

# Uncovering mortality patterns and hospital effects in COVID-19 heart failure patients: a novel multilevel logistic cluster-weighted modeling approach

Luca Caldera<sup>1,\*</sup>, Chiara Masci<sup>2</sup>, Andrea Cappozzo<sup>3</sup>, Marco Forlani<sup>4</sup>, Barbara Antonelli<sup>5</sup>,  
Olivia Leoni<sup>6</sup>, Francesca Ieva<sup>1,7</sup>

<sup>1</sup>MOX, Department of Mathematics, Politecnico di Milano, Via Bonardi 9, Milan 20133, Italy, <sup>2</sup>Department of Economics, Management, and Quantitative Methods, University of Milan, Via Festa del Perdono 7, Milan 20122, Italy, <sup>3</sup>Department of Statistical Sciences, Università Cattolica del Sacro Cuore, Largo Gemelli 1, Milan 20123, Italy, <sup>4</sup>Regione Lombardia, Informatica SPA, Piazza Città di Lombardia 1, Milan 20124, Italy, <sup>5</sup>Regione Lombardia, Divisione Servizi per il Welfare Regionale, Piazza Città di Lombardia 1, Milan 20124, Italy, <sup>6</sup>U.O. Osservatorio Epidemiologico, DG Welfare, Regione Lombardia, Piazza Città di Lombardia 1, Milan 20124, Italy, <sup>7</sup>Health Data Science Centre, Human Technopole, Viale Rita Levi-Montalcini 1, Milan 20157, Italy

\*Corresponding author: Luca Caldera, MOX, Department of Mathematics, Politecnico di Milano, Via Bonardi 9, Milan 20133, Italy ([luca.caldera@polimi.it](mailto:luca.caldera@polimi.it)).

## ABSTRACT

Evaluating hospital performance and its relationship to patients' characteristics is of utmost importance to ensure timely, effective, and optimal treatment. This is particularly relevant in areas and situations where the healthcare system must deal with an unexpected surge in hospitalizations, such as heart failure patients in the Lombardy Region of Italy during the COVID-19 pandemic. Motivated by this issue, the paper introduces a novel multilevel logistic cluster-weighted model for predicting 45-day mortality following hospitalization due to COVID-19. The methodology flexibly accommodates dependence patterns among continuous and dichotomous variables; effectively accounting for group-specific effects in distinct subgroups showing different attributes. A tailored classification expectation-maximization algorithm is developed for parameter estimation, and extensive simulation studies are conducted to evaluate its performance against competing models. The novel approach is applied to administrative data from the Lombardy Region, with the aim of profiling heart failure patients hospitalized for COVID-19 and investigating the hospital-level impact on their overall mortality. A scenario analysis demonstrates the model's efficacy in managing multiple sources of heterogeneity, thereby yielding promising results in aiding healthcare providers and policymakers in the identification of patient-specific treatment pathways.

**KEYWORDS:** cluster-weighted models; expectation-maximization algorithm; healthcare system; hierarchical data; Ising model; multilevel models.

## 1 INTRODUCTION

Medical data often present a hierarchical structure and encompass hidden clusters of patients. To evaluate the effectiveness of the healthcare system or assess the impact of a treatment or pathology on patient outcomes, classical techniques are often incapable of jointly taking into account the heterogeneity given by the admitting facility and the intrinsic characteristics of the patient. This holds significance as failing to separate the hospital effect from the individual patient's factors limits the policy's practicability. Consequently, there is a growing need to adopt advanced analytical pipelines that can effectively address these complexities. Utilizing such methodologies enables the medical field to derive deeper insights, capturing the intricate interplay of variables that collectively influence patient outcomes (Ammenwerth et al., 2003; Committee et al., 2012; Berta et al., 2016).

Driven by the potential of healthcare administrative data, this study aims to profile and categorize heart failure (HF) patients

from the Lombardy Region of Italy who were hospitalized for COVID-19. The importance of studying this population arises from the observed correlation between COVID-19 and HF, which was associated with a significantly high mortality rate during the pandemic. The coexistence of these conditions often complicates the recovery process for patients infected with the virus (Rey et al., 2020; Bader et al., 2021). Moreover, the virus itself can also trigger HF in patients who have contracted it (Adeghate et al., 2021). Consequently, the objective of the present work is multifold. First, it aims to disentangle the combined influence of hospital features and individual patient factors to improve patient profiling, particularly for vulnerable populations such as those affected by COVID-19, and to assess hospital effects accordingly. Second, the study seeks to examine and analyze the impact of respiratory illnesses on the clinical outcomes and overall well-being of patients within this demographic. This understanding will be crucial in developing tailored interventions to improve patient outcomes.

From a methodological perspective, the proposed approach originates from the general framework of mixture modeling, specifically within the realm of cluster-weighted models (CWM, Gershensfeld, 1997; Ingrassia et al., 2012, 2014). CWMs are employed to represent the joint distribution of a random vector, including a response variable and a set of covariates, in situations in which the data can be naturally divided into clusters.

We introduce ML-CWMD, a novel methodology that integrates and extends CWM with logistic mixed-effects models.

This innovative approach is designed to accommodate multiple types of dependence, capturing not only relationships among observations within the same cluster and group but also enabling the inclusion of dichotomous dependent covariates through the Ising model (Cheng et al., 2014; Ghosal and Mukherjee, 2020). This aspect is particularly impactful in the medical domain as it facilitates the analysis of relationships between the presence or absence of various diseases (referred to as comorbidities) and distinct risk factors observed in hospital-grouped patients. The novelty introduced in the paper is therefore 2-fold: First, it generalizes the multilevel logistic cluster-weighted framework (Berta and Vinciotti, 2019), contributing to the methodological advancement of hierarchical modeling in a complex setup; second, it addresses a critical public health challenge by uncovering group-specific mortality patterns that account for dependent comorbidities in COVID-19 HF patients. Notice that the efficacy of multilevel models in this context is reinforced by their successful application to diverse COVID-19 datasets, as demonstrated in recent studies (Verbeek et al., 2023; Berta et al., 2024).

The remainder of the paper is structured as follows. In Section 2, we introduce the ML-CWMD model. In Section 3, we outline the algorithm necessary for model estimation, discussing inferential aspects and model selection criteria. In Section 4, we showcase the application of our methodology to the Lombardy Region healthcare administrative data, including a comprehensive scenario analysis. A discussion follows in Section 5. Additionally, an extensive simulation study has been conducted and is available in Web Appendix A.

The proposed model has been implemented in R (Team, 2024), and the developed routines are freely available. In addition, we have developed a Shiny app that allows clinicians to access model estimates by uploading their patients' data. The application automatically generates visualizations of results and predictions, enhancing usability and insights while serving as a promising off-the-shelf tool for seamlessly integrating the ML-CWMD modeling approach into clinical practice. For more details on the developed routines and the Shiny app, refer to the Supplementary Materials.

## 2 METHODOLOGY

Consider a dichotomous response variable  $\mathbf{Y}$  and a set of covariates  $\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{D})$ , where  $\mathbf{U}$  represents a  $p$ -dimensional vector of continuous variables,  $\mathbf{V}$  represents a  $q$ -dimensional vector of categorical variables, and  $\mathbf{D}$  represents an  $h$ -dimensional vector of dichotomous variables that may possess some degree of dependence. We consider data points grouped in a known

2-level hierarchy, where units  $i$ , for  $i = 1, \dots, n_j$ , are nested within groups  $j$ , for  $j = 1, \dots, J$ .  $\mathbf{X}$  and  $\mathbf{Y}$  are defined in a finite space  $\Omega$  with values in  $\mathbb{R}^{(p+q+h)} \times \{0, 1\}$ , which is assumed to be divided into  $C$  clusters, irrespective of the known hierarchy, denoted as  $\Omega_1, \dots, \Omega_C$ . Given this setup, we assume that the joint probability across the clusters can be factorized as follows:

$$p((\mathbf{x}, \mathbf{y})|\theta) = \sum_{c=1}^C p(\mathbf{y}|\mathbf{x}, \xi_c)\phi(\mathbf{u}|\mu_c, \Sigma_c) \psi(\mathbf{v}|\lambda_c)\zeta(\mathbf{d}|\Gamma_c, \nu_c)w_c. \quad (1)$$

In detail,  $\theta$  denotes the vector encompassing all model parameters. The mixing proportion of observations within cluster  $c$  is indicated with  $w_c$  ( $w_c > 0, \forall c = 1, \dots, C, \sum_{c=1}^C w_c = 1$ ). Continuous covariates  $\mathbf{U}$  are modeled as a multivariate normal density  $\phi(\cdot|\mu_c, \Sigma_c)$ , with cluster-wise different mean vectors  $\mu_c$  and covariance matrices  $\Sigma_c$ . The categorical covariates  $\mathbf{V}$  are assumed independent and distributed as a  $q$ -variate multinomial distribution  $\psi(\cdot|\lambda_c)$  with cluster-wise different parameter vectors  $\lambda_c$ . The dichotomous dependent covariates  $\mathbf{D}$  follow the Ising model presented in Ghosal and Mukherjee (2020), characterized by distribution  $\zeta(\cdot|\Gamma_c, \nu_c)$ , with cluster-wise different threshold vectors  $\nu_c$  and interaction matrices  $\Gamma_c$ .  $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{D}$  are assumed locally independent within clusters. For each  $c$ ,  $p(\mathbf{y}|\mathbf{x}, \xi_c)$  represents a multilevel logistic regression model where  $\xi_c$  are the cluster-specific parameters for both fixed and random effects. Without loss of generality, we hereafter discuss the setting in which covariates in the conditional and marginal distributions of (1) are identical, although in practical scenarios they might vary partially or entirely (see Section 4). This variation corresponds to what is commonly referred to in the literature as mixtures of regression models with concomitant variables (Dayton and Macready, 1988). In each cluster, for every observation  $i$  within group  $j$ , the conditional distribution takes the form  $p(y_{ij}|\mathbf{x}_{ij}, \xi_c) = [\pi_{ij}]^{y_{ij}} [1 - \pi_{ij}]^{1-y_{ij}}$ . Let  $\boldsymbol{\pi}_j = (\pi_{1j}, \dots, \pi_{n_jj})$  be the vector of conditional probabilities relative to group  $j$ . The mixed effects logistic regression model takes on the following form:

$$\text{logit}(\boldsymbol{\pi}_j) = \mathbf{F}_j \boldsymbol{\beta}_c + \mathbf{R}_j \mathbf{b}_{j,c},$$

$$\boldsymbol{\pi}_j = \frac{\exp\{\mathbf{F}_j \boldsymbol{\beta}_c + \mathbf{R}_j \mathbf{b}_{j,c}\}}{1 + \exp\{\mathbf{F}_j \boldsymbol{\beta}_c + \mathbf{R}_j \mathbf{b}_{j,c}\}}, \quad (2)$$

where  $\text{logit}(\boldsymbol{\pi}_j)$  represents the odds associated with group  $j$ , indicating the probability of an event occurring divided by the probability of not occurring. The terms  $\mathbf{b}_{j,c} \sim \mathcal{N}(\mathbf{0}, \mathbf{D}_c)$ ,  $j = 1, \dots, J$ , represent the random effects for group  $j$  in cluster  $c$ . The vector  $\boldsymbol{\beta}_c$  comprises the coefficients corresponding to the  $c$ -th cluster fixed effects. Matrices  $\mathbf{F}_j$  and  $\mathbf{R}_j$  represent the design matrices associated with fixed effects and random effects in group  $j$ , respectively. We denote the design matrix encompassing all types of covariates as  $\mathbf{X}_j = (\mathbf{F}_j \cup \mathbf{R}_j)$ . All specifications regarding the other terms outlined in Equation 1 will be detailed in subsequent sections.

### 2.1 Model for categorical covariates $V$

Consider  $q$  categorical covariates indexed by  $r = 1, \dots, q$ . In each cluster, we model the categorical covariates  $\mathbf{V}$  with  $q$  independent multinomial distributions of parameter  $\lambda_{cr} = (\lambda_{cr1}, \dots, \lambda_{crs}, \dots, \lambda_{crk_r})$  ( $c = 1, \dots, C, r = 1, \dots, q$ ), identifying the vector of probabilities for the  $k_r$  categories of the  $r$ th covariate. Here, we are making the assumption that every categorical variable takes on the  $k_r$  categories across clusters  $c = 1, \dots, C$ . In this case, we can represent each of these covariates using a binary vector  $\mathbf{v}^r = (v^{r1}, \dots, v^{rk_r})$  such that  $v^{rs} = 1$  if  $v^r = s$  (where  $s \in \{1, \dots, k_r\}$ ), and 0 otherwise (Ingrassia et al., 2015). Therefore, the density  $\psi$ , given by the product of  $q$  conditionally independent multinomial distributions, can be written as

$$\psi(\mathbf{v}|\lambda_c) = \prod_{r=1}^q \prod_{s=1}^{k_r} \lambda_{crs}^{v^{rs}}, \quad c = 1, \dots, C,$$

where  $\lambda_c = (\lambda_{c1}, \dots, \lambda_{cr}, \dots, \lambda_{cq})$ . (3)

### 2.2 Model for dichotomous dependent covariates $D$

Consider  $h$  dichotomous variables indexed by  $l = 1, \dots, h$ , which may possess some degree of dependence, and let  $\mathbf{d} = (d_1, \dots, d_h)$  be an  $h$ -dimensional binary random vector. The corresponding probability is specified by the Ising model with cluster-specific parameters  $\Gamma_c$  and  $\mathbf{v}_c$ :

$$\zeta(\mathbf{d}|\Gamma_c, \mathbf{v}_c) = \frac{1}{S(\Gamma_c, \mathbf{v}_c)} \exp\left(\frac{1}{2} \mathbf{d}^T \Gamma_c \mathbf{d} + \mathbf{d}^T \mathbf{v}_c\right),$$

$c = 1, \dots, C,$  (4)

$$\mathbf{v}_c = [v_{c1}, \dots, v_{ch}], \quad \Gamma_c = \begin{bmatrix} 0 & \gamma_{c12} & \dots & \gamma_{c1h} \\ \gamma_{c21} & 0 & \dots & \gamma_{c2h} \\ \vdots & \vdots & \ddots & \vdots \\ \gamma_{ch1} & \gamma_{ch2} & \dots & 0 \end{bmatrix}.$$

$\Gamma_c \in \mathbb{R}^{h \times h}$  is a symmetric matrix, and  $S$  is the normalizing constant. The interaction parameter  $\Gamma_c$  depicts the interrelationships between all pairs of binary variables. Meanwhile, the threshold parameter  $\mathbf{v}_c$  highlights the propensity of a variable to lean toward one state or the other when all interaction parameters related to that particular variable are equal to zero. The Ising model has formulations in 2 domains ( $\{0,1\}$  and  $\{-1,1\}$ ), each offering a different interpretation of interaction parameters. Within the  $\{0, 1\}$  domain, augmenting the interaction parameter between 2 variables (eg,  $\gamma_{c12}$ ) results in an increased likelihood of the state (1, 1) when compared to the other feasible states: (0, 0), (0, 1), and (1, 0). Further details on the formulation of the Ising model can be found in [Web Appendix B](#).

### 2.3 Likelihood of the model

Consider a sample of  $N = n_1 + \dots + n_j$  observation pairs  $\{(\mathbf{x}_{ij}, y_{ij})\}_{j=1, \dots, J, i=1, \dots, n_j}$  drawn from model (1). To express the likelihood, we introduce a latent variable  $\mathbf{z}$ , where  $z_{ijc} = 1$  if  $(\mathbf{x}_{ij}, y_{ij})$  belongs to the  $c$ th cluster, and 0 otherwise. Also, we make the assumption of independence between observations belonging to different clusters and independence between observations within the same cluster but belonging to different groups (highest-level units). With a slight abuse of notation, we use the

vector  $\mathbf{z}_{jc}$  to indicate that only observations within group  $j$  that are part of cluster  $c$  contribute to the likelihood associated with cluster  $c$  and thus we express the complete likelihood and log-likelihood, respectively, as

$$L((\mathbf{x}, \mathbf{y}, \mathbf{z}); \theta) = \prod_{c=1}^C \prod_{j=1}^J \left[ p(\mathbf{y}_j|\mathbf{x}_j, \xi_c) \phi(\mathbf{u}_j|\mu_c, \Sigma_c) \psi(\mathbf{v}_j|\lambda_c) \zeta(\mathbf{d}_j|\Gamma_c, \mathbf{v}_c) w_c \right]^{z_{jc}},$$

$$l((\mathbf{x}, \mathbf{y}, \mathbf{z}); \theta) = \sum_{c=1}^C \sum_{j=1}^J z_{jc} \log(p(\mathbf{y}_j|\mathbf{x}_j, \xi_c) \phi(\mathbf{u}_j|\mu_c, \Sigma_c) \psi(\mathbf{v}_j|\lambda_c) \zeta(\mathbf{d}_j|\Gamma_c, \mathbf{v}_c) w_c). \quad (6)$$

The log-likelihood can be factorized as

$$l((\mathbf{x}, \mathbf{y}, \mathbf{z}); \theta) = \sum_{c=1}^C \sum_{j=1}^J z_{jc} \log(p(\mathbf{y}_j|\mathbf{x}_j, \xi_c))$$

$$+ \sum_{c=1}^C \sum_{j=1}^J \sum_{i=1}^{n_j} z_{ijc} \log(\phi(\mathbf{u}_{ij}|\mu_c, \Sigma_c))$$

$$+ \sum_{c=1}^C \sum_{j=1}^J \sum_{i=1}^{n_j} z_{ijc} \log(\psi(\mathbf{v}_{ij}|\lambda_c))$$

$$+ \sum_{c=1}^C \sum_{j=1}^J \sum_{i=1}^{n_j} z_{ijc} \log(\zeta(\mathbf{d}_{ij}|\Gamma_c, \mathbf{v}_c))$$

$$+ \sum_{c=1}^C \sum_{j=1}^J \sum_{i=1}^{n_j} z_{ijc} \log(w_c), \quad (7)$$

since we take into account the dependence between observations in the same cluster and group only in the regression component. In particular, for the regression component, we assume independence

- between observations originating from different clusters;
- between observations originating from the same cluster but belonging to different groups.

The first assumption offers the advantage of simplifying the log-likelihood form, facilitating its maximization as elaborated in the next section. However, this assumption comes with a drawback, that is, it precludes the model from capturing potential associations among observations belonging to the same group (eg, hospital) but to different clusters. In the healthcare context, while it may seem reasonable to assume that individuals treated in the same hospital are interconnected, it is important to recognize that patients assigned to different clusters often exhibit significantly different health statuses (see Section 4), reducing the relevance of hospital-induced dependence. Conversely, the associations among individuals treated at the same hospital and within the same cluster are effectively captured through the random effects incorporated into the regression component of the model.

### 3 MODEL ESTIMATION

To estimate the model parameters, we propose a tailored classification expectation-maximization (CEM) algorithm whose theoretical guarantees of the monotonic increase of the objective function at each iteration and its convergence to a stationary point have been proved in Celeux and Govaert (1992). In the following, we present the E, C, and M steps. All computational details concerning the derivation of the updates are provided in Web Appendix C. Lastly, additional details regarding the estimation procedure such as initialization, convergence, model selection, and prediction are discussed in Section 3.4.

$$\begin{aligned} \mathbb{E}[z_{ijc} | (\mathbf{x}, \mathbf{y}), \hat{\boldsymbol{\theta}}^{(k)}] &= \mathbb{P}[z_{ijc} = 1 | (\mathbf{x}_{ij}, y_{ij}), \hat{\boldsymbol{\theta}}^{(k)}] = \hat{\tau}_{ijc}^{(k+1)}((\mathbf{x}_{ij}, y_{ij}), \hat{\boldsymbol{\theta}}^{(k)}) \\ &= \frac{p(y_{ij} | \mathbf{x}_{ij}, \hat{\boldsymbol{\xi}}_c^{(k)}) \phi(\mathbf{u}_{ij} | \hat{\boldsymbol{\mu}}_c^{(k)}, \hat{\boldsymbol{\Sigma}}_c^{(k)}) \psi(\mathbf{v}_{ij} | \hat{\boldsymbol{\lambda}}_c^{(k)}) \zeta(\mathbf{d}_{ij} | \hat{\boldsymbol{\Gamma}}_c^{(k)}, \hat{\mathbf{v}}_c^{(k)}) \hat{w}_c^{(k)}}{\sum_{c=1}^C p(y_{ij} | \mathbf{x}_{ij}, \hat{\boldsymbol{\xi}}_c^{(k)}) \phi(\mathbf{u}_{ij} | \hat{\boldsymbol{\mu}}_c^{(k)}, \hat{\boldsymbol{\Sigma}}_c^{(k)}) \psi(\mathbf{v}_{ij} | \hat{\boldsymbol{\lambda}}_c^{(k)}) \zeta(\mathbf{d}_{ij} | \hat{\boldsymbol{\Gamma}}_c^{(k)}, \hat{\mathbf{v}}_c^{(k)}) \hat{w}_c^{(k)}}, \end{aligned} \quad (8)$$

which corresponds to the posterior probability that observation  $(\mathbf{x}_{ij}, y_{ij})$  belongs to the  $c$ -th cluster, using the current value of  $\boldsymbol{\theta}$ , that is,  $\hat{\boldsymbol{\theta}}^{(k)}$  (Ingrassia et al., 2015).

#### 3.2 C-step

The classification step involves *hard-assigning* each observation to a cluster according to the expectation calculated in Equation 8. We compute the indicator variables  $z_{ijc}^{(k+1)}$  as follows:

$$z_{ijc}^{(k+1)} = \begin{cases} 1 & \text{if } c = \operatorname{argmax}_{t \in \{1, \dots, C\}} \hat{\tau}_{ijt}^{(k+1)}((\mathbf{x}_{ij}, y_{ij}), \hat{\boldsymbol{\theta}}^{(k)}), \\ 0 & \text{otherwise} \end{cases}, \quad (9)$$

where each observation is uniquely allocated to the cluster with the highest estimated posterior probability.

#### 3.3 M-step

During the M-step, in the  $(k+1)$ th iteration, the goal is to maximize the conditional expectation of  $l((\mathbf{x}, \mathbf{y}, \mathbf{z}); \boldsymbol{\theta})$  given the observed data, in which  $z_{ijc}$  is replaced by  $z_{ijc}^{(k+1)}$ . See Web Appendix C for the derivation of all formulas in the M-step.

#### 3.4 Additional details on the estimation procedure

Initialization is a crucial factor for any deterministic algorithm. Although various methods have been suggested in the literature (for an in-depth discussion, see, eg, Biernacki et al., 2003; Karlis and Xekalaki, 2003), empirical evaluation suggests that using multiple random initializations is both reliable and not overly computationally expensive in our context. Thus, we initialize the CEM algorithm by assigning observations to latent clusters randomly, repeating this process several times, and selecting the result with the highest maximized log-likelihood for each potential number of clusters  $C \in \{1, \dots, C_{max}\}$ , with  $C_{max}$  being a pre-specified upper bound for the number of sought clusters. Convergence is monitored by looking at the relative increase of the objective function into 2 consecutive iterations. Lastly, model selection is performed using the Bayesian information criterion (BIC, Schwarz, 1978) to identify the best  $C$  in a data-driven fashion.

#### 3.1 E-step

At the  $(k+1)$ th iteration, the E-step involves computing the expected value of the log-likelihood  $l((\mathbf{x}, \mathbf{y}, \mathbf{z}); \boldsymbol{\theta})$  defined in Equation 7, given the observed data and the previous estimate  $\hat{\boldsymbol{\theta}}^{(k)}$ , obtained from iteration  $k$ . Therefore, the E-step entails calculating the expectation of the random variables  $Z_{ijc}$  associated to  $z_{ijc}$ . In particular, for  $j = 1, \dots, J, i = 1, \dots, n_j, c = 1, \dots, C$ , we have

In detail, the BIC for ML-CWMD is obtained as

$$\begin{aligned} \text{BIC} &= -2 \cdot \ln(L) + k \cdot \ln(N), \\ k &= \underbrace{C \cdot [1 + m]}_{k_{\text{reg}}} + \underbrace{C \cdot \left[ \frac{p(p+3)}{2} \right]}_{k_{\text{cont}}} \\ &+ \underbrace{C \cdot \sum_{r=1}^q (k_r - 1)}_{k_{\text{cat}}} + \underbrace{C \cdot \left[ \frac{l(l+1)}{2} \right]}_{k_{\text{dich}}} + \underbrace{C - 1}_{k_{\text{weights}}}, \end{aligned} \quad (10)$$

where  $\ln(L)$  is the maximized log-likelihood of the model, and  $m$  is the number of regression parameters. Additionally, the formula for generating a prediction for a new observation with our model, whose complete derivation is detailed in Web Appendix B, is as follows:

$$\begin{aligned} \mathbb{E}[y | \mathbf{x}; \boldsymbol{\theta}] &= p(y = 1 | \mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{\sum_{c=1}^C \phi(\mathbf{u} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \psi(\mathbf{v} | \boldsymbol{\lambda}_c) \zeta(\mathbf{d} | \boldsymbol{\Gamma}_c, \mathbf{v}_c) w_c \cdot p(y = 1 | \mathbf{x}, \boldsymbol{\xi}_c)}{\sum_{c=1}^C \phi(\mathbf{u} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \psi(\mathbf{v} | \boldsymbol{\lambda}_c) \zeta(\mathbf{d} | \boldsymbol{\Gamma}_c, \mathbf{v}_c) w_c}. \end{aligned} \quad (11)$$

### 4 ML-CWMD FOR PROFILING COVID-19 HF PATIENTS AND EVALUATING HEALTH PROVIDERS

The proposed ML-CWMD model is now employed to profile HF patients affected by COVID-19 and evaluate the heterogeneity of healthcare providers effects on different risk groups. A patient is classified as HF if their records indicate hospitalizations or contacts with emergency room under diagnosis-related group (DRG) code 127. The DRG 127 corresponds to the definition “Heart Failure and Shock” according to the version of the DRG grouping system used in the administrative database of the Lombardy Region, which follows the MS-DRG v40 version. This code is assigned to hospital discharges where the primary or secondary diagnoses (ICD-9-CM) indicate HF (ICD-9-CM: 428.\*) or related conditions such as hypertension with HF (ICD-9-CM: 402.01, 402.11, and 402.91), according to the coding criteria adopted by the region. We focus on HF patients who

have been hospitalized due to COVID-19 and who appear in the healthcare administrative database of the Lombardy Region in Italy. The purpose of this study is 2-fold. On one hand, our aim is to enhance comprehension of the characteristics and requirements of HF patients in the region, particularly amidst the exceptional circumstances following the onset of the COVID-19 pandemic. On the other hand, we aim to investigate the capabilities of the proposed model in capturing heterogeneous dynamics within hierarchical data. Specifically, by applying the ML-CWMD model, we explore the presence of latent HF COVID-19 patient subpopulations, hereafter called *patient profiles*, that display unique sets of attributes and intrinsic characteristics. Next, we evaluate and rank hospitals based on their performance in treating patients with different profiles. Finally, we deepen our understanding on the association between respiratory illnesses and HF COVID-19 patients death risk across different patient profiles. This might offer valuable insights into the complex interplay between patient characteristics, hospital dynamics, and the impact of respiratory illnesses in the context of COVID-19 hospitalizations, ultimately resulting in actionable margins for defining healthcare interventions to enhance the management of patients with HF.

#### 4.1 Dataset description

From the healthcare administrative records of the [Lombardy Region](#), we identify patients who have been diagnosed with HF (as described in the previous paragraph) in the period between January 1, 2018 and December 31, 2020. Among them, we select the ones who have been hospitalized with COVID-19 in the Lombardy Region between January 31, 2020 and June 18, 2021. The hospitalization for COVID-19 is identified through an explicit “COVID-19 flag” in the database. This flag is based on the coding criteria established by the Lombardy Region during the pandemic. The sample contains 3193 HF COVID-19 patients within 32 distinct hospitals. Hospitals with less than 50 cases were excluded. For each patient, we observe personal and clinical characteristics, such as

- Age: HF patient’s age at the moment of COVID-19 hospitalization.
- Sex: gender of the patient.
- Respiratory diseases: binary variables indicating whether the patient suffers or not from chronic obstructive pulmonary disease (COPD), pneumonia (PNA), respiratory failure (RF), and bronchitis (BRH) at the moment of admission.
- Modified multisource-comorbidity score (MCS): The aim of this scoring system is to offer a succinct snapshot of the patient’s clinical condition as described in Corrao et al. (2017). We customized the score to the current setting (indicating it as “Modified MCS”) excluding comorbidities related to respiratory system. In doing so, we aim at evaluating the influence of these particular diseases on a standalone basis with the adjustment in the regression model. See [Web Appendix D](#) for a detailed breakdown of patients’ Modified MCS across 6 distinct categories, reflecting varying degrees of health severity.

[Web Appendix D](#) provides a comprehensive overview of both personal and clinic variables in the dataset. Additionally, it includes a description of the clinical variables along with their corresponding ICD-9-CM codes (Romano et al., 1993). The endpoint is the overall mortality within a specific timeframe  $t$  following hospitalization for COVID-19. We identify the 45-day threshold as adequate, aligning with our research objectives. We selected this timeframe for 2 main reasons. Firstly, 89.7% of the deceased patients passed away within 45 days. Secondly, this timeframe enables us to evaluate both the initial severity of the illness and the potential influence of hospital treatment on patient outcomes. A shorter duration might not offer ample time for the effects of hospital treatment to manifest, while a longer duration could introduce more confounding variables or uncertainties. Therefore, we created a response encoded as a dichotomous variable that takes the value 1 if patient  $i$  passed away within  $t = 45$  days starting from the moment of hospital admission and 0 otherwise. More generally, given the considered time frame  $t$  from admission, the regression model will return predictions about the likelihood of a patient’s mortality within this specific time frame. This enables the creation of a risk score for each individual patient.

#### 4.2 Model setting and results

We make use of age, MCS, gender, COPD, and BRH to effectively discern latent patient profiles by including them in the set of random covariates in the CMW framework. Concerning respiratory diseases, we incorporate PNA and RF exclusively into the regression term of the model, since COPD and BRH are primary diseases, whereas PNA and RF often result from other underlying conditions. Consequently, we utilize COPD and BRH as a driver for patient profiles while examining the impact of PNA and RF on the cluster-wise different probability of patient mortality. Regarding the distributions of the covariates, we model age and MCS as continuous variables using a multivariate normal distribution, while, gender, COPD, and BRH are modeled as dichotomous dependent covariates using the novel Ising model specification introduced in the proposed ML-CWMD framework. In detail, for each patient  $i$  in hospital  $j$  and for each cluster, we consider

$$\text{logit}(\pi_{ij}) = \alpha_c + \beta_{1c} \cdot \text{PNA}_{ij} + \beta_{2c} \cdot \text{RF}_{ij} + b_{cj}, \quad (12)$$

where  $\pi_{ij}$  is the probability of that patient to die within 45 days from admission. Following a similar procedure as in Section 3.4, we fit ML-CWMD models varying  $C \in \{1, 2, 3, 4\}$ . The optimal result in terms of BIC is associated with  $C = 3$  clusters (for the complete comparison table and the visualization of the 3 clusters in the space Age-MCS, see [Web Appendix D](#)).

Results are subsequently reported in terms of cluster descriptions, fixed and random effects interpretation, and model predictions. The first cluster comprises 2722 individuals (85.2%) aged 82.52 on average with a mean MCS of 9.87, that is, older individuals with relatively limited occurrence of comorbidities. We label this profile as the cohort of elderly patients who are in a state of good health. Within this cluster, there is gender balance (males 54.6% versus females 46.4%), while patients without COPD (67.8%) and BRH (83.2%) predominate. As we can see from Figure 1A, the interaction parameters, which capture

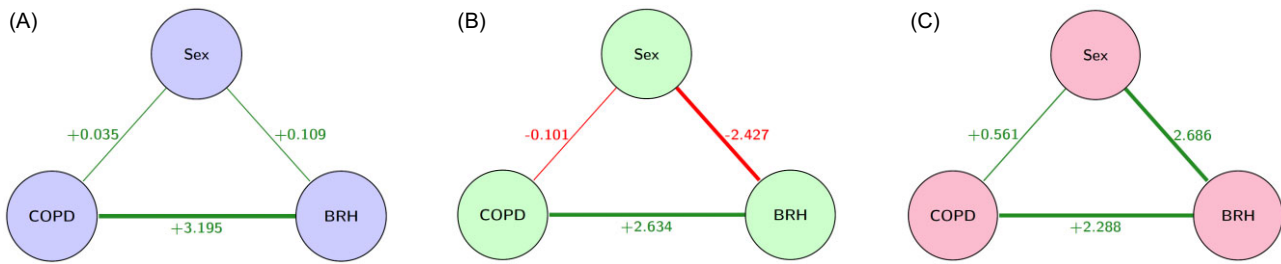


FIGURE 1 (A) Ising model interaction parameters  $\Gamma_1$  obtained from the novel multilevel logistic cluster-weighted model (ML-CWMD) estimation for the first cluster. (B) Ising model interaction parameters  $\Gamma_2$  obtained from ML-CWMD estimation for the second cluster. (C) Ising model interaction parameters  $\Gamma_3$  obtained from ML-CWMD estimation for the third cluster.

TABLE 1 Parameter values acquired through the novel multilevel logistic cluster-weighted model (ML-CWMD) estimation across the 3 clusters.

Parameter	Cluster 1		Cluster 2		Cluster 3	
$\mu$	(82.52; 9.87)		(59.67; 4.59)		(74.65; 29.81)	
$\Sigma$	$\begin{bmatrix} 51.3 & -3.8 \\ -3.8 & 36.9 \end{bmatrix}$		$\begin{bmatrix} 91.6 & -6.3 \\ -6.3 & 19.2 \end{bmatrix}$		$\begin{bmatrix} 98.6 & 14.7 \\ 14.7 & 35.9 \end{bmatrix}$	
$\nu$	(0.11; -3.41; -1.37)		(1.80; -2.10; -1.90)		(0.28; -4.86; -1.40)	
Variable	Estimate	$\Pr(>  z )$	Estimate	$\Pr(>  z )$	Estimate	$\Pr(>  z )$
PNA	+0.18	0.09	-0.08	0.21	+0.27	0.47
RF	+0.10	0.37	+1.51	0.018*	-0.90	0.046*

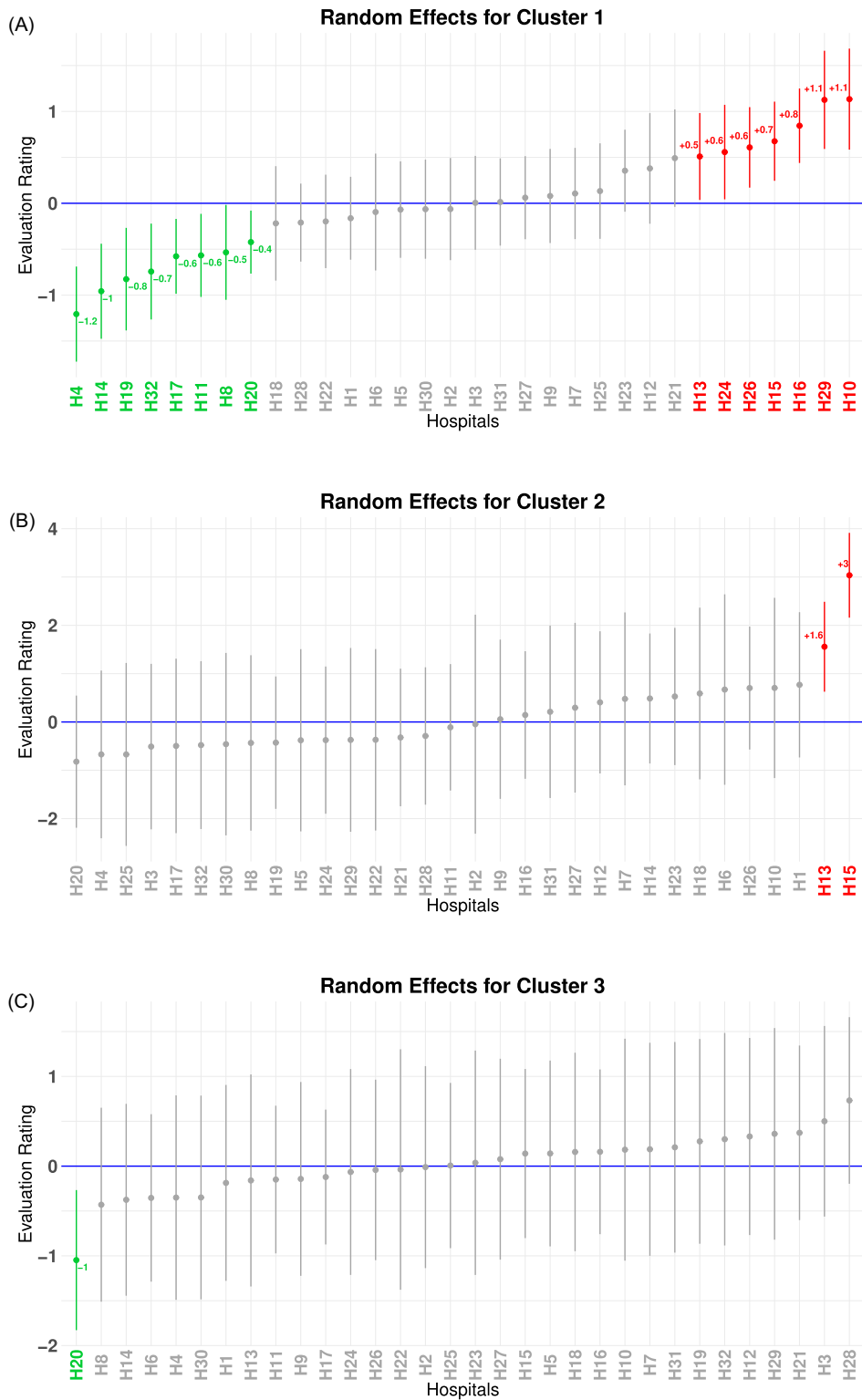
The upper section displays the values of the parameters  $\mu$ ,  $\Sigma$ , and  $\nu$ . The lower section shows the estimate of fixed effects parameters  $\beta_c$  ( $c = 1, \dots, C$ ) and their corresponding  $P$ -values.

the relationships between pairs of binary variables, reveal a positive correlation. The interaction between sex and the 2 respiratory diseases remains minimal, whereas the correlation between COPD and BRH emerges as very high. In this cluster, 35.6% of individuals passed away within 45 days of being admitted. The fixed effects estimates, reported in Table 1, reveal that pneumonia is statistically significant, indicating a higher probability of mortality among patients with this condition. When exploring the random effects related to different hospitals (Figure 2A), we discover that 15 out of 32 hospitals demonstrate a significant influence, either in reducing or increasing the likelihood of death among patients in this cluster, with respect to the average. This phenomenon can be explained since these patient profiles comprise elder individuals who are still in good health, making the appropriateness of hospital instruments and techniques particularly influential. Among the hospitals, those highlighted in green in Figure 2A are the ones that significantly lower the probability of death, and thus can be regarded as the most proficient in treating patients belonging to Cluster 1; conversely, hospitals highlighted in red are the ones that significantly increase the probability of death.

The second cluster comprises 321 patients aged 59.67 on average with a mean MCS score of 4.59. This points toward a patient profile composed of young and middle-aged individuals who exhibit favorable health conditions. Within this cluster, there is a prevalence of males (82.2%) and a prevalence of patients without COPD (84.4%) and BRH (93.5%). Figure 1B illustrates the interaction parameters, indicating a notable positive correlation between COPD and BRH, alongside a significant negative correlation between BRH and sex. This suggests a higher prevalence of females with BRH. Within a 45-day period post-hospitalization, the mortality rate among individuals in this cluster

stands at 19.3%. Upon examining the fixed effects (Table 1), we observe that the presence of RF in patients within this cluster substantially increases their risk of mortality. This correlation likely arises from the fact that RF implies impaired lung function, resulting in diminished oxygen levels in the bloodstream. Given that COVID-19 primarily targets the respiratory system, the co-existence of these 2 conditions further exacerbates respiratory dysfunction, making it challenging for the body to maintain adequate oxygen levels. Consequently, despite their overall good health, the severe oxygen deficiency can lead to exceptionally severe complications. In Figure 2B, we observe that only 2 hospitals exhibit significant associations with an increased likelihood of death in these patients. The limited impact of hospital variability on Cluster 2 is likely attributable to the fact that patients within this profile are generally younger and consistently healthier than the rest of the cohort, making them less susceptible to the overall quality of care provided by the treatment facility.

The third cluster comprises a total of 150 patients aged 74.65 on average along with a mean MCS score of 29.81. This patient profile is predominantly composed of individuals characterized by advanced age and a significant amount of comorbidities. Within this cluster, there is a prevalence of males (67.5%) and a higher occurrence of patients without COPD (62.6%) and BRH (79.3%). Figure 1C depicts interaction parameters, revealing a substantial positive correlation between COPD and BRH, as well as a significant positive correlation between BRH and gender. This implies a greater prevalence of BRH among males within this cluster. For this cluster, the percentage of patients who died within a 45-day period following hospitalization amounts to 43.3%. Interestingly, the presence of RF appears to significantly reduce the probability of death in patients within this cluster (Table 1). This somewhat counterintuitive finding



**FIGURE 2** Random effects of the 32 hospitals in the first (A), second (B) and third (C) cluster. Hospitals increasing the probability of death are those with a confidence interval entirely above (highlighted in red), while hospitals that decrease the probability of death are those with a confidence interval entirely below 0 (highlighted in green).

TABLE 2 Mean and standard deviation for all metrics across the 20 held-out data. Bold font denotes the best results.

Evaluation type	ML-CWMD	GLMM	GLM
Train Prediction Accuracy	<b>0.660 ± 0.01</b>	0.639 ± 0.01	0.588 ± 0.003
Test Prediction Accuracy	<b>0.620 ± 0.02</b>	0.610 ± 0.01	0.588 ± 0.01
ECE	<b>0.045 ± 0.009</b>	0.048 ± 0.012	0.056 ± 0.012
Brier score	<b>0.209 ± 0.005</b>	0.212 ± 0.004	0.225 ± 0.003

is frequently observed among patients of this profile (West et al., 2014). RF acts as a protective factor; its occurrence in patients already in critical condition prompts medical professionals to prioritize their monitoring and treatment, thereby enhancing their chances of survival. Regarding the random effects (Figure 2C), only 1 hospital shows a significant decrease in the probability of death for this profile. This further suggests that the health conditions of patients in this cluster are so critical that the specific hospital where they are treated has relatively little or no impact on their outcomes.

When examining the random effects across clusters, a noteworthy trend emerges. For instance, Hospital H20 consistently retrieves positive evaluations in both Clusters 1 and 3, whereas Hospitals H13 and H15 receive unfavorable ratings in both Clusters 1 and 2. This discovery has important implications for policymaking, allowing suitable monitoring and second-level analysis of “out of control” situations.

Finally, we evaluate the predictive accuracy and calibration of the proposed model comparing it with the generalized linear mixed effects model (GLMM) and the generalized linear model (GLM) under the same regression settings using held-out data. Specifically, we generate 20 training-test splits with an 80-20 proportion stratified by hospitals. The models are trained on the training data and then used to predict death outcomes for patients in the test set. The results, including training and test predictive accuracy, expected calibration error (ECE), and Brier score, are summarized in Table 2, reporting the mean and standard deviation across the 20 splits. The proposed model consistently outperforms the alternatives, achieving superior performance on all metrics.

### 4.3 Scenario analysis

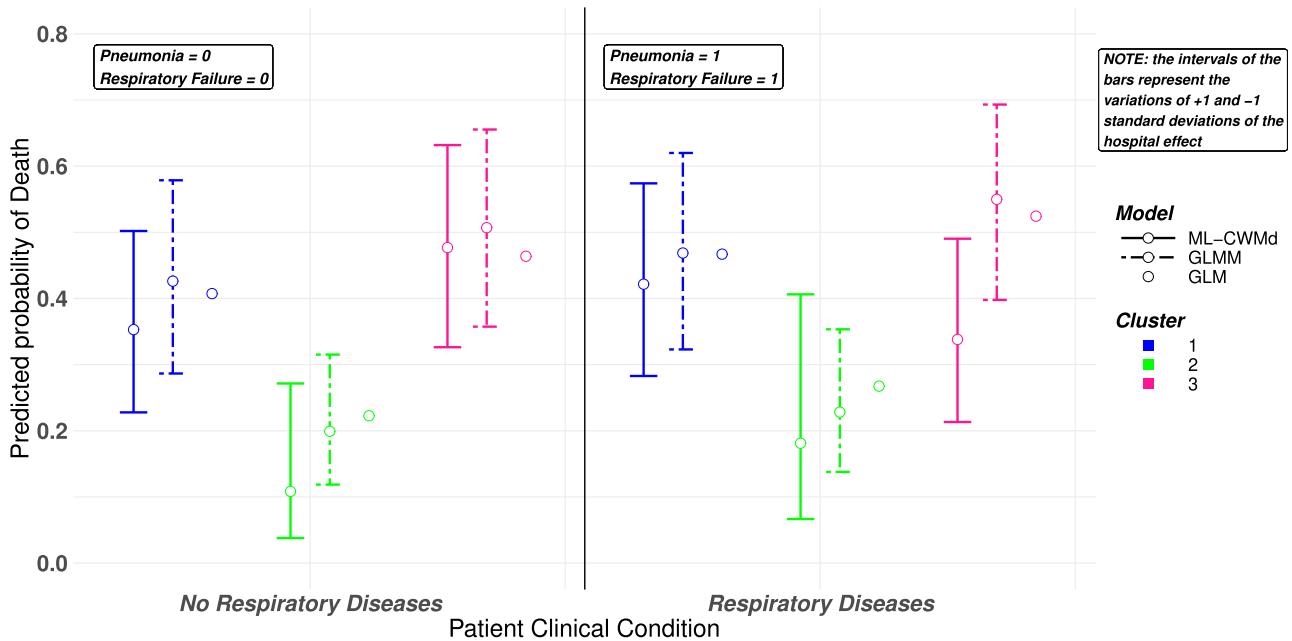
In this section, we implement a scenario analysis to emphasize the importance of considering both the cluster factor and the hospital effect to capture the real-world complexity. We define 3 new COVID-19 HF patients, with characteristics closely mirroring those identified within the 3 clusters. From each of these patients, we select one with respiratory diseases (PNA and RF) and one without them, resulting in 6 possible profiles. For the complete table, including patients’ characteristics, see [Web Appendix D](#). The objective is to illustrate the difference between predicting the likelihood of mortality for these new patients using the proposed model (Equation 1), a GLMM that neglects the consideration of latent patient clusters, and a GLM that neglects both the consideration of latent patient cluster and hospital effect. For each patient, we delineate 3 distinct predictions, corresponding, respectively, to random effects equivalent to  $-\hat{\sigma}_{bc}$  (hospital with commendable performance concerning that spe-

cific patient profile), 0 (hospital with negligible influence on the outcome), and  $+\hat{\sigma}_{bc}$  (hospital with poor performance).

In Figure 3, we show the predictions obtained with the 3 models. For the GLM, we have only a single value since it does not take into account the hospital effect. The bar intervals on the graph show how the predicted probabilities of death change when the hospital is considered beneficial for patient characteristics (at  $-\hat{\sigma}_{bc}$ ) versus when it is considered detrimental (at  $+\hat{\sigma}_{bc}$ ). The GLMM generates highly similar predictions across the same patient with and without respiratory diseases because it does not identify any respiratory diseases as significant covariates. Notably, only age, sex, and MCS emerge as significant variables in the fixed regression component of the GLMM. Consequently, GLMM fails to capture the variability attributed to respiratory diseases, which could significantly influence a patient probability of death, particularly among COVID-19 HF patients. Instead, using the ML-CWMD enables us to take into consideration the respiratory diseases effects observing a significant change in the probability of death based on the presence or absence of these diseases in the ML-CWMD predictions. For instance, for the patients in the green cluster (young individuals with low MCS), the probability of death in a hospital characterized as detrimental for the patient characteristics is approximately 0.28 if the patient has no respiratory diseases. This probability increases to around 0.4 if the patient has both respiratory diseases. In conclusion, the ML-CWMD provides more precise estimations of patient survival probabilities across different risk categories. Additionally, it offers a valid tool for the assessment and monitoring of hospital facilities utilizing administrative databases.

## 5 DISCUSSION

This paper has introduced a novel methodology for simultaneously risk-stratifying patients and conducting cluster-specific hospital evaluations. In detail, a novel ML-CWMD is devised, extending previous works with the ability to effectively capture the dependence among observations within the same cluster and hierarchy; as well as among dichotomous variables through the inclusion of the Ising model contribution. Resorting to maximum likelihood estimation, we have implemented a tailored CEM algorithm to perform model fitting, testing its performance in a simulated setting, and comparing it with state-of-the-art alternatives. Our proposal has demonstrated promising results when dealing with complex scenarios encompassing latent clusters of observations, group effects, and the interdependence among binary covariates. This research has been motivated by the challenge of developing tailored models needed for accommodating diverse patient profiles and hospital-specific effects. Specif-



**FIGURE 3** Predicted probabilities of mortality for the 6 patients across the 3 competing models. Bar intervals illustrate the shift in predicted probabilities when the hospital is deemed advantageous versus detrimental. Solid bars represent the novel multilevel logistic cluster-weighted model (ML-CWMD) introduced in the paper, dotted bars denote GLMM, and single points indicate GLM.

ically, we have applied our proposal to a real-world administrative dataset of Lombardy Region, Italy, including information about HF patients hospitalized for COVID-19. The analysis has revealed the existence of 3 distinct patient profiles, each characterized by cluster-wise different survival patterns and comorbidities. On top of this, the model setting has allowed for valuable insights into the ways respiratory diseases and hospitals impact individual profiles of patients. The analysis has thus demonstrated promising results in terms of actionable margins for defining healthcare interventions to enhance the territorial management of patients with HF through the planning of optimal care pathways, thereby reducing adverse clinical outcomes and improving system efficiency.

The devised methodology also possesses limitations. The independence assumption among observations belonging to the same known hierarchy but different latent clusters may not always be tenable, as the grouping effect could potentially exhibit shared patterns across clusters. Furthermore, the restriction of modeling dependence among dichotomous covariates only was solely motivated by the type of covariates available within the Lombardy Region database. Specifically, the considered Ising model represents the simplest form of Markov random fields (Kindermann and Snell, 1980), and future methodological advancements could indeed extend the current proposal to allow for more general higher-order interactions. In addition, tackling the well-known over-parameterization issue associated with multiple continuous covariates could improve both the flexibility and adaptability of the proposed approach. Solutions to this issue are proposed in Banfield and Raftery (1993), Celeux and Govaert (1995), and Murphy and Murphy (2020), among others. Lastly, one aspect not addressed in this work is the evaluation of uncertainty associated with parameter estimates. While non-

differentiability issues may prevent the implementation of an information matrix-based approach in our setting, sampling-based methods provide an alternative and generalizable solution for estimating standard errors. Jackknife and bootstrap techniques could be adapted to our problem, similar to the approaches proposed by O'Hagan et al. (2019) and Berta and Vinciotti (2019) in the context of Gaussian mixtures and ML-CWMDs, respectively. Future research will extend the current proposal to address its limitations and enhance its ability to model time-to-event outcomes, with several proposals already under study.

## ACKNOWLEDGMENTS

This work is part of the ENHANCE-HEART project: Efficacy evaluation of the therapeutic-care pathways, of the healthcare providers effects, and of the risk stratification in patients suffering from HEART failure. The authors thank the "Unit" Organizzazione Osservatorio Epidemiologico Regionale and ARIA S.p.A for providing data and technological support. The authors gratefully acknowledge the support from the Department of Mathematics of Politecnico di Milano, which facilitated this research as part of the department's activities of "Dipartimento di Eccellenza 2023-2027."

## SUPPLEMENTARY MATERIALS

Supplementary material is available at *Biometrics* online.

**Web Appendices** referenced in Sections 1-4 are available with this paper at the Biometrics website on Oxford Academic. These include (A) a comprehensive simulation study, (B) additional details on the Ising Model, (C) complete mathematical derivations for updates in the M step, and (D) further details on the

ML-CWMD for COVID-19 heart failure profiling. The proposed model has been implemented in R (Team, 2024), and the developed routines are freely available at the Biometrics website on Oxford Academic and at the GitHub repository <https://github.com/luca2245/ML-CWMD>. In addition, we have developed a Shiny app that is available at <https://mcdhmq-luca-caidera.shinyapps.io/app-ml-cwmd/>.

## FUNDING

Chiara Masci acknowledges financial support from the Italian Ministry of University and Research (MUR) under the Department of Excellence 2023-2027 grant agreement 'Centre of Excellence in Economics and Data Science' (CEEDS).

## CONFLICT OF INTEREST

None declared.

## DATA AVAILABILITY

The Lombardy Region dataset analyzed in Section 4 contains sensitive information and cannot be shared due to privacy and confidentiality considerations.

## REFERENCES

- Adeghate, E. A., Eid, N. and Singh, J. (2021). Mechanisms of COVID-19-induced heart failure: a short review. *Heart Failure Reviews*, 26, 363–369.
- Ammenwerth, E., Gräber, S., Herrmann, G., Bürkle, T. and König, J. (2003). Evaluation of health information systems—problems and challenges. *International Journal of Medical Informatics*, 71, 125–135.
- Bader, F., Manla, Y., Atallah, B. and Starling, R. C. (2021). Heart failure and COVID-19. *Heart Failure Reviews*, 26, 1–10.
- Banfield, J. D. and Raftery, A. E. (1993). Model-based Gaussian and non-Gaussian clustering. *Biometrics*, 49, 803.
- Berta, P., Ingrassia, S., Punzo, A. and Vittadini, G. (2016). Multilevel cluster-weighted models for the evaluation of hospitals. *Metron*, 74, 275–292.
- Berta, P., Ingrassia, S., Vittadini, G. and Spinelli, D. (2024). Latent heterogeneity in COVID-19 hospitalisations: a cluster-weighted approach to analyse mortality. *Australian & New Zealand Journal of Statistics*, 66, 1–20.
- Berta, P. and Vinciotti, V. (2019). Multilevel logistic cluster-weighted model for outcome evaluation in health care. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 12, 434–443.
- Biernacki, C., Celeux, G. and Govaert, G. (2003). Choosing starting values for the EM algorithm for getting the highest likelihood in multivariate Gaussian mixture models. *Computational Statistics & Data Analysis*, 41, 561–575.
- Celeux, G. and Govaert, G. (1992). A classification EM algorithm for clustering and two stochastic versions. *Computational Statistics & Data Analysis*, 14, 315–332.
- Celeux, G. and Govaert, G. (1995). Gaussian parsimonious clustering models. *Pattern Recognition*, 28, 781–793.
- Cheng, J., Levina, E., Wang, P. and Zhu, J. (2014). A sparse Ising model with covariates. *Biometrics*, 70, 943–953.
- Committee, C.-C. W. P. et al. (2012). Statistical issues in assessing hospital performance. *The Committee of the Presidents of Statistical Societies*.
- Corrao, G., Rea, F., Di Martino, M., De Palma, R., Scodotto, S., Fusco, D. et al. (2017). Developing and validating a novel multisource comorbidity score from administrative data: a large population-based cohort study from Italy. *BMJ Open*, 7, e019503.
- Dayton, C. M. and Macready, G. B. (1988). Concomitant-variable latent-class models. *Journal of the American Statistical Association*, 83, 173–178.
- Gershensfeld, N. (1997). Nonlinear inference and cluster-weighted modeling. *Annals of the New York Academy of Sciences*, 808, 18–24.
- Ghosal, P. and Mukherjee, S. (2020). Joint estimation of parameters in Ising model. *The Annals of Statistics*, 48, 785–810.
- Ingrassia, S., Minotti, S. C. and Punzo, A. (2014). Model-based clustering via linear cluster-weighted models. *Computational Statistics & Data Analysis*, 71, 159–182.
- Ingrassia, S., Minotti, S. C. and Vittadini, G. (2012). Local statistical modeling via a cluster-weighted approach with elliptical distributions. *Journal of Classification*, 29, 363–401.
- Ingrassia, S., Punzo, A., Vittadini, G. and Minotti, S. C. (2015). The generalized linear mixed cluster-weighted model. *Journal of Classification*, 32, 85–113.
- Karlis, D. and Xekalaki, E. (2003). Choosing initial values for the EM algorithm for finite mixtures. *Computational Statistics & Data Analysis*, 41, 577–590.
- Kindermann, R. and Snell, J. L. (1980). *Markov Random Fields and their Applications*, Vol. 1, Providence, RI American Mathematical Society.
- Murphy, K. and Murphy, T. B. (2020). Gaussian parsimonious clustering models with covariates and a noise component. *Advances in Data Analysis and Classification*, 14, 293–325.
- O'Hagan, A., Murphy, T. B., Scrucca, L. and Gormley, I. C. (2019). Investigation of parameter uncertainty in clustering using a Gaussian mixture model via jackknife, bootstrap and weighted likelihood bootstrap. *Computational Statistics*, 34, 1779–1813.
- Rey, J. R., Caro-Codón, J., Rosillo, S. O., Iniesta, Á. M., Castrejón-Castrejón, S., Marco-Clement, I. et al. (2020). Heart failure in COVID-19 patients: prevalence, incidence and prognostic implications. *European Journal of Heart Failure*, 22, 2205–2215.
- Romano, P. S., Roost, L. L. and Jollis, J. G. (1993). Further evidence concerning the use of a clinical comorbidity index with ICD-9-CM administrative data. *Journal of Clinical Epidemiology*, 46, 1085–1090.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6, 461–464.
- Team R. C. (2024). *R: A Language and Environment for Statistical Computing*.
- Verbeeck, J., Faes, C., Neyens, T., Hens, N., Verbeke, G., Deboosere, P. et al. (2023). A linear mixed model to estimate COVID-19-induced excess mortality. *Biometrics*, 79, 417–425.
- West, E., Barron, D. N., Harrison, D., Rafferty, A. M., Rowan, K. and Sanderson, C. (2014). Nurse staffing, medical staffing and mortality in intensive care: an observational study. *International Journal of Nursing Studies*, 51, 781–794.