

# Computational Statistics II

Unit C.1: Missing data problems, Gibbs sampling and the EM algorithm

**Tommaso Rigon**

**University of Milano-Bicocca**

Ph.D. in Economics, Statistics and Data Science



# Unit C.1

## Main concepts

- Missing data problems;
- Data augmentation and Gibbs sampling;
- The EM algorithm and generalizations;
- Minorize maximize (MM) algorithms.

## Main references

- Bishop, C. M. (2006). Pattern Recognition and Machine Learning, Chapter 9. Springer.
- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *JRSS-B*, **39**(1), 1–38.
- Hunter, D. R., and Lange, K. (2004). A Tutorial on MM Algorithms. *The American Statistician*, **58**(1), 30–37.
- McLachlan, G. J. and Krishnan, T. (1998). The EM Algorithm and Extensions. Wiley.
- Robert, C. P., and Casella, G. (2009). Introducing Monte Carlo methods with R. Springer.

# Missing data problems

- In this unit we will take advantage of specific structures of the model to facilitate both frequentist and Bayesian computations via the EM and Gibbs sampling.
- In most cases, this will involve the introduction of **hidden features** of the model, sometimes called **latent variables**.
- Depending on the context, these latent quantities will have a precise meaning or they will be regarded as purely abstract objects.
- An obvious examples of latent components with a precise interpretation is the case of **missing** or **censored observations**.
- **Key idea**. If the complete data were available, computations would be easier. Besides, imputing the missing values could be interesting on its own.

# Example: survival analysis with an exponential model

- Let  $\mathbf{z} = (z_1, \dots, z_n)^\top$  be iid exponential random variables with rate parameter  $\theta > 0$ .
- If the prior  $\theta \sim \text{Ga}(a, b)$ , then thanks to conjugacy we get the following posterior

$$(\theta \mid \mathbf{z}) \sim \text{Ga} \left( a + n, b + \sum_{i=1}^n z_i \right).$$

- However, in many cases observations are **censored**, as in **Unit A.1**. In fact, we observe the values  $\mathbf{t} = (t_1, \dots, t_n)^\top$  which are either complete ( $t_i = z_i$ ) or censored ( $t_i \leq z_i$ ).
- If the observations were all **complete**, then inference would be straightforward.
- Intuitively, we aim at **sampling** or imputing the **missing information** from the appropriate conditional distribution, in order to make inference about  $\theta$ .

# Data augmentation

- Let  $\mathbf{X}$  be the **observed** data, following some distribution  $\pi(\mathbf{X} \mid \theta)$ , i.e. the **likelihood**, with  $\theta \in \Theta \subseteq \mathbb{R}^p$  being an unknown set of parameters.
- Let  $\pi(\theta)$  be the prior distribution associated to  $\theta$  and let  $\pi(\theta \mid \mathbf{X})$  be the posterior.
- Let  $\mathbf{z} \in \mathcal{Z} \subseteq \mathbb{R}^q$  be a vector of **latent variables**, which are not observed.
- We assume that the likelihood function  $\pi(\mathbf{X} \mid \theta)$  can be written as the marginal distribution of a **complete likelihood**, namely

$$\pi(\mathbf{X} \mid \theta) = \int_{\mathcal{Z}} \pi(\mathbf{X}, \mathbf{z} \mid \theta) d\mathbf{z}.$$

- **Remark.** We focus on continuous densities w.r.t. the Lebesgue measure for the sake of notational simplicity, but these ideas apply in general.

# Data augmentation

- The quantity  $\pi(\mathbf{X}, \mathbf{z} \mid \theta)$  is the **complete** or **augmented** likelihood.
- Within the Bayesian framework, we treat the latent variables  $\mathbf{z}$  as if they were an additional set of unknown parameters, leading to the **augmented posterior**

$$\pi(\theta, \mathbf{z} \mid \mathbf{X}) \propto \pi(\mathbf{X}, \mathbf{z} \mid \theta)\pi(\theta).$$

- In other words, we aim at sampling  $(\theta^{(r)}, \mathbf{z}^{(r)})$  using MCMC from the joint posterior  $\pi(\theta, \mathbf{z} \mid \mathbf{X})$ , which can be performed using any of the strategies we have described.
- If one is interested only in the original parameters  $\theta$  or in the latent dimensions  $\mathbf{z}$ , then it suffices to **ignore** the other set of parameters.
- We sample from  $\pi(\theta, \mathbf{z} \mid \mathbf{X})$  and then discard  $\mathbf{z}$  rather than directly targeting  $\pi(\theta \mid \mathbf{X})$  because the **augmented likelihood** is typically **more tractable** than the original one.

# Data augmentation schemes

- Unfortunately, there are **no general recipes** for finding useful data augmentation schemes. We will see proposals in the probit and logit case in **unit C.2**.
- In principle, whenever the likelihood can be expressed in an integral form, this leads to a potential data augmentation mechanism.
- However, the resulting augmented likelihood must be tractable, otherwise the whole procedure is of little practical utility.
- **Mixture models** greatly benefit from data-augmentation schemes, but we do not discuss them here because they would deserve an entire course on their own.

# Data augmentation and Gibbs sampling

- Although in principle any MCMC strategy could be used to target  $\pi(\theta, \mathbf{z} \mid \mathbf{X})$ , the Gibbs sampling is a natural choice in this setting.
- In fact, it is often the case that the following **full conditional distributions** are available in closed form. Moreover, they also have a nice interpretation.

- **Step 1.** Sample from the “posterior” of  $\theta$  based on the complete likelihood, namely

$$\pi(\theta \mid \mathbf{X}, \mathbf{z}) \propto \pi(\mathbf{X}, \mathbf{z} \mid \theta)\pi(\theta).$$

- **Step 2.** Impute the missing observations  $\mathbf{z}$  by sampling from the full conditional

$$\pi(\mathbf{z} \mid \mathbf{X}, \theta) \propto \pi(\mathbf{X}, \mathbf{z} \mid \theta).$$

- Obviously, we are allowed to split  $\theta$  and  $\mathbf{z}$  into blocks of parameters if this facilitate the Gibbs sampling.



# Example: survival analysis with an exponential model

- Recall the exponential model example with censored data  $\mathbf{t}$  and censorship indicators  $\mathbf{d} = (d_1, \dots, d_n)^\top$ . The **original likelihood** is therefore equal to

$$\pi(\mathbf{t}, \mathbf{d} \mid \theta) = \theta^{n_c} \exp \left\{ -\theta \sum_{i=1}^n t_i \right\}, \quad n_c = \sum_{i=1}^n d_i.$$

- **Remark.** This is a toy example whose purpose is fixing ideas. Indeed, under a Gamma prior, the posterior distribution of  $\theta$  using this likelihood is also available.
- In this setting, the latent variables  $\mathbf{z}$  represent the complete survival times having exponential distribution, so that the **complete likelihood** is

$$\pi(\mathbf{z} \mid \theta) = \theta^n \exp \left\{ -\theta \sum_{i=1}^n z_i \right\}.$$

- The **Gibbs sampling** alternates between the Gamma full conditional  $\pi(\theta \mid \mathbf{z})$  and a sampling step from  $\pi(\mathbf{z} \mid \mathbf{t}, \theta)$ . Note that  $(z_i - t_i \mid t_i, d_i, \theta) \stackrel{\text{ind}}{\sim} \text{Exp}(\theta)$  when  $d_i = 0$ .

# The EM algorithm

- A Gibbs sampling based on data augmentation strategies is strongly connected with the so-called **expectation-maximization** (EM) algorithm.
- The EM is a deterministic algorithm that aims at **maximizing** the likelihood (MLE) or the posterior distribution (MAP), namely at finding

$$\arg \max_{\theta \in \Theta} \pi(\theta \mid \mathbf{X}) = \arg \max_{\theta \in \Theta} \pi(\mathbf{X} \mid \theta) \pi(\theta).$$

- The EM is widely used both within the frequentist and the Bayesian framework. The MLE case is recovered whenever  $\pi(\theta) \propto 1$ .
- Compared to other gradient-based maximizers, it leads to a **monotonic sequence**. The target function always increases during the procedure, thus being more stable.
- On the other hand, the EM **requires** a (tractable) **augmented likelihood**. Moreover, the EM could be slower than other algorithms to reach convergence.

# The EM algorithm

- The EM algorithm alternates between the following steps, which are reminiscent of those of the Gibbs sampling, as they involve similar quantities.
- Initialize the algorithm at a reasonable  $\theta^{(0)}$ . The generic iteration proceeds as follows.
- **Step 1 (Expectation)**. Let  $\theta^{(r)}$  be the current value of the maximization procedure, then obtain the function

$$Q(\theta \mid \theta^{(r)}) = \mathbb{E}\{\log \pi(\mathbf{X}, \mathbf{z} \mid \theta)\},$$

where the expectation is taken with respect to the conditional law  $\pi(\mathbf{z} \mid \mathbf{X}, \theta^{(r)})$ .

- **Step 2 (Maximization)**. The new value of the procedure  $\theta^{(r+1)}$  is obtained by maximizing the function

$$\arg \max_{\theta \in \Theta} Q(\theta \mid \theta^{(r)}) + \log \pi(\theta).$$

- In many cases, the E-step amounts at calculating  $\mathbb{E}(\mathbf{z})$  and then plugging-in this quantity in the augmented log-likelihood. Indeed,  $\log \pi(\mathbf{X}, \mathbf{z} \mid \theta)$  is often linear in  $\mathbf{z}$ .

# Example: survival analysis with an exponential model

- Recall that in the exponential model example, we have that  $(z_i - t_i \mid t_i, d_i, \theta) \stackrel{\text{ind}}{\sim} \text{Exp}(\theta)$  when  $d_i = 0$  and the augmented likelihood is  $\pi(\mathbf{z} \mid \theta) = \theta^n \exp\{-\theta \sum_{i=1}^n z_i\}$ .
- Let us focus on the maximum likelihood, so that  $\pi(\theta) \propto 1$ .

- **Step 1 (Expectation)**. Let  $\theta^{(r)}$  be the current value of the procedure, then

$$\mathcal{Q}(\theta \mid \theta^{(r)}) = n \log \theta - \theta \sum_{i=1}^n \mathbb{E}(z_i) = n \log \theta - \theta \sum_{i=1}^n \{t_i + (1 - d_i)\theta^{(r)}\},$$

where the expectation is taken with respect to the conditional law  $\pi(\mathbf{z} \mid \mathbf{t}, \mathbf{d}, \theta^{(r)})$ .

- **Step 2 (Maximization)**. The new value of the procedure  $\theta^{(r+1)}$  is obtained by considering the maximum of  $\mathcal{Q}(\theta \mid \theta^{(r)})$ , thus obtaining

$$\theta^{(r+1)} = \left( \frac{1}{n} \sum_{i=1}^n t_i + \frac{n - n_c}{n} \theta^{(r)} \right)^{-1}.$$

# Why does the EM work?

## Theorem (monotonic EM sequence)

The EM sequence for finding the MLE satisfies the following inequality

$$\pi(\mathbf{X} \mid \boldsymbol{\theta}^{(r+1)}) \geq \pi(\mathbf{X} \mid \boldsymbol{\theta}^{(r)}).$$

Similarly, the EM sequence for finding the MAP satisfies the following inequality

$$\pi(\boldsymbol{\theta}^{(r+1)} \mid \mathbf{X}) \geq \pi(\boldsymbol{\theta}^{(r)} \mid \mathbf{X}).$$

- With some further continuity assumptions w.r.t.  $\boldsymbol{\theta}$ , this theorem implies that the EM is guaranteed to reach a **stationary point**.
- If the posterior / likelihood function is concave, the stationary point will be also the global maximum.
- In general, as in any maximization procedure, it is recommended to initialize the algorithm at different starting points.

# Sketch of the proof

- In first place, recognize that the following identity holds true (do it as an exercise!)

$$\log \pi(\boldsymbol{\theta} \mid \mathbf{X}) = \log \pi(\mathbf{X}, \mathbf{z} \mid \boldsymbol{\theta}) + \log \pi(\boldsymbol{\theta}) - \log \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}) - \log \pi(\mathbf{X}),$$

Consequently, one gets the following identity

$$\log \pi(\boldsymbol{\theta} \mid \mathbf{X}) = \mathcal{Q}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^r) + \log \pi(\boldsymbol{\theta}) - \mathbb{E}\{\log \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta})\} - \log \pi(\mathbf{X}),$$

after taking the expectation w.r.t.  $\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^r)$ .

- Let  $\boldsymbol{\theta}^{(r)}$  and  $\boldsymbol{\theta}^{(r+1)}$  be subsequent steps in the EM procedure. Then necessarily it holds that

$$\mathcal{Q}(\boldsymbol{\theta}^{(r+1)} \mid \boldsymbol{\theta}^{(r)}) + \log \pi(\boldsymbol{\theta}^{(r+1)}) \geq \mathcal{Q}(\boldsymbol{\theta}^{(r)} \mid \boldsymbol{\theta}^{(r)}) + \log \pi(\boldsymbol{\theta}^{(r)}),$$

as the value  $\boldsymbol{\theta}^{(r+1)}$  is indeed maximizing the left-hand-side. Furthermore note that because of Jensen's inequality we get

$$\mathbb{E} \left\{ \log \frac{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r+1)})}{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})} \right\} \leq \log \mathbb{E} \left\{ \frac{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r+1)})}{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})} \right\} = 0,$$

expectations being taken w.r.t. to  $\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})$ . This implies that

$$-\mathbb{E}\{\log \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r+1)})\} \geq -\mathbb{E}\{\log \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})\}.$$

- The proof follows by combining the above results, after noting that

$$\begin{aligned} \log \pi(\boldsymbol{\theta}^{(r+1)} \mid \mathbf{X}) &= \mathcal{Q}(\boldsymbol{\theta}^{(r+1)} \mid \boldsymbol{\theta}^{(r)}) + \log \pi(\boldsymbol{\theta}^{(r+1)}) - \mathbb{E}\{\log \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r+1)})\} - \log \pi(\mathbf{X}) \geq \\ &\geq \mathcal{Q}(\boldsymbol{\theta}^{(r)} \mid \boldsymbol{\theta}^{(r)}) + \log \pi(\boldsymbol{\theta}^{(r)}) - \mathbb{E}\{\log \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})\} - \log \pi(\mathbf{X}) = \log \pi(\boldsymbol{\theta}^{(r)} \mid \mathbf{X}). \end{aligned}$$

# An alternative derivation of the EM

- There exists an alternative derivation of the EM purely based on **maximization**.
- Albeit less common, this way of thinking leads to a more **elegant proof** and puts the basis for variational Bayes (VB) procedures **unit D.1**.
- Let  $q(\mathbf{z}) \in \mathbb{Q}$  be a generic density of the latent variables  $\mathbf{z}$  and define

$$\mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \boldsymbol{\theta}\} = \mathbb{E}_q \left( \log \frac{\pi(\mathbf{X}, \mathbf{z} \mid \boldsymbol{\theta})}{q(\mathbf{z})} \right),$$

where the expectations are taken w.r.t.  $q(\mathbf{z})$ .

- Moreover, define the **Kullback-Leibler divergence**

$$\text{KL}\{q(\mathbf{z}) \parallel \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta})\} = -\mathbb{E}_q \left( \log \frac{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{z})} \right).$$

# A maximization / maximization procedure

- Let us focus on the MLE case for notational simplicity. The MAP case is recovered with some minor adjustments (do it as an exercise!)
- For any  $q \in \mathbb{Q}$  the following identity holds true

$$\log \pi(\mathbf{X} \mid \theta) = \mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \theta\} + \text{KL}\{q(\mathbf{z}) \parallel \pi(\mathbf{z} \mid \mathbf{X}, \theta)\}.$$

- Since the Kullback-Leibler divergence  $\text{KL}\{q(\mathbf{z}) \parallel \pi(\mathbf{z} \mid \mathbf{X}, \theta)\} \geq 0$ , then we will have

$$\mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \theta\} \leq \log \pi(\mathbf{X} \mid \theta),$$

meaning that  $\mathcal{L}\{q(\mathbf{z}) \mid \theta, \mathbf{X}\}$  is the **lower bound** of the log-likelihood.

- This suggests that the MLE can be found **maximizing the lower bound**, since

$$\arg \max_{\theta \in \Theta} \log \pi(\mathbf{X} \mid \theta) = \arg \max_{\theta \in \Theta} \max_{q \in \mathbb{Q}} \mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \theta\}.$$

- Indeed, the value  $q(\mathbf{z}) = \pi(\mathbf{z} \mid \mathbf{X}, \theta)$  is the maximum of  $\mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \theta\}$ , because

$$\mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \theta\} = \log \pi(\mathbf{X} \mid \theta) - \underbrace{\text{KL}\{q(\mathbf{z}) \parallel \pi(\mathbf{z} \mid \mathbf{X}, \theta)\}}_{=0} = \log \pi(\mathbf{X} \mid \theta).$$



# A maximization / maximization procedure

- Consequently, the MLE can be obtained by iteratively maximizing  $\mathcal{L}\{q(\mathbf{z}) \mid \boldsymbol{\theta}, \mathbf{X}\}$  over  $q(\mathbf{z})$  for a given value of  $\boldsymbol{\theta}$  and then over  $\boldsymbol{\theta}$  for a given  $q(\mathbf{z})$ .
- Let  $\boldsymbol{\theta}^{(r)}$  be the current value of the procedure.

- **Step 1 (Maximization over  $q$ )**. Given the fixed value  $\boldsymbol{\theta}^{(r)}$ , obtain

$$\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)}) = \arg \max_{q \in \mathcal{Q}} \mathcal{L}\{q(\mathbf{z}) \mid \mathbf{X}, \boldsymbol{\theta}^{(r)}\} = \arg \min_{q \in \mathcal{Q}} \text{KL}\{q(\mathbf{z}) \parallel \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})\}.$$

- **Step 2 (Maximization over  $\boldsymbol{\theta}$ )**. Given the locally optimal value  $q(\mathbf{z}) = \pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)})$ , obtain the new value  $\boldsymbol{\theta}^{(r+1)}$  as the maximizer

$$\boldsymbol{\theta}^{(r+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}\{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)}) \mid \mathbf{X}, \boldsymbol{\theta}\} = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{Q}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(r)}).$$

- These are the steps of the EM, which therefore has an alternative interpretation.
- Moreover, recalling that  $\mathcal{L}\{\pi(\mathbf{z} \mid \mathbf{X}, \boldsymbol{\theta}^{(r)}) \mid \mathbf{X}, \boldsymbol{\theta}^{(r)}\} = \log \pi(\mathbf{X} \mid \boldsymbol{\theta}^{(r)})$ , the monotonicity property of the EM is obvious.

# Generalizations of the EM

- Sometimes the **maximization** of  $Q(\theta \mid \theta^{(r)}) + \log \pi(\theta)$ , namely the maximization step, could be **difficult**.
- Thus, an obvious generalization of the EM algorithm that preserves the monotonicity of the procedure is considering some value  $\theta^{(r+1)}$  such that

$$Q(\theta^{(r+1)} \mid \theta^{(r)}) + \log \pi(\theta^{(r+1)}) \geq Q(\theta^{(r)} \mid \theta^{(r)}) + \log \pi(\theta^{(r)})$$

that is,  $\theta^{(r+1)}$  **increases the function** rather maximizing it.

- An example is the **expectation conditional maximization** (ECM) of Meng and Rubin (1993), where the parameters are partitioned into sub-groups and iteratively maximized.
- Similar ideas can be applied to generalize the expectation step by doing a “partial” update in the maximization of  $q$ .

- We finally consider a large class of optimization methods called **minorize maximize** (MM) that includes the EM as special case.
- The MM methods do not involve missing data or data augmentations, but they rather rely on general **convexity** arguments.
- The MM is used to optimize a  $\ell(\theta; \mathbf{X})$  of the parameters  $\theta$  and the data  $\mathbf{X}$ , with  $f(\cdot)$  being the posterior distribution, the likelihood, or a general loss function.
- Let  $\theta^{(r)}$  be the current value of the iterative maximization procedure. We are seeking for a **minorization function**  $g(\theta \mid \theta^{(r)})$ , such that

$$g(\theta \mid \theta^{(r)}) \leq \ell(\theta; \mathbf{X}), \quad \text{for any } \theta \in \Theta,$$

and satisfying  $g(\theta \mid \theta) = \ell(\theta; \mathbf{X})$ .

- In MM algorithms we iteratively maximize the **lower bound**  $g(\theta; \theta^{(r)}, \mathbf{X})$ , so that

$$\theta^{(r+1)} = \arg \max_{\theta \in \Theta} g(\theta \mid \theta^{(r)})$$

- MM leads to **monotonic sequences**, since

$$\ell(\theta^{(r+1)}; \mathbf{X}) \geq g(\theta^{(r+1)} \mid \theta^{(r)}) \geq g(\theta^{(r)} \mid \theta^{(r)}) = \ell(\theta^{(r)}; \mathbf{X}).$$

- This property ensures remarkable numerical stability, but does not provide any hint about the actual construction of  $g(\theta \mid \theta^{(r)})$ .

- The EM is indeed a **special case** of this framework, recovered in the MLE case by defining

$$g(\theta \mid \theta^{(r)}) = \mathcal{L}\{\pi(\mathbf{z} \mid \mathbf{X}, \theta^{(r)}) \mid \mathbf{X}, \theta\} \leq \log \pi(\mathbf{X} \mid \theta).$$

and recalling that  $g(\theta \mid \theta^{(r)}) = \mathcal{Q}(\theta \mid \theta^{(r)}) + \text{const}$ , and that  $g(\theta \mid \theta) = \log \pi(\mathbf{X} \mid \theta)$ .

- We will see an example in **unit C.2** for the logistic regression case.