MATH 829: Introduction to Data Mining and
Analysis
The EM algorithm

Dominique Guillot

Departments of Mathematical Sciences
University of Delaware

April 18, 2016

## Missing values in data

Missing data is a common problem in statistics.

- No measurement for a given individual/time/location, etc.
- Device failed.
- Error in data entry.
- Data was not disclosed for privacy reasons.
- etc.



Missing data in the titanic passenger dataset.

How can we deal with missing values?

- Many possible procedures.
- The choice of the procedure can significantly impact the conclusions of a study.

## Some strategies for dealing with missing values

Some options for dealing with missing values:

- **Deletion** (delete observations, remove variable, etc.).
  Solves the problem, but ignores some of the data (can be significant). May lead to ignoring an entire "category" of observations. Can generate significant bias.

- **Interpolation.**
  Sometimes it is possible to interpolate missing values (e.g. timeseries). However, we need enough data to be able to produce a good interpolation. In some problems, interpolation is not an option (e.g. age in the titanic passenger data).

- **Replace missing value with mean.**
  May introduce bias. Only valid for numerical observations.

- **Imputation with the EM algorithm.**
  Replace missing values by the *most likely values*. Account for all information available. Much more rigorous. However, requires a model. Can be computationally intensive.

## Missing data mechanism

"Types" of missing data:

1. **Missing completely at random** (MCAR): The events that lead to a missing value are independent both of *observable variables* and of the *unobservable parameters* of interest, and occur entirely at random. (Rarely the case in practice.)

2. **Missing at random** (MAR): missingness is not random, but can be fully accounted for by *observed values*.

3. **Missing not at random** (MNAR): neither MAR nor MCAR.

**Example:** a study about people's weight. We measure (weight, sex).

- Some respondent may not answer the survey for no particular reason. MCAR

- Maybe women are less likely to answer than male (independently of their weight). MAR

- Heavy or light people may be less likely to disclose their weight. MNAR.

## Example

- Suppose we have **independent** observations of a *discrete* random vector $X = (X_1, X_2, X_3, X_4)$ taking values in $\{0, 1, 2, 3\}$.

| $X_1$ | $X_2$ | $X_3$ | $X_4$ |
|---|---|---|---|
| 2 | 0 | 2 | 3 |
| 3 | NA | 1 | 1 |
| 1 | 3 | NA | NA |
| 2 | NA | 1 | NA |

- Let $p(x_1, x_2, x_3, x_4) = P(X_1 = x_1, X_2 = x_2, X_3 = x_3, X_4 = x_4)$ be the pmf of $X$.

- Ignoring the missing data mechanism, we have

$$p(x_1, \mathrm{NA}, x_3, x_4) = \sum_{x=0}^{3} p(x_1, x, x_2, x_3).$$

## Example (cont.)

- Suppose the data comes from a parametric model $p(x_1, x_2, x_3, x_4; \theta)$ where $\theta \in \Theta$ is unknown.

| $X_1$ | $X_2$ | $X_3$ | $X_4$ |
|---|---|---|---|
| 2 | 0 | 2 | 3 |
| 3 | NA | 1 | 1 |
| 1 | 3 | NA | NA |
| 2 | NA | 1 | NA |

- We compute the *likelihood* of the data:

$$L(\theta) = p(2, 0, 2, 3) \times p_{1,3,4}(3, 1, 1) \times p_{1,2}(1, 3) \times p_{1,3}(2, 1),$$

where $p_{1,3,4}(x_1, x_3, x_4) = \sum_{x_2=0}^{3} p(x_1, x_2, x_3, x_4)$, $p_{1,2}(x_1, x_2) = \sum_{x_3=0}^{3} \sum_{x_4=0}^{3} p(x_1, x_2, x_3, x_4)$, and $p_{1,3}(x_1, x_3) = \sum_{x_2=0}^{3} \sum_{x_4=0}^{3} p(x_1, x_2, x_3, x_4)$ denote *marginals* of $p$.

- The *likelihood* can now be maximized as a function of $\theta$.

## Imputing the missing values

- Recall that $f(x) = E(Y | X = x)$ has the following optimality property:

$$E(Y | X = x) = \operatorname*{argmin}_{c \in \mathbb{R}} E(Y - c)^2$$

where $c$ is some function of $x$.

- So $E(Y | X = x)$ is the "best prediction" of $Y$ given $X$ in the mean squared error sense.

- As a result, once $p(x; \theta)$ is known (after estimating $\theta$ by maximum likelihood for example), we can *impute* missing values using:

$$\hat{x}_{\mathrm{miss}} = E(x_{\mathrm{miss}} | x_{\mathrm{observed}}).$$

For example, if $x = (1, 3, \mathrm{NA}, \mathrm{NA})$ then:

$$(\hat{x}_3, \hat{x}_4) = E((X_3, X_4) | X_1 = 1, X_2 = 3),$$

where $E$ is computed with respect to $p(x_1, x_2, x_3, x_4; \theta)$.

## Summary

In summary, given a family of probability models $p(x; \theta)$ for the data, under MAR, we can:

1. Compute the likelihood of $\theta$ by *marginalizing* over the missing values.
2. Estimate the parameter $\theta$ by maximum likelihood.
3. Impute missing values using $\hat{x}_{\mathrm{miss}} = E_{\theta}(x_{\mathrm{miss}} | x_{\mathrm{observed}})$. where $E_{\theta}$ denotes the expected value with respect to the probability distribution $p_{\theta}$.

**Remark:** We assumed above that the variables are discrete, and the observations are independent for simplicity. The same procedure applied in the general case.

- The methodology described so far solves our missing data problem in principle.
- However, explicitly finding the maximum of the likelihood function can be very difficult.

The **Expectation-Maximization** (EM) algorithm of *Dempster, Laird, and Rubin, 1977* provides a more efficient way of solving the problem.

The EM algorithm leverages the fact the the likelihood is often easy to maximize if there is no missing values.

For simplicity, we will assume our observations are independent and the random variables are discrete.

Some notation:

- We have a random vector $W$ taking values in $\mathbb{R}^p$.
- The distribution of the vector is $p(w; \theta)$.
- We want to estimate $\theta$.
- We only observe a part of the vector

$$(x^{(i)}, z^{(i)}) \in \mathbb{R}^{p_i} \times \mathbb{R}^{p - p_i} \qquad (i = 1, \dots, n).$$

- So $x^{(i)}$ is the **observed** part and $z^{(i)}$ is the **unobserved** part.
- The log-likelihood function is given by

$$l(\theta) = \sum_{i=1}^{n} \log p(x^{(i)}; \theta) = \sum_{i=1}^{n} \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta).$$

(the second sum is over all the possible values of $z^{(i)}$).

- We would like to maximize that function over $\theta$ (generally difficult).

Instead of trying to maximize the log-likelihood directly, the EM algorithm constructs a sequence of approximations $\theta^{(i)}$ of $\theta$.

- Let $\theta^{(0)}$ be an **initial guess** for $\theta$.
- Given the current estimate $\theta^{(i)}$ of $\theta$, compute

$$Q(\theta | \theta^{(i)}) := E_{z|x; \theta^{(i)}} \ \log p(x, z; \theta)$$

$$= \sum_{i=1}^{n} E_{z^{(i)} | x^{(i)}; \theta^{(i)}} \left( \log p(x^{(i)}, z^{(i)}; \theta) \right) \qquad \text{(E step)}$$

(In other words, we average the missing values according to their distribution after observing the observed values.)

- We then optimize $Q(\theta | \theta^{(i)})$ with respect to $\theta$:

$$\theta^{(i+1)} := \operatorname*{argmax}_{\theta} Q(\theta | \theta^{(i)}) \qquad \text{(M step)}.$$

**Theorem:** The sequence $\theta^{(i)}$ constructed by the EM algorithm satisfies:

$$l(\theta^{(i+1)}) \geq l(\theta^{(i)}).$$

**Remark:** There is no guarantee that the EM algorithm will find the global max of the likelihood. It may only find a local max.