

Mixture of experts models for multilevel data: Modeling framework and approximation theory

Tsz Chai Fung^a, Spark C. Tseung^b

^a Maurice R. Greenberg School of Risk Science, Georgia State University, 35 Broad Street NW, Atlanta, GA 30303, United States

^b Department of Statistical Sciences, University of Toronto, Ontario Power Building, 700 University Avenue, 9th Floor, Toronto, ON M5G 1Z5, Canada

ARTICLE INFO

Communicated by Y. Lai

Keywords:

Artificial neural network
Crossed and nested random effects
Denseness
Mixed effects models
Universal approximation theorem

ABSTRACT

Multilevel data are prevalent in many real-world applications. However, it remains an open research problem to identify and justify a class of models that flexibly capture a wide range of multilevel data. Motivated by the versatility of the mixture of experts (MoE) models in fitting regression data, in this article we extend upon the MoE and study a class of mixed MoE (MMoE) models for multilevel data. Under some regularity conditions, we prove that the MMoE is dense in the space of any continuous mixed effects models in the sense of weak convergence. As a result, the MMoE has a potential to accurately resemble almost all characteristics inherited in multilevel data, including the marginal distributions, dependence structures, regression links, random intercepts and random slopes. In a particular case where the multilevel data is hierarchical, we further show that a nested version of the MMoE universally approximates a broad range of dependence structures of the random effects among different factor levels.

1. Introduction

1.1. Background and literature review

Mixture of experts (MoE) model, which is first introduced by Jacobs et al. [1] (see also, e.g., Jordan and Jacobs [2], McLachlan and Peel [3] for details), is a probabilistic version of neural network architecture useful for flexible regression, classification and distribution modeling, with applications to various areas including healthcare, business, social and environmental science. Readers may refer to Yuksel et al. [4], Masoudnia and Ebrahimpour [5], Nguyen and Chamroukhi [6] for the literature reviews on both the theories and applications of the MoE.

The model structure of the MoE is as follows. Suppose that we have N observations $(y, \mathbf{x}) = \{(y_i, \mathbf{x}_i)\}_{i=1, \dots, N}$, where $\mathbf{y}_i = (y_{i1}, \dots, y_{iK})$ is a K -dimensional response variable with output space $\mathcal{Y} \subseteq \mathbb{R}^K$ and $\mathbf{x}_i = (x_{i1}, \dots, x_{iP})$ are the P covariates or features with input space $\mathcal{X} \subseteq \mathbb{R}^P$. Under the MoE framework, the conditional distribution function of y_i given \mathbf{x}_i is

$$F(y_i; \boldsymbol{\alpha}, \boldsymbol{\psi}, g | \mathbf{x}_i) = \sum_{j=1}^g \pi_j(\mathbf{x}_i; \boldsymbol{\alpha}) F_0(y_i; \boldsymbol{\psi}_j | \mathbf{x}_i), \quad (1)$$

where g is the number of latent classes. Here, $\pi_j(\mathbf{x}_i; \boldsymbol{\alpha}) > 0$ is called the gating function with $\sum_{j=1}^g \pi_j(\mathbf{x}_i; \boldsymbol{\alpha}) = 1$ and parameters $\boldsymbol{\alpha}$. The left panel of Fig. 2 graphically illustrates the MoE model architecture. While the most commonly used gating function is the logit-linear or softmax

gating [1], expressed as $\pi_j(\mathbf{x}_i; \boldsymbol{\alpha}) = \exp\{\alpha_{j,0} + \boldsymbol{\alpha}_j^T \mathbf{x}_i\} / \sum_{j'=1}^g \exp\{\alpha_{j',0} + \boldsymbol{\alpha}_{j'}^T \mathbf{x}_i\}$ with $\boldsymbol{\alpha} = \{\alpha_{j,0}, \boldsymbol{\alpha}_j : j = 1, \dots, g\}$, various alternative gating functions have also been explored in the literature. These include Gaussian gating [7], student-t gating [8], probit gating [9], as well as more advanced “sparse” gating functions designed to reduce the complexity of MoE models, such as the Top- k gate [10] and the DSelect- k gate [11]. Also, $F_0(y_i; \boldsymbol{\psi}_j | \mathbf{x}_i)$ is a probability distribution called the expert function with parameters $\boldsymbol{\psi} := \{\boldsymbol{\psi}_j : j = 1, \dots, g\}$. While a common choice for the expert function is a Gaussian distribution [12], there has been substantial developments on alternative choices of expert functions to cater for various distributional characteristics such as heavy-tailedness (Laplace by Nguyen and McLachlan [13], t-distribution by Chamroukhi [14], skewed t by Chamroukhi [15] and transformed gamma by Fung et al. [16]) and discrete distributions (Poisson by Grun and Leisch [17] and Erlang Count by Fung et al. [18]).

Model flexibility is a crucial desirable property for the class of MoE, and there are extensive research work on the approximation theory for the MoE. Zeevi et al. [19] shows that the mean function of a univariate ($K = 1$) logit-gated MoE can approximate any Sobolev class functions. This result is extended by Jiang and Tanner [20], who considers the transformed Sobolev class. Without considering the convergence rate, Nguyen et al. [21] proves that the MoE mean function is dense in the class of any continuous functions without the restriction of the

* Corresponding author.

E-mail addresses: tfung@gsu.edu (T.C. Fung), spark.tseung@mail.utoronto.ca (S.C. Tseung).

Sobolev class, and Nguyen et al. [22] shows similar denseness results using a multivariate ($K > 1$) Gaussian-gated MoE.

Apart from studying the mean functions, some research works focus on conditional density approximation with respect to the Hellinger distance, Kullback–Leibler (KL) divergence, or Lebesgue space. Jiang and Tanner [23], Mendes and Jiang [24] generalize the results of Jiang and Tanner [20] by demonstrating the approximation capability of the MoE to any exponential family non-linear regression models. Norets et al. [25] shows that the logit-gated MoE with Gaussian expert function can approximate any conditional densities. Similar results are proved by Nguyen et al. [22], Norets and Pelenis [26], who consider the Gaussian gating functions. Recently, Nguyen et al. [27] proves that the class of MoE is dense in the Lebesgue space.

Another stream of distribution approximation theorems studies the denseness in the sense of Prohorov metric of weak convergence. Extending upon Tijms [28], Breuer and Baum [29] who explore denseness of finite mixture and phase-type distributions in the space of any probability distributions, Fung et al. [30] formulates the concept of “denseness” in regression settings and shows that the class of MoE is dense in the space of any regression distributions, subject to some regularity conditions such as Lipschitz continuity and distribution tightness. In contrast to other existing works on distribution approximations, the results of Fung et al. [30] are very general as: (i) they hold under a wide range of choices of expert functions (not restricted to Gaussian or other symmetric expert functions); (ii) the target distribution is not restricted to a special class (e.g., exponential family regression models).

Despite its model flexibility, the aforementioned framework implicitly assumes that input–output pairs are independent across observations. However, this assumption does not hold for multilevel data [31]. In addition to the inputs x_i , the dataset includes L levels of factors, each corresponding to a categorical variable used to group or cluster data into different units or categories. Observations that share the same unit within a factor level are likely to have unobserved common characteristics, creating clear interdependencies among them. Overlooking these dependencies can result in spurious, misleading, or biased clustering and prediction outcomes [31].

Multilevel data are prevalent across many applications. The most classic one is the school problem [32–34], where “school” and “classroom within a school” act as two levels of factors affecting the performance of a student (as an observation). Multilevel data structure can also be caused by repeated measurements collected in longitudinal studies. This is common across various areas including health (e.g., Molenberghs et al. [35]) and business (e.g., Boucher and Denuit [36]). For instance, the medical outcome of a patient or the amount of insurance claims by a policyholder are measured or collected repeatedly over time. A remarkable special case of multilevel data is the hierarchical (or nested) data, where the L factor levels can be ranked from high to low. The school problem is a clear example of hierarchical data.

A popular model to account for the interdependencies among observations is the generalized linear mixed effects model (GLMM) [33,37], which assumes that the output y_i depends on the sum of fixed effects (i.e., the impact of the inputs x_i) and random effects (i.e., the impact of the factors θ_i). To improve model flexibility or achieve specific clustering purposes, the GLMM framework is extended to a non-linear setting [38,39], formulated in a neural network structure [40] or integrated to a finite-mixture modeling framework [41,42]. Despite of the desirable properties of the MoE models leading to extensive applications, the research works of mixed effects models in the context of MoE framework are relatively scarce. Yau et al. [43] first proposes a two-component logit-gated Gaussian-expert MoE with random effects incorporated in both gating and expert functions. Ng and McLachlan [44] then formulates a general g -component mixed effect MoE with the use of logistic expert functions for binary classifications. Ng and McLachlan [45] considers a similar framework with random effects only incorporated to the expert functions. Nonetheless, all the aforementioned mixed effect MoE models only deal with a single level of random effect (i.e., $L = 1$).

1.2. Contributions of this paper

Driven by the growing popularity of MoE models, the widespread use of multilevel data, and the need to formally establish model flexibility through approximation theories, this paper introduces several innovative contributions to the theoretical modeling framework and approximation theory in the field of mixture of experts (MoE) models for multilevel data analysis.

First, we introduce the mixed MoE (MMoE) model for multilevel regression data, addressing a gap in the literature on MoE models for multilevel data. Our model incorporates multiple levels of random effects, providing a more comprehensive approach to multilevel data compared to the limited scope of existing works by Yau et al. [43], Ng and McLachlan [44], and Ng and McLachlan [45], which only consider a single level of random effects. Such a restriction is inadequate for handling complex multilevel data structures, such as the hierarchical structure seen in the school problem discussed in Section 1.1. Our model also features a reduced structure aimed at preserving parsimony and interpretability without sacrificing flexibility or approximation capability. Specifically, our model is simplified in two key ways: (i) we assume that random effects in the gating functions are shared across mixture latent classes and are independent, unlike the approach in Ng and McLachlan [44], which allows for varying and dependent random effects across latent classes; and (ii) we exclude the effects of observed inputs and random effects from the expert functions, as specified by Ng and McLachlan [44]. These reductions significantly decrease the number of model parameters and, as demonstrated in our subsequent work [46], are essential for feasible and effective model estimation, particularly with large datasets.

Second, we define the concept of denseness for mixed effects models in terms of weak convergence and demonstrate that the proposed MMoE class is dense in the space of continuous mixed effects models under certain mild regularity conditions. While universal approximation theories for standard input–output models, including MoE and neural network models, have been extensively explored in the statistics and machine learning literature (see Section 1.1), this paper represents the first attempt to formulate and prove approximation theories for multilevel models. Existing studies on mixed effects models within the MoE framework (e.g., Ng and McLachlan [44]) do not address the theoretical approximation capability of their models. This paper extends the setting in Fung et al. [30], which formulated denseness for regression distributions, to accommodate multilevel data. Our novel denseness theory not only illustrates the flexibility of the proposed model in capturing various characteristics of multilevel data, such as joint distributions, regression patterns, random intercepts, and random slopes, but also suggests that our model is parsimonious with a reduced structure. Compared to Fung et al. [30], this paper also relaxes several assumptions required for the denseness theorem, such as Lipschitz continuity and distribution tightness. Consequently, the proof techniques employed in this paper differ significantly from those in Fung et al. [30].

Third, in the context of hierarchical data, we prove that a nested version of the MMoE can effectively approximate a wide range of dependence structures between upper and lower-level factors, even when the MMoE is simplified to include only independent random effects across levels. This result is particularly relevant for applications where dependencies exist between factors, such as the impact of a classroom on a student’s performance being influenced by the school the student attends.

The focus of this paper is to formulate the MMoE model for multilevel data and theoretically justify its versatility. In a subsequent paper [46], we will address the estimation and application problems under the proposed MMoE. A stochastic variational ECM algorithm is proposed to efficiently estimate the model parameters. Also, the MMoE is applied to an automobile insurance dataset, demonstrating its ability to reasonably predict policyholders’ future claims based on their past claim histories.

Table 1
Hypothetical school problem dataset with $N = 12$ students, including scores, gender, and income.

Student	School ID	Classroom ID	Score (year 1)	Score (year 2)	Gender (Male)	Income
i	$c_1(i)$	$c_2(i)$	y_{i1}	y_{i2}	x_{i1}	x_{i2}
1	1	1	79.48	80.18	0	10.82
2	1	1	66.16	60.96	1	11.29
3	1	2	77.28	67.48	0	12.21
4	2	3	77.30	53.53	1	10.31
5	2	3	85.91	80.29	0	11.16
6	2	4	75.65	73.99	1	10.13
7	2	5	85.02	74.35	0	11.41
8	2	5	72.84	82.08	1	11.92
9	3	6	78.21	83.41	0	11.00
10	3	7	57.64	65.50	0	11.61
11	4	8	64.47	66.21	1	10.60
12	4	8	73.38	52.63	1	11.70

1.3. Structure of the paper

This paper is structured as follows. Section 2 mathematically defines a multilevel data with an example. Section 3 defines a generalized class of mixed effect models for multilevel data, which includes nearly all mixed effect models in the literature. In Section 4, we introduce the MME as a candidate class of mixed effect models to flexibly capture multilevel data. Interpretation and visualization of the proposed model are also provided. Section 5 defines “denseness” in the context of mixed effect models and proves that the MME is a universal approximator of most mixed effect models subject to some mild conditions. In Section 6, we discuss the model formulation and denseness property in a special case where the dataset is hierarchical with nested random effects. The intuitive explanations on the proof idea of the theoretical results are demonstrated in Section 7, accompanied by a numerical illustration in Section 8. The findings are summarized in Section 9, accompanying some limitations of the denseness theory in justifying the approximation capability of the proposed MME.

2. Multilevel data structure

This section mathematically defines a multilevel data and provides a hypothetical example to illustrate how multilevel dataset is structured.

The observed multilevel data structure can be mathematically represented as $D := \{(y_i, x_i, c(i))\}_{i=1, \dots, N}$, where y_i , x_i , and N are defined in Section 1.1, and $c(\cdot)$ is a function that maps each observation i to the corresponding units across all L levels of factors. Specifically, we denote $c(i) := (c_1(i), \dots, c_L(i))$, where $c_l(\cdot) : \{1, \dots, N\} \mapsto \{1, \dots, S_l\}$ is a known function that maps observation i to its corresponding level- l unit, with S_l being the number of possible level- l units for $l = 1, \dots, L$. In other words, $c_l(i)$ identifies the level- l unit or category to which observation i belongs. The box in the top left of Fig. 1 provides a visual representation of the multilevel data structure.

Example 1. We construct a small hypothetical multilevel dataset in the context of the “school problem”, consisting of $N = 12$ students, as shown in Table 1. The output of this dataset, $y_i = (y_{i1}, y_{i2})$, is bivariate ($K = 2$), with y_{i1} and y_{i2} representing the average exam scores of student i in year 1 and year 2, respectively. The dataset also includes $P = 2$ covariates $x_i = (x_{i1}, x_{i2})$ as inputs, where x_{i1} is a binary indicator for gender (1 for male) and x_{i2} is the log of household income. There are $L = 2$ levels of factors: “school” and “classroom”. The second and third columns of the table, $c_1(i)$ and $c_2(i)$, represent the identifiers for the school and classroom to which each student belongs. For instance, observation $i = 7$ corresponds to a student in classroom 5 (a level-2 unit of the factor “classroom”), which is in school 2 (a level-1 unit of the factor “school”). This dataset includes a total of $S_1 = 4$ schools and $S_2 = 8$ classrooms.

Remark 1. One might wonder why not treat these factors as categorical variables, include them in the input x_i using one-hot encoding, and then fit a standard input–output model (e.g., regression or neural network), rather than organizing the dataset in a multilevel format and using mixed effects models. Although this approach is theoretically possible, it faces practical challenges. Typically, the number of possible level- l units, S_l , grows with the sample size at a rate of $O(N)$. When N is large, as in the case of the automobile insurance dataset discussed by Tseung et al. [46], the dimensionality of x_i would become computationally prohibitive, and the model would be highly susceptible to overfitting if the factors were treated as one-hot encoded categorical variables. Moreover, this approach ignores the interdependencies between outputs across different observations, leading to the problems discussed in Section 1.1. Additionally, standard input–output models have limitations in prediction when new units within a factor emerge (e.g., a new school or classroom not represented in the training data). Mixed effects models, which we will discuss in the next section, avoid this issue and can effectively predict future observations with previously unseen factor units.

3. Mixed effect models for multilevel data

Datasets with a multilevel structure are typically modeled using a mixed effects framework. In this approach, the known inputs x_i are considered the “fixed effects” or “hard parameter sharing”. Additionally, it is assumed that there are L unobserved variables $\theta_i = (\theta_{i1}, \dots, \theta_{iL})$ that jointly influence the output y_i , where θ_{il} represents the effect of the level- l unit on observation i . We define $\theta_l^{(s)} := \theta_{il} = \theta_{i'l}$ if $c_l(i) = c_l(i') = s$ for $s = 1, \dots, S_l$ and $l = 1, \dots, L$. This means that if two observations i and i' belong to the same level- l unit, they share the same level- l unobserved variable $\theta_{il} = \theta_{i'l}$, making the observations statistically dependent. The unobserved variables θ_i are treated as random and are specified by a probability distribution, thus θ_i is regarded as a “random effect” or “soft parameter sharing”.

Example 2. For the hypothetical “school problem” data in Table 1, the variable $\theta_1^{(s)}$ represents the unobserved characteristics shared by students from the same school, so $\theta_{i1} = \theta_{i'1}$ if students i and i' attend school s . Similarly, $\theta_2^{(s)}$ represents the shared characteristics of students in the same classroom, meaning $\theta_{i2} = \theta_{i'2}$ if they belong to classroom s . For instance, $\theta_{4,1}^{(2)} = \theta_{6,1}^{(2)} = \theta_{6,1}$ because students 4 and 6 attend the same school (#2), but $\theta_{4,2}^{(2)} \neq \theta_{6,2}^{(2)}$ since they are in different classrooms. The variables $\theta_1^{(s)}$ and $\theta_2^{(s)}$ can be interpreted as unobserved factors, such as the quality of teachers and the learning environment, that influence the exam performance of all students within the same school or classroom.

In this section, we will discuss some technical details regarding a generalized framework of the mixed effects models. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be the probability space and suppose that $\theta_l^{(s)}$ is a \mathcal{F} -measurable map from (Ω, \mathcal{F}) to $(\Theta_l, \mathcal{Q}_l)$ for every $l = 1, \dots, L$ and $s = 1, \dots, S_l$, where Θ_l is the space of $\theta_l^{(s)}$ or θ_{il} . $\theta_l^{(s)}$ is a random variable if $(\Theta_l, \mathcal{Q}_l) = (\mathbb{R}, \mathcal{R})$ where

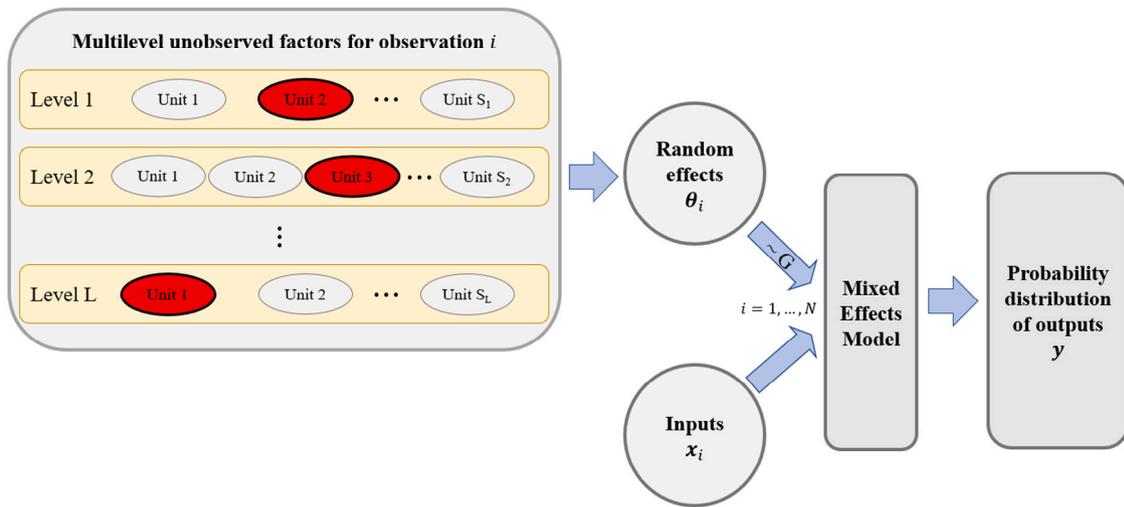


Fig. 1. Multilevel data structure and mixed effects modeling framework. The box in the top left illustrates the level- l unit to which observation i belongs across all factor levels $l = 1, \dots, L$, corresponding to the mapping function $c(i)$. In this example, observation i is associated with the second level-1 unit, third level-2 unit, and first level- L unit. These factor unit classifications govern the random effects θ_i for observation i . For all observations $i = 1, \dots, N$, the random effects θ_i , drawn from a distribution G , along with the inputs x_i , are used in the mixed effects model H as defined in Eq. (2). This model generates the joint probability distribution of outputs y , denoted by $\tilde{H}(y|x)$ in Eq. (4).

\mathcal{R} is a Borel set, but we do not want to impose such a restriction. It is because the random effects $\theta_l^{(s)}$ are unobserved and we are unsure if they can be quantified as a real number. In the context of the school problem, summarizing the aforementioned unobserved characteristics of schools or classrooms into a single number is not feasible. The space Θ_l also varies among different mixed effects models in literature. For example, $\Theta_l = \mathbb{R}$ for most GLMMs [37], $\Theta_l = \mathbb{R}^{2g-1}$ for the GLMM MoE model by Ng and McLachlan [44], and $\Theta_l = \mathbb{R}^{gK}$ for the mixture of random effects models by Ng and McLachlan [45], where g and K are those defined in Section 1.

Under the generalized mixed effects model, we assume that y_i depends only on x_i and θ_i . Conditioned on x_i and θ_i , we further assume that $\{y_i\}_{i=1, \dots, N}$ are mutually independent. Specifically, we denote $y_i | x_i, \theta_i \stackrel{\text{ind}}{\sim} H(\cdot | x_i, \theta_i)$, $i = 1, \dots, N$, (2)

where $H := H(\cdot | x_i, \theta_i)$ can be any probability distributions on y_i given x_i and θ_i . With these assumptions, the joint distribution of y given x is given by

$$\tilde{H}(y|x) = \int_{\Omega} \left[\prod_{i=1}^N H(y_i | x_i, \theta_i) \right] d\mathbb{P}. \quad (3)$$

Suppose that each measurable space $(\Theta_l, \mathcal{Q}_l)$ is also equipped by a probability measure G_l , corresponding to the “distribution” of $\theta_l^{(s)}$. Denote further $(\tilde{\Theta}, \tilde{\mathcal{Q}}, G)$ as the product of probability spaces $\{(\Theta_l, \mathcal{Q}_l, G_l)\}_{l=1, \dots, L; s=1, \dots, S_l}$. Then, the joint distribution in Eq. (3) can be re-written as

$$\tilde{H}(y|x) = \int_{\tilde{\Theta}} \left[\prod_{i=1}^N H(y_i | x_i, \theta_i) \right] dG(\tilde{\theta}), \quad (4)$$

where $\tilde{\theta} = \{\theta_l^{(s)}\}_{l=1, \dots, L; s=1, \dots, S_l}$. Fig. 1 graphically summarizes the modeling framework. The model framework above encompasses a broad range of models designed for multilevel or hierarchical data, including generalized linear mixed models (GLMM) and nonlinear mixed effects models (see, e.g., Goldstein [33], Davidian and Gallant [38]). Since there are no constraints on the functional forms of H and G , this structure is highly flexible. It naturally incorporates any potential joint distributions of $y_i | x_i, \theta_i$, regression links between x_i and y_i (nonlinear effects and interactions among covariates), the influence of random effects θ_i on y_i (random intercepts), and the interactions between θ_i and x_i (random slopes).

The following assumption has been implicitly made on $\theta_l^{(s)}$:

Assumption 1. $\{\theta_l^{(s)}\}_{l=1, \dots, L; s=1, \dots, S_l}$ are mutually independent.

The above assumption implies that $(\tilde{\Theta}, \tilde{\mathcal{Q}}, G)$ is a product probability space, so $G(\tilde{\theta})$ can be written as

$$G(\tilde{\theta}) = \prod_{l=1}^L \prod_{s=1}^{S_l} G_l(\theta_l^{(s)}) \quad \text{or} \quad dG(\tilde{\theta}) = \prod_{l=1}^L \prod_{s=1}^{S_l} G_l(d\theta_l^{(s)}). \quad (5)$$

Note that the prior independence assumption across factors within a level (i.e., $s = 1, \dots, S_l$) is very natural for most data structures especially for those involving repeated measurements (see, e.g., Goldstein [33], Boucher and Denuit [36], Ng et al. [41], Yau et al. [43], and Ng and McLachlan [44]). The prior independence across levels (i.e., $l = 1, \dots, L$) is also often assumed for datasets with multilevel structure (see, e.g., Goldstein [33], McGilchrist [37]).

4. Mixture of experts model with random effect

Despite of the generality of the above mixed effect model (Eq. (3)), it is essential to appropriately specify the functional forms of H and G to model a multilevel dataset. However, this is challenging, especially when the space $\tilde{\Theta}$ of the latent random effects $\tilde{\theta}$ defined in Section 3 is not observed from the dataset. Recall that a multilevel dataset only provides information on how each observation is classified into one of the factors for each level l , but not on what the factors are or how to quantify these factors.

In this section, we introduce the mixture of experts (MoE) model with random effects, called the mixed MoE (MMoE), as a candidate regression model to cater for multilevel data structure. The justification of the proposed model, which analyzes the ability of the MMoE to accurately approximate the generalized form of the mixed effect model (Eq. (3)), will be presented in the next section.

4.1. Model set-up

Under the MMoE, we assume that each observation i is equipped by L levels of random effects, denoted as an L -vector $w_i = (w_{i1}, \dots, w_{iL})$. Similar to the mapping of the unobserved factors introduced in Section 3, we also have $w_{il} = w_{i'l} := w_l^{(s)}$ if $c_l(i) = c_l(i') = s$. Hence, $w_l^{(s)}$ represents the level- l random effect associated with all observations that correspond to the s th level- l factor. The only difference between θ_i and w_i is that we restrict $w_{il} \in \mathbb{R}$ into a Euclidean space instead of a general space $\tilde{\Theta}$, which is unknown and hard to specify. Similar to

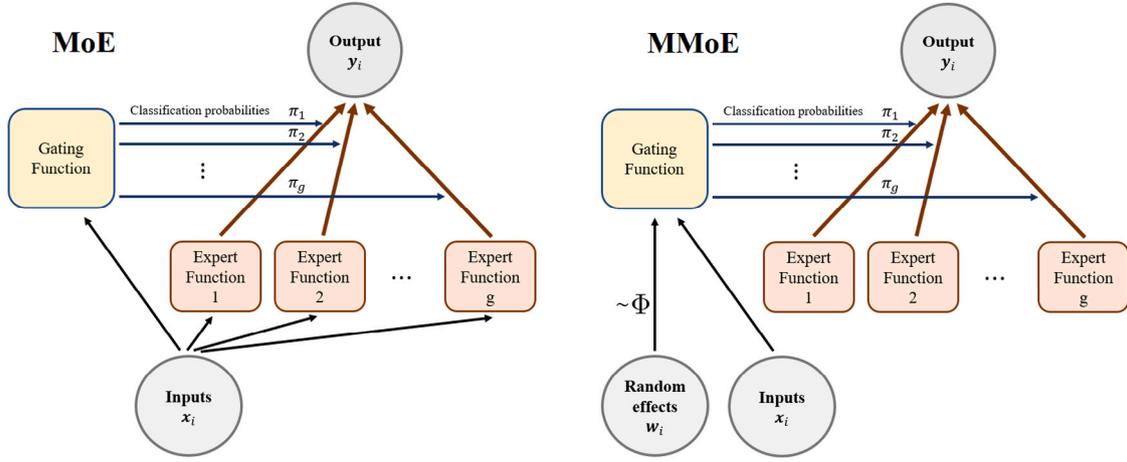


Fig. 2. Model architectures for the MoE (left panel) and MMoE (right panel) models. In the MoE model, inputs x_i are processed by both the gating and expert functions. The gating function determines the classification probabilities $\pi_j := \pi_j(x_i; \alpha)$ ($j = 1, \dots, g$) for each of the g expert functions. In the proposed MMoE model, the random effects w_i , drawn from a distribution Φ , together with the inputs x_i , are fed exclusively into the gating functions.

θ , we also define $w = \{w_l^{(s)}\}_{l=1, \dots, L; s=1, \dots, S_l}$ as the random effects across all levels and factors.

The distribution function of y_i conditional on x_i and w_i is given by

$$F(y_i; \alpha, \beta, \Psi, g | x_i, w_i) = \sum_{j=1}^g \pi_j(x_i, w_i; \alpha, \beta) F_0(y_i; \Psi_j), \quad (6)$$

where g is the number of latent classes, $\pi_j(x_i, w_i; \alpha, \beta)$ is the mixing weight for the j th class (called the gating function) parameterized by α and β , and $F_0(y_i; \Psi_j)$ is the multivariate distribution function of y_i for the j th class (called the expert function) parameterized by Ψ_j . Here, we denote $\alpha = \{\alpha_{j0}, \alpha_j : j = 1, \dots, g\} \in \mathcal{A}$ as the regression parameters of the gating function, $\beta = \{\beta_j : j = 1, \dots, g\} \in \mathcal{B}$ as the coefficients of the random effects, and $\Psi = \{\Psi_j : j = 1, \dots, g\} \in \mathcal{P}$ as the parameters of the expert functions, where \mathcal{A} , \mathcal{B} and \mathcal{P} are defined respectively as the parameter spaces for α , β and Ψ . Moreover, we specify $\pi_j(x_i, w_i; \alpha, \beta)$ as a logit linear gating function, given by

$$\pi_j(x_i, w_i; \alpha, \beta) = \frac{\exp\{\alpha_{j0} + \alpha_j^T x_i + \beta_j^T w_i\}}{\sum_{j'=1}^g \exp\{\alpha_{j'0} + \alpha_{j'}^T x_i + \beta_{j'}^T w_i\}}, \quad j = 1, 2, \dots, g. \quad (7)$$

Apart from that, the random effects $\{w_l^{(s)}\}_{l=1, \dots, L; s=1, \dots, S_l}$ are assumed to be independent across l and s , and $w_l^{(s)}$ follows a fixed pre-specified distribution Φ_l with no extra parameters in it. Based on the above model specification, the joint distribution of y given x is

$$\tilde{F}(y; x) := \tilde{F}(y; \alpha, \beta, \Psi, g | x) = \int \prod_{i=1}^N F(y_i; \alpha, \beta, \Psi | x_i, w_i) d\Phi(w), \quad (8)$$

where Φ is the joint distribution of w , given by

$$\Phi(w) = \prod_{l=1}^L \prod_{s=1}^{S_l} \Phi_l(w_l^{(s)}) \quad \text{or} \quad d\Phi(w) = \prod_{l=1}^L \prod_{s=1}^{S_l} \Phi_l(dw_l^{(s)}). \quad (9)$$

The model can be interpreted as follows with a visual illustration displayed in the right panel of Fig. 2. Each observation is assigned into one of the g homogeneous subgroups with classification probabilities $\pi_j(x_i, w_i; \alpha, \beta)$. The classification probabilities vary among observations as they depend on both inputs x_i and unobserved factors w_i . Conditioned on the subgroup observation i belongs to, the outputs y_i are governed by a homogeneous probability distribution $F_0(y_i; \Psi_j)$ independent of x_i and w_i .

Example 3. The hypothetical data in Table 1 was generated using the proposed MMoE model with two subgroups ($g = 2$). The parameters are set as follows: $\alpha_{1,0} = -5.5$, $\alpha_1 = (0, 0.5)$, $\beta_1 = (0.5, 0.5)$, $\alpha_{2,0} = 0$, $\alpha_2 = (0, 0)$, and $\beta_2 = (0, 0)$. Both Φ_1 and Φ_2 are assumed to follow a standard normal distribution. The distributions of the outputs for the

subgroups are modeled as $F_0(y_i; \Psi_1) \sim \text{MVN}(\mu = (80, 80), \Sigma = \text{diag}(5, 5))$ and $F_0(y_i; \Psi_2) \sim \text{MVN}(\mu = (65, 65), \Sigma = \text{diag}(10, 10))$, where ‘‘MVN’’ denotes a multivariate normal distribution, μ is the mean vector, Σ is the covariance matrix, and ‘‘diag’’ indicates a diagonal matrix. The model can be interpreted as follows:

1. The two subgroups likely represent ‘‘good students’’ (subgroup 1) and ‘‘bad students’’ (subgroup 2), as indicated by the higher average score in subgroup 1 (80) compared to subgroup 2 (65). If it is known that student i is ‘‘good’’ (resp. ‘‘bad’’), then $F_0(y_i; \Psi_1)$ (resp. $F_0(y_i; \Psi_2)$) represents the probability distribution of student i 's exam scores.
2. The parameters α indicate how gender and income influence the likelihood of a student being classified as ‘‘good’’ or ‘‘bad’’. The values $\alpha_1 = (0, 0.5)$ suggest that gender does not influence the probability of a student being classified as ‘‘good’’ or ‘‘bad’’. However, students from wealthier families are more likely to be categorized as ‘‘good students’’ (subgroup 1).
3. Students attending the same school (level-1 random effect w_{i1}) or classroom (level-2 random effect w_{i2}) share the same random effects, making them more likely to be classified into the same category. For instance, both students 11 and 12, who belong to the same school and classroom, are classified as ‘‘bad students’’ and have lower exam scores. The magnitudes of β govern the likelihood that students from the same school or classroom will be categorized within the same group.

4.2. Comparisons to the literature

Our proposed MMoE model extends the standard MoE model (e.g., Jacobs et al. [1]) described in Eq. (1) to account for multilevel data by incorporating random effects w_i into the gating function. If the coefficients of the random effects in Eq. (6) are set to $\beta = 0$, the proposed MMoE model simplifies back to the standard MoE model.

Our proposed MMoE model builds on the limited literature concerning mixed effects MoE models, such as Yau et al. [43], Ng and McLachlan [44], and Ng and McLachlan [45], which only account for a single level of random effects (i.e., $L = 1$). By incorporating multiple levels of random effects, our model addresses the complex multilevel data structures seen in contexts like the school problem. Additionally, our approach diverges from existing mixed effects MoE models in two key ways that simplify the model structure. First, we eliminate the regression relationship in the expert functions, meaning they do not depend on the inputs x_i (see also Fig. 2), following the reduced MoE

(RMoE) model concept introduced by Fung et al. [30]. Second, our model assumes that the level- l random effect w_{ij} is consistent across all g gating functions (i.e., w_{ij} does not vary with j), contrasting with Ng and McLachlan [44], which allows for varying and independent level- l random effects across gating functions.

These model simplifications offer several advantages. First, they allow for a broader selection of probability distributions as the expert function F_0 , including more complex non-exponential distributions (e.g., phase-type distributions) where regression modeling may be impractical or computationally intensive. Second, the simplified structure enhances interpretability, as our model can cluster observations into homogeneous subgroups and explain the variability of each level- l factor through a single source (i.e., $w_{ij} \in \mathbb{R}$) rather than the multiple sources considered by Ng and McLachlan [44]. Third, the reduction in model parameters significantly decreases the computational burden during parameter estimation. Specifically, our follow-up work [46], which employs a variational ECM algorithm for parameter estimation, identifies that the primary computational challenge lies in simulating and computing the posterior samples of the random effects w_i given the observed data. By reducing the dimension of random effects from \mathbb{R}^{2g-1} [44] or \mathbb{R}^{gK} [45] to \mathbb{R} for each level, our estimation algorithm becomes significantly more computationally efficient and scalable to large datasets.

The remaining issue is: how does such a reduced structure affect its model flexibility? To justify the proposed model, we will demonstrate the denseness property in the following section, meaning that the MMoE model structure of Eq. (8) can approximate any generalized form of mixed effect models expressed by Eq. (3). This will provide evidences suggesting that our proposed model is parsimonious. In other words, the MMoE has the simplest structure without harming its representation capability.

5. Denseness theory

This section studies the approximation capability of the class of MMoE models. Our goal is to show that the proposed MMoE is versatile enough to approximate any mixed effects models under mild regularity conditions, even if the MMoE is constructed in a reduced form: (i) the gating function is restricted to be a logit linear gating; (ii) regression link is removed in the expert functions; (iii) the random effects are restricted to follow some fixed pre-determined distributions. Before that, we need to technically formulate a class of mixed effects models and define “denseness” for mixed effects models. These definitions are the extensions of Fung et al. [30], who defines “regression distributions” and “denseness” in the regression settings without considering random effects.

Following the notations defined in Section 3, for $l = 1, \dots, L$, let \mathcal{T}_l denote a collection of some spaces of level- l random effects Θ_l , and \mathcal{G}_l denote a collection of probability measures G_l on $\theta_l^{(s)}$. We also denote $\mathcal{T} := \mathcal{T}_1 \times \dots \times \mathcal{T}_L$ as a collection of product spaces $\Theta := \Theta_1 \times \dots \times \Theta_L$, \mathcal{H} as a collection of some distribution functions H on $(y_i|x_i, \theta_i)$, and $\mathcal{G} := \mathcal{G}_1 \times \dots \times \mathcal{G}_L$ as a collection of product probability measures G on $\tilde{\theta}$. Also, let \mathcal{C} be a set containing all possible mappings $c(\cdot)$, and define a vector $\mathcal{S} = (S_1, \dots, S_L)$ with $S = \mathbb{N}^L$. “A class of mixed effects models” and “mixed effects distributions” are first defined as follows:

Definition 1. A class of mixed effects models $\mathcal{M}_L(\mathcal{X}; \mathcal{T}, \mathcal{H}, \mathcal{G}) := \{\tilde{H}(\cdot; \mathcal{X}; \Theta, H, G) : \Theta_l \in \mathcal{T}_l, H \in \mathcal{H}, G_l \in \mathcal{G}_l, l = 1, \dots, L\}$ is a collection of mixed effects distributions $\tilde{H}(\cdot; \mathcal{X}; \Theta, H, G)$, where each mixed effects distribution $\tilde{H}(\cdot; \mathcal{X}; \Theta, H, G) := \{\tilde{H}(y|x) := \tilde{H}(y|x; \Theta, H, G) = \int_{\tilde{\theta}} \left[\prod_{i=1}^N H(y_i|x_i, \theta_i) \right] dG(\tilde{\theta}) : x_i \in \mathcal{X}, i \in \{1, \dots, N\}, N \in \mathbb{N}, \mathcal{S} \in \mathcal{S}, c \in \mathcal{C}\}$ is itself a collection of joint probability distributions.

In simpler terms, a “mixed effects distribution” refers to the collection of all possible joint probability distributions, as described by Eq. (4), that can be generated by a mixed effects model outlined in

Section 3, based on a specified H and a specified G . It is important to note that a single mixed effects model can produce various possible joint probability distributions since Eq. (4) depends on the covariates (x_1, \dots, x_N) , sample size N , and the mapping $c(\cdot)$. “A class of mixed effects models” refers to the set of all possible “mixed effects distributions” that can be generated by varying Θ , H , and G .

In the spirit of Fung et al. [30], denseness is defined in the sense of weak convergence of probability distributions. Therefore, before defining denseness, we need to define weak convergence of mixed effects distributions. Let $\{\Theta^{(n)}\}_{n=1,2,\dots}$ denote a sequence of (product) spaces of random effects, $\{H^{(n)}\}_{n=1,2,\dots}$ denote a sequence of distribution functions on y_i given x_i and θ_i , and $\{G^{(n)}\}_{n=1,2,\dots}$ denote a sequence of product probability measures on $\tilde{\theta}$. The definition is as follows:

Definition 2. Consider a sequence of mixed effects distributions $\tilde{H}^{(n)} := \tilde{H}(\cdot; \mathcal{X}; \Theta^{(n)}, H^{(n)}, G^{(n)})$ and a target mixed effects distribution $\tilde{H} := \tilde{H}(\cdot; \mathcal{X}; \Theta, H, G)$. We say that $\{\tilde{H}^{(n)}\}_{n=1,2,\dots}$ weakly converges to \tilde{H} if and only if, for every $x_i \in \mathcal{X}$ (for all $i \in \{1, \dots, N\}$), $N \in \mathbb{N}$, $\mathcal{S} \in \mathcal{S}$, and $c \in \mathcal{C}$, we have $\tilde{H}(\cdot|x; \Theta^{(n)}, H^{(n)}, G^{(n)}) \xrightarrow{D} \tilde{H}(\cdot|x; \Theta, H, G)$ as $n \rightarrow \infty$, where \xrightarrow{D} denotes weak convergence or convergence in distribution. If the convergence in distribution is uniform across any compact input space $\tilde{\mathcal{X}} \subseteq \mathcal{X}$, meaning that $\tilde{H}(y|x; \Theta^{(n)}, H^{(n)}, G^{(n)}) \rightarrow \tilde{H}(y|x; \Theta, H, G)$ uniformly on (y, x) with $x_i \in \tilde{\mathcal{X}}$ (for all $i \in \{1, \dots, N\}$), then we say that $\{\tilde{H}^{(n)}\}_{n=1,2,\dots}$ weakly converges to \tilde{H} compactly.

In essence, the definition above requires that all distributions produced by the mixed effects models converge in order to assert that the “mixed effects distributions” themselves converge. We are now able to extend the formalism of Fung et al. [30] and define denseness in the settings of mixed effects models. Similar as above, denote \mathcal{T}' , \mathcal{H}' and \mathcal{G}' , respectively, as the collections of Θ , H and G . The definition is as follows:

Definition 3. Consider two classes of mixed effects models $\mathcal{M}_L := \mathcal{M}_L(\mathcal{X}; \mathcal{T}, \mathcal{H}, \mathcal{G})$ and $\mathcal{M}'_L := \mathcal{M}_L(\mathcal{X}; \mathcal{T}', \mathcal{H}', \mathcal{G}')$. \mathcal{M}_L is dense in \mathcal{M}'_L if and only if for all $(\Theta, H, G) \in \mathcal{T}' \times \mathcal{H}' \times \mathcal{G}'$, there exists a sequence of $\{(\Theta^{(n)}, H^{(n)}, G^{(n)})\}_{n=1,2,\dots}$ with $(\Theta^{(n)}, H^{(n)}, G^{(n)}) \in \mathcal{T} \times \mathcal{H} \times \mathcal{G}$ such that the mixed effects distributions $\{\tilde{H}^{(n)} := \tilde{H}(\cdot; \mathcal{X}; \Theta^{(n)}, H^{(n)}, G^{(n)})\}_{n=1,2,\dots}$ weakly converge to $\tilde{H} := \tilde{H}(\cdot; \mathcal{X}; \Theta, H, G)$. If $\{\tilde{H}^{(n)}\}_{n=1,2,\dots}$ weakly converge to \tilde{H} compactly, then \mathcal{M}_L is said to be compactly dense in \mathcal{M}'_L .

It is evident that if \mathcal{M}_L includes \mathcal{M}'_L , then \mathcal{M}_L is dense in \mathcal{M}'_L . However, if \mathcal{M}_L is a smaller class than \mathcal{M}'_L , the definition of denseness implies that any mixed effects model in \mathcal{M}'_L can be represented or closely approximated by models within the more extensive class \mathcal{M}_L . Put more simply, \mathcal{M}_L can be seen as a model class that is, in practice, at least as rich or flexible as \mathcal{M}'_L even if it is a theoretically smaller set.

With all the necessary definitions established, we now introduce a class of generalized mixed effects models denoted by $\mathcal{M}_L^{\text{gen}}(\mathcal{X}) := \mathcal{M}_L(\mathcal{X}; \mathcal{T}^{\text{gen}}, \mathcal{H}^{\text{gen}}, \mathcal{G}^{\text{gen}})$. Here, $\mathcal{T}^{\text{gen}} := \mathcal{T}_1^{\text{gen}} \times \dots \times \mathcal{T}_L^{\text{gen}}$, $\mathcal{T}_l^{\text{gen}}$ ($l = 1, \dots, L$), \mathcal{H}^{gen} , and \mathcal{G}^{gen} represents the collections of all possible Θ , Θ_l , H , and G that satisfy the following two mild technical assumptions:

Assumption 2. Each space $\Theta_l \in \mathcal{T}_l^{\text{gen}}$ is equipped with a complete separable metric d_{Θ_l} .

Assumption 3. For every probability distribution functions $H \in \mathcal{H}^{\text{gen}}$, $H(y_i|x_i, \theta_i)$ is continuous with respect to (y_i, x_i, θ_i) .

Assumption 2 is a mathematical framework designed to ensure the rigor of the theoretical results and has no practical implications or constraints. Although **Assumption 3** requires H to be a continuous distribution, it can be readily adapted to accommodate discrete distributions. For further details, see **Remark 2**.

Next, we consider the class of MMoE models with a predetermined expert function F_0 and joint distribution of random effects Φ

defined in Section 4. We denote this class as $\mathcal{M}_L^{\text{MMoE}}(\mathcal{X}; F_0, \Phi) := \mathcal{M}_L(\mathcal{X}; \mathcal{T}^{\text{MMoE}}, \mathcal{H}^{\text{MMoE}}, \mathcal{G}^{\text{MMoE}})$. Here, $\mathcal{T}^{\text{MMoE}} := \mathcal{T}_1^{\text{MMoE}} \times \dots \times \mathcal{T}_L^{\text{MMoE}}$ represents the set of all possible spaces for MMoE random effects \mathbf{w} , where $\mathcal{T}_l^{\text{MMoE}}$ contains the possible spaces for $w_l^{(s)}$ for level $l \in \{1, \dots, L\}$. The set $\mathcal{H}^{\text{MMoE}}$ contains all conditional distributions given by Eq. (6), and $\mathcal{G}^{\text{MMoE}}$ encompasses all possible distributions of the random effects \mathbf{w} . Given that the proposed MMoE model assumes a scalar random effect for each level, with \mathbb{R} as its domain, it follows that $\mathcal{T}_l^{\text{MMoE}} = \{\mathbb{R}\}$ and $\mathcal{T}^{\text{MMoE}} = \{\mathbb{R}^L\}$ each contain only a single element. Additionally, since the distribution of \mathbf{w} is fixed as Φ according to Eq. (9), we have $\mathcal{G}^{\text{MMoE}} = \{\Phi\}$, also containing just one element. Moreover, we can express $\mathcal{H}^{\text{MMoE}} = \{F(y; \alpha, \beta, \psi, g | x_i, w_i) : \alpha \in \mathcal{A}, \beta \in \mathcal{B}, \psi \in \mathcal{P}, g \in \mathbb{N}\}$, where $F(y; \alpha, \beta, \psi, g | x_i, w_i)$ takes the form given in Eq. (6). The following two conditions on the choices of F_0 and Φ are crucial for establishing the denseness property of the MMoE model class:

Assumption 4. F_0 satisfies the denseness condition outlined by Proposition 3.1 of Fung et al. [30], meaning that for every $q \in \mathbb{R}^K$, there exists a sequence of parameters $\{\psi^{(n)}(q)\}_{n=1,2,\dots}$ such that $F_0(\cdot; \psi^{(n)}(q)) \xrightarrow{D} q$ as $n \rightarrow \infty$.

Assumption 5. Φ_l is a continuous distribution function for every $l = 1, \dots, L$.

The above two assumptions are not necessarily mild. As discussed in Fung et al. [30], some common distributions, such as Pareto and exponential distributions, do not satisfy the denseness condition under Assumption 4. Moreover, Assumption 5 does not hold whenever we choose any discrete distributions for Φ_l . Nonetheless, note that the expert function F_0 and random effect distribution Φ_l are both predetermined, so we have the control to choose suitable functions that fulfill Assumptions 4 and 5 before modeling a multilevel dataset via the MMoE. For example, one may choose Gamma, Weibull, log-normal, or inverse-Burr distributions as the expert function F_0 [30], and select a normal distribution as the random effect distribution Φ_l .

We now introduce the main result that validates the representational power of the proposed MMoE model class. The detailed technical proof is provided in Appendix A, while Section 7 offers intuitive explanations of the key steps in the proof to aid in understanding the core concepts.

Theorem 1. Suppose that Assumptions 1 to 5 are satisfied. Then, $\mathcal{M}_L^{\text{MMoE}}(\mathcal{X}; F_0, \Phi)$ is compactly dense in $\mathcal{M}_L^{\text{gen}}(\mathcal{X})$.

Remark 2. It is also important to investigate into the above theorem for discrete distributions (see, e.g., Jiang and Tanner [23], Fung et al. [30]) instead of continuous distributions. As discussed by Norets et al. [25], any discrete distribution can be represented by a continuous latent distribution. Then, it is obvious that Theorem 1 still holds for discrete distributions if the denseness condition in Assumption 4 is changed from $q \in \mathbb{R}^K$ to $q \in \mathbb{N}^K$.

Theorem 1 demonstrates that the proposed MMoE model, described in Eq. (8), can closely approximate any mixed effects models in the form of Eq. (4), which has not been addressed by the existing literature on multilevel MoE models, including Ng and McLachlan [44, 45]. This denseness property offers a theoretical foundation for the MMoE model's flexibility in capturing a wide range of model features, including joint distributions (such as multimodal distributions and dependencies among outputs), regression patterns (like non-linear links and interactions among covariates), random intercepts (reflecting unique effects of unobserved factors on the outputs), and random slopes (accounting for interactions between covariates and random effects). As discussed in Section 4, our proposed model features a reduced structure. Therefore, the result in Theorem 1, which provides approximations within a narrower class of models, is remarkably stronger

compared to proving the result based on a broader class of unreduced MMoE models, such as those proposed by Ng and McLachlan [44, 45]. Additionally, the denseness property suggests the parsimony of the proposed model.

6. Nested mixed effects models for hierarchical data

A nested mixed effects model is a specific case of the model discussed in Section 3, involving multiple random effects. In nested models, the L levels of factors in a multilevel dataset must be arranged hierarchically, with the first level representing the highest level and the L th level representing the lowest. In this structure, observations that share the same unit at a lower level must also share the same unit at the higher level. If the factors across levels are not nested, the resulting model is referred to as a ‘‘crossed’’ mixed effects model.

Before defining a nested model, it is helpful to introduce an alternative notation for identifying observations. Denote $i = (i_1, i_2, \dots, i_{L+1})$ as an identifier of an observation, and $i_l = (i_1, i_2, \dots, i_l)$ as an identifier up to level l . Here, i_1 indicates the label of the level-1 unit to which the observation belongs. For $l = 2, 3, \dots, L$, i_l serves as the label for the level- l unit, given that the labels for the first $(l - 1)$ levels are i_{l-1} . Finally, i_{L+1} represents the label of the observation itself, given that the L units corresponding to the L levels of factors are labeled by i_L .

A classic example of nested factors is the school problem (see, e.g., Aitkin and Longford [32]), where two levels of factors influence a student's performance: school and classroom. Here, since students in the same classroom must belong to the same school, we have $L = 2$, with ‘‘school’’ as the first level and ‘‘classroom’’ as the second. For example, in the school problem with $L = 2$, $i = (2, 3, 5)$ refers to the fifth student in the third classroom of the second school.

Furthermore, denote N_0 as the number of possible level-1 factors, so that the support of i_1 is given by $\mathcal{I}_1 := \{1, \dots, N_0\}$. For $l = 2, 3, \dots, L$, define $N_{i_{l-1}}$ as the number of possible level- l factors where the first $(l - 1)$ factor levels are labeled as i_{l-1} , such that the support of i_l is $\mathcal{I}_l := \{i_l : i_{l-1} \in \mathcal{I}_{l-1}, i_l = 1, \dots, N_{i_{l-1}}\}$. Similarly, N_{i_L} is the number of observations having i_L as the label of the L factors. Then, the support of i is $\mathcal{I} := \{i : i_L \in \mathcal{I}_L, i_{L+1} = 1, \dots, N_{i_L}\}$. Also, the total number of observations is given by $N = \sum_{i_1=1}^{N_0} \sum_{i_2=1}^{N_{i_1}} \dots \sum_{i_L=1}^{N_{i_{L-1}}} N_{i_L}$. Hence, the nested multilevel dataset is given by $(y, \mathbf{x}) := \{(y_i, \mathbf{x}_i)\}_{i \in \mathcal{I}}$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$ are respectively the input and output of an observation labeled as $i \in \mathcal{I}$.

Example 4. The school problem serves as an example of nested factors. As seen in Table 1, students in the same classroom must also belong to the same school. Therefore, we have $L = 2$, where ‘‘school’’ and ‘‘classroom’’ are the first and second levels of factors, respectively. In this dataset, there are $N_0 = 4$ schools, and the number of classrooms in each school is $(N_1, N_2, N_3, N_4) = (3, 5, 2, 2)$. Observations are relabeled from i to i to reflect this structure. For instance, student $i = 7$ is relabeled as $i = (2, 3, 1)$, indicating the first student in the third classroom (classroom #5) of the second school (school #2). The support set for this data is $\mathcal{I} = \{(1, 1, 1), (1, 1, 2), (1, 2, 1), (2, 1, 1), (2, 1, 2), (2, 2, 1), (2, 3, 1), (2, 3, 2), (3, 1, 1), (3, 2, 1), (4, 1, 1), (4, 1, 2)\}$. To summarize the above descriptions, Fig. 3 provides a visual representation of the hierarchical data structure within the context of the school problem.

In this section, we will define the generalized class of nested mixed effects models, formulate the proposed MMoE model with nested random effects, and construct the denseness theory for the class of nested MMoE.

6.1. Generalized nested mixed effect model

Similar to the MMoE defined in Section 3, under the nested MMoE, the response y_i depends on its covariates \mathbf{x}_i and L levels of latent random effects denoted as $\theta_{i_1}, \dots, \theta_{i_L}$. The dependence assumption among the latent factors is stated as follows:

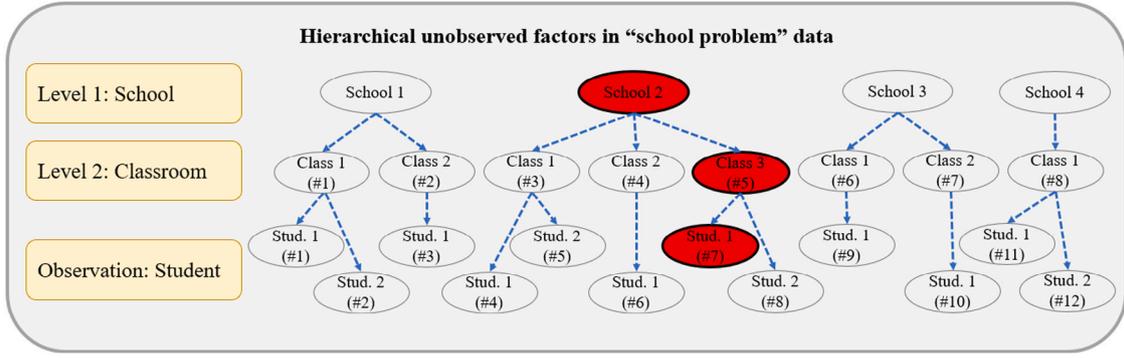


Fig. 3. Hierarchical data structure for a hypothetical “school problem” dataset, illustrating two nested levels of factors: school and classroom. Note that “Stud.” means “Student”. The units shaded in darker color (red) indicate the specific school and classroom to which a student belongs. In this example, the highlighted units correspond to student #7, who is the first student in the third classroom (classroom #5) within the second school (school #2). The dashed arrows depict the hierarchical organization of the data, indicating that students from the same classroom must also belong to the same school.

Assumption 6. Conditioned on $\theta_{i_1}, \dots, \theta_{i_L}$, the N observations are independent. The set of parents of θ_{i_l} is given by $\text{pa}(\theta_{i_l}) = (\theta_{i_1}, \dots, \theta_{i_{l-1}})$ for $l = 1, \dots, L$, meaning that θ_{i_l} may depend only on $(\theta_{i_1}, \dots, \theta_{i_{l-1}})$.

In other words, the lower level factor only depends directly on its corresponding upper level factors. Under a nested hierarchical data structure, we are able to relax the independence assumption in [Assumption 1](#) by allowing the dependence of lower level random effects on their parents (upper level effects). For $l = 1, \dots, L$, we also denote G_l as the distribution of the level- l factor conditioned on $\text{pa}(\theta_{i_l})$.

The joint distribution of \mathbf{y} given \mathbf{x} is given by

$$\tilde{H}(\mathbf{y}|\mathbf{x}) = \int_{\tilde{\theta}} \left[\prod_{i \in I} H(\mathbf{y}_i; \mathbf{x}_i | \theta_{i_l}) \right] dG(\tilde{\theta}) \quad (10)$$

with

$$G(\tilde{\theta}) = \prod_{i_1=1}^{N_0} G_1(\theta_{i_1}) \prod_{i_2=1}^{N_{i_1}} G_2(\theta_{i_2} | \theta_{i_1}) \cdots \prod_{i_L=1}^{N_{i_{L-1}}} G_L(\theta_{i_L} | \theta_{i_1}, \dots, \theta_{i_{L-1}}), \quad (11)$$

where $\tilde{\theta} = \{\theta_{i_l} : i_l \in I_l, l = 1, \dots, L\}$.

6.2. Nested mixed mixture of experts model

Analogous to the MMoE introduced in [Section 4](#), we construct a nested MMoE for hierarchical data. Denote $\mathbf{w}_i = (w_{i_1}, \dots, w_{i_L}) \in \mathbb{R}^L$ as the L levels of random effects of observation i . Also, define $\mathbf{w} = \{w_{i_l}\}_{i_l \in I_l, l=1, \dots, L}$ as all the random effects aggregated across all observations. With a slight change of notations from [Eqs. \(6\) and \(7\)](#), the distribution function of \mathbf{y}_i conditional on \mathbf{x}_i and \mathbf{w}_i is then

$$F(\mathbf{y}_i; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\Psi}, g | \mathbf{x}_i, \mathbf{w}_i) = \sum_{j=1}^g \pi_j(\mathbf{x}_i, \mathbf{w}_i; \boldsymbol{\alpha}, \boldsymbol{\beta}) F_0(\mathbf{y}_i; \boldsymbol{\Psi}_j), \quad (12)$$

where the gating function given by

$$\pi_j(\mathbf{x}_i, \mathbf{w}_i; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\exp\{\alpha_{j0} + \boldsymbol{\alpha}_j^T \mathbf{x}_i + \boldsymbol{\beta}_j^T \mathbf{w}_i\}}{\sum_{j'=1}^g \exp\{\alpha_{j'0} + \boldsymbol{\alpha}_{j'}^T \mathbf{x}_i + \boldsymbol{\beta}_{j'}^T \mathbf{w}_i\}}, \quad j = 1, 2, \dots, g. \quad (13)$$

Similar to [Section 4](#), the random effects $\{w_{i_l}\}_{l=1, \dots, L; i_l=1, \dots, N_{i_{l-1}}}$ are constructed to be independent across l and i_l with $w_{i_l} \sim \Phi_l$. This construction is simplified from [Assumption 6](#) of the generalized nested mixed models where the random effects may depend on their parents. Adapting from [Eqs. \(8\) and \(9\)](#), the joint distribution of \mathbf{y} given \mathbf{x} is

$$\tilde{F}(\mathbf{y}; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\Psi}, g | \mathbf{x}) = \int \prod_{i \in I} F(\mathbf{y}_i; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\Psi} | \mathbf{x}_i, \mathbf{w}_i) d\Phi(\mathbf{w}), \quad (14)$$

where Φ is the joint distribution of \mathbf{w} given by

$$\Phi(\mathbf{w}) = \prod_{l=1}^L \prod_{i_l \in I_l} \Phi_l(w_{i_l}) \quad \text{or} \quad d\Phi(\mathbf{w}) = \prod_{l=1}^L \prod_{i_l \in I_l} \Phi_l(dw_{i_l}). \quad (15)$$

From [Eq. \(15\)](#) above, the random effects are still assumed to be independent under the nested MMoE. However, we will show in the

following subsection that such a specification suffices to approximate the dependence of lower level random effects on their parents under a hierarchical data structure.

6.3. Denseness theory for nested MMoE

Analogous to [Section 5](#), it is desirable to develop an approximation theory for the nested MMoE in the space of the generalized nested mixed effect models. Denote $\mathcal{N} = (N_0, \{N_{i_1}\}_{i_1 \in I_1}, \{N_{i_2}\}_{i_2 \in I_2}, \dots, \{N_{i_L}\}_{i_L \in I_L})$ as the number of factors belonging to each parent factors for each level with $N_0 \in \mathbb{N}$ and $N_{i_l} \in \mathbb{N}$ for $l = 1, \dots, L$, and \mathcal{N} contains all combinations of possible \mathcal{N} . Other notations, unless specified otherwise, are consistent to those defined by [Section 5](#). The equivalent definitions analogous to [Section 5](#) for hierarchical data structure are listed as follows:

Definition 4. A class of nested mixed effects models $\mathcal{M}_L(\mathcal{X}; \mathcal{T}, \mathcal{H}, \mathcal{G}) := \{\tilde{H}(\cdot; \mathcal{X}; \boldsymbol{\theta}, H, G) : \boldsymbol{\theta}_l \in \mathcal{T}_l, H \in \mathcal{H}, G_l \in \mathcal{G}_l, l = 1, \dots, L\}$ is a collection of nested mixed effects distributions $\tilde{H}(\cdot; \mathcal{X}; \boldsymbol{\theta}, H, G)$, where each nested mixed effects distribution $\tilde{H}(\cdot; \mathcal{X}; \boldsymbol{\theta}, H, G) := \{\tilde{H}(\mathbf{y}|\mathbf{x}) := \tilde{H}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}, H, G) = \int_{\tilde{\theta}} [\prod_{i \in I} H(\mathbf{y}_i | \mathbf{x}_i, \theta_{i_l})] dG(\tilde{\theta}) : \mathbf{x}_i \in \mathcal{X}, i \in I, \mathcal{N} \in \mathcal{N}\}$ is itself a collection of joint probability distributions.

Definition 5. Consider a sequence of nested mixed effects distributions $\tilde{H}^{(n)} := \tilde{H}(\cdot; \mathcal{X}; \boldsymbol{\theta}^{(n)}, H^{(n)}, G^{(n)})$ and a target nested mixed effects distribution $\tilde{H} := \tilde{H}(\cdot; \mathcal{X}; \boldsymbol{\theta}, H, G)$. We say that $\{\tilde{H}^{(n)}\}_{n=1,2,\dots}$ weakly converge to \tilde{H} if and only if for every given $\mathbf{x}_i \in \mathcal{X}$ (for all $i \in I$), $\mathcal{N} \in \mathcal{N}$, we have $\tilde{H}(\cdot | \mathbf{x}_i; \boldsymbol{\theta}^{(n)}, H^{(n)}, G^{(n)}) \xrightarrow{D} \tilde{H}(\cdot | \mathbf{x}_i; \boldsymbol{\theta}, H, G)$ as $n \rightarrow \infty$, where \xrightarrow{D} represents a weak convergence or convergence in distribution. If the distributional convergence is uniform across any compact input space $\bar{\mathcal{X}} \subseteq \mathcal{X}$, i.e. $\tilde{H}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}^{(n)}, H^{(n)}, G^{(n)}) \rightarrow \tilde{H}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}, H, G)$ uniformly on (\mathbf{y}, \mathbf{x}) with $\mathbf{x}_i \in \bar{\mathcal{X}}$ (for all $i \in I$), then we say that $\{\tilde{H}^{(n)}\}_{n=1,2,\dots}$ weakly converge to \tilde{H} compactly.

The denseness definition for hierarchical data structure (nested mixed effects models) is exactly the same as [Definition 3](#). Define $\mathcal{M}_L^{\text{gen}}(\mathcal{X})$ as the class of generalized nested mixed effects models expressed in [Eq. \(10\)](#), subject to [Assumptions 2 and 3](#). Also denote $\mathcal{M}_L^{\text{MMoE}}(\mathcal{X}; F_0, \Phi)$ as the class of nested MMoE given by [Eq. \(14\)](#). We have the following denseness theorem:

Theorem 2. Suppose that [Assumptions 2 to 6](#) hold. Then, $\mathcal{M}_L^{\text{MMoE}}(\mathcal{X}; F_0, \Phi)$ is compactly dense in $\mathcal{M}_L^{\text{gen}}(\mathcal{X})$.

The proof is leveraged to [Appendix B](#). [Theorem 2](#) suggests that the nested MMoE has a potential to approximate any generalized nested mixed effect models arbitrarily accurately, even if the random effects are restricted to be independent under the nested MMoE while the

random effects under the generalized nested mixed effect models can be dependent (Assumption 6). In contrast, under Theorem 1, the MMoE can only approximate mixed effect models with independent random effects (Assumption 1). Therefore, given that the data structure is hierarchical, Theorem 2 is a stronger theoretical result than Theorem 1.

7. Proof idea of Theorem 1

In this section, we provide a simplified illustration and intuitive explanation of the key steps involved in proving Theorem 1. The detailed technical proofs can be found in Appendix A. To demonstrate that the proposed MMoE class is dense in the space of generalized mixed effects models, we need to identify an MoE distribution $F(y_i; \alpha, \beta, \Psi | x_i, w_i)$ in Eq. (6) such that the MMoE joint distribution $\tilde{F}(y; x)$ described in Eq. (8) closely approximates the joint distribution $\tilde{H}(y|x)$ given by Eq. (4), i.e.,

$$\int \left[\prod_{i=1}^N F(y_i; \alpha, \beta, \Psi | x_i, w_i) \right] d\Phi(w) \approx \int_{\tilde{\Theta}} \left[\prod_{i=1}^N H(y_i | x_i, \theta_i) \right] dG(\tilde{\theta}) \quad (16)$$

The main challenge in finding a suitable $F(y_i; \alpha, \beta, \Psi | x_i, w_i)$ that satisfies Eq. (16) is that the two sides of the formula integrate two different spaces. Specifically, the left side of Eq. (16) involves a Lebesgue integration on Euclidean space, while the right side integrates with respect to an arbitrary abstract measure. One potential approach to address this issue is to apply an integral transformation on the right side. To ensure that Eq. (16) holds, we construct a mapping $M(\cdot) : \tilde{\Theta} \rightarrow \tilde{\mathcal{W}}$ and identify an appropriate MoE distribution such that $F(y_i; \alpha, \beta, \Psi | x_i, w_i) \approx H(y_i | x_i, \theta_i)$ for all $i = 1, \dots, N$ whenever $M(\tilde{\theta}) = w$ with $\tilde{\theta} \in \tilde{\Theta}$ and $w \in \tilde{\mathcal{W}}$, where $\tilde{\mathcal{W}}$ represents the space of w . The technical details correspond to Appendix A.2 (Step 2). For simplicity, we consider only a single level of factors ($L = 1$), where the random effect in the generalized mixed effects models is $\theta_i = \theta_{i1}$, and the MMoE random effect is $w_i = w_{i1}$. Here, Θ_1 and \mathcal{W}_1 represent the spaces for θ_{i1} and w_{i1} respectively. This result can be extended to cases with multiple levels of factors, with rigorous derivations provided in Appendices A.1 and A.2 (Steps 1 and 2).

To construct a mapping $M(\cdot)$ that effectively links the two spaces Θ_1 and \mathcal{W}_1 , we propose discretizing both Θ_1 and \mathcal{W}_1 into D_1 subspaces. Let $\{\Theta_{1,d_1}\}_{d_1=1,\dots,D_1}$ denote the disjoint partitions on Θ_1 , and $\{\mathcal{W}_{1,d_1}\}_{d_1=1,\dots,D_1}$ denote the disjoint partitions on \mathcal{W}_1 . These partitions are selected so that the probabilities match for each partition, i.e., $G_1(\Theta_{1,d_1}) = \Phi_1(\mathcal{W}_{1,d_1})$ for $d_1 = 1, 2, \dots, D_1$. For each subspace Θ_{1,d_1} , we select a representative point $\theta_{d_1}^* \in \Theta_{1,d_1}$ such that $H(y_i | x_i, \theta_{d_1}^*)$ serves as a reasonable approximation of $H(y_i | x_i, \theta_i)$ for any $\theta_i \in \Theta_{1,d_1}$. Rigorous technical details regarding the space partitioning procedures are presented in Appendices A.1 and A.2 (Steps 1 and 2). We then approximate $H(y_i | x_i, \theta_i)$ using an MoE model with a representation analogous to $F(y_i; \alpha, \beta, \Psi | x_i, w_i)$ in Eq. (6), denoted as

$$F^{**(\omega)}(y_i | x_i, w_i) = \sum_{d_1=1}^{D_1} \pi_{d_1}(x_i, w_i; \alpha, \beta) H(y_i | x_i, \theta_{d_1}^*), \quad (17)$$

where $\pi_{d_1}(x_i, w_i; \alpha, \beta)$ is given by Eq. (7). Note that Eq. (17) corresponds exactly to $F^{**(\omega)}(y_i | x_i, w_i)$ in Eq. (22) of the appendix, using simplified notations. The only distinction between $F^{**(\omega)}(y_i | x_i, w_i)$ and $F(y_i; \alpha, \beta, \Psi | x_i, w_i)$ lies in the expert functions they employ: $F^{**(\omega)}(y_i | x_i, w_i)$ uses $H(y_i | x_i, \theta_{d_1}^*)$, which depends on x_i , while $F(y_i; \alpha, \beta, \Psi | x_i, w_i)$ uses $F_0(y_i; \Psi_j)$, which is independent of x_i . Nevertheless, following the approach outlined in Fung et al. [30], it is shown in Appendices A.3 and A.4 (Steps 3 and 4) that $H(y_i | x_i, \theta_{d_1}^*)$ can itself be approximated by an MoE distribution, with $F_0(y_i; \Psi_j)$ serving as the expert function. This approximation process is not illustrated here to avoid unnecessary complexity, as it is unrelated to random effects and has already been demonstrated by Fung et al. [30]. Therefore, if $H(y_i | x_i, \theta_i)$ can be well approximated by $F^{**(\omega)}(y_i | x_i, w_i)$, it follows that it can also be effectively approximated by $F(y_i; \alpha, \beta, \Psi | x_i, w_i)$.

To perform approximations, we first observe that the softmax gating function $\pi_{d_1}(x_i, w_i; \alpha, \beta)$ in Eq. (17) is shown by Lemma 3.2 of Fung et al. [30] to be “fully flexible”, allowing an observation i to be assigned to any of the D_1 mixture components based on the random effects w_i . Specifically, by expressing $\beta = u\tilde{\beta}$, where $u > 0$ is a tuning parameter that controls how “hard” the softmax gating function assigns observations to the mixture components, one can carefully select parameters $\tilde{\beta}$ such that for any $w_i \in \mathcal{W}_{1,d_1}$ and sufficiently large u , we have $\pi_{d_1}(x_i, w_i; \alpha, \beta) \approx 1$ and $\pi_{d'_1}(x_i, w_i; \alpha, \beta) \approx 0$ for any $d'_1 \neq d_1$, effectively assigning the observation to mixture component d_1 when $w_i \in \mathcal{W}_{1,d_1}$. This process interprets Lemma 2 in the appendix. Hence, we have $F^{**(\omega)}(y_i | x_i, w_i) \approx H(y_i | x_i, \theta_{d_1}^*)$ from Eq. (17). If the partition of spaces Θ_1 and \mathcal{W}_1 are sufficiently fine by selecting a large enough D_1 , then each subspace Θ_{1,d_1} will become sufficiently small, ensuring that $H(y_i | x_i, \theta_{d_1}^*) \approx H(y_i | x_i, \theta_i)$. Summarizing, we have

$$F(y_i; \alpha, \beta, \Psi | x_i, w_i) \approx F^{**(\omega)}(y_i | x_i, w_i) \approx H(y_i | x_i, \theta_{d_1}^*) \approx H(y_i | x_i, \theta_i), \quad (18)$$

and thus Eq. (16) holds, leading to the desired result.

8. Numerical illustration

We provide a numerical example to illustrate the denseness property of the proposed MMoE model, using the approximation procedures discussed in Section 7. For simplicity, the output considered in this example is one-dimensional, meaning $y_i := y_i \in \mathbb{R}$ with $K = 1$. Additionally, we exclude the covariates x_i from this illustration, i.e., $P = 1$, so that $\tilde{\beta}$, as defined in Section 7, can be expressed as $\tilde{\beta} = \{(\tilde{\beta}_{d_1,0}, \tilde{\beta}_{d_1,1})\}_{d_1=1,2,\dots,D_1}$. Here, $\tilde{\beta}_{d_1,0} := u\tilde{\beta}_{d_1,0}$ and $\tilde{\beta}_{d_1,1} := u\tilde{\beta}_{d_1,1}$ correspond, respectively, to the intercept and random effect coefficients of the d_1 -th latent class within $\pi_{d_1}(x_i, w_i; \alpha, \beta)$ in Eq. (17), which can be notationally simplified as $\pi_{d_1}(x_i, w_i; \alpha, \beta) := \exp\{\beta_{d_1,0} + \beta_{d_1,1}w_{i1}\} / \sum_{d'_1=1}^{D_1} \exp\{\beta_{d'_1,0} + \beta_{d'_1,1}w_{i1}\}$ in this example. Considering only $L = 1$, we choose $\Theta_1 = [-2, 2] \subseteq \mathbb{R}$ as the space for generalized random effect $\theta_i := \theta_{i1}$. We construct θ_{i1} to follow a uniform distribution over the range -2 to 2 , so that the distribution function of θ_{i1} is $(2 + \theta_{i1})/4$. The conditional distribution of y_i given θ_{i1} is assumed to be normal with a mean of θ_{i1} and a standard deviation of 1, such that the density of y_i given θ_{i1} is $\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_i - \theta_{i1})^2}$.

In this context, many of the notations introduced in Section 7 can be greatly simplified. Specifically, we express $F^{**(\omega)}(y_i | x_i, w_i)$ as $F^{**(\omega)}(y_i | w_{i1})$, $H(y_i | x_i, \theta_{i1})$ as $H(y_i | \theta_{i1})$, and $H(y_i | x_i, \theta_{d_1}^*)$ as $H(y_i | \theta_{1,d_1}^*)$, where θ_{1,d_1}^* is a scalar representative point in Θ_{1,d_1} as defined in the previous section. To improve clarity, we restate the approximation relations in Eq. (18) using these simplified notations as follows:

$$F^{**(\omega)}(y_i | w_{i1}) \stackrel{u \rightarrow \infty}{\approx} H(y_i | \theta_{1,d_1}^*) \stackrel{D_1 \rightarrow \infty}{\approx} H(y_i | \theta_{i1}). \quad (19)$$

As a concrete example, Table 2 illustrates the procedure described in Section 7 with $D_1 = 5$ partitions for Θ_1 and \mathcal{W}_1 . We select the midpoint θ_{1,d_1}^* as the representative value for the subspace Θ_{1,d_1} . The values of $\{(\tilde{\beta}_{d_1,0}, \tilde{\beta}_{d_1,1})\}_{d_1=1,2,\dots,D_1}$ are carefully chosen following Lemma 2 in the appendix (or see Section 7 for brief description). This example is visually represented in the second row of Fig. 4. Specifically, in Fig. 4, each vertical slice of the blue 2D images plots $F^{**(\omega)}(y_i | w_{i1})$ against y_i (vertical axis) and w_{i1} (horizontal axis) for $u = 1, 10, 100$ and 1000 from left to right and for $D_1 = 2, 5, 10$ and 100 from top to bottom, while the green 2D images in the rightmost panel plot $H(y_i | \theta_{1,d_1}^*)$ against y_i (vertical axis) and θ_{i1} for $D_1 = 2, 5, 10$ and 100 from top to bottom, where θ_{1,d_1}^* is the representative point of the subspace Θ_{1,d_1} to which θ_{i1} belongs. For example, under Table 2, if $\theta_{i1} = -0.6$, then $\theta_{i1} \in \Theta_{1,2}$ and thus $\theta_{1,d_1}^* = \theta_{1,2}^* = -0.8$. Two key observations can be made:

- As u increases, i.e., moving from the left panels to the right panels, $F^{**(\omega)}(y_i | w_{i1})$ visually converges towards $H(y_i | \theta_{1,d_1}^*)$ as

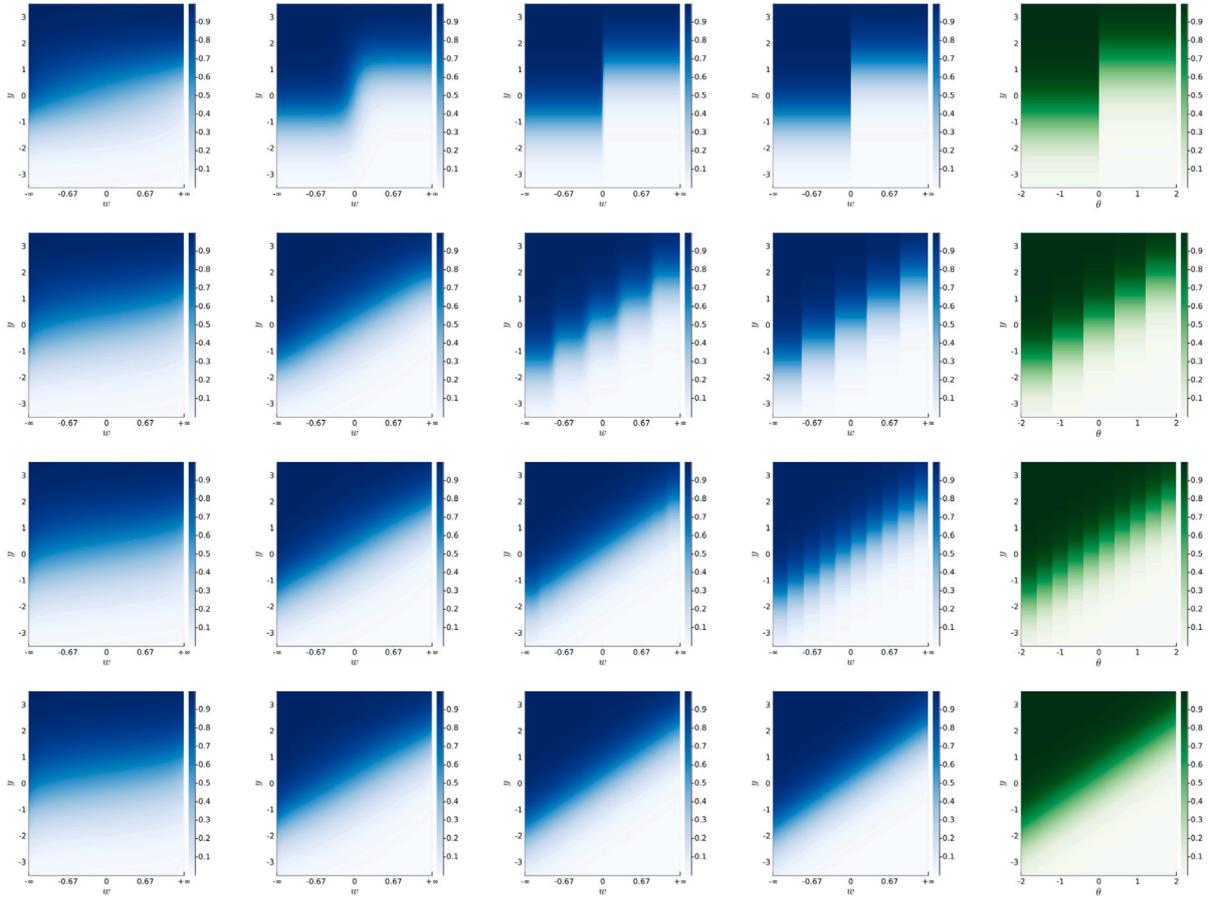


Fig. 4. Illustration of using $F^{***(u)}(y_i|w_{i1})$ (blue) to approximate $H(y_i|\theta_{1,d_1}^*)$ (green) for different numbers of partitions of Θ_1 and \mathcal{W}_1 . Each vertical slice corresponds to the conditional density of $F^{***(u)}(y_i|w_{i1})$ or $H(y_i|\theta_{1,d_1}^*)$. From left to right, the control parameter u ranges from 1, 10, 100 to 1000 in the blue plots, and eventually $F^{***(u)}(y_i|w_{i1}) = H(y_i|\theta_{1,d_1}^*)$ for all $w_i \in \mathcal{W}_{1,d_1}$ as $u \rightarrow \infty$. From top to bottom, the number of partition D_1 ranges from 2, 5, 10 to 100. When the partition becomes infinitely fine and $u \rightarrow \infty$, each $H(y_i|\theta_{1,d_1}^*)$ is approximated arbitrarily well by some $F^{***(\infty)}(y_i|w_{i1})$, e.g., $H(y_i|\theta_{1,d_1}^* = 1) = F^{***(\infty)}(y_i|w_{i1} = 0.67)$.

Table 2
Example of approximation with five partitions of Θ_1 and \mathcal{W}_1 .

d_1	1	2	3	4	5
Θ_{1,d_1}	$[-2.00, -1.20]$	$(-1.20, -0.40]$	$(-0.40, 0.40]$	$(0.40, 1.20]$	$(1.20, 2.00]$
\mathcal{W}_{1,d_1}	$(-\infty, -0.84]$	$(-0.84, -0.25]$	$(-0.25, 0.25]$	$(0.25, 0.84]$	$(0.84, +\infty)$
θ_{1,d_1}^*	-1.60	-0.80	0.00	0.80	1.60
$\tilde{\beta}_{d_1,0}$	0.00	0.23	0.29	0.23	0.00
$\tilde{\beta}_{d_1,1}$	0.00	0.25	0.50	0.75	1.00

shown in the rightmost panel. This provides empirical support for the convergence $F^{***(u)}(y_i|w_{i1}) \rightarrow H(y_i|\theta_{1,d_1}^*)$ as $u \rightarrow \infty$.

- As D_1 increases, making the partitions finer (moving from top to bottom in the rightmost panel), $H(y_i|\theta_{1,d_1}^*)$ transitions from resembling a step function to becoming a smooth function of θ_{i1} . The bottom-right panel closely resembles $H(y_i|\theta_{i1})$, which is a continuous function of (y_i, θ_{i1}) . This supports the convergence $H(y_i|\theta_{1,d_1}^*) \rightarrow H(y_i|\theta_{i1})$ as $D_1 \rightarrow \infty$.

In summary, this numerical example supports the approximation relationship given by Eq. (19), thereby demonstrating the denseness property of the proposed MMoE model as outlined in Section 7.

9. Discussions

In this paper, we introduce a class of the mixed mixture of experts models (MMoE) for multilevel regression data. We prove that the MMoE is dense in the space of generalized mixed effects models, which

is a rich class containing almost all models in the literature having independent random effects, in the sense of weak convergence. We further study a special case where the data is hierarchical in structure. In this case, the proposed nested MMoE is shown to be dense in the space of generalized nested mixed effects models where the random effects can possibly be dependent. The two denseness theorems justify the versatility of the MMoE in catering for a broad range of multilevel data characteristics, including the marginal distribution, dependence, regression link, random intercept and random slope.

This paper aims to prove the most general results imposing minimal assumptions. The only practical restriction is that the expert function $F_0(y_i; \psi_j)$ in Eq. (6) needs to approximate any degenerate distributions (Assumption 4). This assumption is much weaker than those applied to the existing approximation theorems (see, e.g., Nguyen et al. [22], Norets and Pelenis [26]), which require that the expert density function is a scalable symmetric function (Equation (3.1) of Norets and Pelenis [26]), and that the target density function does not change abruptly w.r.t. y_i and x_i (Equation (3.2) of Norets and Pelenis [26]). On the other hand, there are several limitations of the denseness theorems formulated in this paper. Firstly, weak convergence does not guarantee the approximation capability in terms of moments (e.g., the mean function studied by Nguyen et al. [21]) or some distance metrics (e.g., the KL divergence studied by Jiang and Tanner [20], Nguyen et al. [22], Norets et al. [25]). To establish such denseness theorems, one needs to further assume that the moments of the MMoE expert functions and the target distributions are finite, and that the conditions indicated by Equations (3.1) and (3.2) of Norets and Pelenis [26] are fulfilled. Secondly, as described in Section 3.4 of Fung et al. [30], the denseness theorems

do not provide a convergence rate, so there is no control on the mixture components g to approximate any generalized mixed effects distributions at a desired level of accuracy. To establish the rate results, one needs to impose further assumptions on the target distribution $\tilde{H}(y|\mathbf{x})$ in Eq. (4) and the MMoE distribution $\tilde{F}(y;\mathbf{x})$ in Eq. (8). We leave these technical establishments as a future research direction.

CRedit authorship contribution statement

Tsz Chai Fung: Conceptualization, Formal analysis, Investigation, Methodology, Project administration, Supervision, Validation, Visualization, Writing – original draft, Writing – review & editing. **Spark C. Tseung:** Formal analysis, Visualization, Writing – original draft, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Appendix A. Proof of Theorem 1

We begin by introducing the following functions, with detailed notation to be clarified later in the proof.

$$\tilde{H}(y|\mathbf{x}) = \int_{\tilde{\theta}} \left[\prod_{i=1}^N H(y_i|\mathbf{x}_i, \theta_i) \right] dG(\tilde{\theta}), \quad (20)$$

$$\tilde{H}^*(y|\mathbf{x}) = \sum_{\tilde{d} \in \tilde{D}} \prod_{i=1}^N H(y_i|\mathbf{x}_i, \theta_{d^{(c(i))}}^*) G(\tilde{\theta}_{\tilde{d}}), \quad (21)$$

$$\tilde{F}^{**(\mathbf{u})}(y|\mathbf{x}) = \int_{\mathbf{w}} \prod_{i=1}^N F^{**(\mathbf{u})}(y_i|\mathbf{x}_i, \mathbf{w}_i) d\Phi(\mathbf{w}) \quad (22)$$

$$\text{with } F^{**(\mathbf{u})}(y_i|\mathbf{x}_i, \mathbf{w}_i) = \sum_{d \in D^+} \xi_d^{(\mathbf{u})}(\mathbf{w}_i) H(y_i|\mathbf{x}_i, \theta_d^*),$$

$$\tilde{F}^{*(M,Q,t,\mathbf{u})}(y|\mathbf{x}) = \int_{\mathbf{w}} \prod_{i=1}^N F^{*(M,Q,t,\mathbf{u})}(y_i|\mathbf{x}_i, \mathbf{w}_i) d\Phi(\mathbf{w}) \quad (23)$$

$$\text{with } F^{*(M,Q,t,\mathbf{u})}(y_i|\mathbf{x}_i, \mathbf{w}_i) = \sum_{m \in \mathcal{M}} \sum_{q \in \mathcal{Q}} \sum_{d \in D^+} \pi_j^{(t,\mathbf{u})}(\mathbf{x}_i, \mathbf{w}_i; \tilde{\alpha}_q^Q, \tilde{\beta}_d) 1\{y_i \geq |y|_m^M\},$$

$$\tilde{F}^{(M,Q,t,\mathbf{u},v)}(y|\mathbf{x}) = \int_{\mathbf{w}} \prod_{i=1}^N F^{(M,Q,t,\mathbf{u},v)}(y_i|\mathbf{x}_i, \mathbf{w}_i) d\Phi(\mathbf{w}) \quad (24)$$

with $F^{(M,Q,t,\mathbf{u},v)}(y_i|\mathbf{x}_i, \mathbf{w}_i) = \sum_{m \in \mathcal{M}} \sum_{q \in \mathcal{Q}} \sum_{d \in D^+} \pi_j^{(t,\mathbf{u})}(\mathbf{x}_i, \mathbf{w}_i; \tilde{\alpha}_q^Q, \tilde{\beta}_d) F_0(y_i; \psi_m^{M(v)})$. Note that Eq. (20) is in the form of the generalized mixed effects models under Eq. (4), while Eq. (24) is in the MMoE framework under Eq. (8). The main proof idea is to bound the approximation errors between any two consecutive functions from Eqs. (20) to (24):

A.1. Step 1: Approximating Eq. (20) by Eq. (21)

As the metric space (Θ_l, d_{Θ_l}) is complete separable (Assumption 2), the probability measure G_l on $\theta_l^{(s)}$ is tight, i.e. $\forall \epsilon_1 > 0, \exists \tilde{\Theta}_l \subseteq \Theta_l$ compact such that $G_l(\tilde{\Theta}_l) := \mathbb{P}(\theta_l^{(s)} \in \tilde{\Theta}_l) \geq 1 - \epsilon_1$. Since $\tilde{\Theta}_l$ is compact, for any $\delta_1 > 0$ we can find subspaces $\{\Theta_{l,d_l}\}_{d_l=1,\dots,D_l}$ (d_l is called the subspace index) and points $\{\theta_{l,d_l}^*\}_{d_l=1,\dots,D_l}$ such that $\cup_{d_l=1,\dots,D_l} \Theta_{l,d_l} = \tilde{\Theta}_l$, $\theta_{l,d_l}^* \in \Theta_{l,d_l}$ for every $d_l = 1, \dots, D_l$ and Θ_{l,d_l} is covered by the ball with radius δ_1/L centered at θ_{l,d_l}^* . Due to the uniform continuity of H w.r.t. θ_i on $\tilde{\Theta}_l$ (Assumption 3 implies uniform convergence in compact space), for any $\epsilon_2 > 0$ we can choose sufficiently small δ_1 with the aforementioned subspace partitioning such that $|H(y_i|\mathbf{x}_i, \theta_i) -$

$H(y_i|\mathbf{x}_i, \theta_{d^{(c(i))}}^*)| \leq \epsilon_2$ if $\theta_{il} \in \Theta_{l,d^{(c(i))}}$ for every $l = 1, \dots, L$, where $\theta_{d^{(c(i))}}^* = \{\theta_{1,d_l^{(c(i))}}^*\}_{l=1,\dots,L}$ and $d^{(c(i))} = \{d_l^{(c(i))}\}_{l=1,\dots,L}$ with $d_l^{(c(i))} \in \{1, \dots, D_l\}$. Note here the superscript $(c(i))$ of $d_l^{(c(i))}$ represents the subspace index d_l corresponding to the i th observation.

Define $\tilde{H}^*(y|\mathbf{x})$ in the form of Eq. (21), where $\tilde{D} = \prod_{l=1}^L \prod_{s=1}^{S_l} \mathcal{D}_l^{(s)}$ with $\mathcal{D}_l^{(s)} = \{1, \dots, D_l\}$, $\tilde{d} = \{d_l^{(s)}\}_{l=1,\dots,L; s=1,\dots,S_l}$ and $\tilde{\theta}_{\tilde{d}} = \prod_{l=1}^L \prod_{s=1}^{S_l} \Theta_{l,d_l^{(s)}}$. We first introduce the following technical lemma which can be easily proved by induction:

Lemma 1. Given $0 \leq a_i, b_i \leq 1$ and $|a_i - b_i| \leq \epsilon$ for all $i = 1, \dots, N$, then $|\prod_{i=1}^N a_i - \prod_{i=1}^N b_i| \leq N\epsilon$.

Since $\mathbb{P}(\tilde{\theta} \notin \cup_{\tilde{d} \in \tilde{D}} \tilde{\theta}_{\tilde{d}}) \leq \sum_{l=1,\dots,L; s=1,\dots,S_l} \mathbb{P}(\theta_l^{(s)} \notin \tilde{\Theta}_l) \leq NL\epsilon_1$, we can now rewrite $\tilde{H}(y|\mathbf{x})$ as

$$\tilde{H}(y|\mathbf{x}) = \sum_{\tilde{d} \in \tilde{D}} \int_{\tilde{\theta}_{\tilde{d}}} \prod_{i=1}^N H(y_i|\mathbf{x}_i, \theta_i) dG(\tilde{\theta}) + NL\mathcal{O}_1(\epsilon_1), \quad (25)$$

and $\tilde{H}^*(y|\mathbf{x})$ as

$$\tilde{H}^*(y|\mathbf{x}) = \sum_{\tilde{d} \in \tilde{D}} \int_{\tilde{\theta}_{\tilde{d}}} \prod_{i=1}^N H(y_i|\mathbf{x}_i, \theta_{d^{(c(i))}}^*) dG(\tilde{\theta}), \quad (26)$$

where $0 \leq \mathcal{O}_1(\epsilon_1) \leq \epsilon_1$, and we have the following approximation error bound between $\tilde{H}(y|\mathbf{x})$ and $\tilde{H}^*(y|\mathbf{x})$:

$$|\tilde{H}^*(y|\mathbf{x}) - \tilde{H}(y|\mathbf{x})| \leq \sum_{\tilde{d} \in \tilde{D}} \int_{\tilde{\theta}_{\tilde{d}}} \left| \prod_{i=1}^N H(y_i|\mathbf{x}_i, \theta_{d^{(c(i))}}^*) - \prod_{i=1}^N H(y_i|\mathbf{x}_i, \theta_i) \right| dG(\tilde{\theta}) + NL\epsilon_1$$

$$\leq N\epsilon_2 + NL\epsilon_1, \quad (27)$$

where the second inequality is resulted from Lemma 1.

A.2. Step 2: Approximating Eq. (21) by Eq. (22)

Partition the space of w_{il} (i.e. \mathbb{R}) into $D_l + 2$ adjacent half open half closed intervals $\mathcal{W}_{l,d_l} = (w_{l,d_l-1}^*, w_{l,d_l}^*]$ (for $d_l = 1, \dots, D_l$), $\mathcal{W}_{l,0} = (-\infty, w_{l,0}^*]$ and $\mathcal{W}_{l,D_l+1} = (w_{l,D_l}^*, \infty)$ such that $\Phi_l(\mathcal{W}_{l,d_l}) = G_l(\Theta_{l,d_l})$ for every $d_l = 1, \dots, D_l$. Note that such a partitioning always exists as Φ_l is a continuous distribution. Also denote the domain of (the L -vector) $\mathbf{d} = (d_1, \dots, d_L)$ as $\mathcal{D} = \prod_{l=1}^L \{1, \dots, D_l\}$ and the corresponding extended domain $\mathcal{D}^+ = \prod_{l=1}^L \{0, 1, \dots, D_l + 1\}$.

Denote $\tilde{F}^{**(\mathbf{u})}(y|\mathbf{x})$ as the form of Eq. (22) with $\xi_d^{(\mathbf{u})}(\mathbf{w}_i)$ given by $\xi_d^{(\mathbf{u})}(\mathbf{w}_i) = \exp\{u(\tilde{\beta}_{d,0} + \tilde{\beta}_d^T \mathbf{w}_i)\} / \sum_{d' \in \mathcal{D}^+} \exp\{u(\tilde{\beta}_{d',0} + \tilde{\beta}_{d'}^T \mathbf{w}_i)\}$. (28)

We have the following technical lemma to help us choose suitable parameters $\{(\tilde{\beta}_{d,0}, \tilde{\beta}_d)\}_{d \in \mathcal{D}^+}$ of $\xi_d^{(\mathbf{u})}(\mathbf{w}_i)$:

Lemma 2. There exists parameters $\{(\tilde{\beta}_{d,0}, \tilde{\beta}_d)\}_{d \in \mathcal{D}^+}$ of $\xi_d^{(\mathbf{u})}(\mathbf{w}_i)$ such that $\xi_d^{(\mathbf{u})}(\mathbf{w}_i) \xrightarrow{u \rightarrow \infty} \prod_{l=1}^L 1_{w_{il}}^*(\mathcal{W}_{l,d_l})$ for every $\mathbf{d} \in \mathcal{D}^+$, where $1_w^*(\mathcal{W})$ is an indicator which equals to 1 if w falls inside the interior of \mathcal{W} , equals to a constant $c \in [0, 1]$ if w is on the boundary of \mathcal{W} , and equals to 0 otherwise.

Proof. With a slight (notational) variant of Lemma 3.1 of Fung et al. [30], one can easily show the existence of parameters $\{(\tilde{\beta}_{d,0}, \tilde{\beta}_d)\}_{d \in \mathcal{D}^+}$ such that $\tilde{\beta}_{d,0} + \tilde{\beta}_d^T \mathbf{w}_i > \max_{d' \neq d, d' \in \mathcal{D}^+} \tilde{\beta}_{d',0} + \tilde{\beta}_{d'}^T \mathbf{w}_i$ if and only if $w_{il} \in \mathcal{W}_{l,d_l}^*$ for every $\mathbf{d} \in \mathcal{D}^+$, where \mathcal{W}_{l,d_l}^* is the interior of \mathcal{W}_{l,d_l} . Under a slight variation of Lemma 3.2 of Fung et al. [30], we have $\xi_d^{(\mathbf{u})}(\mathbf{w}_i) \xrightarrow{u \rightarrow \infty} 1\{\tilde{\beta}_{d,0} + \tilde{\beta}_d^T \mathbf{w}_i > \max_{d' \neq d, d' \in \mathcal{D}^+} \tilde{\beta}_{d',0} + \tilde{\beta}_{d'}^T \mathbf{w}_i\} + \mathcal{O}_1(1) \times 1\{\tilde{\beta}_{d,0} + \tilde{\beta}_d^T \mathbf{w}_i = \max_{d' \neq d, d' \in \mathcal{D}^+} \tilde{\beta}_{d',0} + \tilde{\beta}_{d'}^T \mathbf{w}_i\} = \prod_{l=1}^L 1_{w_{il}}^*(\mathcal{W}_{l,d_l})$, where $0 \leq \mathcal{O}_1(1) \leq 1$. \square

Choosing the parameters indicated by the above lemma, we have the following approximation result between $\tilde{H}^*(y|\mathbf{x})$ and $\tilde{F}^{**(\mathbf{u})}(y|\mathbf{x})$:

$$\tilde{F}^{**(\mathbf{u})}(y|\mathbf{x}) \xrightarrow{u \rightarrow \infty} \int_{\mathbf{w}} \sum_{\tilde{d} \in \tilde{D}^+} \left(\prod_{l=1}^L 1_{w_{il}}^*(\mathcal{W}_{l,d_l}) \right) H(y_i|\mathbf{x}_i, \theta_{d^{(c(i))}}^*) d\Phi(\mathbf{w})$$

$$\begin{aligned} & \pi_j^{(t,u)}(\mathbf{x}_i, \mathbf{w}_i; \tilde{\alpha}_q^O, \tilde{\beta}_d) \\ &= \frac{\exp\{\log H(\mathcal{Y}_m^M | \mathbf{x}_q^{*Q}, \theta_d^*) + t(\tilde{\alpha}_{q,0}^O + \tilde{\alpha}_q^{OT} \mathbf{x}_i) + u(\tilde{\beta}_{d,0} + \tilde{\beta}_d^T \mathbf{w}_i)\}}{\sum_{m' \in \mathcal{M}} \sum_{q' \in \mathcal{Q}} \sum_{d' \in \mathcal{D}} \exp\{\log H(\mathcal{Y}_{m'}^M | \mathbf{x}_{q'}^{*Q}, \theta_{d'}^*) + t(\tilde{\alpha}_{q',0}^O + \tilde{\alpha}_{q'}^{OT} \mathbf{x}_i) + u(\tilde{\beta}_{d',0} + \tilde{\beta}_{d'}^T \mathbf{w}_i)\}} \end{aligned} \quad (30)$$

$$\begin{aligned} &= \frac{\exp\{t(\tilde{\alpha}_{q,0}^O + \tilde{\alpha}_q^{OT} \mathbf{x}_i)\}}{\sum_{q' \in \mathcal{Q}} \exp\{t(\tilde{\alpha}_{q',0}^O + \tilde{\alpha}_{q'}^{OT} \mathbf{x}_i)\}} \frac{\exp\{u(\tilde{\beta}_{d,0} + \tilde{\beta}_d^T \mathbf{w}_i)\}}{\sum_{d' \in \mathcal{D}} \exp\{u(\tilde{\beta}_{d',0} + \tilde{\beta}_{d'}^T \mathbf{w}_i)\}} H(\mathcal{Y}_m^M | \mathbf{x}_q^{*Q}, \theta_d^*) \\ &:= \gamma_q^{(t)}(\mathbf{x}_i) \times \xi_d^{(u)}(\mathbf{w}_i) \times H(\mathcal{Y}_m^M | \mathbf{x}_q^{*Q}, \theta_d^*), \end{aligned} \quad (31)$$

Box 1.

$$\begin{aligned} &= \int \tilde{\mathcal{W}} \prod_{i=1}^N \sum_{d \in \mathcal{D}^+} \left(\prod_{i=1}^L 1_{\mathcal{W}_{i,d}}^*(\mathcal{W}_{i,d}) \right) H(y_i | \mathbf{x}_i, \theta_d^*) d\Phi(\mathbf{w}) + \mathcal{O}_1(\epsilon_1) \\ &= \sum_{d \in \mathcal{D}} \int \mathcal{W}_d \prod_{i=1}^N \sum_{d \in \mathcal{D}} \left(\prod_{i=1}^L 1_{\mathcal{W}_{i,d}}^*(\mathcal{W}_{i,d}) \right) H(y_i | \mathbf{x}_i, \theta_d^*) d\Phi(\mathbf{w}) + \mathcal{O}_1(\epsilon_1) \\ &= \sum_{d \in \mathcal{D}} \int \mathcal{W}_d \prod_{i=1}^N H(y_i | \mathbf{x}_i, \theta_{d^{(c(i))}}^*) d\Phi(\mathbf{w}) + \mathcal{O}_1(\epsilon_1) \\ &= \sum_{d \in \mathcal{D}} \prod_{i=1}^N H(y_i | \mathbf{x}_i, \theta_{d^{(c(i))}}^*) G(\tilde{\Theta}_d) + \mathcal{O}_1(\epsilon_1) = \tilde{H}^*(\mathbf{y} | \mathbf{x}) + \mathcal{O}_1(\epsilon_1), \end{aligned} \quad (29)$$

where $\tilde{\mathcal{W}}_l^{(s)} = \cup_{d=1}^{D_l} \mathcal{W}_{l,d}$, $\tilde{\mathcal{W}} = \prod_{l=1}^L \prod_{s=1}^{S_l} \tilde{\mathcal{W}}_l^{(s)}$ and $\mathcal{W}_d = \prod_{l=1}^L \prod_{s=1}^{S_l} \tilde{\mathcal{W}}_{l,d}^{(s)}$. The convergence is resulted from the Dominated Convergence Theorem (DCT), which is obviously uniform on (\mathbf{y}, \mathbf{x}) as $\xi_d^{(u)}(\mathbf{w}_i)$ (the only term in $\tilde{F}^{**(\mathbf{y} | \mathbf{x})}$ which depends on u) does not depend on (\mathbf{y}, \mathbf{x}) and $H(y_i | \mathbf{x}_i, \theta_d^*)$ is bounded above at 1. The third last equality holds as the events of the indicator functions only (jointly) hold if $d = d^{(c(i))}$ when $\mathbf{w} \in \mathcal{W}_d$. The second last equality holds as the integrand is constant on $\mathbf{w} \in \mathcal{W}_d$ and $\Phi(\mathcal{W}_d) = \prod_{l=1}^L \prod_{s=1}^{S_l} \Phi_l(\mathcal{W}_{l,d}^{(s)}) = \prod_{l=1}^L \prod_{s=1}^{S_l} G_l(\Theta_{l,d}^{(s)}) = G(\tilde{\Theta}_d)$.

A.3. Step 3: Approximating Eq. (22) by Eq. (23)

Partition the space of $y_{ik} \in \mathbb{R}$ into adjacent half open half closed intervals $\{\mathcal{Y}_{k,m}^M = ((m-1/2)h_M^y, (m+1/2)h_M^y]\}_{m=-M, \dots, M-1, M}$. Define $\mathcal{Y}^M := \prod_{k=1}^K \cup_{m=-M}^M \mathcal{Y}_{k,m}^M := \prod_{k=1}^K \mathcal{Y}_k^M = (-M+1/2)h_M^y, (M+1/2)h_M^y]^K$, $\mathcal{Y}_{k,-(M+1)}^M = (-\infty, -(M+1/2)h_M^y]$ and $\mathcal{Y}_{k,M+1}^M = ((M+1/2)h_M^y, \infty)$. Choose h_M^y such that $h_M^y \downarrow 0$ and $Mh_M^y \uparrow \infty$ as $M \rightarrow \infty$. Also denote $\mathcal{Y}_m^M = \mathcal{Y}_{1,m_1}^M \times \dots \times \mathcal{Y}_{K,m_K}^M$ as a hypercube with $\mathbf{m} := (m_1, \dots, m_K) \in \mathcal{M} := \{-M+1, \dots, M+1\}^K$.

Similarly, partition the space of $x_{ip} \in \mathbb{R}$ into adjacent half open half closed intervals $\{\mathcal{X}_{p,q}^Q = ((q-1/2)h_Q^x, (q+1/2)h_Q^x]\}_{q=-Q, \dots, Q-1, Q}$. Define $\mathcal{X}^Q := \prod_{p=1}^P \cup_{q=-Q}^Q \mathcal{X}_{p,q}^Q := \prod_{p=1}^P \mathcal{X}_p^Q = (-Q+1/2)h_Q^x, (Q+1/2)h_Q^x]^P$, $\mathcal{X}_{p,-(Q+1)}^Q = (-\infty, -(Q+1/2)h_Q^x]$ and $\mathcal{X}_{p,Q+1}^Q = ((Q+1/2)h_Q^x, \infty)$. Choose h_Q^x such that $h_Q^x \downarrow 0$ and $Qh_Q^x \uparrow \infty$ as $Q \rightarrow \infty$. Also denote $\mathcal{X}_q^Q = \mathcal{X}_{1,q_1}^Q \times \dots \times \mathcal{X}_{P,q_P}^Q$ as a hypercube with $\mathbf{q} := (q_1, \dots, q_P) \in \mathcal{Q} := \{-Q+1, \dots, Q+1\}^P$.

Denote $\tilde{F}^{*(M,Q,t,u)}(\mathbf{y} | \mathbf{x})$ as the form of Eq. (23), where $[\mathbf{y}]_m^M := ([y]_{m_1}^M, \dots, [y]_{m_K}^M)$ (inside the expression of $\tilde{F}^{*(M,Q,t,u)}(\mathbf{y} | \mathbf{x}_i, \mathbf{w}_i)$) represents the leftmost vertex of the hypercube \mathcal{Y}_m^M . Also, the function $\pi_j^{(t,u)}(\mathbf{x}_i, \mathbf{w}_i; \tilde{\alpha}_q^O, \tilde{\beta}_d)$ is a logit linear gating function given by Eq. (31) (see Box 1), where $\mathbf{x}_q^{*Q} = (x_{1,q_1}^{*Q}, \dots, x_{P,q_P}^{*Q})$ is the mid-point of the hypercube \mathcal{X}_q^Q . Here, "mid-points" for $\mathcal{X}_{p,-(Q+1)}^Q$ and $\mathcal{X}_{p,(Q+1)}^Q$ are respectively defined as $-(Q+1)h_Q^x$ and $(Q+1)h_Q^x$.

Further denote $R^Q(\mathbf{x}_i) = \{q : \forall p \text{ we have } |p_j^Q(x_{ip}) - x_{p,q_p}^{*Q}| \leq h_Q^x\}$ and $R'^Q(\mathbf{x}_i) = \{q : \exists p \text{ such that } |p_j^Q(x_{ip}) - x_{p,q_p}^{*Q}| > h_Q^x\}$, where $p_j^Q(x_{ip})$ is the projection of x_{ip} on the interval \mathcal{X}_p^Q (i.e. $p_j^Q(x_{ip})$ is the point in \mathcal{X}_p^Q where x_{ip} is closest to). To derive the approximation results, we first introduce the following two technical lemmas.

Lemma 3. Given any fixed Q and h_Q^x , there exists parameters $\{\tilde{\alpha}_{q,0}^O, \tilde{\alpha}_q^{OT}\}_{q \in \mathcal{Q}}$, such that for any $\epsilon_3 > 0$, we have $\sum_{q \in R^Q(\mathbf{x}_i)} \gamma_q^{(t)}(\mathbf{x}_i) \leq \epsilon_3$ for all $\mathbf{x}_i \in \mathcal{X}$ with sufficiently large t .

Proof. This follows directly from the proof of Theorem 3.2 of Fung et al. [30]. \square

Lemma 4. The probability distributions $\{H(y_i | \mathbf{x}_i, \theta_i)\}_{\mathbf{x}_i \in \tilde{\mathcal{X}}, \theta_i \in \tilde{\Theta}}$ are tight for any compact spaces $\tilde{\mathcal{X}}$ and $\tilde{\Theta}$.

Proof. Divide the compact space $\tilde{\mathcal{Z}} := \tilde{\mathcal{X}} \times \tilde{\Theta}$ into D subspaces $\{\mathcal{Z}_d\}_{d=1, \dots, D}$, where each subspace \mathcal{Z}_d is small enough to be covered by a ball with radius δ . Define $\mathbf{z}_d^* := (\mathbf{x}_d^*, \theta_d^*)$ as an arbitrary interior point of \mathcal{Z}_d for $d = 1, \dots, D$. For each $d = 1, \dots, D$, we choose a response space $\tilde{\mathcal{Y}}_d \in \mathcal{Y}$ such that $H(\tilde{\mathcal{Y}}_d | \mathbf{x}_d^*, \theta_d^*) \geq 1 - \epsilon/2$. Select a compact space $\tilde{\mathcal{Y}}^*$ which covers all $\{\tilde{\mathcal{Y}}_d\}_{d=1, \dots, D}$, and we have $H(\tilde{\mathcal{Y}}^* | \mathbf{x}_d^*, \theta_d^*) \geq 1 - \epsilon/2$ true for all $d = 1, \dots, D$. Uniform continuity of H on any compact space implies that $|H(\tilde{\mathcal{Y}}^* | \mathbf{x}_d^*, \theta_d^*) - H(\tilde{\mathcal{Y}}^* | \mathbf{x}_i, \theta_i)| \leq \epsilon/2$ for sufficient small δ if $(\mathbf{x}_i, \theta_i) \in \mathcal{Z}_d$. Overall, we have $H(\tilde{\mathcal{Y}}^* | \mathbf{x}_i, \theta_i) \geq 1 - \epsilon$, and hence the result follows. \square

As a result, for any compact covariates space of $\tilde{\mathcal{X}}$ and any $\epsilon_4 > 0$, we can use Lemma 4 to select a rectangular output space $\tilde{\mathcal{Y}}$ such that $H(\tilde{\mathcal{Y}} | \mathbf{x}_i, \theta_i) \geq 1 - \epsilon_4$ for any $\mathbf{x}_i \in \tilde{\mathcal{X}}$ and $\theta_i \in \tilde{\Theta}$, where $\tilde{\Theta} = \prod_{l=1}^L \tilde{\Theta}_l$ and note that $\tilde{\Theta}_l$ is defined in Appendix A.1. Then we have for all $\mathbf{x}_i, \mathbf{x}_q^{*Q} \in \tilde{\mathcal{X}}$ and $\theta_d^* \in \tilde{\Theta}$:

$$|H([\mathbf{y}]_m^M | \mathbf{x}_q^{*Q}, \theta_d^*) - H(p_j^M([\mathbf{y}]_m^M) | \mathbf{x}_q^{*Q}, \theta_d^*)| \leq \epsilon_4, \quad (32)$$

and

$$|H(y_i | \mathbf{x}_i, \theta_d^*) - H(p_j^M(y_i) | \mathbf{x}_i, \theta_d^*)| \leq \epsilon_4, \quad (33)$$

where $p_j^M(y_i)$ is the projection of y_i on $\tilde{\mathcal{Y}}$, and $[\mathbf{y}]_m^M := ([y]_{m_1}^M, \dots, [y]_{m_K}^M)$ is the rightmost vertex of hypercube \mathcal{Y}_m^M . Further, because of the uniform continuity of $H(y_i | \mathbf{x}_i, \theta_d^*)$ on $(\mathbf{x}_i, \theta_d^*)$ within a compact support, for any $\epsilon_5 > 0$, we can choose sufficient large Q (to make h_Q^x small while \mathcal{X}^Q covers $\tilde{\mathcal{X}}$) and M (to make h_M^y small while \mathcal{Y}^M covers $\tilde{\mathcal{Y}}$) such that for any $y_i \in \mathcal{Y}_m^M$ and $q \in R^Q(\mathbf{x}_i)$, we have:

$$|H(p_j^M(y_i) | \mathbf{x}_i, \theta_d^*) - H(p_j^M([\mathbf{y}]_m^M) | \mathbf{x}_q^{*Q}, \theta_d^*)| \leq \epsilon_5. \quad (34)$$

Summarizing the above three equations, we have:

$$|H(y_i | \mathbf{x}_i, \theta_d^*) - H([\mathbf{y}]_m^M | \mathbf{x}_q^{*Q}, \theta_d^*)| \leq 2\epsilon_4 + \epsilon_5. \quad (35)$$

Note that $\tilde{F}^{*(M,Q,t,u)}(\mathbf{y} | \mathbf{x}_i, \mathbf{w}_i)$ can be re-written as

$$\begin{aligned} \tilde{F}^{*(M,Q,t,u)}(\mathbf{y} | \mathbf{x}_i, \mathbf{w}_i) &= \sum_{\mathbf{m} \in \mathcal{M}} \sum_{q \in \mathcal{Q}} \sum_{d \in \mathcal{D}^+} \gamma_q^{(t)}(\mathbf{x}_i) \xi_d^{(u)}(\mathbf{w}_i) H([\mathbf{y}]_m^M | \mathbf{x}_q^{*Q}, \theta_d^*) 1\{\mathbf{y}_i \in \mathcal{Y}_m^M\} \\ &= \sum_{\mathbf{m} \in \mathcal{M}} \sum_{q \in R^Q(\mathbf{x}_i)} \sum_{d \in \mathcal{D}^+} \gamma_q^{(t)}(\mathbf{x}_i) \xi_d^{(u)}(\mathbf{w}_i) H([\mathbf{y}]_m^M | \mathbf{x}_q^{*Q}, \theta_d^*) 1\{\mathbf{y}_i \in \mathcal{Y}_m^M\} + \mathcal{O}_1(\epsilon_3), \end{aligned} \quad (36)$$

where $0 \leq \mathcal{O}_1(\epsilon_3) \leq \epsilon_3$. The last equality is resulted from Lemma 3. The approximation result is:

$$|\tilde{F}^{*(M,Q,t,u)}(\mathbf{y} | \mathbf{x}) - \tilde{F}^{**(\mathbf{y} | \mathbf{x})}|$$

$$\begin{aligned}
 &\leq \int \left| \prod_{i=1}^N F^{*(M,Q,t,u)}(y_i|x_i, \mathbf{w}_i) - \prod_{i=1}^N F^{*(u)}(y_i|x_i, \mathbf{w}_i) \right| d\Phi(\mathbf{w}) \\
 &\leq N \int \left| F^{*(M,Q,t,u)}(y_i|x_i, \mathbf{w}_i) - F^{*(u)}(y_i|x_i, \mathbf{w}_i) \right| d\Phi(\mathbf{w}) \\
 &\leq N \int \sum_{m \in \mathcal{M}} \sum_{q \in R^Q(x_i)} \sum_{d \in D^+} \gamma_q^{(i)}(x_i) \xi_d^{(u)}(\mathbf{w}_i) \\
 &\quad \times \left| H(y_i|x_i, \theta_d^*) - H(\lfloor y_i \rfloor_m^M | x_q^{*Q}, \theta_d^*) \right| 1\{y_i \in \mathcal{Y}_m^M\} d\Phi(\mathbf{w}) + 2\epsilon_3 \\
 &\leq N(2\epsilon_4 + \epsilon_5 + 2\epsilon_3), \tag{37}
 \end{aligned}$$

where the second inequality is resulted from Lemma 1, the third and last inequalities are respectively resulted from Eqs. (35) and (36).

A.4. Step 4: Approximating Eq. (23) by Eq. (24)

Write $\tilde{F}^{(M,Q,t,u,v)}(y|x)$ as in Eq. (24), where $F_0(y_i; \Psi_m^{M(v)})$ in $F^{(M,Q,t,u,v)}$ ($y_i|x_i, \mathbf{w}_i$) is chosen in the way that $F_0(y_i; \Psi_m^{M(v)}) \xrightarrow{D} 1\{y_i \geq \lfloor y_i \rfloor_m^M\}$ as $v \rightarrow \infty$. Due to the distributional convergence as well as \mathcal{M} is a finite set, for any $\epsilon_6 > 0$ we can find a sufficient large v such that for every $m \in \mathcal{M}$, we have:

$$|F_0(y_i; \Psi_m^{M(v)}) - 1\{y_i \geq \lfloor y_i \rfloor_m^M\}| \leq \epsilon_6 + \sum_{k=1}^K 1\{y_{ik} \in \mathcal{L}_k^{\delta^*}(\lfloor y_i \rfloor_m^M)\}, \tag{38}$$

where δ^* is chosen to be $0 < \delta^* < h_M^y/2$, and $\mathcal{L}_k^{\delta^*}(\lfloor y_i \rfloor_m^M)$ represents non-overlapping intervals for $m_k = -(M+1), \dots, (M+1)$ with

$$\mathcal{L}_k^{\delta^*}(\lfloor y_i \rfloor_m^M) = \begin{cases} [\lfloor y_i \rfloor_m^M - \delta^*, \lfloor y_i \rfloor_m^M + \delta^*], & \text{if } m_k > -(M+1) \\ (-\infty, \lfloor y_i \rfloor_m^M - 2\delta^*], & \text{if } m_k = -(M+1) \end{cases} \tag{39}$$

Note that the rightmost term in Eq. (38) is to control for the fact that the weak convergence of F_0 to the indicator is not uniform when y_{ik} is close to $\lfloor y_i \rfloor_m^M$. Consider the bound

$$\begin{aligned}
 &|F^{(M,Q,t,u,v)}(y_i|x_i, \mathbf{w}_i) - F^{*(M,Q,t,u)}(y_i|x_i, \mathbf{w}_i)| \\
 &\leq \sum_{m \in \mathcal{M}} \sum_{q \in Q} \sum_{d \in D^+} \gamma_q^{(i)}(x_i) \xi_d^{(u)}(\mathbf{w}_i) H(\mathcal{Y}_m^M | x_q^{*Q}, \theta_d^*) F_0(y_i; \Psi_m^{M(v)}) - 1\{y_i \geq \lfloor y_i \rfloor_m^M\} \\
 &\leq \left\{ \max_{q \in Q, d \in D^+} \sum_{m \in \mathcal{M}} \sum_{k=1}^K H(\mathcal{Y}_m^M | x_q^{*Q}, \theta_d^*) 1\{y_{ik} \in \mathcal{L}_k^{\delta^*}(\lfloor y_i \rfloor_m^M)\} \right\} + \epsilon_6 \\
 &= \sum_{k=1}^K \left\{ \max_{q \in Q, d \in D^+} \sum_{m_i = -(M+1), \dots, (M+1)} H_k(\mathcal{Y}_{k,m_i} | x_q^{*Q}, \theta_d^*) 1\{y_{ik} \in \mathcal{L}_k^{\delta^*}(\lfloor y_i \rfloor_m^M)\} \right\} + \epsilon_6. \tag{40}
 \end{aligned}$$

Since $\mathcal{L}_k^{\delta^*}(\lfloor y_i \rfloor_m^M)$ is non-overlapping for $m_k = -(M+1), \dots, (M+1)$, only one term in the summation of Eq. (40) is non-zero. Since H_k is a continuous distribution, for any $\epsilon_7 > 0$, we have $H_k(\mathcal{Y}_{k,m_i} | x_q^{*Q}, \theta_d^*) \leq \epsilon_7$ for any $q \in Q, d \in D^+$ and $m_k = -(M+1), \dots, (M+1)$ given that M is sufficiently large. Finally, using the same proof idea as Eq. (37), we have

$$|\tilde{F}^{(M,Q,t,u,v)}(y|x) - \tilde{F}^{*(M,Q,t,u)}(y|x)| \leq N(K\epsilon_7 + \epsilon_6). \tag{41}$$

In summary, based on Eqs. (27), (29) (37) and (41), Theorem 1 holds because for sufficiently large M, Q, t and v , the following inequality holds uniformly for each x_i falling into a compact covariates space:

$$|\tilde{F}^{(M,Q,t,u \rightarrow \infty,v)}(y|x) - \tilde{H}(y|x)| \leq (N\epsilon_2 + \epsilon_1) + \mathcal{O}_1(\epsilon_1) + N(2\epsilon_4 + \epsilon_5 + 2\epsilon_3) + N(K\epsilon_7 + \epsilon_6), \tag{42}$$

where ϵ_1 to ϵ_7 can be chosen to be arbitrarily small, and any parameters chosen in Steps 1 to 4 are independent of $N, \mathcal{S} = (\mathcal{S}_1, \dots, \mathcal{S}_L)$ and $c(\cdot)$.

Appendix B. Proof of Theorem 2

Other than the notational differences, the proof ideas of Theorems 1 and 2 are substantially similar. Precisely, the 4-step framework used to prove Theorem 1 also applies to prove Theorem 2. As a result, we

only present a sketch proof of Theorem 2, with an emphasis on the key differences of proof techniques between the two theorems. Unless specified otherwise, the notations adopted in this proof section is the same as those defined in Appendix A. Analogous to Eqs. (20) and (21), in Step 1 we examine an approximation bound between the following two equations:

$$\tilde{H}(y|x) = \int_{\tilde{\Theta}} \left[\prod_{i \in \mathcal{I}} H(y_i|x_i, \theta_i) \right] dG(\tilde{\Theta}), \tag{43}$$

and

$$\tilde{H}^*(y|x) = \sum_{\tilde{d} \in \tilde{D}} \prod_{i \in \mathcal{I}} H(y_i|x_i, \theta_{d^{(i)}}^*) G(\tilde{\Theta}_{\tilde{d}}). \tag{44}$$

Similar to Appendix A.1, we choose compact spaces $\{\tilde{\Theta}_l \in \Theta\}_{l=1, \dots, L}$ such that $\mathbb{P}(\cap_{l=1}^L \cap_{i \in \mathcal{I}_l} \{\theta_{li} \in \tilde{\Theta}_l\}) \geq 1 - NL\epsilon_1$, where $\mathcal{I}_l = \{i : i_1 = 1, \dots, N_0; \dots; i_l = 1, \dots, N_{l-1}\}$. Also, partition granular subspaces $\{\Theta_{l,d_l}\}_{d_l=1, \dots, D_l}$ of $\tilde{\Theta}_l$ and define the interior points $\{\theta_{l,d_l}^*\}_{d_l=1, \dots, D_l}$ analogous to Appendix A.1. Define $\tilde{D} = \prod_{l=1}^L \prod_{i \in \mathcal{I}_l} \mathcal{D}_l^{(i)}$ with $\mathcal{D}_l^{(i)} = \{1, \dots, D_l\}$, $\tilde{d} = \{d_l^{(i)}\}_{l=1, \dots, L; i \in \mathcal{I}_l}$ with $d_l^{(i)} \in \mathcal{D}_l^{(i)}$, $d^{(L)} = \{d_l^{(i)}\}_{l=1, \dots, L}$, and $\tilde{\Theta}_{\tilde{d}} = \prod_{l=1}^L \prod_{i \in \mathcal{I}_l} \Theta_{l,d_l^{(i)}}^*$. Then, using the same idea as Eqs. (25) to (27), we obtain an arbitrarily small approximation error bound between Eqs. (43) and (44).

In Step 2, we define the following function analogous to Eq. (22) and derive its error bound for approximating Eq. (44):

$$\tilde{F}^{*(u)}(y|x) = \int_{i \in \mathcal{I}} \prod_{i \in \mathcal{I}} F^{*(u)}(y_i|x_i, \mathbf{w}_i) d\Phi(\mathbf{w}) \tag{45}$$

with $F^{*(u)}(y_i|x_i, \mathbf{w}_i) = \sum_{d \in D^+} \xi_d^{(u)}(\mathbf{w}_i) H(y_i|x_i, \theta_d^*)$. Here, we choose $\xi_d^{(u)}(\mathbf{w}_i)$ in a different way as Eq. (28), which is crucial to cater for the dependencies of random effects across levels, as follows:

$$\xi_d^{(u)}(\mathbf{w}_i) = \exp\left\{ \sum_{l=1}^L u^{1/l} (\tilde{\beta}_{d_l,0} + \tilde{\beta}_{d_l,1} w_{il}) \right\} / \sum_{d' \in D^+} \exp\left\{ \sum_{l=1}^L u^{1/l} (\tilde{\beta}_{d'_l,0} + \tilde{\beta}_{d'_l,1} w_{il}) \right\}, \tag{46}$$

where $D = \prod_{l=1}^L \{1, \dots, D_l\}$ and $D^+ = \prod_{l=1}^L \{0, 1, \dots, D_l + 1\}$ are exactly the same as those defined in Appendix A.2, and $d_l = (d_1, \dots, d_l)$ and $d'_l = (d'_1, \dots, d'_l)$ with $\mathbf{d} = \mathbf{d}_L$ (and $\mathbf{d}' = \mathbf{d}'_L$).

In contrast to Appendix A.2, we construct intervals \mathcal{W}_{l,d_l} such that $\cup_{d_l=0}^{D_l+1} \mathcal{W}_{l,d_l} = \mathbb{R}$ for any $d_{l-1} \in \prod_{l=1}^{l-1} \{1, \dots, D_{l'}\}$ and $\Phi_l(\mathcal{W}_{l,d_l}) = G_l(\Theta_{l,d_l} | \Theta_{1,d_1}, \dots, \Theta_{l-1,d_{l-1}})$, which represents the probability of level- l random effect belongs to Θ_{l,d_l} conditioned on the corresponding upper level random effects fall into $(\Theta_{1,d_1}, \dots, \Theta_{l-1,d_{l-1}})$. The following lemma is analogous to Lemma 2:

Lemma 5. *There exists parameters $\{(\tilde{\beta}_{d_l,0}, \tilde{\beta}_{d_l,1})\}_{d \in D^+, l=1, \dots, L}$ of $\xi_d^{(u)}(\mathbf{w}_i)$ such that $\xi_d^{(u)}(\mathbf{w}_i) \xrightarrow{u \rightarrow \infty} \prod_{l=1}^L 1_{w_{il}}(\mathcal{W}_{l,d_l})$ for every $d \in D^+$.*

Proof. Similar to the proof of Lemma 2, we choose suitable parameters such that $\tilde{\beta}_{d_l,0} + \tilde{\beta}_{d_l,1} w_{il} > \max_{d'_l \neq d_l} \tilde{\beta}_{d'_l,0} + \tilde{\beta}_{d'_l,1} w_{il}$ if and only if $w_{il} \in \mathcal{W}_{l,d_l}^*$ for every $d \in D^+$ and $l = 1, \dots, L$, where \mathcal{W}_{l,d_l}^* is the interior of \mathcal{W}_{l,d_l} and $d_l^* = (d_1, \dots, d_{l-1}, d_l^*)$. Observe the expression $\sum_{l=1}^L u^{1/l} (\tilde{\beta}_{d_l,0} + \tilde{\beta}_{d_l,1} w_{il})$ in Eq. (46) where the term corresponding to a higher level factor dominates when u is large. Therefore, for sufficient large u , we have $\sum_{l=1}^L u^{1/l} (\tilde{\beta}_{d_l,0} + \tilde{\beta}_{d_l,1} w_{il}) > \max_{d' \neq d} \sum_{l=1}^L u^{1/l} (\tilde{\beta}_{d'_l,0} + \tilde{\beta}_{d'_l,1} w_{il})$ for \mathbf{w}_i satisfying $w_{il} \in \mathcal{W}_{l,d_l}^*$ for every $l = 1, \dots, L$. Following the same logic as Lemma 2, the result follows. \square

The approximation bound between Eqs. (44) and (45) can be obtained using the same logic as that outlined by Eq. (29), where we further note that $\Phi(\mathcal{W}_{\tilde{d}}) = \prod_{l=1}^L \prod_{i \in \mathcal{I}_l} \Phi_l(\mathcal{W}_{l,d_l^{(i)}}) = \prod_{l=1}^L \prod_{i \in \mathcal{I}_l} G_l(\Theta_{l,d_l^{(i)}} | \Theta_{1,d_1^{(i_1)}}, \dots, \Theta_{l-1,d_{l-1}^{(i_{l-1})}}) = G(\tilde{\Theta}_{\tilde{d}})$ with $\mathcal{W}_{\tilde{d}} = \prod_{l=1}^L \prod_{i \in \mathcal{I}_l} \mathcal{W}_{l,d_l^{(i)}}^*$.

Step 3 and 4 involve evaluations of the following expressions in analogous to Eqs. (23) and (24):

$$\tilde{F}^{*(M,Q,t,u)}(y|x) = \int_{i \in \mathcal{I}} \prod_{i \in \mathcal{I}} F^{*(M,Q,t,u)}(y_i|x_i, \mathbf{w}_i) d\Phi(\mathbf{w}) \tag{47}$$

$$\text{with } F^{*(M,Q,t,u)}(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}_i) = \sum_{m \in \mathcal{M}} \sum_{q \in \mathcal{Q}} \sum_{d \in \mathcal{D}^+} \pi_j^{(t,u)}(\mathbf{x}_i, \mathbf{w}_i; \tilde{\alpha}_q^Q, \tilde{\beta}_d) 1_{\{\mathbf{y}_i \geq \lfloor \mathbf{y} \rfloor_m^M\}},$$

$$\tilde{F}^{(M,Q,t,u,v)}(\mathbf{y}|\mathbf{x}) = \int \prod_{i \in \mathcal{I}} F^{(M,Q,t,u,v)}(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}_i) d\Phi(\mathbf{w}) \quad (48)$$

with $F^{(M,Q,t,u,v)}(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}_i) = \sum_{m \in \mathcal{M}} \sum_{q \in \mathcal{Q}} \sum_{d \in \mathcal{D}^+} \pi_j^{(t,u)}(\mathbf{x}_i, \mathbf{w}_i; \tilde{\alpha}_q^Q, \tilde{\beta}_d) F_0(\mathbf{y}_i; \boldsymbol{\psi}_m^{M(v)})$. The derivation techniques here are exactly the same as those presented in [Appendices A.3](#) and [A.4](#), so this part of the proof is omitted.

Data availability

No data was used for the research described in the article.

References

- [1] R.A. Jacobs, M.I. Jordan, S.J. Nowlan, G.E. Hinton, Adaptive mixtures of local experts, *Neural Comput.* 3 (1) (1991) 79–87.
- [2] M.I. Jordan, R.A. Jacobs, Hierarchical mixtures of experts and the EM algorithm, *Neural Comput.* 6 (2) (1994) 181–214.
- [3] G.J. McLachlan, D. Peel, *Finite Mixture Models*, John Wiley & Sons, 2000.
- [4] S.E. Yuksel, J.N. Wilson, P.D. Gader, Twenty years of mixture of experts, *IEEE Trans. Neural Netw. Learn. Syst.* 23 (8) (2012) 1177–1193.
- [5] S. Masoudnia, R. Ebrahimpour, Mixture of experts: A literature survey, *Artif. Intell. Rev.* 42 (2) (2014) 275–293.
- [6] H.D. Nguyen, F. Chamroukhi, Practical and theoretical aspects of mixture-of-experts modeling: An overview, *Wiley Interdiscip. Reviews: Data Min. Knowl. Discov.* 8 (4) (2018) e1246.
- [7] L. Xu, M.I. Jordan, G.E. Hinton, An alternative model for mixtures of experts, *Adv. Neural Inf. Process. Syst.* (1995) 633–640.
- [8] S. Ingrassia, S.C. Minotti, G. Vittadini, Local statistical modeling via a cluster-weighted approach with elliptical distributions, *J. Classification* 29 (3) (2012) 363–401.
- [9] J. Geweke, M. Keane, Smoothly mixing regressions, *J. Econometrics* 138 (1) (2007) 252–290.
- [10] C. Riquelme, J. Puigcerver, B. Mustafa, M. Neumann, R. Jenatton, A. Susano Pinto, D. Keysers, N. Houlsby, Scaling vision with sparse mixture of experts, *Adv. Neural Inf. Process. Syst.* 34 (2021) 8583–8595.
- [11] H. Hazimeh, Z. Zhao, A. Chowdhery, M. Sathiamoorthy, Y. Chen, R. Mazumder, L. Hong, E. Chi, Dselect-k: Differentiable selection in the mixture of experts with applications to multi-task learning, *Adv. Neural Inf. Process. Syst.* 34 (2021) 29335–29347.
- [12] M.I. Jordan, R.A. Jacobs, Hierarchies of adaptive experts, in: *Advances in Neural Information Processing Systems*, 1992, pp. 985–992.
- [13] H.D. Nguyen, G.J. McLachlan, Laplace mixture of linear experts, *Comput. Statist. Data Anal.* 93 (2016) 177–191.
- [14] F. Chamroukhi, Robust mixture of experts modeling using the t distribution, *Neural Netw.* 79 (2016) 20–36.
- [15] F. Chamroukhi, Skew t mixture of experts, *Neurocomputing* 266 (2017) 390–408.
- [16] T.C. Fung, A.L. Badescu, X.S. Lin, A new class of severity regression models with an application to IBNR prediction, *N. Am. Actuar. J.* 25 (2) (2021) 206–231.
- [17] B. Grun, F. Leisch, FlexMix version 2: Finite mixtures with concomitant variables and varying and constant parameters, *J. Stat. Softw.* 28 (4) (2008) 1–35.
- [18] T.C. Fung, A.L. Badescu, X.S. Lin, A class of mixture of experts models for general insurance: Application to correlated claim frequencies, *ASTIN Bull.: J. IAA* 49 (3) (2019) 647–688.
- [19] A.J. Zeevi, R. Meir, V. Maiorov, Error bounds for functional approximation and estimation using mixtures of experts, *IEEE Trans. Inform. Theory* 44 (3) (1998) 1010–1025.
- [20] W. Jiang, M.A. Tanner, On the approximation rate of hierarchical mixtures-of-experts for generalized linear models, *Neural Comput.* 11 (5) (1999) 1183–1198.
- [21] H.D. Nguyen, L.R. Lloyd-Jones, G.J. McLachlan, A universal approximation theorem for mixture-of-experts models, *Neural Comput.* 28 (12) (2016) 2585–2593.
- [22] H.D. Nguyen, F. Chamroukhi, F. Forbes, Approximation results regarding the multiple-output Gaussian gated mixture of linear experts model, *Neurocomputing* 366 (2019) 208–214.
- [23] W. Jiang, M.A. Tanner, Hierarchical mixtures-of-experts for exponential family regression models: Approximation and maximum likelihood estimation, *Ann. Statist.* (1999) 987–1011.
- [24] E.F. Mendes, W. Jiang, On convergence rates of mixtures of polynomial experts, *Neural Comput.* 24 (11) (2012) 3025–3051.
- [25] A. Norets, et al., Approximation of conditional densities by smooth mixtures of regressions, *Ann. Statist.* 38 (3) (2010) 1733–1766.
- [26] A. Norets, J. Pelenis, Posterior consistency in conditional density estimation by covariate dependent mixtures, *Econometric Theory* 30 (3) (2014) 606–646.
- [27] H.D. Nguyen, T. Nguyen, F. Chamroukhi, G.J. McLachlan, Approximations of conditional probability density functions in lebesgue spaces via mixture of experts models, *J. Stat. Distrib. Appl.* 8 (1) (2021) 1–15.
- [28] H. Tijms, *Stochastic Models: An Algorithmic Approach*, John Wiley, 1994.
- [29] L. Breuer, D. Baum, *An Introduction to Queueing Theory and Matrix-Analytic Methods*, Springer Science & Business Media, 2005.
- [30] T.C. Fung, A.L. Badescu, X.S. Lin, A class of mixture of experts models for general insurance: Theoretical developments, *Insurance Math. Econom.* 89 (2019) 111–127.
- [31] H. Goldstein, *Multilevel Statistical Models*, vol. 922, John Wiley & Sons, 2011.
- [32] M. Aitkin, N. Longford, Statistical modelling issues in school effectiveness studies, *J. R. Stat. Soc. Ser. A (General)* 149 (1) (1986) 1–26.
- [33] H. Goldstein, Multilevel mixed linear model analysis using iterative generalized least squares, *Biometrika* 73 (1) (1986) 43–56.
- [34] E.W. Frees, J.-S. Kim, Multilevel model prediction, *Psychometrika* 71 (1) (2006) 79–104.
- [35] G. Molenberghs, G. Verbeke, C.G. Demétrio, A.M. Vieira, A family of generalized linear models for repeated measures with normal and conjugate random effects, *Statist. Sci.* 25 (3) (2010) 325–347.
- [36] J.-P. Boucher, M. Denuit, Fixed versus random effects in Poisson regression models for claim counts: A case study with motor insurance, *ASTIN Bull. J. IAA* 36 (1) (2006) 285–301.
- [37] C. McGilchrist, Estimation in generalized mixed models, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 56 (1) (1994) 61–69.
- [38] M. Davidian, A.R. Gallant, The nonlinear mixed effects model with a smooth random effects density, *Biometrika* 80 (3) (1993) 475–488.
- [39] T.G. Gregoire, O. Schabenberger, A non-linear mixed-effects model to predict cumulative bole volume of standing trees, *J. Appl. Stat.* 23 (2–3) (1996) 257–272.
- [40] B. Bakker, T. Heskes, Task clustering and gating for bayesian multitask learning, *J. Mach. Learn. Res.* 4 (May) (2003) 83–99.
- [41] S.-K. Ng, G. McLachlan, K.K. Yau, A.H. Lee, Modelling the distribution of ischaemic stroke-specific survival time using an EM-based mixture approach with random effects adjustment, *Stat. Med.* 23 (17) (2004) 2729–2744.
- [42] S.-K. Ng, G.J. McLachlan, K. Wang, L. Ben-Tovim Jones, S.-W. Ng, A mixture model with random-effects components for clustering correlated gene-expression profiles, *Bioinformatics* 22 (14) (2006) 1745–1752.
- [43] K.K. Yau, A.H. Lee, A.S. Ng, Finite mixture regression model with random effects: Application to neonatal hospital length of stay, *Comput. Statist. Data Anal.* 41 (3–4) (2003) 359–366.
- [44] S.-K. Ng, G.J. McLachlan, Extension of mixture-of-experts networks for binary classification of hierarchical data, *Artif. Intell. Med.* 41 (1) (2007) 57–67.
- [45] S.-K. Ng, G.J. McLachlan, Mixture models for clustering multilevel growth trajectories, *Comput. Statist. Data Anal.* 71 (2014) 43–51.
- [46] S.C. Tseung, I.W. Chan, T.C. Fung, A.L. Badescu, X.S. Lin, Improving risk classification and ratemaking using mixture-of-experts models with random effects, *J. Risk Insurance* 90 (3) (2023) 789–820.



Tszy Chai Fung is an Assistant Professor in the Maurice R. Greenberg School of Risk Science, Georgia State University (GSU). Prior to joining GSU, he worked as a Postdoctoral Researcher at ETH Zurich during 2020–2021 and earned a Ph.D. Statistics degree at the University of Toronto in 2020. His current research interests include probabilistic neural network, universal approximation theory, statistical modeling, and inference techniques with applications in business and economics problems, including actuarial science and risk management.



Spark C. Tseung is a Ph.D. student in actuarial science at the Department of Statistical Sciences, University of Toronto (U of T). He is interested in applying novel statistical and machine learning methods to classical actuarial problems such as pricing and claims reserving. Prior to joining U of T, he obtained a bachelor's degree in actuarial science from the University of Hong Kong.