



**HAL**  
open science

## A macro-DAG structure based mixture model

Bernard Chalmond

► **To cite this version:**

| Bernard Chalmond. A macro-DAG structure based mixture model. 2013. hal-00947454

**HAL Id: hal-00947454**

**<https://hal.science/hal-00947454>**

Preprint submitted on 16 Feb 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A macro-DAG structure based mixture model

BERNARD CHALMOND

University of Cergy-Pontoise, France and CMLA, Ecole Normale Supérieure de Cachan, France \*

**Abstract-** In the context of unsupervised classification of multidimensional data, we revisit the classical mixture model in the case where the dependencies among the random variables are described by a DAG structure. The structure is considered at two levels, the original DAG and its macro-representation. This two-level representation is the main base of the proposed mixture model. To perform unsupervised classification, we propose a dedicated algorithm called EM-mDAG, which extends the classical EM algorithm. In the Gaussian case, we show that this algorithm can be efficiently implemented. The experiments reveal that this method favors the selection of a small number of classes.

*Keywords:* Mixture model, DAG structure, Bayesian network, EM algorithm

## 1. Introduction

Let  $\mathbf{X}$  be a random vector with values in  $\mathbf{R}^n$  for which we have a  $N$ -sample  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  with  $n < N$ . Our goal is the clustering of  $\mathcal{X}$ . This task is approached through a mixture model but with a particular constraint that makes the specificity of our contribution.

The dependency structure among the  $n$  components  $X^j$  of  $\mathbf{X}$  is subject to a structure represented by a DAG, in other words  $\mathbf{X}$  is a Bayesian network. This structure induces a partition of  $\mathbf{X}$  into  $M + 1$  random vectors called *macro-variables* :  $\mathbf{X} = \uplus_{m=0}^M X^{J_m}$ , where  $X^{J_m} = (X^{j_1}, \dots, X^{j_m})$  when  $J_m = \{j_1, \dots, j_m\}$ . Fig.1 depicts an example with  $M = 3$  and  $J_0 = \{1\}$ ,  $J_1 = \{2, 3\}$ ,  $J_2 = \{4, 5\}$ ,  $J_3 = \{6, 7, 8\}$ .

Each macro-variable  $X^{J_m}$  is dependent on a hidden class variable  $C^m$  with values in  $\mathcal{K}_m = \{1, 2, \dots, \nu_m\}$ . Each occurrence in  $\mathcal{K}_m$  is the number of a class called *elementary class*. Therefore  $\mathbf{X}$  is dependent on the hidden multi-class variable  $\mathbf{C} = (C^1, \dots, C^M)$  whose values are in  $\mathcal{K} = \otimes_{m=0}^M \mathcal{K}_m$ . Each  $(M + 1)$ -tuple of  $\mathcal{K}$  refers to a set of elementary classes called *composite class*. The  $(M + 1)$ -tuples can be interpreted as pathways connecting the elementary classes through the macro-variables as it is illustrated in Table 1. The objective is to find the most probable pathways. We consider the mixture model

$$p_{\bar{\theta}}(\mathbf{x}) = \sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x} | \mathbf{k}),$$

where the probability distribution  $p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x} | \mathbf{k})$  is that of the Bayesian network conditional on the composite class  $\mathbf{k}$  and  $\bar{\theta}_{\mathbf{k}}$  denotes the set of parameters defining this distribution.

---

\* E-mail : bernard.chalmond@cmla.ens-cachan.fr

**Table 1.** Composite class numbering for  $M = 3$  and  $\nu_0 = 1, \nu_1 = \nu_2 = 2, \nu_3 = 4$ . This table gives the exhaustive list of the 16 composite classes  $\mathcal{K}$ , where each column is an  $(M + 1)$ -tuple  $(1, \mathbf{k})$  with  $\mathbf{k} \in \mathcal{K}$ .

m=0 :	1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
m=1 :	1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2
m=2 :	1, 1, 1, 1, 2, 2, 2, 2, 1, 1, 1, 1, 2, 2, 2
m=3 :	1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4

In this paper we describe this mixture model and we give a version of the EM algorithm, called EM-mDAG, for performing unsupervised classification. One of the main role of the EM-mDAG algorithm is to reveal probabilistic relationships among hidden elementary classes. Its implementation is done in the Gaussian case. Simulations illustrate the method and reveal a specific property. The EM-mDAG algorithm can select a small number of significant composite classes in  $\mathcal{K}$ .

## 2. Models and Method

### 2.1. Basic knowledge

- *Conventional mixture model for non supervised classification.*

Let a random vector  $\mathbf{X} = (X^1, \dots, X^j, \dots, X^n)$  with values in  $\mathbf{R}^n$ . We assume that its probability distribution  $p_\phi(\mathbf{x})$  is a mixture of  $\nu$  distributions  $\{p_{\theta_k}(\mathbf{x})\}$  as follows :

$$p_\phi(\mathbf{x}) = \sum_{k=1}^{\nu} \alpha_k p_{\theta_k}(\mathbf{x}) \quad \text{with} \quad \sum_{k=1}^{\nu} \alpha_k = 1. \quad (2.1)$$

$p_{\theta_k}(\mathbf{x})$  is defined by a parametric law of parameters  $\theta_k$ , as for instance the Gaussian law. The parameter set is denoted  $\phi = \{\alpha, \theta\}$  where  $\alpha = \{\alpha_k\}$  and  $\theta = \{\theta_k\}$ . This mixture model can be interpreted in the context of unsupervised data classification. Let  $C$  be the hidden variable, which is an indicator variable of classes with values in  $\{1, \dots, \nu\}$ . Then, (2.1) is rewritten as

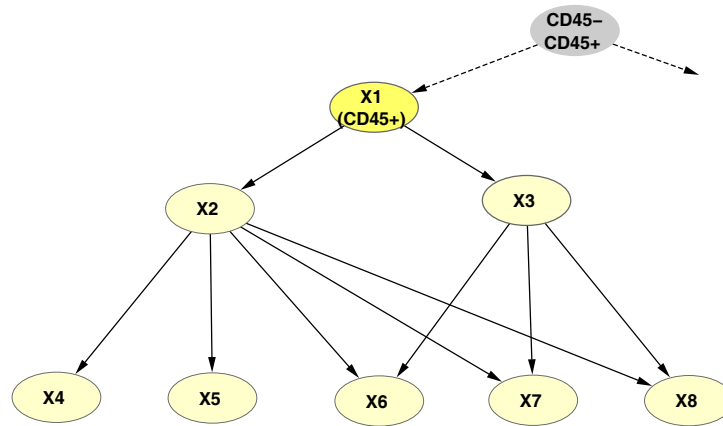
$$p_\phi(\mathbf{x}) = \sum_{k=1}^{\nu} P(C = k) p_{\theta_k}(\mathbf{x} | C = k). \quad (2.2)$$

The classification is to assign a class to every observation  $\mathbf{x}$ ,<sup>1</sup>. When  $\phi$  is given, the MAP decision rule consists to choose the class

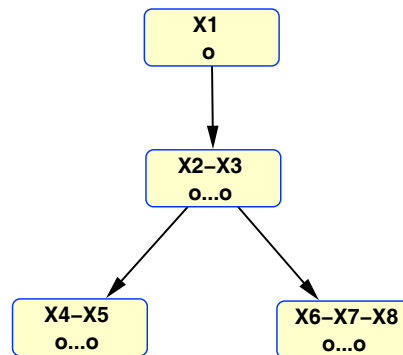
$$\hat{k}(\mathbf{x}) = \arg \max_k P_\phi(C = k | \mathbf{x}). \quad (2.3)$$

Otherwise, things are more complicated because  $k(\mathbf{x})$  and  $\phi$  have to be simultaneously estimated. On the basis of maximum likelihood, the EM algorithm allows this estimation from a sample  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  of  $\mathbf{X}$ .

<sup>1</sup>A class is defined by its number and its parameters. More often, we confound "class" and "class number".



(a)



(b)

**Figure 1.** Two-level structure. (a) DAG structure. (b) Macro-DAG structure with its macro-variables  $X^{J_1} = (X^2, X^3)$ ,  $X^{J_2} = (X^4, X^5)$  and  $X^{J_3} = (X^6, X^7, X^8)$ ; the small circles depict the elementary classes.

The general formulation of the EM algorithm, which is also valid for our particular case, reads as follows. If  $\phi(\ell)$  is an estimation of  $\phi$ , then an updated estimation is :

$$\begin{aligned}\phi(\ell + 1) &= \arg \max_{\phi} Q(\phi | \phi(\ell)) , \\ Q(\phi | \phi(\ell)) &= \mathbb{E}_{\mathbf{C} | \mathcal{X}} [\log p_{\phi}(\mathcal{X}, \mathbf{C}) | \phi(\ell)] ,\end{aligned}\tag{2.4}$$

where  $\mathbf{C} = \{C_1, \dots, C_N\}$  is a series of i.i.d. variables related to  $C$ .  $Q$  is an expected log-likelihood with respect to  $p_{\phi(\ell)}(\mathbf{C} | \mathcal{X})$ . The EM algorithm is an iterative procedure. From an initial estimate  $\phi(0)$ , it computes successively  $\phi(0) \rightarrow \dots \rightarrow \phi(\ell) \rightarrow \dots$ . The marginal likelihood series  $\{p_{\phi(\ell)}(\mathcal{X}), \ell = 0, 1, \dots\}$  is non-decreasing.

- *Bayesian network.*

The previous classical formalism is the primal version for mixture modeling in the context of classification [4]. The EM algorithm also applies to more complex situations such as those where the  $\mathbf{X}_i$  are not i.i.d. variables, but are dependent through hidden variables  $C_i$  governed by a Markov chain [2] or a Markov random field [3]. In this article, we remain in the case where  $\mathcal{X}$  is a sample from i.i.d. variables, but we consider a Markov structure for the dependence of the components  $X^j$ . This Markovian structure is based on a DAG denoted  $G = (V, E)$ .  $V = \{1, \dots, j, \dots, n\}$  denotes the variable numbers. The edges  $E \subset V \times V$  are directed :  $(j', j) \in E$  is denoted  $j' \rightarrow j$ . The set  $\bar{j} = \{j' : j' \rightarrow j\}$  denotes the parents of the node  $j$ . The DAG structure has a fundamental property due to its acyclic nature : there is a numbering of the nodes such that  $\bar{j} \subset \{1, 2, \dots, j - 1\}$ . We assume that the nodes have been ordered in this way. With this property and that of Markov, we get the factorization

$$p(\mathbf{x}) = \prod_j p(x^j | x^{\bar{j}}).\tag{2.5}$$

The set  $B = (\mathbf{X}, G, \{p(x^j | x^{\bar{j}})\})$  is called Bayesian network. When the distribution  $p(\mathbf{x})$  is non homogeneous, a mixture model as (2.2) can be considered in which  $p_{\theta_k}(\mathbf{x} | C = k)$  denotes a Bayesian network conditional on the hidden class  $C$ . This mixture model has been investigated in [6] with a particular interest for DAG structure estimation.

## 2.2. Mixture model, composite class and Bayesian network

### 2.2.1. Composite class model

Let a partition of  $V$  composed of  $M + 1$  macro-nodes :  $V = J_0 \uplus J_1 \uplus \dots \uplus J_M$ , built from the DAG structure :  $J_m$  is a macro-node if all its nodes have the same parents (In Fig.1,  $M = 3$ , and  $J_0 = \{1\}$ ,  $J_1 = \{2, 3\}$ ,  $J_2 = \{4, 5\}$ ,  $J_3 = \{6, 7, 8\}$ ).  $J_0$  is the root of the tree and most often is a single node <sup>2</sup>. Let  $J_{\bar{1}}, \dots, J_{\bar{M}}$  be the parents of  $J_1, \dots, J_M$ , respectively. Given the definition of macro-nodes, each  $J_{\bar{m}}$  is composed of a single macro-node (In Fig.1,  $J_{\bar{1}} = J_0$ ,  $J_{\bar{2}} = J_1$ ,  $J_{\bar{3}} = J_1$ ). The macro-nodes  $\mathcal{V} = \{J_m\}$  and their connexions  $\mathcal{E}$  induced by  $\{J_{\bar{m}}\}$  define a new directed acyclic graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  called macro-DAG.

<sup>2</sup>  $J_0$  has only one class and therefore  $\nu_0 = 1$ .

Given a set of specifications  $\{p(x^{J_m} | x^{\bar{J}_m})\}$  for  $B$ , a Bayesian network  $\mathcal{B} = (\mathbf{X}, \mathcal{G}, \{p(x^{J_m} | x^{\bar{J}_m})\}_{m=0}^M)$  can be defined for the macro-variables  $\{X_m^J\}_{m=0}^M$ . The difference with  $B$  is essentially that  $\mathcal{B}$  is a vectorial process whose factorization formula is written as

$$p(\mathbf{x}) = \prod_m p(x^{J_m} | x^{\bar{J}_m}). \quad (2.6)$$

The factorization (2.6) assumes that the probability distribution is homogeneous, whereas it is not the case in our context. The distribution is depending on a hidden class variable  $\mathbf{C}$ , which implies that  $p(\mathbf{x})$  is a mixture of distributions.

Firstly we assume that each macro-variable  $X_m^J$  is characterized by  $\nu_m$  classes, called *elementary classes*, whose parameters are denoted  $\theta^m = \{\theta_1^m, \dots, \theta_{\nu_m}^m\}$ . If we forget for a while the DAG structure, then each variable taken independently of the others, is defined by a mixture model for which (2.1) is rewritten as

$$p(x^{J_m}) = \sum_{k=1}^{\nu_m} P(C^m = k) p_{\theta_k^m}(x^{J_m} | C^m = k). \quad (2.7)$$

Secondly, we consider the indicator variable of *composite classes*  $\mathbf{C} = (C^1, \dots, C^M)$  with values in the set of  $M$ -tuples  $\mathcal{K} = \{\mathbf{k} = (k_1, \dots, k_M)\}$  where  $k_m \in \{1, \dots, \nu_m\}$ , as represented in Table 1. The classification is to assign a composite class to each observation  $\mathbf{x}$ . This involves selecting an elementary class  $k_m$  for each macro-variable. An immediate solution would be to perform  $M$  independent classifications, based on (2.7) but this approach would have the disadvantage of not considering the DAG structure. Therefore we must address the classification as a whole.

Considering the DAG structure, a composite class  $\mathbf{k}$  is not only defined by the parameters  $\theta_{\mathbf{k}} = \{\theta_{k_m}^m\}_{m=1}^M$  of its elementary classes, but also by the dependency parameters  $\bar{\theta}_{\mathbf{k}} = \{\bar{\theta}_{k_m}^m\}_{m=1}^M$  that define the specifications of the Bayesian network  $\mathbf{X}$  conditionally to  $\mathbf{C} = \mathbf{k}$ <sup>3</sup>. These parameters are related to the parameters  $\theta_{\mathbf{k}}$ . For each composite class, the factorization formula (2.6) based on the macro-DAG is written as

$$p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x} | \mathbf{k}) = p(x^{J_0}) \prod_{m=1}^M p_{\bar{\theta}_{\mathbf{k}}}^m(x^{J_m} | x^{\bar{J}_m}, k_m, \bar{k}_m), \quad (2.8)$$

where  $\bar{k}_m$  denotes the class number associated to  $x^{\bar{J}_m}$  and appearing in  $\mathbf{k}$ . In the notation  $p_{\bar{\theta}_{\mathbf{k}}}^m$ , only the classes  $k_m$  and  $\bar{k}_m$  of  $\mathbf{k}$  are active. Finally the mixture model is written as

$$p_{\bar{\theta}}(\mathbf{x}) = \sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x} | \mathbf{k}). \quad (2.9)$$

Initially in (2.7) the definition of elementary classes has been made independently for each macro-variables. Now, the Markov dependence (2.8) introduces dependencies among these classes. The parameter setting of the mixture model (2.9) differs from the classical mixture model (2.1).

<sup>3</sup>In this paper, the notation  $\bar{\cdot}$  is reserved to parameters associated to the DAG dependencies.

Two  $M$ -tuples may have common components. For example, all components of  $(1, k_2, \dots, k_M)$  and  $(2, k_2, \dots, k_M)$  are identical, except the first. Thus, since two  $M$ -tuples may have common components, two components of the mixture may have common parameters<sup>4</sup>. In fact, there is one parameter setting per class, totaling  $\sum_m \nu_m$  settings, while there are  $|\mathcal{K}| = \prod_m \nu_m$  composite classes.

### 2.2.2. EM-mDAG algorithm

The ultimate objective is to assign a composite class to every observation  $\mathbf{x}$  :

$$\mathbf{x} \rightarrow \widehat{\mathbf{k}}(\mathbf{x}) = \arg \max_{\mathbf{k} \in \mathcal{K}} P_\phi(\mathbf{C} = \mathbf{k} \mid \mathbf{x}) .$$

Therefore, it is necessary to estimate  $\phi = (\alpha, \bar{\theta})$ . In an equivalent manner to (2.4), the estimation of  $\phi$  is based on the log-likelihood by maximizing the Lagrangian function

$$\begin{aligned} \mathcal{L}(\alpha, \bar{\theta}) &= \sum_{i=1}^N \log \left[ \sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x}_i \mid \mathbf{k}) \right] + \lambda \left[ \sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} - 1 \right] , \\ &= \sum_{i=1}^N \log \left[ \sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} \prod_{m=1}^M p_{\bar{\theta}_{\mathbf{k}}^m}(x_i^{J_m} \mid x_i^{J_{\overline{m}}}, k_m, k_{\overline{m}}) \right] + \lambda \left[ \sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} - 1 \right] , \end{aligned} \quad (2.10)$$

where  $\lambda$  denotes the Lagrangian parameter associated to the constraint  $\sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}} = 1$ . At the iteration  $\ell$  of the EM algorithm, the re-estimation formula of  $\alpha$  is written as in the classical case :

$$\alpha_{\mathbf{k}}(\ell + 1) = \frac{1}{N} \sum_i p_{\phi(\ell)}(\mathbf{k} \mid \mathbf{x}_i) , \quad (2.11)$$

where the a posteriori probability of the composite class  $\mathbf{k}$  is defined by

$$p_{\phi(\ell)}(\mathbf{k} \mid \mathbf{x}_i) = \frac{\alpha_{\mathbf{k}}(\ell) p_{\bar{\theta}_{\mathbf{k}}(\ell)}(\mathbf{x}_i \mid \mathbf{k})}{p_{\phi(\ell)}(\mathbf{x}_i)} = \frac{\alpha_{\mathbf{k}}(\ell) p_{\bar{\theta}_{\mathbf{k}}(\ell)}(\mathbf{x}_i \mid \mathbf{k})}{\sum_{\mathbf{k} \in \mathcal{K}} \alpha_{\mathbf{k}}(\ell) p_{\bar{\theta}_{\mathbf{k}}(\ell)}(\mathbf{x}_i \mid \mathbf{k})} . \quad (2.12)$$

As we said above, the peculiarity of this variant of the EM algorithm is the fact that a same parameter  $\bar{\theta}_{k_m}^m$  can be present in several composite classes. In the classical case (2.1), the gradient of the Lagrangian function with respect to  $\theta_k$  concerns only  $p_{\theta_k}$  while in (2.10), the gradient with respect to  $\bar{\theta}_{k_m}^m$  relates to several components  $p_{\bar{\theta}_{\mathbf{k}}}$ .

**Proposition 1.** *The re-estimation formula of  $\bar{\theta}$  is given by the solution  $\bar{\theta}(\ell + 1)$  of the linear system*

$$\sum_{i=1}^N \sum_{\substack{\mathbf{k}=(k_1, \dots, k_M): \\ k_m=\tau_m}} p_{\phi(\ell)}(\mathbf{k} \mid \mathbf{x}_i) \frac{\partial}{\partial \bar{\theta}_{\tau_m}^m} \log p_{\bar{\theta}_{\tau_m}^m}(x_i^{J_m} \mid x_i^{J_{\overline{m}}}, \tau_m, k_{\overline{m}}) \Big|_{\bar{\theta}=\bar{\theta}(\ell+1)} = 0 , \quad (2.13)$$

$$\tau_m = 1, \dots, \nu_m , \quad m = 1, \dots, M .$$

<sup>4</sup>In the classical case, several components can also be concerned by a same parameter, for instance the same variance in the Gaussian case, but it is not a constraint contrary to our mixture model where many parameters are necessarily shared.

**Proof.** Let's focus on  $\bar{\theta}_{\tau_m}^m$  where  $\tau_m$  is a particular class number in  $\{1, \dots, \nu_m\}$ .

$$\begin{aligned}
\frac{\partial \mathcal{L}(\alpha, \bar{\theta})}{\partial \theta_{\tau_m}^m} &= \sum_{i=1}^N \frac{1}{p_\phi(\mathbf{x}_i)} \left[ \sum_{\substack{\mathbf{k}=(k_1, \dots, k_M): \\ k_m=\tau_m}} \alpha_{\mathbf{k}} \frac{\partial}{\partial \theta_{\tau_m}^m} p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x}_i | \mathbf{k}) \right], \\
&= \sum_{i=1}^N \frac{1}{p_\phi(\mathbf{x}_i)} \left[ \sum_{\substack{\mathbf{k}=(k_1, \dots, k_M): \\ k_m=\tau_m}} \alpha_{\mathbf{k}} \frac{p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x}_i | \mathbf{k})}{p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x}_i | \mathbf{k})} \frac{\partial}{\partial \theta_{\tau_m}^m} p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x}_i | \mathbf{k}) \right], \\
&= \sum_{i=1}^N \sum_{\mathbf{k}: k_m=\tau_m} p_\phi(\mathbf{k} | \mathbf{x}_i) \frac{\partial}{\partial \theta_{\tau_m}^m} \log p_{\bar{\theta}_{\mathbf{k}}}(\mathbf{x}_i | \mathbf{k}).
\end{aligned} \tag{2.14}$$

Recalling the factorization formula (2.5), the gradient can be written as

$$\begin{aligned}
\frac{\partial \mathcal{L}(\alpha, \bar{\theta})}{\partial \theta_{\tau_m}^m} &= \sum_{i=1}^N \sum_{\mathbf{k}: k_m=\tau_m} p_\phi(\mathbf{k} | \mathbf{x}_i) \frac{\partial}{\partial \theta_{\tau_m}^m} \left[ \sum_{a=1}^M \log p_{\bar{\theta}_{\mathbf{k}}}^a(x_i^{J_a} | x_i^{J_{\bar{a}}}, k_a, k_{\bar{a}}) \right], \\
&= \sum_{i=1}^N \sum_{\mathbf{k}: k_m=\tau_m} p_\phi(\mathbf{k} | \mathbf{x}_i) \frac{\partial}{\partial \theta_{\tau_m}^m} \left[ \log p_{\bar{\theta}_{\tau_m}^m}(x_i^{J_m} | x_i^{J_{\bar{m}}}, \tau_m, k_{\bar{m}}) \right],
\end{aligned} \tag{2.15}$$

which leads after a shortcut, to the system (2.13).  $\square$

### 2.2.3. Gaussian case, linear dependency model and DAG

#### • Linear dependency model and DAG.

Under the Gaussian assumption, conditionally on the elementary classes, the law of the macro-variables are

$$X^{J_m | k_m} \doteq [X^{J_m} | k_m] \sim \mathcal{N}(\mu_{k_m}^m, \Gamma_{k_m}^m), \tag{2.16}$$

and with respect to the DAG, the transition laws among these variables are

$$[X^{J_m | k_m} | x^{J_{\bar{m}}}, k_{\bar{m}}] = [X^{J_m} | x^{J_{\bar{m}}}, k_m, k_{\bar{m}}] \sim \mathcal{N}(\mu_{k_m | k_{\bar{m}}}^{m, \mathbf{x}}, \Gamma_{k_m | k_{\bar{m}}}^{m, \mathbf{x}}). \tag{2.17}$$

We assume the linear regression model

$$\begin{aligned}
\mu_{k_m | k_{\bar{m}}}^{m, \mathbf{x}} &= A_{k_m, k_{\bar{m}}}^m x^{J_{\bar{m}}} + b_{k_m, k_{\bar{m}}}^m, \\
\Gamma_{k_m | k_{\bar{m}}}^{m, \mathbf{x}} &= \Gamma_{k_m | k_{\bar{m}}}^m.
\end{aligned} \tag{2.18}$$

Therefore, the respective parameter settings of (2.16) and (2.17) are respectively

$$\begin{aligned}
\theta_{k_m}^m &= \{\mu_{k_m}^m, \Gamma_{k_m}^m\}, \\
\bar{\theta}_{k_m}^m &= \{A_{k_m, k_{\bar{m}}}^m, b_{k_m, k_{\bar{m}}}^m, \Gamma_{k_m | k_{\bar{m}}}^m\}.
\end{aligned}$$



Note that the linear regression model (2.18) depends on the direction of the DAG.  $A_{k_m, k_{\overline{m}}}^m$  characterizes the regression of  $X^{J_m}$  on  $x^{J_{\overline{m}}}$  and not the reverse. Note also that the regression model is multidimensional in output since  $X^{J_m|k_m}$  is a random vector in  $\mathbf{R}^{|J_m|}$ .

- *Re-estimation formulas.*

The first approach to obtain these formulas would be to use (2.13) taking into account the Gaussian log-density that for any observation  $\mathbf{x}_i$  can be written as :

$$\begin{aligned} \log p_{\bar{\theta}_{k_m}^m}(x_i^{J_m} | x_i^{J_{\overline{m}}}, k_m, k_{\overline{m}}) &= cst - \frac{1}{2} |\log \Gamma_{k_m|k_{\overline{m}}}^m| \\ &\quad - \frac{1}{2} (x_i^{J_m} - \mu_{k_m|k_{\overline{m}}}^{m, \mathbf{x}_i})' (\Gamma_{k_m|k_{\overline{m}}}^m)^{-1} (x_i^{J_m} - \mu_{k_m|k_{\overline{m}}}^{m, \mathbf{x}_i}). \end{aligned}$$

Establish equation (2.13) requires to derivate with respect to all the components of  $\bar{\theta}_{k_m}^m$ . For instance, the partial derivative with respect to  $A_{\tau_m, k_{\overline{m}}}^m$  is written as

$$\frac{\partial}{\partial A_{\tau_m, k_{\overline{m}}}^m} \log p_{\bar{\theta}_{\tau_m}^m}(x_i^{J_m} | x_i^{J_{\overline{m}}}, \tau_m, k_{\overline{m}}) = (\Gamma_{\tau_m, k_{\overline{m}}}^m)^{-1} x_i^{J_{\overline{m}}} (x_i^{J_m} - \mu_{\tau_m, k_{\overline{m}}}^{m, \mathbf{x}_i})'.$$

We could go on, but there is a more direct way to operate.

**Proposition 2.** *Since (2.18) is a linear regression model, it is easier to use the classical formulas of the maximum likelihood estimation of this model. In our context, these formulas [5] are written as*

$$\begin{aligned} b_{\tau_m, k_{\overline{m}}}^m(\ell + 1) &= \widehat{\mathbb{E}}(X^{J_m|\tau_m}) - A_{\tau_m, k_{\overline{m}}}^m(\ell + 1) \widehat{\mathbb{E}}(X^{J_{\overline{m}}|k_{\overline{m}}}), \\ A_{\tau_m, k_{\overline{m}}}^m(\ell + 1) &= \widehat{\text{Cov}}(X^{J_m|\tau_m}, X^{J_{\overline{m}}|k_{\overline{m}}}) \left[ \widehat{\text{Var}}(X^{J_{\overline{m}}|k_{\overline{m}}}) \right]^{-1}, \end{aligned}$$

where  $\widehat{\cdot}$  denotes an empirical estimate.

However, the empirical estimate of the moments (expectations, covariance matrix and variance-covariance matrix) must be weighted by weights  $w$  derived from the DAG :

$$\begin{aligned} b_{\tau_m, k_{\overline{m}}}^m(\ell + 1) &= \sum_{i=1}^N \sum_{\mathbf{k}: k_m = \tau_m} w_{i, \mathbf{k}}^{\tau_m}(\ell) x_i^{J_m} \\ &\quad - A_{\tau_m, k_{\overline{m}}}^m(\ell + 1) \sum_{i=1}^N \sum_{\mathbf{k}: k_m = \tau_m} w_{i, \mathbf{k}}^{\tau_m}(\ell) x_i^{J_{\overline{m}}}, \end{aligned} \tag{2.19}$$

where the weights at iteration  $\ell$  are

$$w_{i, \mathbf{k}}^{\tau_m}(\ell) = \frac{p_{\phi}(\ell)(\mathbf{k} | \mathbf{x}_i)}{\sum_{i=1}^N \sum_{\mathbf{k}: k_m = \tau_m} p_{\phi}(\ell)(\mathbf{k} | \mathbf{x}_i)}.$$

Similarly, by denoting  $\hat{\mu}^{J_m|\tau_m} = \widehat{E}(X^{J_m|\tau_m})$ , we have

$$A_{\tau_m, k_m}^m(\ell + 1) = \sum_{i=1}^N \sum_{\mathbf{k}: k_m = \tau_m} w_{i, \mathbf{k}}^{\tau_m}(\ell) (x_i^{J_m} - \hat{\mu}^{J_m|\tau_m})(x_i^{J_m} - \hat{\mu}^{J_m|\tau_m})' \\ \times \left[ \sum_{i=1}^N \sum_{\mathbf{k}: k_m = \tau_m} w_{i, \mathbf{k}}^{\tau_m}(\ell) (x_i^{J_m} - \hat{\mu}^{J_m|\tau_m})(x_i^{J_m} - \hat{\mu}^{J_m|\tau_m})' \right]^{-1}. \quad (2.20)$$

Finally, we have also

$$\Gamma_{\tau_m, k_m}^m(\ell + 1) = \sum_{i=1}^N \sum_{\mathbf{k}: k_m = \tau_m} w_{i, \mathbf{k}}^{\tau_m}(\ell) \left( x_i^{J_m} - \mu_{\tau_m, k_m}^m(\ell + 1) \right) \\ \left( x_i^{J_m} - \mu_{\tau_m, k_m}^m(\ell + 1) \right)'. \quad (2.21)$$

Note that the programming of the re-estimation formulas (2.19, 2.20, 2.21) is relatively difficult because two data structures interfere : the dependency structure derived from the DAG, and the list structure of the composite classes (cf. (1) for instance).

- *Back to the elementary class parameters.*

For all  $\mathbf{x}_i$ , the estimated composite class  $\widehat{\mathbf{k}}(\mathbf{x}_i) = (\widehat{k}_1(\mathbf{x}_i), \dots, \widehat{k}_M(\mathbf{x}_i))$  has been computed. The user is also interested by the parameters  $\theta_{\widehat{\mathbf{k}}}$  of the elementary classes  $\{\widehat{k}_m(\mathbf{x}_i)\}_{m=1}^M$  in order to interpret the leaves of the tree. The law (2.16) ignores the DAG, contrary to the law (2.17). However, in the Gaussian case, the parameters  $\theta_{\widehat{\mathbf{k}}}$  are related to the parameters  $\theta_{\mathbf{k}}$  as follows [5] :

$$\mu_{k_m | k_m}^{m, \mathbf{x}} = \mu_{k_m}^m + \Gamma_{k_m, k_m}^m (\Gamma_{k_m}^m)^{-1} (x^{J_m} - \mu_{k_m}^m), \\ \Gamma_{k_m | k_m}^m = \Gamma_{k_m}^m + \Gamma_{k_m, k_m}^m (\Gamma_{k_m}^m)^{-1} \Gamma_{k_m, k_m}^m. \quad (2.22)$$

To avoid the difficulty of solving the system with respect to  $\mu_{k_m}^m$  and  $\Gamma_{k_m}^m$ , we consider

$$\hat{\mu}_{k_m}^m = \frac{1}{N} \sum_i x_i^{J_m} \mathbf{1}_{\widehat{k}_m(\mathbf{x}_i) = k_m}, \\ \hat{\Gamma}_{k_m}^m = \frac{1}{N} \sum_i (x_i^{J_m} - \hat{\mu}_{k_m}^m)(x_i^{J_m} - \hat{\mu}_{k_m}^m)' \mathbf{1}_{\widehat{k}_m(\mathbf{x}_i) = k_m}.$$

- *Initial solution.*

The solution at the first step of the EM-mDAG algorithm is obtained by performing  $M$  independent classifications using the conventional EM algorithm. Therefore, for each macro-variable, we have  $\nu_m$  clusters in  $\mathbf{R}^{|J_m|}$  whose labels are  $\{\widehat{k}_m(x_i^{J_m}), i = 1, \dots, n\}$ . From there, the initial solution  $\hat{\theta}_{k_m, k_m}^m(0)$  at iteration  $\ell = 0$  is computed using ordinary linear regression for every pair of clusters  $(k_m, k_m)$  for which there are observations :  $\{i : \widehat{k}_m(x_i^{J_m}) = k_m, \widehat{k}_m(x_i^{J_m}) = k_m\} \neq \emptyset$ . Starting from this initial solution, the role of the EM-mDAG algorithm is to re-organize the clusters in order to extract from  $\mathcal{K}$  a set of composite classes of high likelihood.

**Table 2.** Expectations  $\mu_{\mathbf{k}}$  of the 5 composite classes  $\mathcal{K}_0$  for data simulation with  $M = 3$ ,  $\nu_0 = 1$ ,  $\nu_1 = 2$ ,  $\nu_2 = 3$ ,  $\nu_3 = 4$ . Here  $c = 1.5$ .

$\mu_{\mathbf{k}}$ :	
$X_1$ :	0, 0, 0, 0, 0
$X_2$ :	-c, -c, +c, +c, +c
$X_3$ :	+c, +c, -c, -c, -c
$X_4$ :	-c, -c, +c, +c, +c
$X_5$ :	-c, -c, -c, +c, +c
$X_6$ :	-c, +c, -c, +c, +c
$X_7$ :	-c, +c, +c, -c, -c
$X_8$ :	-c, +c, -c, +c, +c

### 3. Experiments on simulated data

The random vector  $\mathbf{X}$  of dimension  $n = 8$  consists of  $M = 3$  macro-variables with  $\nu_1 = 2$ ,  $\nu_2 = 3$ ,  $\nu_3 = 4$ . Among the  $|\mathcal{K}| = 24$  potential composite classes, only 5 significant composite classes  $\mathcal{K}_0$  were considered. It implies that for  $\mathbf{k} \notin \mathcal{K}_0$  no data has been generated, or in another words

$$\alpha_{\mathbf{k}} = 0, \forall \mathbf{k} \notin \mathcal{K}_0. \quad (3.1)$$

Table 2 gives the expectation of these 5 composite classes and therefore the expectation of the elementary classes within the 3 macro-variables.

Fig.4 shows the observations of the macro-variables with their labels. This simulation was inspired by the cytometry data analysis domain (see Appendix) but with much more overlapping of the elementary classes. The first macro-variable  $X^{J_1} = (X^2, X^3)$  shows two groups that it is possible to manually split, giving rise to two elementary classes denoted  $X^{2+}$  and  $X^{3+}$ . Each group is a mixture that the other two macro-variables help to identify. The macro-variable  $X^{J_2} = (X^4, X^5)$  highlights the components of the group  $X^{2+}$ , while the macro-variable  $X^{J_3} = (X^6, X^7, X^8)$  highlights the components of the group  $X^{3+}$ . However the overlapping of the mixture components in the groups  $X^{2+}$  and  $X^{3+}$  does not allow a partitioning of these groups as easy as for  $X^{J_1}$ . Therefore we must address the classification as a whole.

- *Data processing.*

Fig.5 shows the initial solution of the EM-mDAG algorithm at step  $\ell = 0$ . This solution results from  $M = 3$  independent classifications by applying the classical EM algorithm on each macro-variable. The class labels  $\hat{k}(x_i^{J_m})$  defined in (2.3) are gathered for making composite class labels by using a table similar to Table 1. This initial solution is unsatisfactory. The macro-variables  $X^{J_2}$  and  $X^{J_3}$  are strongly blurred by several small composite classes that are artefacts. The final solution of the EM-mDAG algorithm is shown in Fig.6. The representation in terms of mixture components is close to the original in Fig.4. The macro-variables  $X^{J_2}$  and  $X^{J_3}$  respectively highlight the components of the group  $X^{2+}$  and  $X^{3+}$ , despite the fifth class that is divided into two neighbour classes.

Fig.7 and Fig.8 show respectively the classifications obtained with the usual EM algorithm successively performed on the basis of 24 classes and 5 classes. With 24 classes, the number of non-empty classes is large and therefore the classification is greatly erroneous. With 5 classes, the classification provided by the EM algorithm does not meet the specificity of the data. The macro-variable  $X^{J_2} = (X^4, X^5)$  does not highlight the mixture components in the group  $X^{2+}$ . This is a major problem that hinders significantly the interpretation of composite classes in terms of macro-variables.

- *Property of the EM-mDAG algorithm.*

The experiments show that the EM-mDAG algorithm has the property to keep a small number of  $\alpha_k$  different from zero when there is a limited number of significant composite classes  $\mathcal{K}_0 \subset \mathcal{K}$ . This selection ability is not so surprising. Firstly,  $\mathbf{X}$  is not observable along  $\mathbf{k}$  when  $\mathbf{k} \notin \mathcal{K}_0$ , which means that its conditional distribution is not defined for these  $\mathbf{k}$ . There exists at least one couple  $(k_m, k_{\overline{m}})$  in  $\mathbf{k}$  whose observability of  $[X^{J_m} | X^{J_{\overline{m}}}, k_m, k_{\overline{m}}]$  is undefined. At every step  $\ell$  of the algorithm, there are several couples  $(k_m, k_{\overline{m}})$  such that no observation  $\mathbf{x}_i$  is simultaneously present in the clusters  $k_m$  and  $k_{\overline{m}}$ :  $\{i : \widehat{k}_m(x_i^{J_m}) = k_m, \widehat{k}_{\overline{m}}(x_i^{J_{\overline{m}}}) = k_{\overline{m}}\} = \emptyset$ . Secondly, the Markovian dependency introduced by the specifications  $p(x_i^{J_m} | x_i^{J_{\overline{m}}}, k_m, k_{\overline{m}})$  has for effect to reorganize the initial clustering while maintaining a well-contrasted partitioning. This is a well-known property of the Markovian approach. These two remarks should be useful for establishing an analytical proof of the selection property.

## 4. Discussion

In this paper, we have presented a mixture model dedicated to the case where the dependencies among the components of the multidimensional random vector are governed by a DAG structure. The mixture model takes advantage of a two-level structure, which is composed by the DAG itself and its macro-representation. A dedicated EM algorithm has been efficiently implemented for the Gaussian case. The experiments show that this algorithm is able to select a small number of composite classes. This selection ability is important because it allows to circumvent the difficulty of choosing the exact number of elementary classes for each macro-variable. In fact, one of the main role of the EM-mDAG algorithm is to reveal significant relationships among the hidden elementary classes, some of them becoming empty during the procedure.

### Appendix : a case study

This section presents a domain in which our method should be useful, as it is currently being undertaken by Xiaoyi Chen at Institut Pasteur (Systems Biology team). A  $N$ -sample  $\mathcal{X}$  of tens and even hundreds of thousands of cells is observed by flow cytometry. For each cell  $i = 1, \dots, N$ , the instrument provides a measurement vector  $\mathbf{x}_i$  of dimension  $n$ . This sample is a mixture of several cell populations. The goal is to group these measurements so that each class corresponds to a well-identified cell type [1].

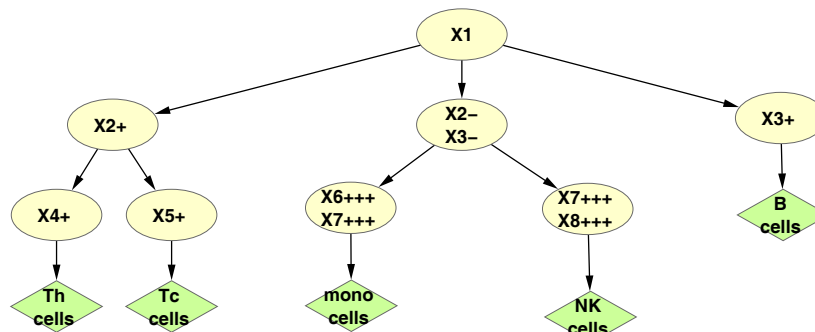


Figure 2. A partial decision tree with biological classes on leaves, (thanks to Xiaoyi Chen, Institut Pasteur).

The analysis, which is based on a dependency tree as illustrated in Fig.1-a, is usually accomplished by sequential manual partitioning (called "Gating") of the sample  $\mathcal{X}$  from the top to the bottom of the tree. Rather than watching simultaneously the  $n$  dimensions, that is to say the cloud  $\mathcal{X}$  in the space  $\mathbf{R}^n$ , the biologist works in subspaces of smaller dimensions, 1, 2 or 3, according to associations of variables  $X^j$ , here called macro-variables, as shown in Fig.1-b.

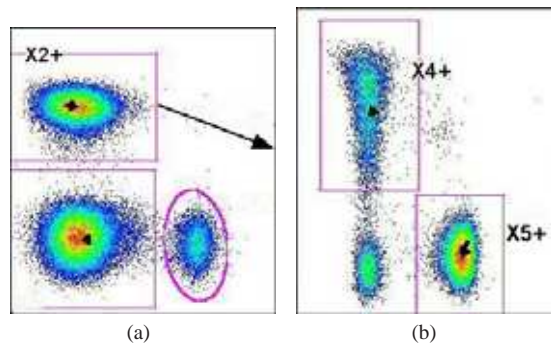
At the top of the tree, only one coordinate of the cloud  $\mathcal{X}$  is analyzed. This is the variable  $X^1$  corresponding to high values  $CD45+$  of the biological variable  $CD45$ . In this example, to simplify, the tree height was reduced by starting the tree with  $X^1 = CD45+$  instead of  $(CD45-, CD45+)$ . To determine the two groups  $CD45-$  et  $CD45+$ , a threshold  $\tau_{CD45}$  is manually selected for separating the small and large values of  $CD45$ . Conditionally on the elementary class  $X^1 = CD45+$ , the procedure continues along the tree structure, as follows.

Three elementary classes are extracted from the 2-D distribution of the sample  $\{(x_i^2, x_i^3)\}_{i=1}^N$  and denoted  $(X^{2+})$ ,  $(X^{2-}, X^{3-})$ ,  $(X^{3+})$  as illustrated in Fig.2 and Fig.3-a. On each group, this operation is repeated on the following macro-variables in dimension 2 for  $(X^4, X^5)$  conditionally on  $(X^{2+})$  as illustrated in Fig.3-b, and in dimension 3 for  $(X^6, X^7, X^8)$  conditionally on  $(X^{2-}, X^{3-})$ .

This conditional and sequential procedure can be represented by a DAG and then modeled by a Bayesian network. The main advantage of using the EM-mDAG is its ability to global classification while keeping the biological dependency structure, which is necessary for identifying the cell types.

## References

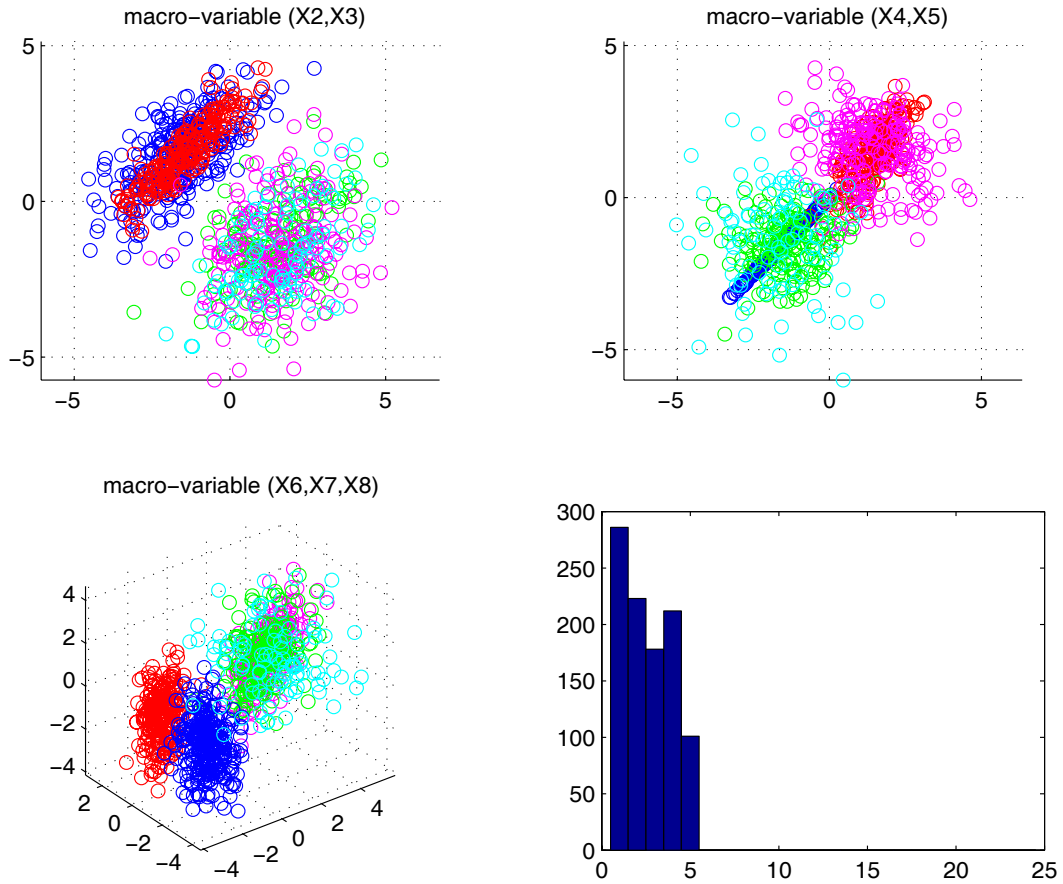
- [1] Nima Aghaeepour, Greg Finak, The FlowCAP Consortium, The DREAM Consortium, Holger Hoos, Tim R Mosmann, Ryan Brinkman, Raphael Gottardo, and Richard H Scheuer-



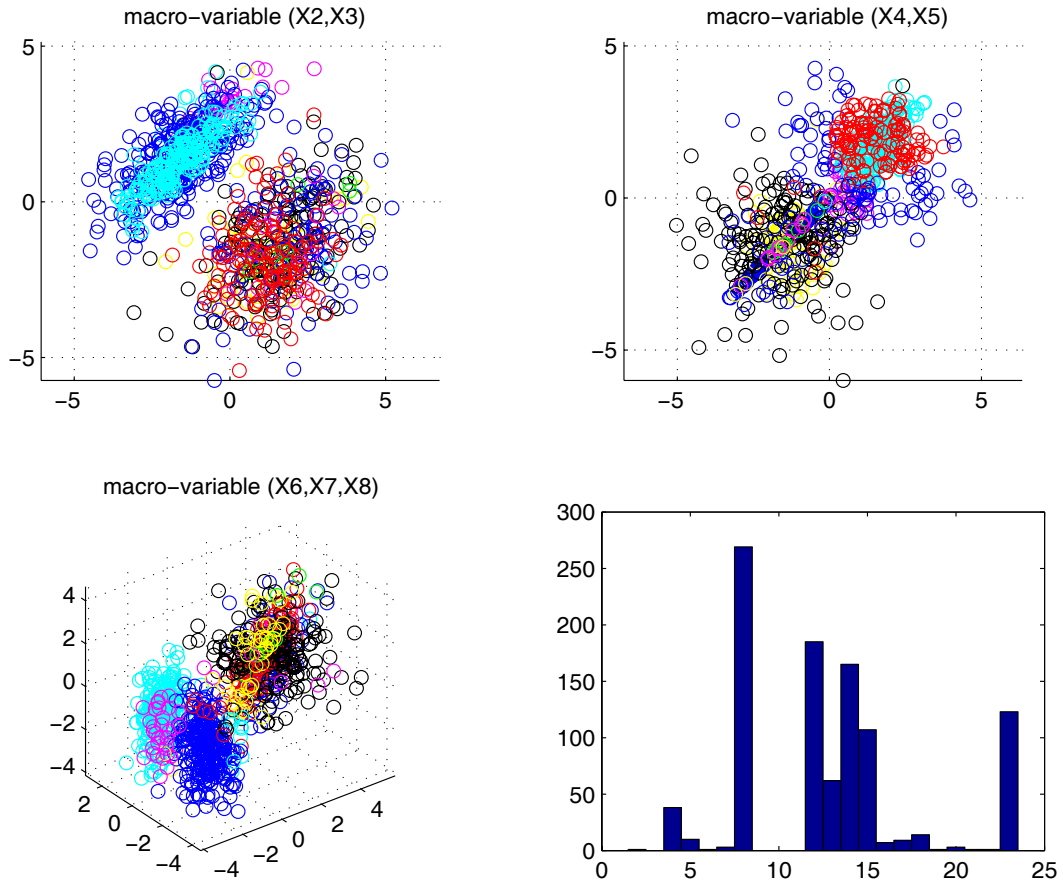
**Figure 3.** Two steps of the sequential procedure leading to the biological classes "Th cells" and "Tc cells" in Fig.2, (thanks to Xiaoyi Chen, Institut Pasteur). From the distribution of the sample  $\{(x_i^2, x_i^3)\}_{i=1}^N$  shown in (a), 3 elementary classes ( $X^{2+}$ ), ( $X^{2-}, X^{3-}$ ), ( $X^{3+}$ ) are manually extracted. (b) shows the distribution of the sample  $\{(x_i^4, x_i^5)\}$  limited to the records  $i$  coming from the class ( $X^{2+}$ ). Conditionally to ( $X^{2+}$ ), 2 new elementary classes ( $X^{4+}$ ) and ( $X^{5+}$ ) are manually extracted.

mann. Critical assessment of automated flow cytometry data analysis techniques. *Nature Methods*, 10(3):228–238, 2013.

- [2] Leonard E. Baum, Ted Petrie, George Soules, and Norman Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *The Annals of Mathematical Statistics*, 41(1):164–171, 1970.
- [3] Bernard Chalmond. An iterative Gibbsian technique for the reconstruction of m-ary images. *Pattern Recognition*, 22:747–761, 1989.
- [4] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39:1–38, 1977.
- [5] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer, 2009.
- [6] Bo Thiesson, Christopher Meek, David Maxwell Chickering, and David Heckerman. Learning mixtures of DAG models. Technical report, Microsoft Research, 1997.

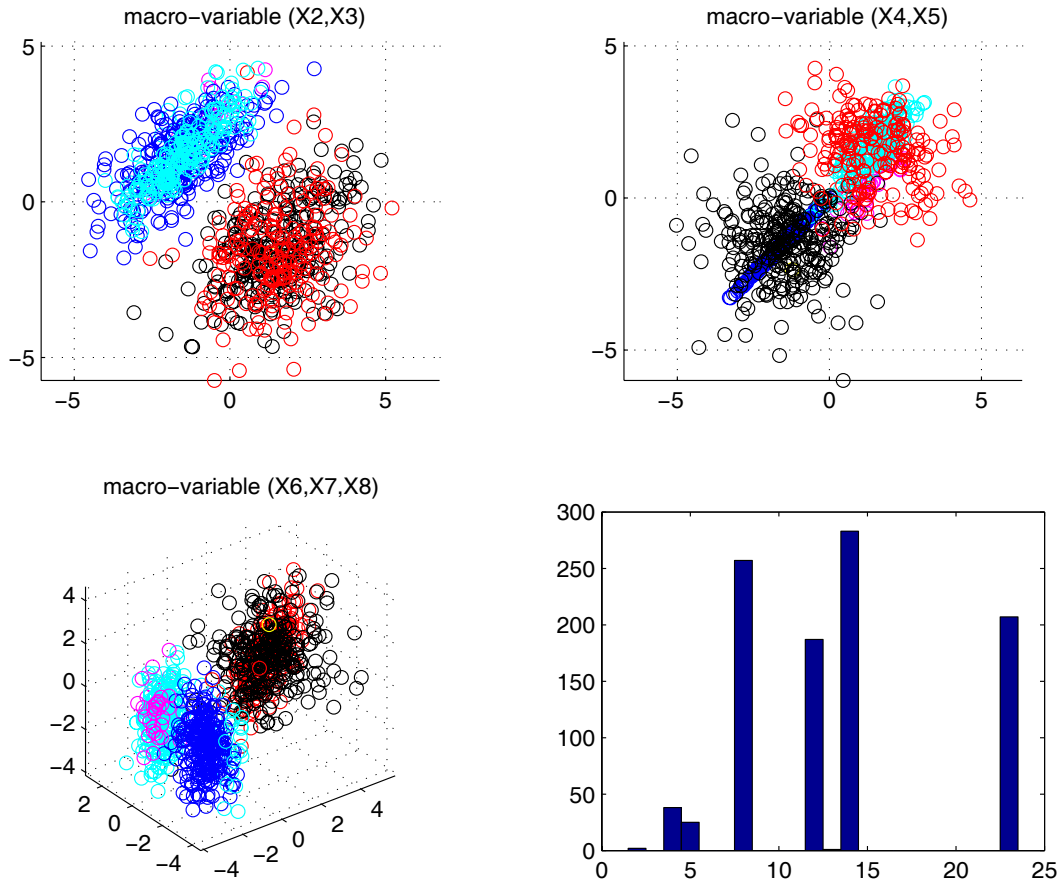


**Figure 4.** True labeling. There are only 5 composite classes. Simulation was performed with  $\nu_1 = 2$ ,  $\nu_2 = 3$ ,  $\nu_3 = 4$  for a sample of size  $N = 1000$ . The graphic "macro-variable (X2,X3)" depicts the coordinates  $\{(x_i^2, x_i^3)\}_{i=1}^N$ , the "macro-variable (X4,X5)" depicts  $\{(x_i^4, x_i^5)\}_{i=1}^N$ , and the "macro-variable (X6,X7,X8)" depicts  $\{(x_i^6, x_i^7, x_i^8)\}_{i=1}^N$ . The histogram gives the empirical distribution of the composite classes.

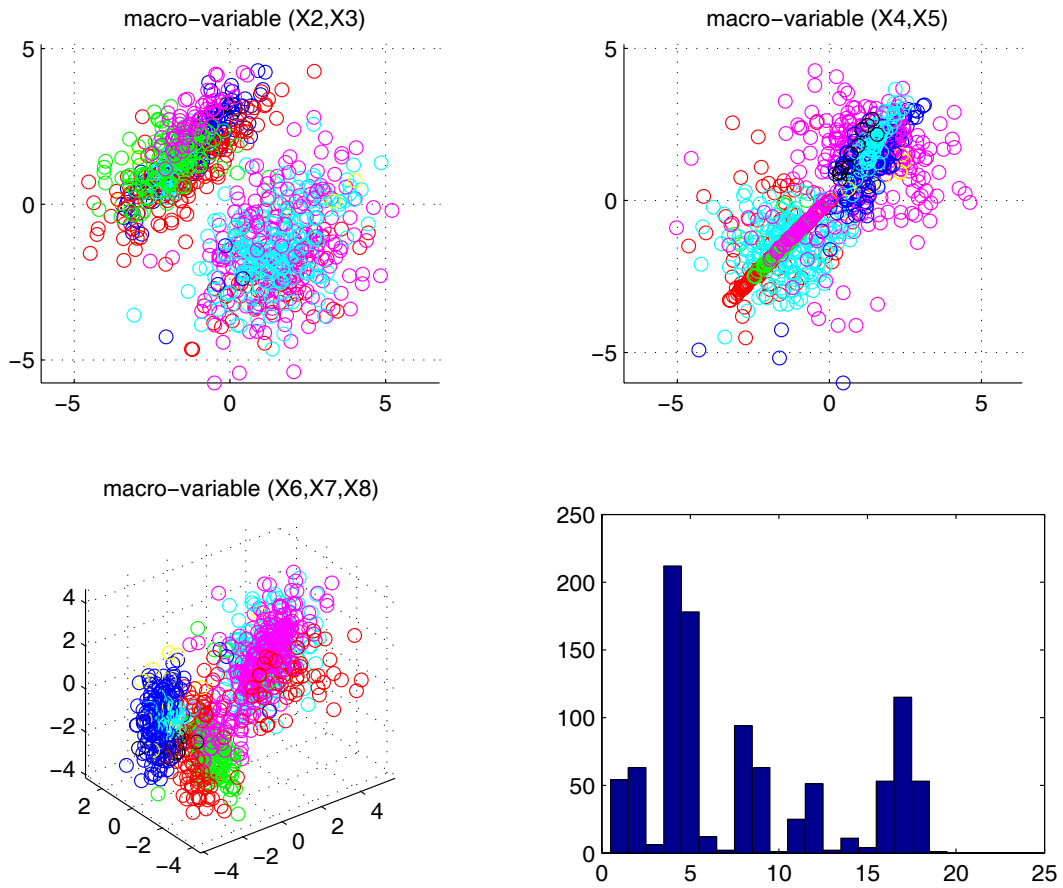


**Figure 5.** Initial solution of the EM-mDAG at step  $\ell = 0$ .  $M = 3$  independent classifications was achieved by applying the classical EM algorithm on each macro-variable. Compared with the ground true in Fig.4, this representation is strongly blurred.

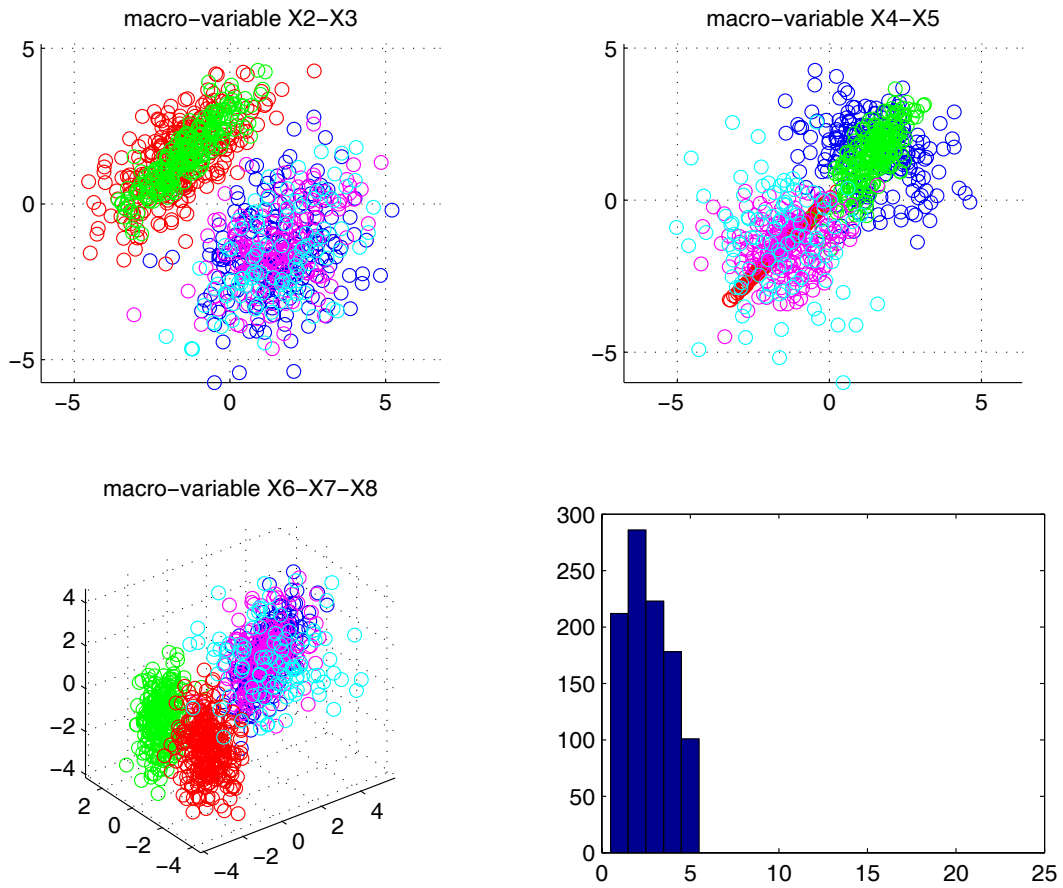




**Figure 6.** EM-mDAG based clustering at step  $\ell = 20$ . As in Fig.4, the macro-variables  $(X^4, X^5)$  and  $(X^6, X^7, X^8)$  respectively highlight the components of the group  $X^{2+}$  and  $X^{3+}$  of the macro-variable  $(X^2, X^3)$ , despite the fifth class that is divided into two neighbour classes, numbered  $k = 4$  and  $k = 5$  in the histogram.



**Figure 7.** Standard EM algorithm for 24 classes. The number of non-empty classes is large and therefore the classification is greatly erroneous.



**Figure 8.** Standard EM algorithm for 5 classes. The macro-variable  $(X^4, X^5)$  does not highlight the mixture components of the group  $X^2$  + of  $(X^2, X^3)$ .