considered here. We use the Expectation-Maximisation (EM) algorithm (Dempster et al. (1977)). At iteration ( $q$ ), a current estimate  $\Psi^{(q-1)}$  is available and the algorithm maximizes the function  $Q$  defined, in a missing data framework, as:

$$Q(\Psi|\Psi^{(q-1)}) \equiv \mathbb{E}[\log P(\mathbf{x}^o, \mathbf{X}^m, \mathbf{Z}|\Psi)|\mathbf{x}^o, \Psi^{(q-1)}] \quad (2)$$

to get updated  $\Psi^{(q)}$ . **It is worth stressing that expectation in equation (2)**

1  
2  
3 is not only taken over unknown labels  $\mathbf{Z}$  (as in the classical fully-observed  
4 data case), but also over missing values  $\mathbf{X}^m$ . An EM algorithm with incom-  
5 plete data has already been studied for *Independent Mixture Model* (IMM)  
6 (see Little and Rubin (2002)), as well as for *Hidden Markov Chain Model* (see  
7 Celeux and Durand (2007) for a recent reference). To our knowledge, it has  
8 never been studied for any HMRF model. Due to the more complex depen-  
9 dence structure, expectation of equation (2) is not explicitly tractable for an  
10 HMRF model, as both the normalizing constant  $W$  of equation (1) and the  
11 conditional probability  $P(\mathbf{z}|\mathbf{x})$  cannot be computed exactly. Approximations  
12 are then required to make the algorithm tractable.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

26 In this paper, we propose to use a mean field-like approximation of the Marko-  
27 vian *a posteriori* distribution  $P_G(\mathbf{z}|\mathbf{x})$  similar to the distribution proposed  
28 in Celeux et al. (2003) in the framework of complete data clustering. It was  
29 shown to be more efficient than the most widely-used clustering approaches on  
30 both simulated and real data (see Celeux et al. (2003) and Vignes and Forbes  
31 (2007)). This suggests good properties of convergence; local convergence of a  
32 very similar algorithm has been proven in Forbes and Fort (2007). The algo-  
33 rithm developed here extends the procedure to the missing data framework.  
34 Informally, the idea of our algorithm when considering a particular site  $i$  is to  
35 neglect the fluctuations of the neighbouring sites by setting them to fixed values  
36  $\tilde{z}_j, j \in N_i$  (means for example). The untractable Markovian distribution  $P_G(\mathbf{z})$   
37 is then approximated by the tractable factorized distribution  $\prod_{i \in \mathcal{S}} P_G(z_i | \tilde{z}_{N_i})$   
38 where  $\tilde{z}_{N_i}$  denotes the set  $\{\tilde{z}_j, j \in N_i\}$ . Due to conditional independence,  
39  $P(\mathbf{x}, \mathbf{z})$  is also approximated as a factorized distribution and equation (2)  
40 becomes tractable. Values  $\tilde{z}_i$  being *a priori* unknown, mean field-like approx-  
41 imations lead to an iterative EM-like algorithm repeating two steps. In what  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

follows, values for the  $\tilde{z}_i$ 's are simulated, as recommended in Celeux et al. (2003) in a complete data framework. More precisely, starting with parameters  $\Psi^{(0)}$ , at iteration  $(q)$ ,

- (1) For each site  $i \in \mathcal{S}$ , simulate from the observed data  $\mathbf{x}_i^{o_i}$  and the current parameter estimate  $\Psi^{(q-1)}$  a configuration  $\tilde{z}_i^{(q)}$  i.e. values for the  $Z_i$ 's.
- (2) Apply one step of the EM algorithm on the factorized model resulting from the mean field-like approximation to get updated estimates  $\Psi^{(q)}$  of the parameters.

This procedure will be referred in what follows as “*SFmiss algorithm*” standing for Simulated Field algorithm with missing values. This algorithm accounts for the Markovian structure of the data while factorizing the distribution on which EM is tractable.

In the E-step, *a posteriori* probabilities are computed for all  $i \in \mathcal{S}$  and  $k \in \llbracket 1, K \rrbracket$  by

$$\tilde{t}_{ik}^{(q)} = P_G(Z_i = k | \mathbf{x}_i^{o_i}, \tilde{z}_{N_i}^{(q)}) \quad (3)$$

The difference with the complete data case of Celeux et al. (2003) is that conditioning in (3) involves the observed data  $\mathbf{x}_i^{o_i} \in \mathbb{R}^{|\mathcal{O}_i|}$ , and not the whole vector  $\mathbf{x}_i \in \mathbb{R}^D$ . As in Celeux et al. (2003), the conditioning also includes neighbours through the  $\tilde{z}_{N_i}^{(q)}$  term.

In the M-step, parameters  $\Psi = (\theta_1, \dots, \theta_K, \beta)$  are updated. The updating of the Markovian parameter  $\beta$  remains unchanged as compared to the complete data case (see Celeux et al. (2003)). No analytical expression is available but the optimal  $\beta^{(q)}$  is unique and can easily be obtained numerically. Unlike  $\beta$ , the updating of the Gaussian class-dependent parameters  $\theta_k = (\mu_k, \Sigma_k)$ ,  $k \in$

$\llbracket 1, K \rrbracket$ , differs from the complete data case. Denote  $\Sigma_k^{o_i o_i} = \{(\Sigma_k)_{st}, s \in o_i, t \in o_i\}$ ,  $\Sigma_k^{o_i m_i} = \{(\Sigma_k)_{st}, s \in o_i, t \in m_i\} = (\Sigma_k^{m_i o_i})^T$  and  $\Sigma_k^{m_i m_i} = \{(\Sigma_k)_{st}, s \in m_i, t \in m_i\}$ . Then  $P(X_i^{m_i} | \mathbf{x}_i^{o_i}, \theta_k)$  is a Gaussian distribution with mean  $\eta_{ik}$  and covariance  $\Gamma_{ik}$  defined as:

$$\begin{aligned} \eta_{ik} &= \mu_k^{m_i} + \Sigma_k^{m_i o_i} (\Sigma_k^{o_i o_i})^{-1} (\mathbf{x}_i^{o_i} - \mu_k^{o_i}) \\ \Gamma_{ik} &= \Sigma_k^{m_i m_i} - \Sigma_k^{m_i o_i} (\Sigma_k^{o_i o_i})^{-1} \Sigma_k^{o_i m_i}. \end{aligned} \quad (4)$$

At iteration  $(q)$  the component  $s \in \llbracket 1, D \rrbracket$  of  $\mu_k$  is updated as:

$$(\mu_k^s)^{(q)} = \frac{\sum_i \tilde{t}_{ik}^{(q)} (r_i^s x_i^s + (1 - r_i^s) \eta_{ik}^s)^{(q)}}{\sum_i \tilde{t}_{ik}^{(q)}} \quad (5)$$

with  $r_i^s = 1$  if variable  $x_i^s$  is observed, 0 otherwise. Compared with the complete data case, equation (5) simply replaces the missing variable  $x_i^s$  by the conditional mean  $(\eta_{ik}^s)^{(q)}$  of the distribution  $P(X_i^{m_i} | \mathbf{x}_i^{o_i}, \theta_k^{(q)})$ .

Similarly, the component  $s, t \in \llbracket 1, D \rrbracket$  of  $\Sigma_k$  is updated as:

$$(\Sigma_k^{st})^{(q)} = \frac{\sum_i t_{ik}^{(q)} (S_{ik}^{st})^{(q)}}{\sum_i t_{ik}^{(q)}}$$

with for all  $i \in S, k \in \llbracket 1, K \rrbracket, s, t \in \llbracket 1, D \rrbracket$ ,

$$\begin{aligned} (S_{ik}^{st})^{(q)} &= r_i^s r_i^t (x_i^s - \mu_k^s)^{(q)} (x_i^t - \mu_k^t)^{(q)} \\ &+ r_i^s (1 - r_i^t) (x_i^s - \mu_k^s)^{(q)} (\eta_{ik}^t)^{(q)} - \mu_k^t)^{(q)} \\ &+ (1 - r_i^s) r_i^t (\eta_{ik}^s)^{(q)} - \mu_k^s)^{(q)} (x_i^t - \mu_k^t)^{(q)} \\ &+ (1 - r_i^s) (1 - r_i^t) \{ (\eta_{ik}^s)^{(q)} - \mu_k^s)^{(q)} (\eta_{ik}^t)^{(q)} - \mu_k^t)^{(q)} + \Gamma_{ik}^{st} \} \end{aligned}$$

It is worth stressing that, because of the  $\Gamma_{ik}^{st}{}^{(q)}$  term in the last factor, this is not equivalent to replacing a missing variable  $x_i^s$  by the mean  $(\eta_{ik}^s)^{(q)}$  of the

conditional distribution  $P(X_i^{m_i} | x_i^{o_i}, \theta_k^{(q)})$ . This is consistent with the remark that mean imputation technique lowers the estimated variance (Little and Rubin (2002)).

### 2.3 A posteriori classification and imputation

Running  $q_{max}$  steps of the SFmiss algorithm leads to estimates  $\Psi^{(q_{max})}$  of the model parameters and to configurations  $\hat{z}_i^{(q_{max})}$ ,  $i \in \mathcal{S}$ , for the mean field-like approximation. These quantities can then be used to both cluster sites and impute missing data. Due to the factorization of  $P(\mathbf{z}|\mathbf{x})$  resulting from mean field-like approximation, MAP (*Maximum A Posteriori*) and MPM (*Maximum Posterior Marginal*) classification rules are equivalent and consist in classifying a site  $i \in \mathcal{S}$  in :

$$\hat{z}_i = \arg \max_{k \in [1, K]} P(Z_i = k | \mathbf{x}_i^{o_i}, \tilde{z}_{N_i}^{(q_{max})}, \beta^{(q_{max})}) = \arg \max_{k \in [1, K]} \tilde{t}_{ik}^{(q_{max})} \quad (6)$$

Classification rule (6) involves therefore (i) the observed data  $\mathbf{x}_i^{o_i} \in \mathbb{R}^{|\mathcal{O}_i|}$  (and not the whole vector  $\mathbf{x}_i \in \mathbb{R}^D$  as in the complete data case) and (ii) the neighbours through the additional  $\tilde{z}_{N_i}^{(q_{max})} = \{\tilde{z}_j^{(q_{max})}, j \in N_i\}$  term. It accounts therefore explicitly for dependencies between sites. This is a clear advantage of our HMRF model over IMM.

Missing observations can also be *a posteriori* reconstructed. MAP (or MPM) rule leads to impute missing observations  $\mathbf{x}_i^{m_i}$  for sites  $i \in \mathcal{S}$  by the most likely values conditionally on observed  $\mathbf{x}_i^{o_i}$  and on class  $\hat{z}_i$ :

$$\hat{\mathbf{x}}_i^{m_i} = \arg \max_{\mathbf{x}_i^{m_i}} P(\mathbf{x}_i^{m_i} | \mathbf{x}_i^{o_i}, \hat{z}_i) = \eta_{i\hat{z}_i}. \quad (7)$$

Equation (7) can be seen as a mean imputation. It differs nevertheless from

the classical pre-processed mean imputation in several points:

- (i) It is performed *a posteriori*, and not as a pre-processing. Therefore it does not artificially biased the parameter estimation (in particular the  $\Sigma_k$ 's, see Little and Rubin (2002)) required for classification.
- (ii) Relationships between sites are taken into account through the classification  $\hat{z}_i$  which, as seen previously (equation (6)), involves the neighborhood structure.
- (iii) The mean is not computed over all sites, but only over sites belonging to the same cluster and therefore sharing information (related biological functions for example, as in Section 4).

### 3 Illustration on synthetic data

The purpose of this section is to illustrate the differences between our method and standard imputation methods, and to emphasize the general aspects of the former with respect to the latter, for both classification and imputation issues. From among several exercises we have performed, we present here some results related to a 4-class synthetic image. Data were obtained as follows. Starting from the synthetic image, a noisy image is generated by considering that observations belonging to the  $k$ -th class (for  $k = 1, \dots, 4$ ) are realizations of a 4-dimensional ( $D = 4$ ) non diagonal Gaussian distribution, with mean  $\mu_k = (k, k, k, k)^T$  and covariance matrix  $\Sigma_k$  with diagonal terms equal to 0.5 and non diagonal terms to 0.2. We then consider two ways of producing missing data. The first one removes randomly a given proportion of data (MCAR case). The second one removes a given proportion of the highest and the lowest data (left and right censorship, NMAR case).

1  
2  
3 The classification results obtained respectively by our method and by the  
4 Markovian Simulated-Field algorithm (referred to as SF, Celeux et al. (2003))  
5 with various prior imputations are shown in Figure 1. Imputation techniques  
6 considered are filling in with zeros (ZERO+SF), with column means (MEAN+SF),  
7 or using standard imputation methods such as K-Nearest Neighbors (Trojan-  
8 skaya et al. (2001), KNN+SF), Bayesian Principal Component Analysis (Oba  
9 et al. (2003), BPCA+SF) or Support Vector Regression (Wang et al. (2006),  
10 SVR+SF). The Local Least Square Impute method (Kim et al. (2005)) gave  
11 poor results on our data and are not reported here. It appears that the SFmiss  
12 algorithm performs better than tested imputation methods, although the un-  
13 derlying model is the same: an HMRF model. To assess the gain in using a  
14 Markovian model, we also compare with the IMM, with parameters estimated  
15 by the EM algorithm with incomplete data (referred to as EMmiss, see Little  
16 and Rubin (2002)).  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

33  
34 [Figure 1 about here.]  
35

36  
37 [Figure 2 about here.]  
38  
39

40  
41 As compared with IMM, it appears that taking dependencies into account  
42 (through the use of Markovian models) improves the results significantly. Fur-  
43 thermore, SFmiss algorithm provides a way of modelling uncertainty over miss-  
44 ing observation values, leading to better classifications and imputations. It can  
45 also be noted that our algorithm performs well even when the correct underly-  
46 ing model is not set as an hypothesis: the synthetic image is not a realization  
47 of a Potts model and censored data are NMAR! The censored data case seems  
48 to be more difficult than MCAR case but SFmiss provides reasonably good  
49 classifications for high percentages of missing values (up to 60%, see Figure 1  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

and visualization on Figure 2).

As mentioned in Section 2.3, in addition to providing a classification, our algorithm has the ability to reconstruct -or impute- missing data. Figure 3 displays imputation errors for the methods mentioned above. These errors are measured by the normalized Root Mean Squared Error (RMSE): if  $\hat{\mathbf{x}}$  is the imputed data matrix, *i.e.*, an estimate of the complete data matrix  $\mathbf{x}$ , the RMSE is defined as:

$$\text{RMSE} = \sqrt{\frac{\text{mean}\{(\mathbf{x} - \hat{\mathbf{x}})^2\}}{\text{mean}\{\mathbf{x}^2\}}},$$

where  $\mathbf{x}^2$ , for example, is component-wise.

[Figure 3 about here.]

Results presented in Figure 3 confirm that SFmiss offers a substantial improvement as concerns imputation issue. Note that, as average measures of error, RMSE for KNN, BPCA and SVR imputation techniques are similar although corresponding imputed component-wise values can be quite different, leading to different classification as reported in Figure 1.

## 4 Experiments on yeast cell-cycle data

### 4.1 Individual Data

Data of Spellman et al. (1998) on *Saccharomyces cerevisiae* that focuses on the identification of cell-cycle regulated genes were used. These data are expression profiles from yeast cultures synchronized by different methods. The full data set consists of a 77 dimensional vector for each of the 6179 genes. The initial



1  
2  
3 dataset has 5% overall missing entries. In the following, we will focus on the  
4  
5 *cdc28* experiment initially performed by Cho et al. (1998) (Spellman et al.  
6  
7 (1998) used this data along with their own for their analysis). Yeast were  
8  
9 synchronized by stopping them in late G1 phase of the cell cycle. 17 time  
10  
11 points (dimensions) were collected every 10 minutes so nearly two cell cycles  
12  
13 have occurred.  
14  
15  
16  
17

#### 18 4.2 Interaction Network

19  
20  
21  
22  
23 Available biological networks contain a significant amount of information that  
24  
25 should not be ignored to provide optimal statistical analysis of the machinery  
26  
27 of the cell. Our aim is to build a graph with a biological entity (gene) at each  
28  
29 node. An edge will stand for a confirmed link between two entities: interac-  
30  
31 tions between genes, gene products, complexes of proteins, families, metabolic  
32  
33 pathways,...

34  
35  
36  
37 For network data, we use the release 7 of STRING (von Mering et al. (2007)),  
38  
39 a consistent database of known and predicted protein-protein interactions. It  
40  
41 gathers information from a wide variety of different sources: genomic context,  
42  
43 literature knowledge, physical interactions, *etc.* The current version contains  
44  
45 401 948 curated interactions for 5611 genes of *Saccharomyces cerevisiae*. Note  
46  
47 that two or more interactions may occur between the same couple of genes,  
48  
49 because different kinds of interactions are considered. We selected the inter-  
50  
51 section between the set of 800 genes identified as cell-cycle regulated in Spell-  
52  
53 man et al. (1998) and those contained in the STRING database. The resulting  
54  
55 graph consists of 612 nodes (genes) and 3530 edges accounting for one or more  
56  
57 interaction(s), which are given equal weight (see Section 5 for a discussion on  
58  
59  
60

1  
2  
3 this aspect).

4  
5  
6 Unlike with synthetic data, the correct number of clusters is unknown. Bayes  
7 Information Criterion (BIC, Schwarz (1978)), a penalized likelihood that ac-  
8 counts for the complexity of the model, is widely used to tackle this issue.  
9  
10 In our HMRF setting, BIC is untractable and we used a mean-field approx-  
11 imation as described in Forbes and Peyrard (2003). We allowed the number  
12 of clusters to range from 2 to 12.  $K = 9$  was the selected number of clusters  
13 (data not shown).  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

#### 24 *4.3 Results and discussion*

25  
26  
27  
28  
29 We performed imputation and 9-class clustering of the yeast data using SFmiss  
30 and several other algorithms for comparison. This section aims at assessing  
31 the quality of the produced clusters, a difficult task as there is no consensus  
32 criterion to rely on. We illustrate the gain in using our approach on some  
33 specific biological features.  
34  
35  
36  
37  
38  
39

40  
41 We first check whether the output clusters of our model are well-suited to sum-  
42 marize biological knowledge compared to other algorithms. A general trend is  
43 that the SFmiss algorithm gathers interacting genes better than other algo-  
44 rithms. More precisely, genes clustered together by SFmiss have more internal  
45 connections than those clustered together by other imputation methods (al-  
46 though they rely on SF algorithm that takes the network into account). It  
47 reveals that the way SFmiss deals with missing observations is certainly more  
48 appropriate. This is consistent with the spatial parameter  $\beta$  of our model:  $\beta$   
49 is estimated to 0.41, which means that the neighbourhood plays a significant  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 role.  
4  
5

6 The clusters reliability can be quantified using *Gene Ontology* (GO, The Gene  
7 Ontology Consortium (2000)) terms representativeness **focusing** on the biolog-  
8 ical process under study: yeast cell cycle. The more GO terms present in the  
9 data set are shared by genes in the same cluster, the more *sensitive* the method  
10 is. The more the 9 clusters isolate different parts of GO, the more *specific* the  
11 method is.  
12  
13  
14  
15  
16  
17  
18

19  
20 [Table 1 about here.]  
21  
22

23 For each GO category, a test is performed to determine whether the category  
24 is over-represented in each cluster. Under-representation can be tested as well  
25 but its analysis is not presented here for brevity reasons. P-values in Table 1 are  
26 computed with the FDR correction of Benjamini and Hochberg (1995) which  
27 is widely used and has proven its efficiency. Very low P-value indicates that  
28 the tested GO term is over-represented in the analyzed cluster, and therefore  
29 that the algorithm successfully grouped genes sharing this biological feature.  
30 For clarity, we only compare in Table 1 our SFmiss algorithm to IMM with  
31 missing data (EMmiss, see Little and Rubin (2002)), and to HMRF model  
32 with prior KNN imputation (SF+KNN). Other tested imputation methods  
33 (mean imputation, BPCA, SVR...) did not give better results. Apart from  
34 few exceptions (as cluster  $k = 8$ ), P-values of Table 1 suggest that SFmiss  
35 algorithm performs better at grouping genes with similar annotations than  
36 other algorithms, *i.e.* is more sensitive.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 Another nice feature of our SFmiss algorithm is that it does not only sum-  
54 marize known biological knowledge but can give directions for putative func-  
55 tions on components of living organisms based on the clustering results. For  
56  
57  
58  
59  
60

1  
2  
3 a detailed example, genes *ydl105w*, *yer111c*, *ykr077w*, *yjl196c*, *ylr212c* and  
4  
5  
6 *ynl082w* are all classified in cluster 1 by SFmiss, whereas there are dispatched  
7  
8 in various clusters by EMmiss. But all of them are brought into play during  
9  
10 cell cycle processes: mitotic spindle complex repair (*ylr212c*) or G1/S tran-  
11  
12 sition of the mitotic cycle (*yer111c*) for example. *ykr077w* is annotated as a  
13  
14 putative transcription activator. Our method suggests that this annotation is  
15  
16 fully coherent and that this gene plays a key role either as a cell cycle regulator  
17  
18 or as a regulated gene of the process.

19  
20  
21 We can also illustrate the advantage of accounting for missing values in a  
22  
23 united fashion as compared to prior filling-in with Troyanskaya et al. (2001)  
24  
25 KNNimpute. Genes *ybl002w*, *ygl093w* and *ypl269w* belong to SFmiss cluster  
26  
27 4 and are dispatched in various clusters by KNN+SF. Their annotations are  
28  
29 making sense when compared with those of their cluster and confirm a possible  
30  
31 functional description of the cluster: chromatin assembly, required for accurate  
32  
33 chromosome segregation localized to the nuclear side of the spindle pole body  
34  
35 and required for cytoplasmic microtubule orientation in yeast (polarization)  
36  
37 respectively.

38  
39  
40  
41 The interpretation of clusters when compared to temporal classes of the cell  
42  
43 cycle (G1, S, S/G2, G2/M and M/G1, Spellman et al. (1998)) emphasizes  
44  
45 the specificity of the SFmiss algorithm, observed to a lesser extent in clusters  
46  
47 resulting from other algorithms. Cluster 0 is almost entirely included in Spell-  
48  
49 man et al. (1998) G2/M group and cluster 1 is in G1 just like cluster 5. Cluster  
50  
51 2 include genes regulated in late G2, M and early G1 phases (quite broad, cer-  
52  
53 tainly a reason why no specific function is highlighted in this cluster). Cluster  
54  
55 3 is similar but with an earlier start. Cluster 8 is focused on M-regulated genes.  
56  
57 Cluster 4 shows its temporal peak in S phase. Lastly, cluster 6 has many genes  
58  
59  
60

1  
2  
3 from early G2 to M. These interpretations are corroborated if we investigate  
4 the expression profiles for each meaningful cluster. Examples of such addi-  
5 tional evidences are given in Figure 4. These profiles are very similar to those  
6 obtained in Figure 4 (C, D) by Cho et al. (1998) for annotated genes.  
7  
8  
9

10  
11  
12  
13  
14 [Figure 4 about here.]  
15  
16  
17

18 Last but not least we would like to present another major advantage of our  
19 approach: it responds with a much greater level of stability than other tested  
20 methods when the number of observed data decreases. This is illustrated in  
21 Figure 5. Additional missing values were generated under MCAR. We then  
22 compared (i) the new classification with the initial one to assess stability of  
23 the classification (Figure 5, left panel), (ii) the new imputed values with the  
24 initial ones to assess stability of the imputation using RMSE (Figure 5, right  
25 panel).  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

37 [Figure 5 about here.]  
38  
39  
40

41 Apart from SFmiss, all algorithms show dramatical instability when the rate  
42 of added missing value increases above 6%. Note that a 11% difference with  
43 the initial classification corresponds to one cluster which is fully lost. This  
44 suggests that these algorithms have an unsatisfactory behaviour when they  
45 are facing datasets even with "as few as" nearly 10% of total missing data in  
46 the favourable MCAR case. On the contrary, the SFmiss algorithm shows an  
47 interesting stability towards the rate of missing value. Its performance impairs  
48 significantly above 25% of overall missing data which is quite acceptable as  
49 regards real encountered situations.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 5 Conclusions and future work

Missing data can bring many difficulties in data analysis simply because most data analysis procedures were not designed for them. This is particularly true in the context of post-genomic data integration. Data absence is usually a nuisance, not the focus of inquiry. We presented a comprehensive integrated statistical tool for modelling individual measurements that have a network-dependant structure. We overcame the conceptual and computational challenges and demonstrated the good features of our method on both synthetic and real biological datasets.

Our results prompt further studies. It would be interesting to analyze a dataset on a whole-genome scale. We restricted our analysis to genes with prior knowledge for validation purpose. Another prospect would be to take into account the missingness mechanism to improve performances on NMAR generated data. A possibility would be to consider the missingness mechanism as a third process and to use the recent triplet Markov field model of Blanchet and Forbes (2008) that would have to be extended to the framework of incomplete observations. A final plan is to account for missing edges as we did for missing individual measurements; biological interaction data are known to be incomplete or noisy. A first step would be to consider confidence levels for interactions as their reliability vary a lot when reported by two-hybrid screening for example. This feature is being developed in our software.

## Acknowledgements

The work was partially funded by the Scottish Government.

## Supplementary Materials

The SpaCEM3 software used in this study is available at <http://spacem3.gforge.inria.fr/> along with datasets.

## References

- Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate—a practical and powerful approach to multiple testing, *Journal of the Royal Statistical Society, Series B* **57**, 289–300.
- Blanchet, J. and Vignes, M. (2007) Combined expression data with missing values and gene interaction network analysis: a Markovian integrated approach, 366–373, in *Proceedings of the 7th IEEE BIBE*.
- Blanchet, J. and Forbes, F. (2008) Triplet Markov fields for supervised classification of complex structure data, *IEEE PAMI*, to appear.
- Bo, T.H., Dysvik, B. and Jonassen, I. (2004) LSImpute: accurate estimation of missing values in microarray data with least square methods, *Nucleic Acids Research* **32**, e34.
- Brock, G.N., Shaffer, J.R., Blakesley, R.E., Lotz, M.J., Tseng, G.C. (2008) Which missing value imputation method to use in expression profiles: a comparative study and two selection schemes, *BMC Bioinformatics*, **9**:12.
- Celeux, G., Forbes, F. and Peyrard, N. (2003) EM procedures using mean-field like approximations for Markov-model based image segmentation, *Pattern Recognition* **36**, 131–144.
- Celeux, G. and Durand, J.B. (2007) Selecting hidden Markov model state number with cross-validated likelihood, *Computational Statistics*, to appear.
- Cho, R.J., Campbell, M.J., Winzeler, E.A., Steinmetz, L., Conway, A., Wod-

- 1  
2  
3 icka,L., Wolfsberg,T.G., Gabrielian,A.E., Landsman,D., Lockhart,D.J. and  
4  
5 Davis,R.W. (1998) A genome-wide transcriptional analysis of the mitotic  
6  
7 cell cycle, *Molecular Cell* **2**, 65–73.  
8  
9  
10 Dempster,A.P., Laird,N.M., Rubin,D.B. (1977) Maximum likelihood from in-  
11  
12 complete data via the EM algorithm, *Journal of the Royal Statistical Soci-*  
13  
14 *ety, Series B* **39**, 1–38.  
15  
16 Forbes,F. and Peyrard,N. (2003) Hidden Markov random field model selection  
17  
18 criteria based on mean field-like approximations, *IEEE Transactions on*  
19  
20 *Pattern Analysis and Machine Intelligence* **25**, 1089–1101.  
21  
22 Forbes, F. and Fort, G. (2007) Combining Monte Carlo and mean-field-like  
23  
24 methods for inference in hidden Markov random fields, *IEEE Transactions*  
25  
26 *on Image Processing* **16**, 824-837.  
27  
28 Hu,J., Li,H., Waterman,M.S. and Zhou,X.J. (2004) Integrative missing value  
29  
30 estimation for microarray data, *BMC Bioinformatics* **7**, 449.  
31  
32 Kim,H., Golub,G.H. and Park,H. (2005) Missing value estimation for DNA  
33  
34 microarray gene expression data: local least squares imputation, *Bioinform-*  
35  
36 *atics* **21**, 187–198.  
37  
38 Kim,D.W., Lee,K.Y., Lee,K.H. and Lee,D. (2007) Towards clustering of in-  
39  
40 complete microarray data without the use of imputation *Bioinformatics*,  
41  
42 **23**, 107–113.  
43  
44 Little,R.J. and Rubin,D.B. (2002) *Statistical analysis with missing data*, New-  
45  
46 York: Wiley, second edition.  
47  
48 McLachlan G.J., Do K.A. and Ambroise C. (2004) *Analyzing microarray gene*  
49  
50 *expression data*, New Jersey: Wiley.  
51  
52 Oba,S., Sato,M., Takemasa,I., Monden,M., Matsubara,K. and Ishii,S. (2003) A  
53  
54 Bayesian Missing value estimation method, *Bioinformatics* **19**, 2088–2096.  
55  
56 Ouyang M., Welsh W.J., Georgopoulos P. (2004) Gaussian mixture clustering  
57  
58  
59  
60



- 1  
2  
3 and imputation of microarray data, *Bioinformatics* **20**, 917–923.  
4  
5 Schwarz G. (1978), Estimating the dimension of a model, *Annals of Statistics*  
6  
7 **6**, 131–134.  
8  
9 Sehgal,M.S., Gondal,I. and Dooley,L.S. (2005) Collateral missing value im-  
10  
11 putation: a new robust missing value estimation algorithm for microarray  
12  
13 data, *Bioinformatics* **21**, 2417–2423.  
14  
15 Spellman,P.T., Sherlock,G., Zhang,M.Q., Iyer,V.R., Anders,K., Eisen,M.B.,  
16  
17 Brown,P.O., Botstein,D. and Futcher,B. (1998) Comprehensive identifica-  
18  
19 tion of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by  
20  
21 microarray hybridization, *Molecular Biology of the Cell* **9**, 3273–3297.  
22  
23 The Gene Ontology Consortium (2000) Gene Ontology: tool for the unification  
24  
25 of biology, *Nature Genetics* **25**, 25–29.  
26  
27 Troyanskaya,O., Cantor,M., Sherlock,G., Brown,P., Hastie,T., Tibshirani,R.,  
28  
29 Botstein,D. and Altman,R.B. (2001) Missing values estimation methods for  
30  
31 DNA microarrays, *Bioinformatics* **17**, 520–525.  
32  
33 Vignes,M. and Forbes,F. (2007) Gene clustering via integrated Markov models  
34  
35 combining individual and pairwise features, *IEEE/ACM Transactions on*  
36  
37 *Computational Biology and Bioinformatics* in press.  
38  
39 von Mering C., Jensen L.J., Kuhn M., Chaffron S., Doerks T., Krüger B., Snel  
40  
41 B., Bork P. (2007) STRING 7 - recent developments in the integration and  
42  
43 prediction of proteins interactions, *Nucleic Acids Research* **35**, D358–62.  
44  
45 Wang,X., Li,A., Jiang,Z. and Feng H. (2006) Missing value estimation for DNA  
46  
47 microarray gene expression data by Support Vector Regression imputation  
48  
49 and orthogonal coding scheme, *BMC Bioinformatics* **7**, 32.  
50  
51 Yeung,K.Y., Fraley,C., Murua,A., Raftery,A. and Ruzzo,L. (2001) Model-  
52  
53 based clustering and data transformations for gene expression data, *Bioin-*  
54  
55 *formatics* **17**, 977–987.  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## List of Figures

- 1 Experiments on image-like simulated data: percentage of misclassified pixels vs. percentage of missing data for (left) randomly missing data (MCAR case) and (right) censored data (NMAR case). 25
- 2 Experiments on image-like simulated data: visualization of the synthetic image results (*i.e.* obtained clusters) for various percent of missing data (30%, 50%, 60% and 70%,) with SFmiss in the NMAR case. 26
- 3 Experiments on image-like simulated data: RMSE vs. percentage of missing data for (left) randomly missing data (MCAR case), (right) censored data (NMAR case) on the synthetic dataset. 27
- 4 Yeast cell-cycle data: examples of expression profiles for three SFmiss clusters (black solid line is the mean profile and dashed lines indicate standard deviation from the mean). 28
- 5 Yeast cell-cycle data. Left: Percentage of error for different algorithms vs. percentage of added (to the inherent approximately 5% in the dataset) missing value. Right: RMSE vs. percentage of added missing value. Algorithms are the same as in Section 3; PREV+SF has prior imputation thanks to an autoregressive model with lag 1 usually suited for time series. 29

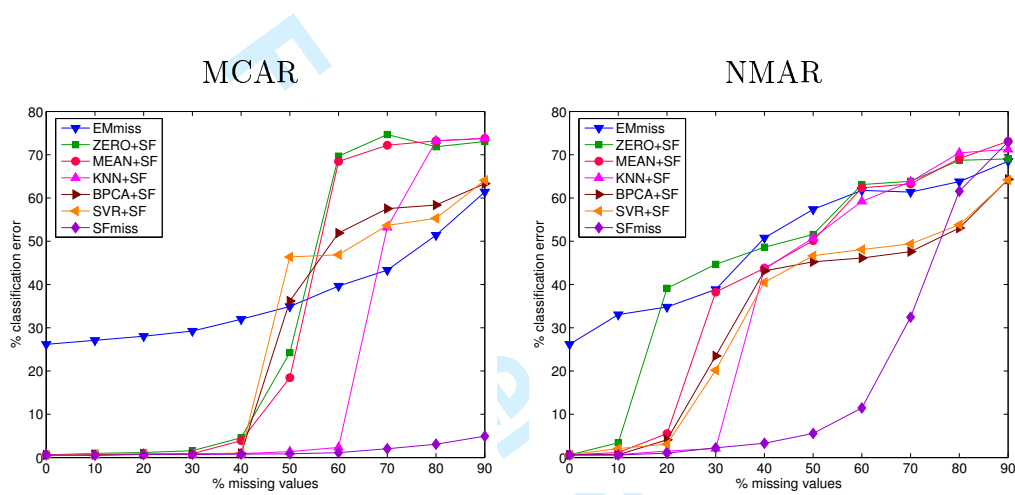


Figure 1. Experiments on image-like simulated data: percentage of misclassified pixels vs. percentage of missing data for (left) randomly missing data (MCAR case) and (right) censored data (NMAR case).

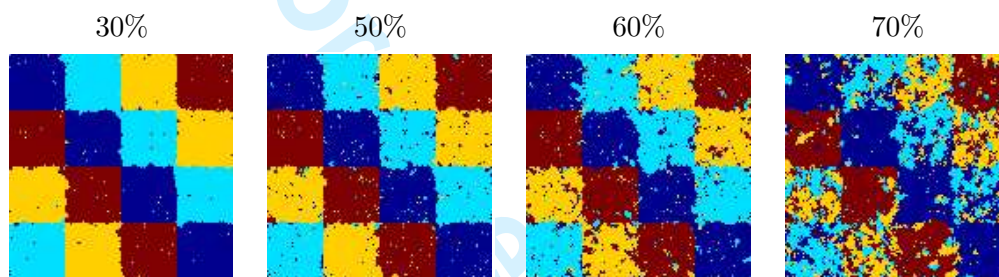


Figure 2. Experiments on image-like simulated data: visualization of the synthetic image results (*i.e.* obtained clusters) for various percent of missing data (30%, 50%, 60% and 70%,) with SFmiss in the NMAR case.

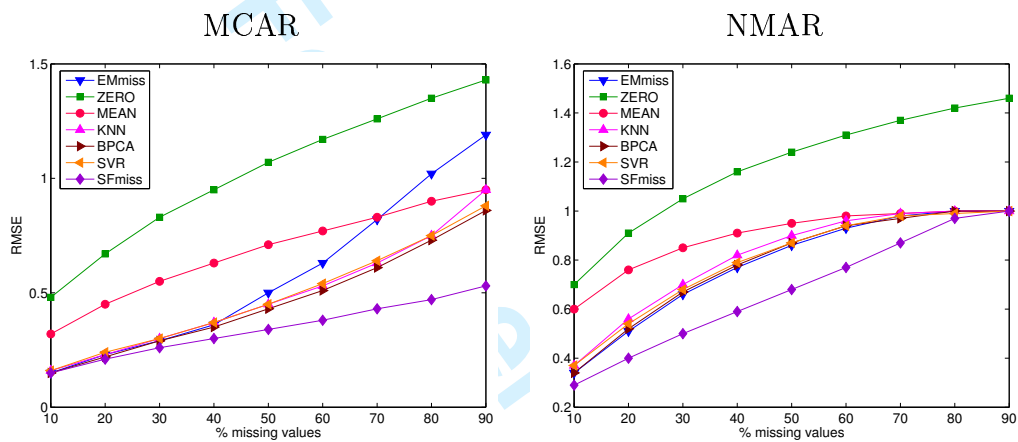
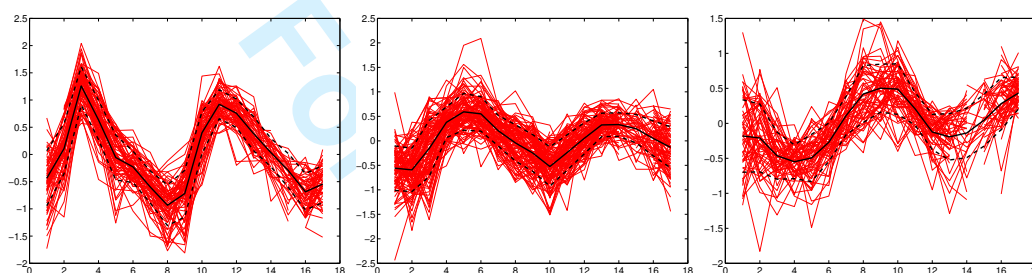


Figure 3. Experiments on image-like simulated data: RMSE vs. percentage of missing data for (left) randomly missing data (MCAR case), (right) censored data (NMAR case) on the synthetic dataset.



SFmiss cluster 1 concerned with G1 phase and DNA replication

SFmiss cluster 4: S phase, chromosome segregation and biosynthesis (*e.g.* S met. proc.)

SFmiss cluster 8 including M phase, polarization and ATP activity

Figure 4. Yeast cell-cycle data: examples of expression profiles for three SFmiss clusters (black solid line is the mean profile and dashed lines indicate standard deviation from the mean).

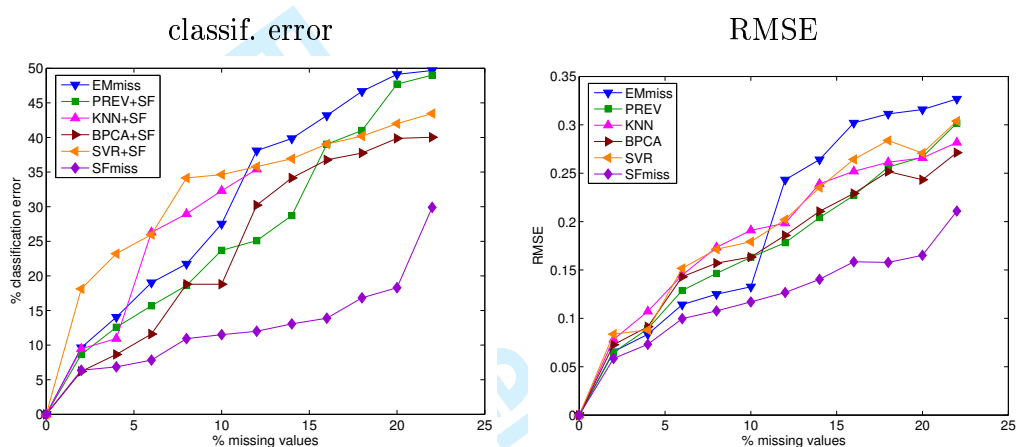


Figure 5. Yeast cell-cycle data. Left: Percentage of error for different algorithms vs. percentage of added (to the inherent approximately 5% in the dataset) missing value. Right: RMSE vs. percentage of added missing value. Algorithms are the same as in Section 3; PREV+SF has prior imputation thanks to an autoregressive model with lag 1 usually suited for time series.

**List of Tables**

- 1 Yeast cell-cycle data: some representative GO terms analysis  
of clusters obtained by tested models and P-values of  
over-representation. The lowest the P-values, the more isolated  
the GO terms. For each GO term (row), the best method is  
indicated by underlined bold P-values. 31

For Peer Review

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



Table 1

Yeast cell-cycle data: some representative GO terms analysis of clusters obtained by tested models and P-values of over-representation. The lowest the P-values, the more isolated the GO terms. For each GO term (row), the best method is indicated by underlined bold P-values.

cluster number	p-values of SFmiss clust.	best p-values among EMmiss clust.	best p-values among KNN+SF clust.
k=3	GO:0006732, coenzyme met. process		
	<u><b><math>1.1 \cdot 10^{-2}</math></b></u>	$> 0.1$	$> 0.1$
k=4	GO:0005819, spindle		
	<u><b><math>4.6 \cdot 10^{-9}</math></b></u>	$6.7 \cdot 10^{-7}$	$2.0 \cdot 10^{-6}$
	GO:0006790, sulf. met. process		
	<u><b><math>1.1 \cdot 10^{-4}</math></b></u>	$2.4 \cdot 10^{-4}$	$8.7 \cdot 10^{-4}$
	GO:0000278, mitotic cell cycle		
	<u><b><math>2.2 \cdot 10^{-3}</math></b></u>	$7.7 \cdot 10^{-3}$	$> 0.1$
	GO:0030472, mit. spin. org. & biogen. in nucleus		
	<u><b><math>5.2 \cdot 10^{-3}</math></b></u>	$8.8 \cdot 10^{-3}$	$2.0 \cdot 10^{-2}$
k=5	GO:0006974, resp. to DNA dam. stim.		
	<u><b><math>1.8 \cdot 10^{-3}</math></b></u>	$3.0 \cdot 10^{-3}$	$8.0 \cdot 10^{-3}$
	GO:0000724, dbl-str. bk rep. via hom. comb.		
	<u><b><math>1.9 \cdot 10^{-2}</math></b></u>	$2.7 \cdot 10^{-2}$	$4.6 \cdot 10^{-2}$
	GO:0000030, mannosyltransf. act.		
	<u><b><math>1.1 \cdot 10^{-2}</math></b></u>	$1.2 \cdot 10^{-2}$	$2.7 \cdot 10^{-2}$
k=8	GO:0042555, MCM cplx		
	<u><b><math>3.4 \cdot 10^{-4}</math></b></u>	$8.3 \cdot 10^{-4}$	$4.0 \cdot 10^{-4}$
	GO:0008026, ATP-dep. helicase act.		
	$5.5 \cdot 10^{-4}$	$1.3 \cdot 10^{-3}$	<u><b><math>4.5 \cdot 10^{-4}</math></b></u>
	GO:0006268, DNA unwind. replic.		
	$2.8 \cdot 10^{-3}$	$6.7 \cdot 10^{-3}$	<u><b><math>1.1 \cdot 10^{-3}</math></b></u>
	GO:0042623, ATPase act. coupl.		
	<u><b><math>4.4 \cdot 10^{-3}</math></b></u>	$1.5 \cdot 10^{-2}$	$4.3 \cdot 10^{-2}$