

Dealing with missing data through mixture models

Vincent Vandewalle^{1,2,3}

Joint work with Christophe Biernacki^{3,4,5}

154th ICB Seminar on "Statistics and clinical practice"
Warsaw
11 may 2017

1



2



3



4



5



Overview

1 Missing data problem

Data (1/2)

Many variables collected :

- Heterogeneous data (of various natures)
- Missing data
- Uncertain data (by interval)

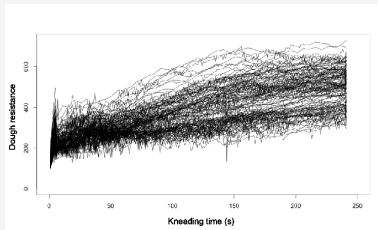
Heterogeneous, missing, uncertain

| Units $x^0 \in \mathcal{X}$ | | | |
|-----------------------------|-------------|--------------|---------|
| ? | 0.5 | ? | 5 |
| 0.3 | 0.1 | green | 3 |
| 0.3 | 0.6 | {red, green} | 3 |
| 0.9 | [0.25 0.45] | red | ? |
| ↓ | ↓ | ↓ | ↓ |
| continuous | continuous | categorical | integer |

Data (2/2)

And also :

- Rank data
- Directional data
- Ordinal data
- Functional data
- Network data
- ...



Imputation of missing data

Motivation for imputing missing values

Most of standard statistical methods cannot cope with missing data \Rightarrow need to complete the dataset

Single imputation

- $E[\mathbf{x}^M]$: average
- $\arg \max p(\mathbf{x}^M)$: mode
- $E[\mathbf{x}^M | \mathbf{x}^O]$: conditional expectation (regression, NPALS, ...)
- Underestimation of the variability related to missing values

Multiple imputation

- Replace missing value by many possible values ideally coming from $p(\mathbf{x}^M | \mathbf{x}^O)$
- Variability inherent to missing data taken into account
- Necessity to consider many versions of the dataset when performing the learning

The fully conditional specification approach (Van Buuren & *al.* (2006)

Presentation

- Specify a conditional model $p(x_j | \mathbf{x}_{-j}, \theta_j)$ for each j (linear regression, logistic regression, ...)
- At each iteration
 - 1 Draw for each variable j the missing value given the observed values, the simulated missing values and the current parameter
 - 2 Update the parameters of the conditional models

Advantage

- Missing data variability taken into account
- Simulated missing values expected to come from $p(\mathbf{x}^M | \mathbf{x}^O)$
- Good performance in practice : implemented in the R package mice.

Supervised classification (1/3)

- **Data** : learning data $\mathcal{D} = (\mathbf{x}^O, \mathbf{z})$
 - n units : $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) = (\mathbf{x}^O, \mathbf{x}^M)$ belonging to the space \mathcal{X}
 - Observed variables \mathbf{x}^O
 - Missing variables \mathbf{x}^M
 - Partition in K clusters $G_1, \dots, G_K : \mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$, $\mathbf{z}_i = (z_{i1}, \dots, z_{iK})'$

$$\mathbf{x}_i \in G_k \Leftrightarrow z_{ik} = 1 \text{ et } \forall h \neq k, z_{ih} = 0$$

- **Goal** : learn a classification rule r on \mathcal{D}

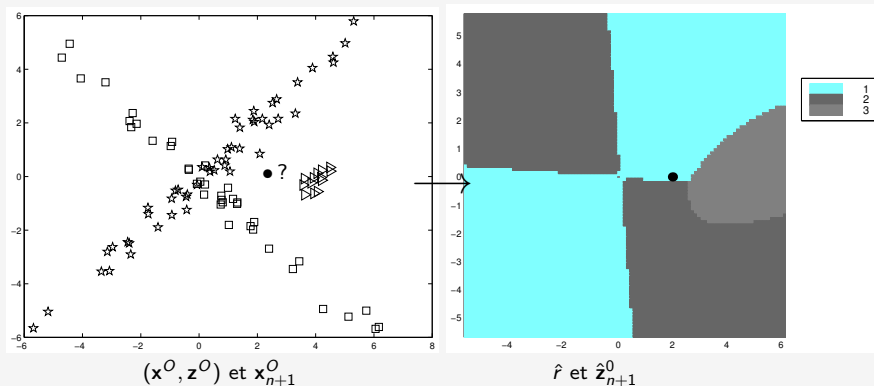
$$\begin{array}{lll} r : & \mathcal{X} & \longrightarrow \{1, \dots, K\} \\ & \mathbf{x}_{n+1}^O & \longmapsto r(\mathbf{x}_{n+1}^O). \end{array}$$

Supervised classification (2/3)

Heterogeneous, missing, uncertain

| Units x^O | | | | Partition z | \Leftrightarrow | Class |
|-------------|-------------|--------------|---------|---------------|-------------------|-------|
| ? | 0.5 | red | 5 | 0 1 0 | \Leftrightarrow | G_2 |
| 0.3 | 0.1 | vert | 3 | 1 0 0 | \Leftrightarrow | G_1 |
| 0.3 | 0.6 | {red, green} | 3 | 1 0 0 | \Leftrightarrow | G_1 |
| 0.9 | [0.25 0.45] | red | ? | 0 0 1 | \Leftrightarrow | G_3 |
| ↓ | ↓ | ↓ | ↓ | | | |
| continuous | continuous | categorical | integer | | | |

Supervised classification (3/3)



Semi-supervised classification (1/3)

- **Data** : learning data $\mathcal{D} = (\mathbf{x}^O, \mathbf{z}^O)$
 - n units : $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) = (\mathbf{x}^O, \mathbf{x}^M)$ belonging to the space \mathcal{X}
 - Observed variables \mathbf{x}^O
 - Missing variables \mathbf{x}^M
 - Partition : $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_n) = (\mathbf{z}^O, \mathbf{z}^M)$
 - Observed class \mathbf{z}^O
 - Missing class \mathbf{z}^M
- **Goal** : estimate the classification rule r from \mathcal{D}

$$r : \begin{array}{ll} \mathcal{X} & \longrightarrow \{1, \dots, K\} \\ \mathbf{x}_{n+1}^O & \longmapsto r(\mathbf{x}_{n+1}^O). \end{array}$$

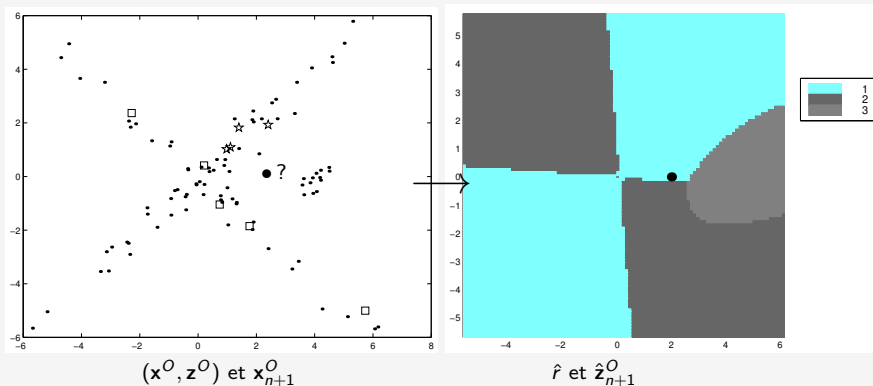
- **Idea** : \mathbf{x} is cheaper than \mathbf{z} thus $n \gg n_\ell$, n_ℓ number of labeled units.

Semi-supervised classification (2/3)

heterogeneous, missing, uncertain

| Units x^O | | | | Partition z^O | \Leftrightarrow | Class |
|-------------|-------------|-------------|---------|-----------------|-------------------|----------------|
| ? | 0.5 | red | 5 | 0 ? ? | \Leftrightarrow | G_2 ou G_3 |
| 0.3 | 0.1 | vert | 3 | 1 0 0 | \Leftrightarrow | G_1 |
| 0.3 | 0.6 | {red,green} | 3 | ? ? ? | \Leftrightarrow | ??? |
| 0.9 | [0.25 0.45] | red | ? | 0 0 1 | \Leftrightarrow | G_3 |
| ↓ | ↓ | ↓ | ↓ | | | |
| continue | continuous | categorical | integer | | | |

Semi-supervised classification (3/3)



Unsupervised classification (1/3)

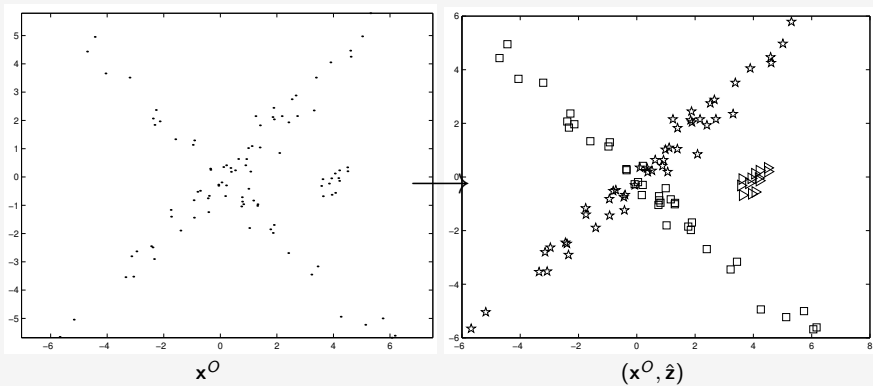
- **Data** : learning data $\mathcal{D} = \mathbf{x}^O$, with $\mathbf{z}^O = \emptyset$
- **Goal** : estimate the partition \mathbf{z} and the number of clusters K
- **Usually usually known as** : clustering

Unsupervised classification (2/3)

heterogeneous, missing, uncertain

| Units x^O | | | | Partition z^O | | | \Leftrightarrow | Cluster |
|-------------|-------------|--------------|---------|-----------------|---|---|-------------------|---------|
| ? | 0.5 | red | 5 | ? | ? | ? | \Leftrightarrow | ??? |
| 0.3 | 0.1 | green | 3 | ? | ? | ? | \Leftrightarrow | ??? |
| 0.3 | 0.6 | {red, green} | 3 | ? | ? | ? | \Leftrightarrow | ??? |
| 0.9 | [0.25 0.45] | red | ? | ? | ? | ? | \Leftrightarrow | ??? |
| ↓ | ↓ | ↓ | ↓ | | | | | |
| continuous | continuous | categorical | integer | | | | | |

Unsupervised classification (3/3)



Usuals solutions (1/2)

Two types of models

- **Generatives models**

- Modeling of $p(\mathbf{x}, \mathbf{z})$
- Which implies a direct modeling of $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$
- Missing data easily taken into account on \mathbf{z} and \mathbf{x}

- **Predictive models**

- Modeling of $p(\mathbf{z}|\mathbf{x})$ eventually only $\mathbf{1}_{\{p(\mathbf{z}|\mathbf{x}) > 1/2\}}$, or only the rank on $p(\mathbf{z}|\mathbf{x})$
- Avoids assumptions on $p(\mathbf{x})$, and thus limits the modeling errors
- Difficulty to take into account missing data on \mathbf{z} and/or \mathbf{x}

Solutions usuelles (2/2)

- **Missing / uncertain data** : multiple imputation
- **Mixed data** : empirical recoding

How to deal with heterogenous, missing, uncertain data in an integrated way ?

Solutions usuelles (2/2)

- **Missing / uncertain data** : multiple imputation
- **Mixed data** : empirical recoding

How to deal with heterogenous, missing, uncertain data in an integrated way?
By using the mixture models.

Mixture models (1/2)

- Rigorous definition of a cluster :

$$\mathbf{x}_1 \in G_k \Leftrightarrow \mathbf{x}_1 \text{ is a realization of } \mathbf{X}_1 \sim p_k(\mathbf{x}_1)$$

- Mixture formulation :

$$\begin{aligned} \mathbf{Z}_1 &\sim \text{Mult}_K(1, \pi_1, \dots, \pi_K) \\ \mathbf{X}_1 |_{Z_{1k}=1} &\sim p_k(\mathbf{x}_1) \end{aligned}$$

- Joint and marginal distribution (mixture) :

$$\begin{aligned} (\mathbf{X}_1, \mathbf{Z}_1) &\sim \prod_{k=1}^K [\pi_k p_k(\mathbf{x}_1)]^{z_{1k}} \\ \mathbf{X}_1 &\sim p(\mathbf{x}_1) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_1) \end{aligned}$$

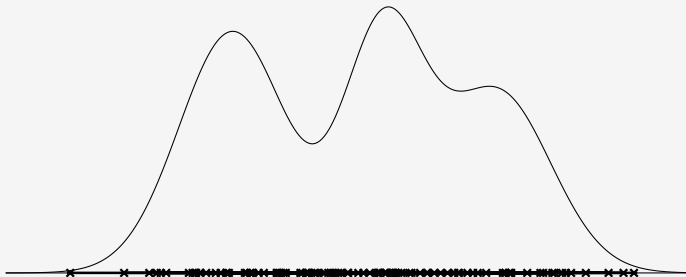
Mixture models (2/2)

$$p(\mathbf{x}_1) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_1).$$



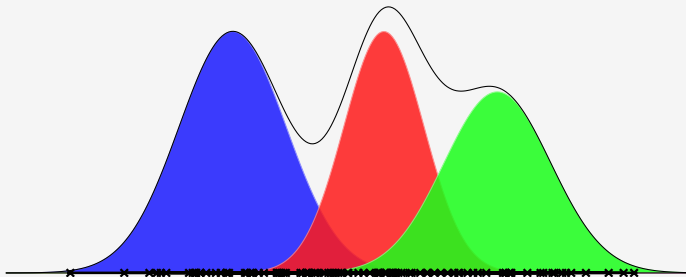
Mixture models (2/2)

$$p(\mathbf{x}_1) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_1).$$



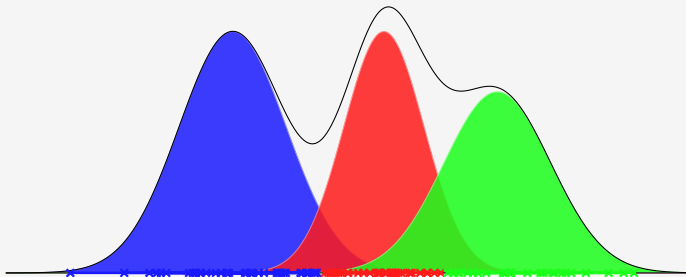
Mixture models (2/2)

$$p(\mathbf{x}_1) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_1)$$



Mixture models (2/2)

$$p(\mathbf{x}_1) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_1)$$



Prostate cancer data (1/5)

- **Units** : patients with prostate cancer classified in grade 3 or 4 of cancer based on advanced study.
- **Variables** : $d = 12$ preliminary measures on each patient :
 - **8 continuous variables** : age, weight, systolic blood pressure, diastolic blood pressure, serum hemoglobin, size of primary tumor, index of tumor stage and histologic grade, serum prostatic acid phosphatase
 - **4 categorical variables** : performance rating, cardiovascular disease history, electrocardiogram code, bone metastases
 - **Missing data**
- **Goal** : predict the grade of the cancer based on preliminary analysis.

| <i>Covariate</i> | <i>Abbreviation</i> | <i>Number of Levels</i> <i>(if categorical)</i> |
|--------------------------------------------|---------------------|----------------------------------------------------|
| Age | Age | |
| Weight | Wi | |
| Performance rating | PF | 4 |
| Cardiovascular disease history | HX | 2 |
| Systolic Blood pressure | SBP | |
| Diastolic blood pressure | DBP | |
| Electrocardiogram code | EKG | 7 |
| Serum haemoglobin | HIG | |
| Size of primary tumour | SZ | |
| Index of tumour stage and histologic grade | SG | |
| Serum prostatic acid phosphatase | AP | |
| Bone metastases | BM | 2 |