# Model-based Inference for Subgroup Analysis

by

Juan Shen

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Statistics)
in The University of Michigan
2014

Doctoral Committee:

Professor Xuming He, Chair
Professor Bhramar Mukherjee
Assistant Professor XuanLong Nguyen
Professor Kerby A. Shedden

To my parents

# ACKNOWLEDGEMENTS

I want to thank many people without whom I would not have had such a wonderful experience as a PhD student.

My deepest gratitude goes to my advisor Xuming He. He is a great advisor with a broad and deep knowledge of statistics, and he provided me with his constant encouragement and advice in many respects. I feel very lucky to have Professor He as my advisor.

I would like to thank my committee members: Professor Bhramar Mukherjee, Professor XuanLong Nguyen, and Professor Kerby Shedden for their valuable time and suggestions on my thesis research and my job search, especially Professor Bhramar Mukherjee for her warm encouragement and support.

I would like to thank the faculty and staff members in the Department of Statistics at the University of Michigan, especially Judy McDonald, George Michailidis, and Elizaveta Levina for their kind help after my transfer to Michigan. Many thanks are given to my fellow students in the department. Particularly, I would like to thank Naveen Narisetty for his support and encouragement in my research, and Huitian Lei, Yiwei Zhang, Jing Ma, Yuan Zhang, Kam Chung Wong, Yingchuan Wang, among others, for their friendship.

I also want to thank the faculty members and fellow students in University of Illinois, including Professor Annie Qu, Professor Feng Liang, Na Cui, and Lu Gan.

Last, but not the least, I would like to thank my beloved parents for their unconditional love and endless support.

# TABLE OF CONTENTS

iv

# LIST OF FIGURES

**Figure**

# LIST OF TABLES

# ABSTRACT

Model-based Inference for Subgroup Analysis

by

Juan Shen

Chair: Professor Xuming He

Subgroup analysis is an important problem in clinical trials. For example, when a new treatment is approved for use, there may be concerns that the efficacy is driven by extreme efficacy in a subgroup only. In recent years, researchers often attempt to identify a potential subgroup with an enhanced treatment effect. In this dissertation, we assume that there exist two potential subgroups in which the subjects react differently to the treatment. We propose a logistic-normal mixture model where the group means as well as the mixing proportions may be covariate-dependent. Testing the existence of subgroups is critical in the mixture model, but requires nonstandard statistical tests. We derive a test based on a small number of $EM$ iterations towards the likelihood, and propose the bootstrap approximation for the critical values of the test. When subgroups exist, the mixture model helps us identify the factors that are associated with the group membership. We apply the proposed method to the Aids Clinical Trials Group 320 study, and demonstrate that the patients with higher values of baseline $CD4$ or $RNA$ tend to benefit significantly more by adding a protease inhibitor to two nucleoside analogues. We also extend our results to the logistic-normal mixture models with unequal variances across subgroups.

# CHAPTER I

# Introduction

## 1.1 The Motivation Example and the Goals

In the Aids Clinical Trials Group 320 study (ACTG320) (*Hammer et al.*, 1997), the efficacy of the treatment of adding a protease inhibitor to two nucleoside analogues to the human immunodeficiency virus type 1 (HIV–1) infection is tested, where the control group receives only the two nucleoside. The goal of the treatment is to increase or to inhibit the decrease of the $CD4$ cell counts, and the outcome is the change of the $CD4$ counts at certain time points. We ask whether there are heterogeneous treatment effects across different subpopulations. Other variables for each subject include age, gender, race, weight, Karnof (Karnofsky performance scale: 100 indicates no evidence of disease, 90 minor symptoms, 80 some symptoms, 70 active work impossible), Ivdrug (IV drug use history: 1 if never, 2 if currently, 3 if previously), Hemophil status, and Priorzdv (months of prior zidovudine therapy, alone or in combination).

The traditional way to assess a new treatment compared to a standard one is based on the summary statistics for the treatment effect over the entire study population. In the ACTG320 study, from *Hammer et al.* (1997), for the response of the $CD4$ cell count changes at the 24th week, the overall estimated mean difference between the treatment group and the control group is 81 cells/mm$^3$, which is statistically highly significant. However, even with the significant mean difference, it is not necessarily

implied that the new treatment works for all patients. In addition, there is a serious concern about protease inhibitor resistance mutations. So we expect a high enough treatment effect within a specified subgroup of patients who receive the new treatment in the future, to compensate for the costs and risks of using the new treatment.

We hope to estimate the treatment effects in different subgroups and their differences simultaneously. Therefore, we propose a mixture model, with the mixing proportions varying through a logistic model on some covariates $X$. Suppose that we have the response $Y$, and we expect that the mean of $Y$ depends on the covariate $Z$, in which the treatment indicator $T$ is included. That is, the density of $Y|X, Z$ is

$$\pi(X^T\gamma)f(Y; Z^T\beta_1, \sigma) + (1 - \pi(X^T\gamma))f(Y; Z^T(\beta_1 + \beta_2), \sigma), \qquad (1.1)$$

where $\pi(a) = \exp(a)/(1 + \exp(a))$, $f(y; \mu, \sigma)$ is the density of the normal distribution with mean $\mu$ and variance $\sigma^2$, and $(\beta_1, \beta_2, \gamma, \sigma) \in \mathbb{R}^{q_1} \times \mathbb{R}^{q_1} \times \mathbb{R}^{q_2} \times \mathbb{R}$ are unknown parameters.

Our goal is to estimate the differential treatment effects in subgroups, that is, the coefficient of the treatment indicator $T$, after we test the existence of the subgroups and reject the null hypothesis : $\beta_2 = 0$ or $\pi(a) \equiv 0$ or 1.

## 1.2 Literature Review

### 1.2.1 Subgroup Analysis in Clinical Trials with Pre-specified Subgroups

In clinical trials, researchers have been interested in the effect of a treatment among a specified subgroup of patients with certain attributes. For example, the treatment effect may not be significant in the whole population but only significant in one or more subgroups, which is very important to discover in clinical trials. In other cases the treatment may result in greater benefits in certain subgroups than others, as often evaluated in a benefit/risk assessment. Often researchers collect $p$-values

about the overall populations and about some subgroups, and use these different $p$-values to claim the existence of subgroups (*Frasure-Smith et al.* (1997) for example). If the subgroups are not pre-specified, there are usually no attempts to account for the issue of multiple testing. Therefore, this kind of evaluation is only considered an exploration without any confirmatory evidence.

In *Song and Chi* (2007), a two-stage testing procedure is provided, where, in the first stage the authors use the combination of the test statistics for the overall effect and the one for the pre-specified subgroup effect to test the null hypothesis of no overall or subgroup treatment effects. If the null hypothesis is rejected, they further test the hypothesis of no treatment effect in the whole population and a subgroup separately. Assume the null hypothesis $H_{01}$ for no treatment effect for the overall population, and the null hypothesis $H_{02}$ for no treatment effect for the targeted subgroup. Let $H_{012} = H_{01} \cap H_{02}$. Let $Z_2$ and $Z_2^*$, independent of $Z_2$, be the standardized test statistics for no treatment effect in the target subgroup and the complimentary one relative to the overall population, respectively. For a given $\alpha$, pre-specify $\alpha_1$ and $\alpha_1^*$ with $\alpha_1 < \alpha < \alpha_1^* \leq 1$, from which the type-1 error at level $\alpha$ is strongly controlled. Their procedure goes as follows.

- Stage 1: test $H_{012}$ at level $\alpha$. Let $p_1$ be the p-value of the test statistic $Z_1 = \sqrt{k}Z_2 + \sqrt{1-k}Z_2^*$ for some $k > 0$.

    - If $p_1 \leq \alpha_1$, $H_{012}$ is rejected.

    - If $p_1 > \alpha_1^*$, $H_{012}$ is not rejected.

    - Otherwise, conduct the subgroup analysis that generates a new p-value $p_2$. Reject $H_{012}$ if and only if $p_2 \leq \alpha_2$.

- Stage 2: if $H_{012}$ is rejected, test $H_{01}$ and $H_{01}$ each at level $\alpha$.

The two-stage test obtains satisfactory power and strongly controls the type-1 error rate.

*Altstein et al.* (2011) uses a mixture of two log-linear models, depending on whether the subject is treatable. The authors assume that the survival time in each group is a log-linear model with a constant population proportion. An *EM* algorithm is used to obtain the estimates. Their simulations show that for the simulated data where the model is well-defined, the results are satisfactory, that is, the coverage of the treatment effect difference parameter is close to the nominal level for the sample sizes of 400 or higher. However, the effect of a covariate, which could be different across subgroups, is not studied. When the subgroups are not distinguishable based on the available covariates, the model is not well-defined. Their studies do not cover this possibility.

Researchers have considered the problem from a Bayesian point of view. Under the assumptions of exchangeability among treatment-covariate interactions and a linear regression model of the response with respect to the treatment, the covariates and their interactions, and with a proper prior distribution, *Dixon and Simon* (1991) derives the posterior distribution of the subset-specific treatment effects. The exchangeability assumption is reasonable for large randomly designed clinical trials. Only binary covariates are included and there is no consideration for the interactions between the covariates. An extension to more general models is discussed in *Simon* (2002), where a proportional hazard model is used to study the treatment-by-gender interaction.

In the aforementioned work, the subgroups are pre-specified by natural factors like gender and certain baseline measurements. However, finding meaningful subgroups is often a critical part of the work.

### 1.2.2   Subgroup Identification in Clinical Trials

*Bonetti and Gelber* (2004) discusses the problem of examining patterns of treatment effects across several overlapping patient subpopulations. According to the value

of a certain covariate, the authors construct overlapping subgroups and estimate the treatment effect within each subgroup. Then, they plot the treatment effect against the covariate to explore the possible interaction. They derive the joint asymptotic distribution of the treatment effects and use it to construct simultaneous confidence bands and to test the null hypothesis of no interaction, that is, all the treatment effects are the same across the overlapping subgroups. The way they divide the data into overlapping subgroups guarantees reasonably large subgroup sizes.

*Song and Pepe* (2004) considers the case where the response is binary and there is a monotone relationship between the treatment effect and a single covariate. The authors propose a procedure to identify a threshold value for that particular variable. The treatment is assigned depending on whether the variable is above the threshold value. Based on this policy, we have an overall mean response rate for each threshold value. The authors introduce a graphical display, called the selection impact curve, that shows the overall mean response rate as a function of the threshold value. The curve is then used to choose the threshold value, and it can also be used to compare the effects of different covariates on the population response rate. The simple dichotomous criteria are discussed in the paper instead of individual decision making; arguing that it is often the case that medical decisions are made by checking whether a variable exceeds a percentile threshold. However, as pointed out in *Cai et al.* (2011), this method is only useful with respect to an overall utility for the whole population but could not provide a treatment choice scheme at a subject-specific level.

*Foster et al.* (2011) proposes the "Virtual Twins" method to identify a subgroup from randomized clinical trial data where the response $Y$ is binary. Let $T$ be the treatment indicator and $X$ be a vector of covariates. Assume that

$$\text{logit}(P(Y = 1|T, X)) = \alpha + \beta T + \gamma h(X) + \theta T \omega(X),$$

where the main term of interest is $\omega(X)$, a large value of which implies an enhanced treatment effect when $\theta > 0$. The "Virtual Twins" method is to predict the response probability for $T = 0, 1$ for each subject by the random forest method. Define $Z$ to be the difference of the probability for $T = 1$ and for $T = 0$. The authors use $Z$ and the covariates $X$ to build a regression or classification tree to define a subgroup $A$ with an enhanced treatment effect. For the regression method, the tree is built with the difference $Z$ as the response directly, and $X$ as the covariates. The subjects with predicted $Z$ above a threshold are defined to be in the estimated subgroup $\hat{A}$. For the classification method, whether the difference $Z$ is above a certain value is used as the response instead of $Z$. A measure $Q(\hat{A})$ for evaluating the performance of $\hat{A}$ is defined to be the difference of the treatment effect in subgroup $\hat{A}$ and the overall treatment effect. To avoid overfitting, methods such as the cross-validation and the bootstrap bias correction are suggested. Drawbacks of this method include that it tends to identify a subgroup when it does not exist and it is not efficient in identifying the important covariates in defining subgroups when the subgroup does exist.

*Cai et al.* (2011) develops a parametric scoring system based on multiple covariates. Suppose that the data contain variables $(Y, T, Z)$, where $Y$ is the response, $T$ is the indicator of treatment $(T = 1)$ or control $(T = 0)$, and $Z$ is a vector of covariates. Let $Y_t$ be the response if a subject is assigned to group $T = t$, $t = 0, 1$. Let $\mu_t(z) = E(Y_t|Z)$, $t = 0, 1$, and the treatment difference $D(Z) = \mu_1(Z) - \mu_0(Z)$. A nonparametric smoothing technique is used for estimating $\mu_1$ and $\mu_0$. The smoothed average treatment effect difference is used for personalized treatment selection, and a global confidence interval is provided for the average treatment effect difference.

*Zhao et al.* (2013) uses the procedure of *Cai et al.* (2011) but without the smoothing step. Suppose that $\hat{D}(Z)$ is an estimator of $D(Z)$. Let $AD(c)$ be the average treatment difference for the subgroup of subjects, that is, $AD(c) = E(Y_1 - Y_0|\hat{D}(Z) \geq c)$. Note that $AD(c)$ can be transformed to a standardized $\tilde{AD}(q) = AD(F^{-1}_{\hat{D}(Z)}(q))$,

where $q$ denotes any quantile levels of $D(Z)$. Let $\hat{AD}(q)$ be an estimator of $\tilde{AD}(q)$ for $q \in (0, 1)$. Then the authors plot $\hat{AD}(q)$ against $q \in (0, 1)$ in a single graph for different scoring systems from different parametric models for estimating the subject-specific treatment differences. To avoid over-fitting, the data are divided into testing data and evaluation data. Some measurements used to compare scoring systems are constructed as a function of $\tilde{AD}(\cdot)$, such as the metric of the area under the curve, which consistently measures the statistic describing the concordance of the true treatment difference and its empirical one. In the work, $\hat{D}(Z)$ is estimated by regression for the treatment group and for the control group separately, followed by a substraction. The regression could be linear or in more general forms. The general procedure is:

1. Build a candidate set of covariates.

2. For each covariate set Z, compute $\hat{D}(Z)$ in the training data, and then calculate the estimates of $\tilde{AD}(q)$ in the evaluation data.

3. Compare the $\tilde{AD}(q)$ curve for each candidate of scoring systems, and choose one that gives the highest curve of $\tilde{AD}(q)$.

4. Based on the chosen scoring system to further determine the subgroup of interest with the score above a threshold value.

The idea of recursively partitioning has been adopted for subgroup analysis. *Su et al.* (2009) introduces an interaction tree (IT) procedure which follows the three major steps of CART (classification and regression trees by Breiman et al., 1984): (1) growing a large initial binary tree by selecting the best split among all the candidate variables and all the possible splitting values of every candidate variable. The criteria is that the resultant interaction is the most significant by $p-$value; (2) Pruning the trees recursively by removing the "weakest link" according to an interaction-complexity measure; (3) determining the best tree size by some validation method.

*Lipkovich et al.* (2011) also use the idea of recursive partitioning to propose a SIDES method (subgroup identification based on differential effect search). Unlike the IT method by *Su et al.* (2009), in each step of SIDES, only the remaining variables are candidate variables for further splitting, and among all the "better" subpopulations by finding the best splitting for each candidate covariates, multiples of them are added to form the "parent subgroup" for the next step. Therefore, the SIDES method refines the subgroup after each step, instead of considering both child nodes. In the later paper of *Lipkovicha and Dmitrienkoa* (2014), a screening step is added after SIDES to better tackle the case when a large number of irrelevant variables are present.

However, in the procedure reviewed above, there is no joint model connecting the response and the grouping structure. Therefore, we introduce the logistic-normal mixture models which allow this connection.

### 1.2.3 Logistic-Normal Mixture Models

The logistic-normal mixture models have been used in various applications. For time series data, *Wong and Li* (2001) proposes the logistic mixture autoregressive with a exogenous variables model (LMARX), which consists of a mixture of two Gaussian transfer function models with the mixing proportions changing over time through a logistic model. Hypothesis testing of the logistic part is carried out by the standard likelihood ratio test when the model is non-degenerated. However, there is no obvious way to perform significance tests about the conditional mean parameters as well as the number of groups.

In *Muthén and Asparouhov* (2009) and *Muthén and Shedden* (1999), models related to logistic mixture models are used for psychological data. In *Muthén and Asparouhov* (2009), a multilevel model is introduced, in which the first level is a logistic-normal mixture model with slopes and intercepts modeled in the second level. The numerical results show that: "level 1 heterogeneity in the form of latent classes

can be mistaken for level 2 heterogeneity in the form of the random effects that are used in conventional two-level regression analysis." Therefore, mixture models allow heterogeneity to be investigated more fully, more correctly attributing different portions of the heterogeneity. In the paper, the authors suggest using the $BIC$ criterion to select the number of groups. Inference under this model can be made with the $EM$ algorithm, as long as the model is not degenerate.

### 1.2.4 Testing the Number of Groups in Mixture Models

It has remained a challenge to avoid over-fitting in mixture models. In *Goeffinet et al.* (1992), for a simple two group normal mixture model with a constant mean in each group and a constant proportion $p$, the null distribution of the likelihood ratio test for equal means are given for each fixed proportion parameter. Assume that

$$g(\boldsymbol{x}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\Sigma}) = p\varphi(\boldsymbol{x}; \boldsymbol{\theta}_1, \boldsymbol{\Sigma}) + (1-p)\varphi(\boldsymbol{x}; \boldsymbol{\theta}_2, \boldsymbol{\Sigma}), \tag{1.2}$$

where $\varphi(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\Sigma})$ is the normal density with mean $\boldsymbol{\theta}$ and covariance matrix $\boldsymbol{\Sigma}$, and $\boldsymbol{x}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2$ are $q \times 1$ vectors. The goal is to test $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$ for a given $p \neq 0$ or 1. When $q = 1$, the limiting distribution of the likelihood ratio test is a $\chi_1^2$ if $p \neq 0.5$ and $\boldsymbol{\Sigma}$ unknown, otherwise it is a mixture of half probability with a point mass at 0 and half probability of $\chi_1^2$. The simulation results suggest that the convergence rate of the likelihood test statistic is poor, especially when $p$ is close to 0.5. For $q = 2$ and known $\Sigma$, the limiting distribution is $0.5\chi_0^2 + M^2$, where $M = M_1 + \sqrt{M_2^2 + M_3^2}$, $M_1$, $M_2$ and $M_3$ are independent standard normal variables. For other cases, no analytical result is known.

Model (1.2) can be extended to any finite number of groups $k$. For $q = 1$, from *Lo et al.* (2001), for testing if the sample is from a $k_0$ component or $k_1$ ($k_1 > k_0$) component normal mixture model, the limiting distribution of the likelihood ratio

test is a weighted sum of chi-squared random variables with one degree of freedom.

In the work of *Naik et al.* (2007), the Akaike Information Criterion (AIC) is extended to this particular problem of deciding the number of components and selecting variables in the mixture models. Due to the "clustering penalty function", the mixture regression criterion (MRC) is shown to yield marked improvement in model selection. There are three terms in the MRC criterion: the first term measures the lack of fit; the second term balances the temptation to add more variables by imposing a penalty for over-fitting; and the third term, the "clustering penalty function", provides a countervailing force to over-clustering. To illustrate the third term, suppose the mixing proportions are equal, then the third term becomes $2n \log(K)$, which increases with $K$, the number of mixture components. The asymptotic efficiency of the MRC is proved, and the criterion performs well in Monte Carlo studies.

Motivated by the Lasso properties (*Tibshirani*, 1996), *Luo et al.* (2008) incorporates both mixture and regression penalties to obtain group and covarite selections (MR-Lasso). The penalty on the mixture is of the form

$$\sum |\beta_j - \beta_k|$$

where $\beta_j$ the coefficient for group $j$. A modified $EM$ algorithm is given to obtain the MR-Lasso estimates.

*Chen* (1995) finds the optimal convergence rate of the densities in finite mixture normal models to be $n^{-1/2}$ when the exact number of components is known, but only $n^{-1/4}$ when unknown. Under the strong identifiability condition, the rate is shown to be attained by some minimum distance estimators.

Then in *Chen and Chen* (2003), for the testing problem of $p = 0$ or $\theta_1 = \theta_2$ in (1.2) with the parameters bounded and $q = 1$, a complicated limiting distribution of the likelihood ratio statistic is given. In *Chen et al.* (2001), a modified likelihood

ratio test is given for this problem in which a penalization is added to restrict $p$ to be bounded away from 0 and 1.

*Zhu and Zhang* (2004, 2006) consider a two group mixture regression model, in which the distribution within each subgroup is any distribution whose Fisher information is positive definite, and the difference of the two groups lies in the difference in the parameter which measures the strength of association that is contributed by some covariates. A resampling method is proposed for the test due to the complicated form of the limiting distribution.

In *Chen and Li* (2009), an *EM* test is given for the same problem as in *Chen and Chen* (2003) but with a much simpler limiting distribution, and the bounded parameter space assumption is not required. For finite choices of $p$, we repeat the EM algorithm for finite steps to calculate the modified likelihood ratio test statistics. The *EM* test uses the maximum of those values. A generalization to more general models is in *Li and Chen* (2010).

In the literature, researchers have considered various forms of mixture models, but have not addressed the testing problems for mixture models with varying proportions as well as covariates-dependent means. But testing the existence of subgroups is critical in the mixture model. We derive a test based on a small number of *EM* iterations towards the likelihood using the same idea in *Chen and Li* (2009), and propose the bootstrap approximation for the critical values of the test. When subgroups exist, the mixture model helps us identify the factors that are associated with the group membership.

### 1.2.5 Mixture of Experts

Model (1.1) can be viewed as a special case of the "mixture of experts" models (*Jordan and Jacobs* (1994), *Yuksel et al.* (2012)) in computer science, where the "gating function" is logistic and the "expert function" is Gaussian. The identifiability

and the statistical properties of the parameter estimates have been studied by *Jiang and Tanner* (1999a,c,b) among others. Bayesian methods for selecting the number of experts have been suggested by *Peng et al.* (1996) and *Ueda and Ghahramani* (2002), whereas *Fritsch et al.* (1997) considered a grow-and-prune strategy of the model for the same purpose. The existing work on the mixture-of-experts models does not cover the results of the present paper for two main reasons. First, our model aims to find predictive variables $\boldsymbol{X}$ for the subgroups that show differential treatment effects in the response $Y$, so each parameter in our model has direct interpretations. In contrast, the mixture-of-experts models aim at the prediction of $Y$ with no distinction between $\boldsymbol{X}$ and $\boldsymbol{Z}$. Second, and more importantly, we consider the hypothesis testing problem with a specific null hypothesis, that is, no predictable subgroups exist for differential treatment effects. As far as we know, no such confirmatory statistical tests are available in literature.

# CHAPTER II

# Logistic–Normal Two-Group Mixture Model

We start from a simple case where the response $Y$ is normally distributed with a covariate $\boldsymbol{Z}$ in each subgroup, where we have two well-defined subgroups, and the group proportion depends on a covariate $\boldsymbol{X}$ through a logistic model. The expectation for this model is that, for a clinical trial, people may react differently to the same treatment depending on the covariates, such as the baseline of some attributes. One interesting case is that, one subgroup shows a desirable treatment effect, while in the other subgroup, the treatment effect is not significant. Our tasks in this chapter are to build a statistical model, identify the covariates that are associated with the subgroup membership, and conduct statistical inferences about the treatment effect in each subgroup.

## 2.1 Statistical Model

Now we specify our model to be the following, for $i = 1, \ldots, n$,

$$
\begin{aligned}
Y_i &= \boldsymbol{Z}_i^T(\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2\delta_i) + \varepsilon_i, \\
P(\delta_i = 1|\boldsymbol{X}_i, \boldsymbol{Z}_i) &= \pi(\boldsymbol{X}_i^T\boldsymbol{\gamma}) \equiv \exp(\boldsymbol{X}_i^T\boldsymbol{\gamma})/(1 + \exp(\boldsymbol{X}_i^T\boldsymbol{\gamma})), \qquad (2.1) \\
P(\delta_i = 0|\boldsymbol{X}_i, \boldsymbol{Z}_i) &= 1 - P(\delta_i = 1|\boldsymbol{X}_i),
\end{aligned}
$$

where $n$ is the sample size, $Y_i \in \mathbb{R}$ is the outcome, $\delta_i \in \{0, 1\}$ is the subgroup indicator, $\boldsymbol{Z}_i \in \mathbb{R}^{q_1}$ is the covariate associated with the subgroup mean, $\boldsymbol{X}_i \in \mathbb{R}^{q_2}$ is the covariate associated with the group membership, $\boldsymbol{\beta}_1 \in \mathbb{R}^{q_1}, \boldsymbol{\beta}_2 \in \mathbb{R}^{q_1}, \boldsymbol{\gamma} \in \mathbb{R}^{q_2}$ are the corresponding coefficients, $\varepsilon_i \sim N(0, \sigma^2)$ for some parameter $\sigma$. The first elements of $\boldsymbol{X}_i$ and $\boldsymbol{Z}_i$ are 1, and the second element of $\boldsymbol{Z}_i$ is the treatment indicator. We can have overlapping variables in the random vectors of $\boldsymbol{X}_i$ and $\boldsymbol{Z}_i$. The overall parameters are $\boldsymbol{\eta}^T = (\boldsymbol{\beta}_1^T, \sigma, \boldsymbol{\beta}_2^T, \boldsymbol{\gamma}^T)$. Write $\boldsymbol{\theta}^T = (\boldsymbol{\beta}_1^T, \sigma, \boldsymbol{\beta}_2^T)$. We observe the data $\{\boldsymbol{W}_i = (Y_i, \boldsymbol{Z}_i^T, \boldsymbol{X}_i^T), i = 1, \ldots n\}$, and $\delta_i$'s are viewed as latent variables. The observations $\boldsymbol{W}_i$'s are independent.

## 2.2 Identifiability

For mixture models, the parameters are not identifiable in the usual sense. As used in *Teicher* (1961, 1963), in Model (2.1), we define the parameters $((\boldsymbol{\theta}^1)^T, (\boldsymbol{\gamma}^1)^T)$ and $((\boldsymbol{\theta}^2)^T, (\boldsymbol{\gamma}^2)^T)$, where $(\boldsymbol{\theta}^1)^T = ((\boldsymbol{\beta}_1^1)^T, \sigma^1, (\boldsymbol{\beta}_2^1)^T)$ and $(\boldsymbol{\theta}^2)^T = ((\boldsymbol{\beta}_1^2)^T, \sigma^2, (\boldsymbol{\beta}_2^2)^T)$, to be in an equivalent class, if and only if $\boldsymbol{\beta}_1^1 = \boldsymbol{\beta}_1^2 + \boldsymbol{\beta}_2^2, \boldsymbol{\beta}_2^1 = -\boldsymbol{\beta}_2^2, \sigma^1 = \sigma^2$, and $\boldsymbol{\gamma}^1 = -\boldsymbol{\gamma}^2$. Then, we define the identifiability of the parameters in Model (2.1) if equal density functions implies that the parameters are from the same equivalent class. In this sense, by Proposition 1 in *Teicher* (1963), the parameters in Model (2.1) are identifiable when the random vectors $\boldsymbol{X}$ and $\boldsymbol{Z}$ are linearly independent.

## 2.3 *EM* Algorithm

Next we derive the *EM* algorithm (*Dempster et al.*, 1977) to get the estimates of the parameters.

Note that the complete data $(\boldsymbol{W}_i, \delta_i)$ has density

$$f(\boldsymbol{W}_i, \delta_i; \boldsymbol{\eta}) = f(Y_i, \delta_i | \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\eta}) f(\boldsymbol{X}_i, \boldsymbol{Z}_i) = f(Y_i | \delta_i, \boldsymbol{Z}_i; \boldsymbol{\theta}) f(\delta_i | \boldsymbol{X}_i; \boldsymbol{\gamma}) f(\boldsymbol{X}_i, \boldsymbol{Z}_i).$$

In one iteration step, suppose that currently we have the parameter $\boldsymbol{\eta}^T = (\boldsymbol{\eta}^{(k)})^T = ((\boldsymbol{\theta}^{(k)})^T, (\boldsymbol{\gamma}^{(k)})^T)$, then at the $(k+1)th$ iteration, we have the following derivations.

In the $E$ step, let

$$
\begin{aligned}
a_i^{(k)} &= P(\delta_i = 1 | Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \eta^{(k)}) \\
&= f(Y_i | \delta_i = 1, \boldsymbol{Z}_i; \theta^{(k)}) P(\delta_i = 1 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k)}) / (f(Y_i | \delta_i = 1, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k)}) P(\delta_i = 1 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k)}) \\
&\quad + f(Y_i | \delta_i = 0, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k)}) P(\delta_i = 0 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k)})),
\end{aligned}
$$

(2.2)

$b_i^{(k)} = 1 - a_i^{(k)}$, $\boldsymbol{a}^{(k)} = (a_1^{(k)}, \ldots, a_n^{(k)})$, and $\boldsymbol{b}^{(k)} = (b_1^{(k)}, \ldots, b_n^{(k)})$.

Then we have that

$$
\begin{aligned}
Q(\boldsymbol{\eta}^{(k+1)} | \boldsymbol{\eta}^{(k)}) &= \mathbb{E}_{\delta_i | \boldsymbol{W}_i; \boldsymbol{\eta}^{(k)}} \sum_{i=1}^n \log f(\delta_i, Y_i | \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\eta}^{(k+1)}) \\
&= \mathbb{E}_{\delta_i | \boldsymbol{W}_i; \boldsymbol{\eta}^{(k)}} \sum_{i=1}^n (\log f(Y_i | \delta_i, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k+1)}) + \log P(\delta_i | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k+1)})) \\
&= \mathbb{E}_{\delta_i | \boldsymbol{W}_i; \boldsymbol{\eta}^{(k)}} \sum_{i=1}^n \log f(Y_i | \delta_i, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k+1)}) + \mathbb{E}_{\delta_i | \boldsymbol{W}_i; \boldsymbol{\eta}^{(k)}} \sum_{i=1}^n \log P(\delta_i | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k+1)}) \\
&= Q(\boldsymbol{\theta}^{(k+1)} | \boldsymbol{\eta}^{(k)}) + Q(\boldsymbol{\gamma}^{(k+1)} | \boldsymbol{\eta}^{(k)}),
\end{aligned}
$$

(2.3)

where

$$
\begin{aligned}
Q(\boldsymbol{\theta}^{(k+1)} | \boldsymbol{\eta}^{(k)}) &= \sum_{i=1}^n \mathbb{E}_{\delta_i | \boldsymbol{W}_i; \boldsymbol{\eta}^{(k)}} \log f(Y_i | \delta_i; \boldsymbol{\theta}^{(k+1)}) \\
&= \sum_{i=1}^n [a_i^{(k)} \log f(Y_i | \delta_i = 1, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k+1)}) + b_i^{(k)} \log f(Y_i | \delta_i = 0, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k+1)})] \\
&= (-n/2) \log(2\pi\sigma^2) - \sum_{i=1}^n a_i (Y_i - \boldsymbol{Z}_i^T (\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2))^2 / (2\sigma^2) \\
&\quad - \sum_{i=1}^n b_i (Y_i - \boldsymbol{Z}_i^T \boldsymbol{\beta}_1)^2 / (2\sigma^2)
\end{aligned}
$$

(2.4)

and

$$
Q(\boldsymbol{\gamma}^{(k+1)} | \boldsymbol{\eta}^{(k)}) = \sum_i [a_i^{(k)} \log P(\delta_i = 1 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k+1)}) + b_i^{(k)} \log P(\delta_i = 0 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k+1)})].
$$

(2.5)

Therefore in the $(k+1)$th step, we obtain the estimates of $\boldsymbol{\theta}^{(k+1)}$ from $Q(\boldsymbol{\theta}^{(k+1)} | \boldsymbol{\eta}^{(k)})$ by the weighted least squares, and the estimates of $\boldsymbol{\gamma}^{(k+1)}$ by maximizing $Q(\boldsymbol{\gamma}^{(k+1)} | \boldsymbol{\eta}^{(k)})$,

15

which is a weighted logistic regression problem.

*Remark* II.1. When the model does have two distinguishable groups, the *EM* algorithm tends to converge quickly. However, there is no guarantee that the solution is the global maximizer. Therefore, we need to carefully choose different starting values to locate a global maximizer in practice.

## 2.4 Covariance Matrix

In the *EM* algorithm, the standard error of the estimators can be calculated in the following way (*Louis*, 1982). In the last step of the *EM* algorithm, suppose that we have $\hat{\boldsymbol{\eta}}^T = (\hat{\boldsymbol{\theta}}^T, \hat{\boldsymbol{\gamma}}^T) = (\hat{\boldsymbol{\beta}}_1^T, \hat{\sigma}, \hat{\boldsymbol{\beta}}_2^T, \hat{\boldsymbol{\gamma}}^T)$. For the $i$-th observation, let $B_i$ and $S_i$ be the individual negative second derivative and first derivative of the complete data log-likelihood. The inverse of the covariance matrix, the observed Fisher information matrix of the parameter estimator is

$$
\begin{aligned}
\boldsymbol{I}_Y \;=\;& \textstyle\sum_{i=1}^n \mathbb{E}_{(\delta_i | \boldsymbol{W}_i; \hat{\boldsymbol{\eta}})} \boldsymbol{B}_i(Y_i, \delta_i; \hat{\boldsymbol{\eta}}) - \sum_{i=1}^n \mathbb{E}_{(\delta_i | \boldsymbol{W}_i; \hat{\boldsymbol{\eta}})} S_i(Y_i, \delta_i; \hat{\boldsymbol{\eta}}) S_i(Y_i, \delta_i | \boldsymbol{Z}_i; \hat{\boldsymbol{\eta}})^T \\
& - \textstyle\sum_{i \neq j}^n (\mathbb{E}_{(\delta_i | \boldsymbol{W}_i; \hat{\boldsymbol{\eta}})} \boldsymbol{S}_i)(\mathbb{E}_{(\delta_j | \boldsymbol{W}_i; \hat{\boldsymbol{\eta}})} \boldsymbol{S}_j).
\end{aligned}
\tag{2.6}
$$

In our setting, let $\varepsilon_i = Y_i - \boldsymbol{Z}_i^T(\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2 \delta_i)$. Then the individual complete log-likelihood is

$$
l(Y_i, \delta_i | \boldsymbol{X}_i, \boldsymbol{Z}_i) = -\log \sigma - \frac{\varepsilon_i^2}{2\sigma^2} + \delta_i \boldsymbol{X}_i^T \gamma - \log(1 + \exp(\boldsymbol{X}_i^T \boldsymbol{\gamma})).
\tag{2.7}
$$

Then

$$
\begin{aligned}
\boldsymbol{S}_i \;=\;& \frac{\partial l}{\partial \boldsymbol{\eta}^T} = (\tfrac{1}{\sigma^2} \varepsilon_i \boldsymbol{Z}_i^T, -\tfrac{1}{\sigma} + \tfrac{1}{\sigma^3} \varepsilon_i^2, \tfrac{1}{\sigma^2} \varepsilon_i \delta_i \boldsymbol{Z}_i^T, \\
& \delta_i \boldsymbol{X}_i^T - \pi(\boldsymbol{X}_i^T \boldsymbol{\gamma}) \boldsymbol{X}_i^T),
\end{aligned}
\tag{2.8}
$$

and $\boldsymbol{B}_i = \mathrm{diag}(\boldsymbol{B}_{i11}, \boldsymbol{B}_{i22})$, where $\boldsymbol{B}_{i11}$ is

$$\frac{1}{\sigma^2} \begin{pmatrix} \boldsymbol{Z}_i \boldsymbol{Z}_i^T & \frac{2\varepsilon_i \boldsymbol{Z}_i}{\sigma} & \delta_i \boldsymbol{Z}_i \boldsymbol{Z}_i^T \\ \frac{2\varepsilon_i \boldsymbol{Z}_i}{\sigma} & -\frac{1}{\sigma^2} + \frac{3\varepsilon^2}{\sigma^4} & \frac{2\varepsilon_i \boldsymbol{Z}_i}{\sigma} \\ \delta_i \boldsymbol{Z}_i \boldsymbol{Z}_i^T & \frac{2\varepsilon_i \boldsymbol{Z}_i}{\sigma} & \delta_i \boldsymbol{Z}_i \boldsymbol{Z}_i^T \end{pmatrix}, \tag{2.9}$$

and

$$\boldsymbol{B}_{i22} = \pi(\boldsymbol{X}_i^T \boldsymbol{\gamma})(1 - \pi(\boldsymbol{X}_i^T \boldsymbol{\gamma}))\boldsymbol{X}_i \boldsymbol{X}_i^T. \tag{2.10}$$

The covariance matrix of the estimators obtained from Section 2.3 can then be computed via (2.6) by substituting with the estimates from the last step in the $EM$ algorithm.

From the asymptotic covariance matrix, we could easily construct confidence intervals for the parameters of interest. The Wald test is feasible for the parameters related to the treatment effects, if the two groups are distinguishable.

## 2.5  Simulations

In Table 2.1 and 2.2 we summarize some simulation results. Data are generated from

$$Y_i = \mu_1 + \nu_1 T_i + \alpha_1 Z_i + (\mu_2 + \nu_2 T_i + \alpha_2 Z_i)\delta_i + \varepsilon_i,$$

$$P(\delta_i = 1 | X_i) = \pi(\gamma_0 + \gamma_1 X_i),$$

for $i = 1, \ldots, n$, where $\varepsilon_i \sim N(0, \sigma^2)$, independent of $X_i, Z_i$. The observations are independent. In Case 1, $X_i, Z_i$ are independent, $X_i \sim N(1,1)$ and $Z_i \sim N(1,1)$. In Case 2, $X_i = Z_i \sim N(1,1)$. In the tables, we have two sets of parameters and sample size $n = 100$. For each case, we collect the means and sample standard deviations of the estimates in 1000 repeated experiments. From the results we observe that as the

17

sample size increases, the estimates become more accurate. Whether $X_i$ and $Z_i$ are independent or not does not play an important role. In the simulation, we add the constriction that $\nu_2$ is positive to guarantee the uniqueness of the parameters. We obtain good estimates of the parameters in both Case 1 and case 2.

Table 2.1: The means and the sample standard deviations of the estimates in 1000 repeated experiments under Case 1 that $X_i$ and $Z_i$ are independent. $n = 100$

| Parameters | True | est | sd | True | est | sd |
|---|---|---|---|---|---|---|
| $\mu_1$ | 2.0 | 2.00 | 0.14 | 2.0 | 2.00 | 0.14 |
| $\nu_1$ | 0.0 | 0.00 | 0.15 | 0.0 | 0.00 | 0.15 |
| $\alpha_1$ | 2.0 | 2.00 | 0.08 | 2.0 | 2.00 | 0.08 |
| $\mu_2$ | 3.0 | 3.00 | 0.19 | 3.0 | 2.99 | 0.19 |
| $\nu_2$ | 3.0 | 3.00 | 0.22 | 8.0 | 8.00 | 0.22 |
| $\alpha_2$ | 5.0 | 5.00 | 0.11 | 5.0 | 5.00 | 0.11 |
| $\gamma_0$ | 1.0 | 1.02 | 0.37 | 1.0 | 1.02 | 0.37 |
| $\gamma_1$ | -1.0 | -1.02 | 0.28 | -1.0 | -1.02 | 0.28 |
| $\sigma$ | 0.5 | 0.48 | 0.04 | 0.5 | 0.48 | 0.04 |

Table 2.2: The means and the sample standard deviations of the estimates in 1000 repeated experiments under Case 2 that $X_i = Z_i$.

| Parameters | True | est | sd | True | est | sd |
|---|---|---|---|---|---|---|
| $\mu_1$ | 2.0 | 2.00 | 0.16 | 2.0 | 2.00 | 0.16 |
| $\nu_1$ | 0.0 | 0.00 | 0.14 | 0.0 | 0.00 | 0.14 |
| $\alpha_1$ | 2.0 | 2.00 | 0.08 | 2.0 | 2.00 | 0.08 |
| $\mu_2$ | 3.0 | 3.01 | 0.19 | 3.0 | 3.01 | 0.19 |
| $\nu_2$ | 3.0 | 3.00 | 0.20 | 8.0 | 8.00 | 0.20 |
| $\alpha_2$ | 5.0 | 5.00 | 0.11 | 5.0 | 5.00 | 0.11 |
| $\gamma_0$ | 1.0 | 1.01 | 0.37 | 1.0 | 1.01 | 0.36 |
| $\gamma_1$ | -1.0 | -1.04 | 0.29 | -1.0 | -1.04 | 0.28 |
| $\sigma$ | 0.5 | 0.49 | 0.04 | 0.5 | 0.49 | 0.04 |

# CHAPTER III

# Hypothesis Testing for the Existence of Subgroups

## 3.1 Maximum Likelihood Ratio Tests for the Existence of Subgroups

### 3.1.1 Maximum Likelihood Ratio Tests

Before we evaluate the likelihood in detail, we summarize a general result of the likelihood ratio test as follows with a univariate parameter $\theta$ and a random variable $W_i$. The results can be extended to high dimensional $\theta$.

For each $i$, let $L(W_i; \theta)$ be the individual likelihood, and $l(W_i; \theta), \dot{l}(W_i; \theta), \ddot{l}(W_i; \theta)$ be the log-likelihood and its first and second derivatives with respect to $\theta$, respectively. In addition, let $L(\theta) = \sum_{i=1}^{n} L(W_i; \theta)$, and similarly let $l(\theta), \dot{l}(\theta)$ and $\ddot{l}(\theta)$ be the summation of $l(W_i; \theta), \dot{l}(W_i; \theta)$ and $\ddot{l}(W_i; \theta)$, respectively. Suppose $\theta_0$ is the true parameter value, $\hat{\theta}_n$ is the maximum likelihood estimator, which satisfies $\dot{l}(\hat{\theta}_n)=0$, and we are testing the null hypothesis of $\theta = \theta_0$. Under the regularity conditions, we have:

$$0 = \sum_{i=1}^{n} \dot{l}(W_i; \hat{\theta}_n) = \sum_{i=1}^{n} \dot{l}(W_i; \theta_0) + \sum_{i=1}^{n} \ddot{l}(W_i; \theta^*)(\hat{\theta}_n - \theta_0),$$

for some $\theta^*$ such that $|\theta^* - \theta_0| \le |\hat{\theta}_n - \theta_0|$. Equivalently, we have

$$\frac{1}{\sqrt{n}}\dot{l}(\theta_0) = \sqrt{n}(\hat{\theta}_n - \theta_0)\{I(\theta_0) - [\frac{1}{n}\sum_{i=1}^n \ddot{l}(W_i; \theta^*) + I(\theta^*)] + [I(\theta^*) - I(\theta_0)]\}.$$

(3.1)

Since $\hat{\theta}_n - \theta_0 = o_p(1)$, $\frac{1}{n}\sum_{i=1}^n \ddot{l}(W_i, \theta^*) + I(\theta^*) = o_p(1)$, $I(\theta)$ is continuous and positive definite at $\theta_0$, then

$$\frac{1}{\sqrt{n}}\dot{l}(\theta_0) = \sqrt{n}(\hat{\theta}_n - \theta_0)(I(\theta_0) + o_p(1)).$$

(3.2)

Because $\frac{1}{\sqrt{n}}\dot{l}(\theta_0)$ is $O_p(1)$ by the central limit theorem, and $I(\theta_0)$ is positive definite, then $\sqrt{n}(\hat{\theta}_n - \theta_0) = O_p(1)$. Therefore, we have

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{1}{\sqrt{n}}I^{-1}(\theta_0)\dot{l}(\theta_0) + o_p(1).$$

(3.3)

Then expansion of the log-likelihood gives that for some $\theta^*$ such that $||\hat{\theta}_n - \theta_0|| \le ||\theta^* - \theta_0||$, we have

$$
\begin{aligned}
l(\hat{\theta}_n) - l(\theta_0) &= \dot{l}(\theta_0)^T(\hat{\theta}_n - \theta_0) + \frac{1}{2}(\hat{\theta}_n - \theta_0)^T\ddot{l}(\theta^*)(\hat{\theta}_n - \theta_0) \\
&= \dot{l}(\theta_0)^T(\hat{\theta}_n - \theta_0) + \frac{1}{2}(\sqrt{n}(\hat{\theta}_n - \theta_0))^T(-I(\theta_0) + (-I(\theta^*) + I(\theta_0)) \\
&\quad + (\ddot{l}(\theta^*)/n + I(\theta^*)))(\sqrt{n}(\hat{\theta}_n - \theta_0)) \\
&= \dot{l}(\theta_0)^T(\hat{\theta}_n - \theta_0) - \frac{1}{2}(\sqrt{n}(\hat{\theta}_n - \theta_0))^T I(\theta_0)(\sqrt{n}(\hat{\theta}_n - \theta_0)) + o_p(1).
\end{aligned}
$$

(3.4)

From Equation (3.3) and (3.4), we get

$$
\begin{aligned}
2(l(\hat{\theta}_n) - l(\theta_0)) &= \sqrt{n}(\hat{\theta}_n - \theta_0)^T I(\theta_0)\sqrt{n}(\hat{\theta}_n - \theta_0) + o_p(1) \\
&= (\frac{1}{\sqrt{n}}\dot{l}(\theta_0))^T I^{-1}(\theta_0)(\frac{1}{\sqrt{n}}\dot{l}(\theta_0)) + o_p(1).
\end{aligned}
$$

(3.5)

Recall that our model is

$$
\begin{aligned}
f(Y, \boldsymbol{Z}, \boldsymbol{X}) &= f(Y|\boldsymbol{X}, \boldsymbol{Z})g(\boldsymbol{X}, \boldsymbol{Z}) \\
&= (\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\varphi(Y - \boldsymbol{Z}^T(\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2), \sigma^2) + (1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma}))\varphi(Y - \boldsymbol{Z}^T\beta_1, \sigma^2)) \\
&\quad g(\boldsymbol{X}, \boldsymbol{Z}),
\end{aligned}
$$

where $\varphi(\mu, \sigma)$ is the density of a normal variable with mean $\mu$ and variance $\sigma^2$.

If we fix the value of $\boldsymbol{\gamma}^T = (\boldsymbol{\gamma}_{-X}, \boldsymbol{\gamma}_X^T)$, the problem of testing $\beta_2 = 0$ is a regular one when $\gamma_X$ is nonzero, which satisfies the regularity conditions, including the condition that the third derivatives are integrable under the null hypothesis. The parameters are identifiable. Assume that the true parameter $\boldsymbol{\theta}_0^T = (\boldsymbol{\beta}_0^T, \sigma_0, 0)$. Under the null hypothesis that $\boldsymbol{\beta}_2 = 0$, we write the MLE as $\hat{\boldsymbol{\theta}}_0$, and under the alternative, the MLE is $\hat{\boldsymbol{\theta}}_n$.

By direct calculation, the score function

$$
\begin{aligned}
\dot{\boldsymbol{l}}_\gamma(Y, \boldsymbol{Z}, \boldsymbol{X}; \theta_0) &= \big(\tfrac{1}{\sigma_0^2}(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)\boldsymbol{Z}^T, -\tfrac{1}{\sigma_0} + \tfrac{1}{\sigma_0^2}(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)^2, \\
&\quad \tfrac{1}{\sigma_0^2}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)\boldsymbol{Z}^T\big)^T,
\end{aligned}
\tag{3.6}
$$

and the Fisher information matrix

$$
I_\gamma(\boldsymbol{\theta}_0) = \frac{1}{\sigma_0^2}
\begin{pmatrix}
\boldsymbol{A} & 0 & \boldsymbol{B}(\boldsymbol{\gamma}) \\
0 & 2 & 0 \\
\boldsymbol{B}(\boldsymbol{\gamma}) & 0 & \boldsymbol{C}(\boldsymbol{\gamma})
\end{pmatrix},
\tag{3.7}
$$

where $\boldsymbol{A} = \mathbb{E}(\boldsymbol{Z}\boldsymbol{Z}^T)$, $\boldsymbol{B}(\boldsymbol{\gamma}) = \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T)$, $\boldsymbol{C}(\boldsymbol{\gamma}) = \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})^2\boldsymbol{Z}\boldsymbol{Z}^T)$.

Partition $I_\gamma(\boldsymbol{\theta}_0)$ to get $I_{\gamma 11}, I_{\gamma 12}, I_{\gamma 21}, I_{\gamma 22}$ and $I_{\gamma 22}$. In Particular,

$$
I_{\gamma 11}(\boldsymbol{\theta}_0) = \frac{1}{\sigma_0^2}
\begin{pmatrix}
\boldsymbol{A} & 0 \\
0 & 2
\end{pmatrix}.
$$

Let

$$\boldsymbol{M}_1(\boldsymbol{\gamma}) = \Big( \sum_{i=1}^n \frac{(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_0)\boldsymbol{Z}_i^T}{\sqrt{n}\sigma_0^2}, \sum_{i=1}^n \frac{1}{\sqrt{n}}\big(-\frac{1}{\sigma_0} + \frac{(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_0)^2}{\sigma_0^3}\big)\Big)^T, \tag{3.8}$$

and

$$\boldsymbol{M}_2(\gamma) = \frac{1}{\sqrt{n}\sigma_0^2} \sum_{i=1}^n \pi(\boldsymbol{X}^T\boldsymbol{\gamma})(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_0)\boldsymbol{Z}_i^T. \tag{3.9}$$

Then,

$$\frac{1}{\sqrt{n}}\dot{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0) = (\boldsymbol{M}_1(\boldsymbol{\gamma}), \boldsymbol{M}_2(\boldsymbol{\gamma})).$$

For the reduced model where $\boldsymbol{\beta}_2 = 0$, we have similar results. Direct calculations give that the information matrix for the reduced model is $\boldsymbol{I}_{\boldsymbol{\gamma}11}(\boldsymbol{\theta}_0)$, the top left submatrix of $\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)$.

Now for the likelihood ratio statistic of testing $\beta_2 = 0$ with a given $\gamma$,

$$T(\boldsymbol{\gamma}) \equiv 2(\boldsymbol{l}_{\boldsymbol{\gamma}}(\hat{\boldsymbol{\theta}}_n) - \boldsymbol{l}_{\boldsymbol{\gamma}}(\theta_0)) - 2(\boldsymbol{l}_{\boldsymbol{\gamma}}(\hat{\theta}_0) - \boldsymbol{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) = T_1(\boldsymbol{\gamma}) - T_2(\boldsymbol{\gamma}),$$

where

$$\begin{aligned}
T_1(\boldsymbol{\gamma}) &= 2(\boldsymbol{l}_{\boldsymbol{\gamma}}(\hat{\boldsymbol{\theta}}_n) - \boldsymbol{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) \\
&= (\boldsymbol{M}_1(\boldsymbol{\gamma}), \boldsymbol{M}_2(\boldsymbol{\gamma}))^T \boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)^{-1}(\boldsymbol{\theta}_0)(\boldsymbol{M}_1(\boldsymbol{\gamma}), \boldsymbol{M}_2(\boldsymbol{\gamma})) + o_p(1), \\
T_2(\boldsymbol{\gamma}) &= 2(\boldsymbol{l}_{\boldsymbol{\gamma}}(\hat{\boldsymbol{\theta}}_0) - \boldsymbol{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) \\
&= \boldsymbol{M}_1(\boldsymbol{\gamma})^T \boldsymbol{I}_{\boldsymbol{\gamma}11}(\boldsymbol{\theta}_0)^{-1}\boldsymbol{M}_1(\boldsymbol{\gamma}) + o_p(1).
\end{aligned}$$

Then, by observing that

$$\begin{pmatrix} \boldsymbol{I}_{11} & \boldsymbol{I}_{12} \\ \boldsymbol{I}_{21} & \boldsymbol{I}_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{I}_{11}^{-1} + \boldsymbol{I}_{11}^{-1}\boldsymbol{I}_{12}\boldsymbol{I}_{22\cdot1}^{-1}\boldsymbol{I}_{21}\boldsymbol{I}_{11}^{-1} & -\boldsymbol{I}_{11}^{-1}\boldsymbol{I}_{12}\boldsymbol{I}_{22\cdot1}^{-1} \\ -\boldsymbol{I}_{22\cdot1}^{-1}\boldsymbol{I}_{21}\boldsymbol{I}_{11}^{-1} & \boldsymbol{I}_{22\cdot1}^{-1} \end{pmatrix},$$

where $\boldsymbol{I}_{\gamma 22\cdot 1} = \boldsymbol{I}_{\gamma 22} - \boldsymbol{I}_{\gamma 21} I_{\gamma 11}^{-1} \boldsymbol{I}_{\gamma 12} = (\boldsymbol{C}(\boldsymbol{\gamma}) - \boldsymbol{B}(\boldsymbol{\gamma})\boldsymbol{A}^{-1}\boldsymbol{B}(\boldsymbol{\gamma}))/\sigma^2$, we have

$$
\begin{aligned}
T(\gamma) &= T_1(\gamma) - T_2(\gamma) \\
&= (\boldsymbol{M}_2(\boldsymbol{\gamma}) - \boldsymbol{I}_{\gamma 21} \boldsymbol{I}_{\gamma 11}^{-1} \boldsymbol{M}_1(\boldsymbol{\gamma}))^T \boldsymbol{I}_{\gamma 22\cdot 1}^{-1} (\boldsymbol{M}_2(\boldsymbol{\gamma}) - \boldsymbol{I}_{\gamma 21} \boldsymbol{I}_{\gamma 11}^{-1} \boldsymbol{M}_1(\boldsymbol{\gamma})) + o_p(1),
\end{aligned}
$$

and let

$$
\begin{aligned}
\boldsymbol{h}(\boldsymbol{\gamma}) &= \boldsymbol{I}_{\gamma 22\cdot 1}^{-1/2} (\boldsymbol{M}_2(\boldsymbol{\gamma}) - \boldsymbol{I}_{\gamma 21} \boldsymbol{I}_{\gamma 11}^{-1} \boldsymbol{M}_1(\boldsymbol{\gamma})) \\
&= \frac{1}{\sqrt{n}\sigma_0^2} \sum_{i=1}^{n} \boldsymbol{I}_{\gamma 22\cdot 1}^{-1/2} \{\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{I}_{q_2} - \boldsymbol{B}(\boldsymbol{\gamma})\boldsymbol{A}^{-1}\}(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_0)\boldsymbol{Z}_i \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \gamma),
\end{aligned}
$$

where

$$
\boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) = \frac{1}{\sigma_0^2} \boldsymbol{I}_{\gamma 22\cdot 1}^{-1/2} \{\pi(\boldsymbol{X}_i^T\boldsymbol{\gamma})\boldsymbol{I}_{q_2} - \boldsymbol{B}(\boldsymbol{\gamma})\boldsymbol{A}^{-1}\}(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_0)\boldsymbol{Z}_i, \qquad (3.10)
$$

then the test statistic $T(\boldsymbol{\gamma}) = ||h(\boldsymbol{\gamma})||^2$. Note that $\mathbb{E}\psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) = 0$.

Therefore, if we have finitely many $\boldsymbol{\gamma}'s$ from $\boldsymbol{\Gamma} = \{\boldsymbol{\gamma}_1, \ldots, \boldsymbol{\gamma}_K\}$ with all nonzero $X$ coefficients, by the central limit theorem, $(\boldsymbol{h}(\boldsymbol{\gamma}_1), \ldots, \boldsymbol{h}(\boldsymbol{\gamma}_K))$ converges to some random variable $(\boldsymbol{H}_1(\boldsymbol{\Gamma}), \cdots, \boldsymbol{H}_K(\boldsymbol{\Gamma}))$, which depends on the pre-specified set $\boldsymbol{\Gamma}$. Our test statistic $\max\{||\boldsymbol{h}(\boldsymbol{\gamma}_1)||^2, \ldots, ||\boldsymbol{h}(\boldsymbol{\gamma}_K)||^2\}$, converges to the random variable $\max\{||\boldsymbol{H}_1(\boldsymbol{\Gamma})||^2, \cdots, ||\boldsymbol{H}_K(\boldsymbol{\Gamma})||^2\}$, which depends on the chosen set of $\boldsymbol{\gamma}$.

We summarize the above derivations in the following theorem.

**Theorem III.1.** *For Model 2.1, if we choose $\boldsymbol{\gamma}_1, \ldots, \boldsymbol{\gamma}_K$ whose $X$ coefficients are nonzero, $K \geq 1$, and for each $\boldsymbol{\gamma}$, calculate the likelihood ratio test statistic $T(\boldsymbol{\gamma}_i)$ for the null hypothesis of $\boldsymbol{\beta}_2 = 0$, $i = 1, \cdots, K$, then the maximum variable $\max_{1 \leq i \leq K} T(\boldsymbol{\gamma}_i)$ converges to a limiting distribution.*

*Remark* III.2. In theory, for $K = 1$, the limiting distribution should be the standard chi-square distribution. However, since for $q_1 = q_2 = 1$ with no covariates, the convergence is very slow for a given constant proportion (*Goeffinet et al.*, 1992), it

is not unexpected that here with the covariates, the convergence rate is very poor in the simulations. Therefore, to have a better finite sample performance, we prefer not to use the limiting distribution. The bootstrap methods of determining the critical values are recommended.

## 3.2  *EM* Tests for the Existence of Subgroups

In the previous section, we considered tests for the existence of a subgroup by choosing a set of $\boldsymbol{\gamma}$'s, computing the MLE for $\boldsymbol{\theta}^T = (\boldsymbol{\beta}_1^T, \sigma, \boldsymbol{\beta}_2^T)$ given each $\boldsymbol{\gamma}$, and then taking the maximum of the likelihood ratio test statistics for the hypothesis $\beta_2 = 0$ for each fixed $\boldsymbol{\gamma}$. To increase power of the tests, we construct an *EM* test in which the *EM* algorithm is used to update $\boldsymbol{\gamma}$. If the underlying parameters satisfies the alternative hypothesis, the *EM* algorithm tends to push $\boldsymbol{\gamma}$ towards the true one, and hence increase the power of the test.

### 3.2.1  *EM* Test Process

To construct the *EM* test, we first choose a compact set

$$\Gamma \equiv \{\boldsymbol{\gamma} = (\boldsymbol{\gamma}_{-x}, \boldsymbol{\gamma}_x^T)^T : c_1 < ||\boldsymbol{\gamma}_x|| < c_2, ||\boldsymbol{\gamma}_{-x}|| < c_3\}, \tag{3.11}$$

where $c_1, c_2, c_3 > 0$ are some constants. We choose a set of $\boldsymbol{\gamma}_j \in \Gamma$, $j = 1, \ldots, J$, for a positive integer $J$ and another positive integer $K$. Here we use two types of indices for $\boldsymbol{\gamma}$: $\boldsymbol{\gamma}_{-x} \in \mathbb{R}$ and $\boldsymbol{\gamma}_x \in \mathbb{R}^{q_2-1}$ represent the intercept and the slope for one $\boldsymbol{\gamma}$; while $\boldsymbol{\gamma}_j$ represents the whole vector in $\mathbb{R}^{q_2}$, and different $j$ indicates different vectors. The parameter $\boldsymbol{\gamma} \in \Gamma$ will be constrained in the *EM* process. For each $j$, let $\boldsymbol{\gamma}_j^{(0)} = \boldsymbol{\gamma}_j$. We use the *EM* algorithm to compute

$$\boldsymbol{\theta}^{(0)} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^{n} \log P(Y_i, \boldsymbol{X}_i, \boldsymbol{Z}_i; \boldsymbol{\theta}, \boldsymbol{\gamma}^{(0)}). \tag{3.12}$$

At the $k$th step, $1 \leq k \leq K - 1$, we use the $E$-step and $M$-step as derived in Section 2.3. In more details, suppose that currently $\boldsymbol{\eta} = \boldsymbol{\eta}^{(k)}$, then in the $E$ step, let

$$
\begin{aligned}
a_i^{(k)} &= P(\delta_i | Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\eta}^{(k)}) \\
&= P(Y_i | \delta_i = 1, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k)}) P(\delta_i = 1 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k)}) / (P(Y_i | \delta_i = 1, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k)}) P(\delta_i = 1 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k)}) \\
&\quad + P(Y_i | \delta_i = 0, \boldsymbol{Z}_i; \boldsymbol{\theta}^{(k)}) P(\delta_i = 0 | \boldsymbol{X}_i; \boldsymbol{\gamma}^{(k)})),
\end{aligned}
\tag{3.13}
$$

$b_i^{(k)} = 1 - a_i^{(k)}$, for $i = 1, \ldots, n$, $\boldsymbol{a}^{(k)} = (a_1^{(k)}, \ldots, a_n^{(k)})$, and $\boldsymbol{b}^{(k)} = (b_1^{(k)}, \ldots, b_n^{(k)})$.

Then in the $(k + 1)$th step, compute

$$
\boldsymbol{\theta}^{(k+1)} = \operatorname*{argmax}_{\boldsymbol{\theta}} (\sum_{i=1}^{n} [a_i^{(k)} \log P(Y_i | \delta_i = 1, \boldsymbol{Z}_i; \boldsymbol{\theta}) + b_i^{(k)} \log P(Y_i | \delta_i = 0, \boldsymbol{Z}_i; \boldsymbol{\theta})]), \tag{3.14}
$$

$$
\boldsymbol{\gamma}_{temp}^{(k+1)} = \operatorname*{argmax}_{\boldsymbol{\gamma}} (\sum_{i} [a_i^{(k)} \log P(\delta_i = 1 | \boldsymbol{X}_i; \boldsymbol{\gamma}) + b_i^{(k)} \log P(\delta_i = 0 | \boldsymbol{X}_i; \boldsymbol{\gamma})]), \tag{3.15}
$$

and let

$$
\boldsymbol{\gamma}^{(k+1)} = \begin{cases} \boldsymbol{\gamma}_{temp}^{(k+1)}, & \text{if } \boldsymbol{\gamma}_{temp}^{(k+1)} \in \Gamma. \\ \boldsymbol{\gamma}^{(k)}, & \text{o.w..} \end{cases} \tag{3.16}
$$

Iterate the above steps (3.13)-(3.16) until $k = K - 1$. In the last step, compute $\boldsymbol{\gamma}^{(K)}$ by (3.15) and (3.16), and let

$$
\boldsymbol{\theta}^{(K)} = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log P(Y_i, \boldsymbol{X}_i, \boldsymbol{Z}_i; \boldsymbol{\theta}, \boldsymbol{\gamma}^{(K)}). \tag{3.17}
$$

Let $(\boldsymbol{\eta}^{(K)})^T = ((\boldsymbol{\theta}^{(K)})^T, (\boldsymbol{\gamma}^{(K)})^T)$.

Let $\hat{\boldsymbol{\theta}}_0^T$ be the MLE of the parameter $\boldsymbol{\theta}^T = (\boldsymbol{\beta}_1^T, \sigma, \boldsymbol{\beta}_2^T)$ under the null hypothesis $\boldsymbol{\beta}_2 = 0$ with any fixed $\boldsymbol{\gamma}$ value $\boldsymbol{\gamma}_0$. That is,

$$
\hat{\boldsymbol{\theta}}_0 = \operatorname*{argmax}_{\boldsymbol{\theta}} l(\boldsymbol{\theta}, \boldsymbol{\gamma}_0), \quad \text{subject to } \boldsymbol{\beta}_2 = 0.
$$

Then for each $j$, define the likelihood ratio test statistic

$$EM_j^{(K)} = 2(l(\boldsymbol{\eta}^{(K)}) - l(\hat{\boldsymbol{\theta}}_0, \boldsymbol{\gamma}_0)).$$

The $EM$ test statistic is then

$$EM^{(K)} = \max\{EM_j^{(K)} : j = 1, \ldots, J\}. \tag{3.18}$$

The critical values are to be determined.

### 3.2.2   Convergence of the $EM$ Test Statistic

Now we discuss the convergence properties of the $EM$ test statistic $EM^{(K)}$ in Section 3.2.1.

Assumption 1: The random vectors $\boldsymbol{X}$ and $\boldsymbol{Z}$ are linearly independent, respectively, that is, $\mathbb{E}[\boldsymbol{Z}\boldsymbol{Z}^T]$ and $\mathbb{E}[\boldsymbol{X}\boldsymbol{X}^T]$ are positive definite.

Under the null hypothesis that there is no subgroup, by the properties of M estimators (*van der Vaart*, 1998), we have $\hat{\boldsymbol{\theta}}^{(j)} = \boldsymbol{\theta}_0 + o_p(1)$ for $j = 1, \ldots, K$. Then, we only need to consider parameter $\boldsymbol{\theta}$ close to $\boldsymbol{\theta}_0$. For some positive constants $c_4, c_5, c_6$ and $c_7$, let $\Theta \equiv \{\theta : |\boldsymbol{\beta}_1 - \boldsymbol{\beta}_0| \leq c_4, |\boldsymbol{\beta}_2| \leq c_5, c_6 \geq \sigma \geq c_7 > 0\}$, such that with probability tending to 1, $\hat{\boldsymbol{\theta}}^{(j)}$ in the process of the $EM$ tests lie in this set $\Theta$.

With these preparations, we have our main theory in the following:

**Theorem III.3.** *Under the null hypothesis and Assumptions 1 , for any finite integers $J > 0$ and $K \geq 0$, the $EM$ test statistic $EM^{(K)}$ converges to a fixed distribution as $n \to \infty$.*

To prove Theorem III.3, we state three lemmas first. Proofs of the lemmas are in the appendix.

**Lemma III.4.** *Under Assumption 1, there exist some constants $0 < c_8, c_9 < \infty$, such that for $\gamma \in \Theta$,*

$$0 < c_8 \leq \inf_{\gamma \in \Gamma} \lambda_{\min}(\boldsymbol{I}_{\gamma}(\boldsymbol{\theta}_0)) \leq \sup_{\gamma \in \Gamma} \lambda_{\max}(\boldsymbol{I}_{\gamma}(\boldsymbol{\theta}_0)) \leq c_9 < \infty, \tag{3.19}$$

*where $\boldsymbol{I}_{\gamma}(\boldsymbol{\theta}_0)$ is defined in Equation (3.7).*

**Lemma III.5.** *Uniformly in $\boldsymbol{\gamma} \in \Gamma$, we have*

$$2(l_{\gamma}(\hat{\boldsymbol{\theta}}_n) - l_{\gamma}(\hat{\boldsymbol{\theta}}_0)) = ||\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})||^2 + o_p(1), \tag{3.20}$$

*where $\hat{\boldsymbol{\theta}}_n$ is the MLE for $\boldsymbol{\theta}$ given $\boldsymbol{\gamma}$, and $\psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) = \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} \{\pi(\boldsymbol{X}_i^T \boldsymbol{\gamma}) \boldsymbol{I}_{q_2} - \boldsymbol{B}(\boldsymbol{\gamma}) \boldsymbol{A}^{-1}\}(Y_i - \boldsymbol{Z}_i^T \boldsymbol{\beta}_0) \boldsymbol{Z}_i / \sigma_0^2$.*

**Lemma III.6.** *Under the null hypothesis that there is no subgroup, we have $\boldsymbol{\gamma}_j^{(K)} = \boldsymbol{\gamma}_j^{(0)} + o_p(1)$, where $\boldsymbol{\gamma}_j^{(K)}$ is obtained in the EM iterations.*

With these lemmas, we are ready to prove Theorem III.3. In the proofs, we will follow the notations in the empirical process theory (*van der Vaart*, 1998; *van der Vaart and Wellner*, 2000) that $\mathbb{P}_n f = \sum_{i=1}^{n} f(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i)/n$, $\mathbb{P}f = \mathbb{E}f(y, \boldsymbol{z}, \boldsymbol{x})$, and $\mathbb{G}_n f = \sqrt{n}(\mathbb{P}_n - \mathbb{P})f$. Given two functions $l$ and $u$, the brackets $[l, u]$ is the set of functions $f$ such that $l \leq f \leq u$. An $\epsilon$-brackets in $L_r(P)$ is a bracket $[l, u]$ such that $P(u - l)^r < \epsilon^r$. The bracketing number $N_{[]}(\epsilon, \mathcal{F}, L_r(P))$ is the minimum number of $\epsilon$-brackets that cover $\mathcal{F}$, and the bracketing integral $J_{[]}(\delta, \mathcal{F}, L_2(P)) = \int_0^{\delta} \sqrt{\log N_{[]}(\epsilon, \mathcal{F}, L_2(P))} d\epsilon$.

*Proof.* From Lemma III.5, we have that uniformly in $\gamma \in \Gamma$,

$$2(l_{\gamma}(\hat{\boldsymbol{\theta}}_n) - l_{\gamma}(\hat{\boldsymbol{\theta}}_0)) = ||\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})||^2 + o_p(1).$$

27

Let

$$\boldsymbol{h}(\boldsymbol{\gamma}) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}).$$

Recall that $\boldsymbol{\psi}(Y, \boldsymbol{Z}, \boldsymbol{X}; \boldsymbol{\gamma}) = \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} \{\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) \boldsymbol{I}_{q_2} - \boldsymbol{B}(\boldsymbol{\gamma}) \boldsymbol{A}^{-1}\}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) \boldsymbol{Z} / \sigma_0^2$, where $\boldsymbol{A} = \mathbb{E}(\boldsymbol{Z}\boldsymbol{Z}^T), \boldsymbol{B}(\boldsymbol{\gamma}) = \mathbb{E}(\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) \boldsymbol{Z}\boldsymbol{Z}^T)$, $\boldsymbol{C}(\boldsymbol{\gamma}) = \mathbb{E}\pi(\boldsymbol{X}^T \boldsymbol{\gamma})^2 \boldsymbol{Z}\boldsymbol{Z}^T)$, and $\boldsymbol{I}_{\gamma 22 \cdot 1} = (\boldsymbol{C}(\boldsymbol{\gamma}) - \boldsymbol{B}(\boldsymbol{\gamma}) \boldsymbol{A}^{-1} \boldsymbol{B}(\boldsymbol{\gamma})) / \sigma_0^2$. Direct calculations give

$$
\begin{aligned}
\sigma_0^2 \boldsymbol{\psi}'(\boldsymbol{\gamma}) &= \frac{d\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2}}{d\boldsymbol{\gamma}} \pi(\boldsymbol{X}^T \boldsymbol{\gamma})(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T - \frac{d\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2}}{d\boldsymbol{\gamma}} \boldsymbol{B}(\boldsymbol{\gamma}) \boldsymbol{A}^{-1}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T \\
&\quad + \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} (\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) - \pi(\boldsymbol{X}^T \boldsymbol{\gamma})^2) \boldsymbol{X}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T \\
&\quad - \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} (\boldsymbol{B}(\boldsymbol{\gamma}) - \boldsymbol{C}(\boldsymbol{\gamma})) \boldsymbol{A}^{-1} \boldsymbol{X}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T.
\end{aligned}
$$

From above $\mathbb{E}\boldsymbol{\psi}'(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) = 0$, then by central limit theorem $||\boldsymbol{h}'(\gamma)|| = O_p(1)$ for each $\gamma \in \Gamma$. Now we need to show that $||\boldsymbol{h}'(\gamma)|| = O_p(1)$ holds uniformly in $\boldsymbol{\gamma} \in \Gamma$.

In the proof of Lemma III.5, we show that $\mathcal{F} = \{\pi(\boldsymbol{X}^T \boldsymbol{\gamma})(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T)^T : \boldsymbol{\gamma} \in \Gamma\}$ is P-Donsker component-wisely. Slight modification gives that $\mathcal{F}_2 = \{(\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) - \pi(\boldsymbol{X}^T \boldsymbol{\gamma})^2) \boldsymbol{X}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T : \boldsymbol{\gamma} \in \Gamma\}$ is P-Donsker component-wisely under Assumption 1. Along with the zero expectation, we have $\mathbb{G}_n \pi(\boldsymbol{X}^T \boldsymbol{\gamma})(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T = O_p(1)$ and $\mathbb{G}_n(\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) - \pi(\boldsymbol{X}^T \boldsymbol{\gamma})^2) \boldsymbol{X}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) Z^T = O_p(1)$ uniformly in $\boldsymbol{\gamma} \in \Gamma$.

Then since

$$
\begin{aligned}
\boldsymbol{h}'(\gamma) &= \frac{d\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2}}{d\boldsymbol{\gamma}} \mathbb{G}_n(\pi(\boldsymbol{X}^T \boldsymbol{\gamma})(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) \boldsymbol{Z}^T) - \frac{d\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2}}{d\boldsymbol{\gamma}} \boldsymbol{B}(\boldsymbol{\gamma}) \boldsymbol{A}^{-1} \mathbb{G}_n((Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) \boldsymbol{Z}^T) \\
&\quad + \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} \mathbb{G}_n(\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) \boldsymbol{X}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) \boldsymbol{Z}^T) - \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} (\boldsymbol{B}(\boldsymbol{\gamma}) - \boldsymbol{C}(\boldsymbol{\gamma})) \boldsymbol{A}^{-1} \\
&\quad \mathbb{G}_n(\boldsymbol{X}(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0) \boldsymbol{Z}^T),
\end{aligned}
$$

and $\mathbb{E}(\boldsymbol{h}'(\boldsymbol{\gamma})) = 0$ since $\mathbb{E}(Y - \boldsymbol{Z}^T \beta_0 | \boldsymbol{X}, \boldsymbol{Z}) = 0$. In addition, for $\boldsymbol{\gamma} \in \Gamma$, the deterministic functions $d\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} / d\boldsymbol{\gamma}$, $\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2}$, and $\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} (\boldsymbol{B}(\boldsymbol{\gamma}) - \boldsymbol{C}(\boldsymbol{\gamma})) \boldsymbol{A}^{-1}$ are bounded. Finally we will have $\boldsymbol{h}'(\boldsymbol{\gamma}) = O_p(1)$ uniformly in $\boldsymbol{\gamma} \in \Gamma$.

Then by expansion of $h(\gamma)$, we get

$$||\boldsymbol{h}(\boldsymbol{\gamma}_j^K) - \boldsymbol{h}(\boldsymbol{\gamma}_j)|| = O_p(||\boldsymbol{\gamma}_j^K - \boldsymbol{\gamma}_j||) = O_p(||\boldsymbol{\gamma}_j^K - \boldsymbol{\gamma}_j^{(0)}||) = o_p(1).$$

Therefore,

$$
\begin{aligned}
EM_j^K &= ||\boldsymbol{h}(\boldsymbol{\gamma}_j^{(K)})||^2 + o_p(1) = ||\boldsymbol{h}(\boldsymbol{\gamma}_j)||^2 + o_p(1) \\
&= ||\tfrac{1}{\sqrt{n}} \textstyle\sum_{i=1}^n \boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}_j)||^2 + o_p(1),
\end{aligned}
$$

for $j = 1, \ldots, J$. Hence, the $EM$ test statistic $\max\{EM_j^K, j = 1, \ldots, J\}$ converges to a limiting distribution. $\qquad\square$

### 3.2.3 Implementation Issues

Although we have characterized the limiting distribution of the proposed $EM$ test under the null hypothesis, we do not suggest using the asymptotic distribution to carry out the test. Even in the simplest case of $q_1 = q_2 = 1$ with no covariates, the convergence to the chi-square distribution is known to be very slow (*Goeffinet et al.*, 1992). When covariates are present, we certainly do not expect the approximation to be good. To conduct the test based on $EM^{(K)}$, we suggest using the bootstrap method. The asymptotic representations given in Section 3.2 imply the validity of the bootstrap method for computing the $p$ values of the proposed test.

For the selection of $J$ and the specific values for $\boldsymbol{\Gamma}$, we recommend a small number of $\boldsymbol{\gamma}_j$ values. If $\gamma$ is $q_2$-dimensional and $q_2$ is small, we recommend using $J = 2^{q_2-1}$, with one positive value and one negative value in each component of $\gamma_X$, so that the points of $\Gamma$ cover all quadrants. The exact values of $\boldsymbol{\gamma}_j$ are not important. Under this choice, a small value of $K = 3$ generally works well. If $q_2$ is large, we may choose a small number of $\boldsymbol{\gamma}_j$ randomly. Our empirical experience shows that higher values of $J$ do not bring sufficient gain in power. The same can be said about the value $K$. The ability to use small values of $J$ and $K$ makes the proposed $EM$ test practically

useful.

### 3.2.4   Local Power

Although the power function of the proposed EM test appears intractable, we obtain in this section the local power of the test with one starting value of $\boldsymbol{\gamma}$ $(J = 1)$. More specifically, we consider, for some $\boldsymbol{h} \in \mathbb{R}^{q_1}$, the parameters under the null hypothesis and the local alternative as $\boldsymbol{\eta}_0 = (\boldsymbol{\beta}_0, \sigma_0, \boldsymbol{0}, \boldsymbol{\gamma}_0)^T$ and $\boldsymbol{\eta}_a = (\boldsymbol{\beta}_0, \sigma_0, n^{-1/2}\boldsymbol{h}^T, \boldsymbol{\gamma}_0)^T$, respectively. In other words, we consider the hypothesis testing problem of

$$
\begin{aligned}
H_0: &\quad \boldsymbol{\beta}_2 = \boldsymbol{0}, \; v.s. \\
H_a: &\quad \boldsymbol{\beta}_2 = n^{-1/2}\boldsymbol{h}.
\end{aligned}
\tag{3.21}
$$

**Theorem III.7.** *Under $H_a$, the test statistic $T_K(\boldsymbol{\gamma}) := EM^{(K)}$, with any value $\boldsymbol{\gamma} \in \tilde{\boldsymbol{\Gamma}}$ and for any positive integer $K$, converges to a noncentral chi-square distribution with the degree of freedom $q_1$ and the noncentrality parameter*

$$
\lambda(\boldsymbol{\gamma}) = \sigma_0^{-2} ||\boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} (\mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}_0))\boldsymbol{h}||^2.
\tag{3.22}
$$

*In particular, when $\boldsymbol{\gamma} = \boldsymbol{\gamma}_0$, we have*

$$
\lambda(\boldsymbol{\gamma}_0) = \sigma_0^{-2}\boldsymbol{h}^T(C(\boldsymbol{\gamma}_0) - B(\boldsymbol{\gamma}_0)A^{-1}B(\boldsymbol{\gamma}_0))\boldsymbol{h}.
\tag{3.23}
$$

Therefore, the power of the test at $H_a$ is $P(\chi_{q_1;\lambda}^2 > \chi_{q_1}^2(1-\alpha))$ where $\chi_{q_1;\lambda}^2$ is the noncentral chi-square variable with the degree of freedom $q_1$ and the noncentrality parameter $\lambda$, and $\chi_{q_1}^2(1-\alpha)$ is the upper $\alpha$th quantile of the $\chi_{q_1}^2$. When the EM test is carried out with $J \geq 2$ values of $\boldsymbol{\gamma}$, the local power no longer has a simple expression, but it relates to the maximum of $J$ correlated noncentral chi-square random variables whose noncentrality parameters are in the form of (3.22).

Remark: we can show that the noncentral parameter $\lambda$ from $\boldsymbol{\gamma}_0$ is larger than

other $\boldsymbol{\gamma} \in \tilde{\boldsymbol{\Gamma}}$.

Note that if a matrix

$$
\begin{pmatrix} A_1 & A_2 \\ A_2^T & A_4 \end{pmatrix} \tag{3.24}
$$

is positive definite, then we also have $A_4 - A_2^T A_1^{-1} A_2 > 0$ by taking a submatrix of the original matrix.

After direct derivation, showing $\lambda(\boldsymbol{\gamma}_0) \geq \lambda(\boldsymbol{\gamma})$ is equivalent to show that the matrix

$$
\begin{aligned}
&C(\boldsymbol{\gamma}_0) - B(\boldsymbol{\gamma}_0)A^{-1}B(\boldsymbol{\gamma}_0) - (\mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}_0))^T \\
&(C(\boldsymbol{\gamma}) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}))^{-1}(\mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}_0))
\end{aligned} \tag{3.25}
$$

is non-negative definite, which suffices to show that the matrix

$$
\begin{pmatrix} C(\boldsymbol{\gamma}) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}) & \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}_0) \\ \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma}_0)A^{-1}B(\boldsymbol{\gamma}) & C(\boldsymbol{\gamma}_0) - B(\boldsymbol{\gamma}_0)A^{-1}B(\boldsymbol{\gamma}_0) \end{pmatrix} \tag{3.26}
$$

is non-negative definite.

For any $\boldsymbol{a}_1, \boldsymbol{a}_2 \in \mathbb{R}^{q_1}$, we have that

$$
\begin{aligned}
&(\boldsymbol{a}_1^T, \boldsymbol{a}_2^T) \\
&\begin{pmatrix} C(\boldsymbol{\gamma}) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}) & \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma}_0) \\ \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{Z}\boldsymbol{Z}^T) - B(\boldsymbol{\gamma}_0)A^{-1}B(\boldsymbol{\gamma}) & C(\boldsymbol{\gamma}_0) - B(\boldsymbol{\gamma}_0)A^{-1}B(\boldsymbol{\gamma}_0) \end{pmatrix} \\
&(\boldsymbol{a}_1^T, \boldsymbol{a}_2^T)^T \\
&= \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{a}_1^T\boldsymbol{Z} + \pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{a}_2^T\boldsymbol{Z})^2 \\
&- \{\mathbb{E}[(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{a}_1^T\boldsymbol{Z} + \pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{a}_2^T\boldsymbol{Z})\boldsymbol{Z}^T]\}A^{-1}\{\mathbb{E}[\boldsymbol{Z}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{a}_1^T\boldsymbol{Z} + \pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{a}_2^T\boldsymbol{Z})]\}. \\
&\geq 0.
\end{aligned} \tag{3.27}
$$

The last inequality comes from the fact that the matrix

$$
\begin{pmatrix}
\mathbb{E}\boldsymbol{Z}\boldsymbol{Z}^T & \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{a}_1^T\boldsymbol{Z} + \pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{a}_2^T\boldsymbol{Z})\boldsymbol{Z} \\
\mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{a}_1^T\boldsymbol{Z} + \pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{a}_2^T\boldsymbol{Z})\boldsymbol{Z}^T & \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{a}_1^T\boldsymbol{Z} + \pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0)\boldsymbol{a}_2^T\boldsymbol{Z})^2
\end{pmatrix}
\tag{3.28}
$$

is non-negative definite.

## 3.3 Simulations

In this section, we use simulation studies to investigate the finite sample performance of the proposed test. We show that the asymptotic distribution of the test statistic is not a good approximation, but with the bootstrap method, the test performs similarly to the likelihood ratio test in an oracle form. In all the empirical studies, we standardize all the covariates (to unit variance) except the treatment indicator and choose $c_1 = 0.2$, $c_2 = c_3 = 5$ for $\tilde{\Gamma}$.

### 3.3.1 Type I Errors

To assess the accuracy of the asymptotic approximation to the proposed test at a given value of $\boldsymbol{\gamma}$, we first report a relatively simple study based on Model (1) with $q_1 = q_2 = 2$, $\boldsymbol{\beta}_1 = (1,2)^T$, $\boldsymbol{\beta}_2 = (0,0)^T$, $\boldsymbol{X} = (1,x)^T$, $\boldsymbol{Z} = (1,z)^T$, where $x$ is distributed as $N(1,1)$, $z$ is independent of $x$, and distributed as either $N(-1,1)$ or Bernoulli with probability 0.5. The error distribution $\varepsilon$ is independent of $(x,z)$ and distributed as $N(0,0.5^2)$.

We fix $\boldsymbol{\Gamma} = \{(0.3,-0.7)^T\}$ with $J = 1$. As we vary the sample size $n$ from 60 to 1000, we report in Table 3.1 the mean value as well as the 0.90 and 0.95 upper quantiles of the test statistic based on a Monte Carlo study with 5000 data sets. They are compared with their counterparts from the limiting distribution of $\chi_2^2$. It is clear that for $n$ up to 300, the asymptotic approximation is unsatisfactory in preserving the significance levels of the test.

Table 3.1: Quality of asymptotic approximation to the the null distribution of the *EM* test statistic in a simple configuration. The column under "Asymptotic" refers to asymptotic values, and "MC" refers to Monte Carlo-based values. The last two columns show the type I errors of the test, if the asymptotic values of the critical values are used.

| $n$ | Expectation | | 5% critical value | | 10% critical value | | Type I error | |
|---|---|---|---|---|---|---|---|---|
| | Asymptotic | MC | Asymptotic | MC | Asymptotic | MC | size 0.05 | size 0.1 |
| | | | | $z \sim N(-1, 1)$ | | | | |
| 60 | 2 | 2.84 | 5.99 | 7.91 | 4.61 | 6.24 | 0.11 | 0.19 |
| 100 | 2 | 2.60 | 5.99 | 7.23 | 4.61 | 5.65 | 0.09 | 0.16 |
| 300 | 2 | 2.33 | 5.99 | 6.66 | 4.61 | 5.21 | 0.07 | 0.14 |
| 600 | 2 | 2.17 | 5.99 | 6.14 | 4.61 | 4.86 | 0.05 | 0.11 |
| 1000 | 2 | 2.15 | 5.99 | 6.31 | 4.61 | 4.85 | 0.06 | 0.11 |

| $n$ | Expectation | | 5% critical value | | 10% critical value | | Type I error | |
|---|---|---|---|---|---|---|---|---|
| | Asymptotic | MC | Asymptotic | MC | Asymptotic | MC | size 0.05 | size 0.1 |
| | | | | $z \sim Bernoulli(1/2)$ | | | | |
| 60 | 2 | 3.07 | 5.99 | 8.32 | 4.61 | 6.64 | 0.13 | 0.23 |
| 100 | 2 | 2.85 | 5.99 | 7.81 | 4.61 | 6.20 | 0.11 | 0.20 |
| 300 | 2 | 2.33 | 5.99 | 6.60 | 4.61 | 5.22 | 0.07 | 0.13 |
| 600 | 2 | 2.25 | 5.99 | 6.74 | 4.61 | 5.25 | 0.07 | 0.13 |
| 1000 | 2 | 2.10 | 5.99 | 6.17 | 4.61 | 4.83 | 0.05 | 0.11 |

We report another study with the data generated from Model (1) with $q_1 = 3, q_2 = 2$, $\boldsymbol{\beta}_1 = (1, 0, 2)^T$, $\boldsymbol{\beta}_2 = (0, 0, 0)^T$, $\boldsymbol{Z} = (1, t, x)^T$, $\boldsymbol{X} = (1, x)^T$, where $t$ resembles a treatment indicator distributed as Bernoulli(0.5), $x$ is independent of $t$ with the distribution $N(-1, 1)$, and the error $\varepsilon$ is white noise $N(0, 0.5^2)$. The $EM$ test uses $\boldsymbol{\Gamma} = \{(1, -2)^T, (1, 2)^T\}$. In Table 3.2, we demonstrate the type I errors of the proposed $EM$ test based on $5,000$ replicates with the critical values determined via the bootstrap method (with the bootstrap sample size 1000). With the sample size as small as 60, the type I errors are quite close to their nominal values, regardless of our choice of $K \in \{0, 3, 9\}$.

Table 3.2: Type I errors of the $EM$ tests with bootstrap approximations.

| $n$ | Nominal level $\alpha$ | $EM^{(0)}$ | $EM^{(3)}$ | $EM^{(9)}$ |
|---|---|---|---|---|
| $n$=60 | 0.01 | 0.012 | 0.013 | 0.012 |
| | 0.05 | 0.045 | 0.053 | 0.055 |
| | 0.10 | 0.094 | 0.098 | 0.102 |
| $n$=100 | 0.01 | 0.011 | 0.012 | 0.013 |
| | 0.05 | 0.050 | 0.053 | 0.053 |
| | 0.10 | 0.010 | 0.100 | 0.098 |

### 3.3.2 Power Comparison

We use the same model and the same $EM$ tests as in the Type I error study associated with Table 3.2, except that $\boldsymbol{\beta}_2 = (1, a, b)^T$, $\boldsymbol{\gamma} = (1, c)^T$ for some non-negative values of $a, b$, and $c$ are now used. The results are given in Table 3.3, with 0.05 as the nominal level of the test.

We also consider the performance of two likelihood ratio tests in some oracle form, where we use the true value $\boldsymbol{\gamma}_0$. In the first variation denoted by "LRT$^{(0)}$", we use the bootstrap method to carry out the likelihood ratio test, where $\boldsymbol{\gamma}_0$ is used as the starting value in maximizing the likelihood. In addition, we include the second variation, denoted by "LRT(Oracle)", where we also use the bootstrap method to

obtain the critical values, but with parameter $\boldsymbol{\gamma} = \boldsymbol{\gamma}_0$ fixed. Note that "LRT$^{(0)}$" is equivalent to $EM^{(+\infty)}$ from an ideal starting value of $\boldsymbol{\gamma}$, while "LRT(Oracle)" is equivalent to $EM^{(0)}$ where the true $\boldsymbol{\gamma}$ is used. Obviously neither variation can be carried out with real data, so they are used here as a benchmark to gauge the performance of the proposed $EM$ test.

We note that in some settings, especially when $\boldsymbol{\Gamma}$ does not contain any value that is close to the true $\boldsymbol{\gamma} = (1, c)^T$, a few $EM$ iterations help power. At $K = 3$, the power of $EM^{(K)}$ is often comparable to those of "LRT$^{(0)}$" and "LRT(Oracle)", that is, the proposed $EM$ test measures up to an oracle form of the likelihood ratio test.

We also note that, in the cases where $c = 0$, the power of the $EM$ test is noticeably lower than the powers of the Oracle tests. This is because the construction of the $EM$ test requires the value of $\gamma_X$ to stay away from zero in the $EM$ iterations. The good news is that the low power of the $EM$ test for the alternatives with $\gamma_X = 0$ is not a concern for us, because those alternatives correspond to the existence of subgroups that cannot be characterized by the covariates $\mathbf{X}$.

Table 3.3: Power (%) of the $EM$ test at the 5% level. The $EM$ test uses $\boldsymbol{\Gamma} = \{(1, 2)^T, (1, -2)^T\}$, with $K = 0, 3, 9$ iterations. The parameters of Model (1) are $\boldsymbol{\beta}_1 = (1, 0, 2)^T$, $\boldsymbol{\beta}_2 = (1, a, b)^T$, and $\boldsymbol{\gamma} = (1, c)^T$.

| $n$ | $a$ | $b$ | $c$ | LRT$^{(0)}$ | LRT(Oracle) | $EM^{(0)}$ | $EM^{(3)}$ | $EM^{(9)}$ |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 60 | 0.5 | 1 | 1 | 76.8 | 80.2 | 73.0 | 75.0 | 77.4 |
| 60 | 0.5 | 0 | 1 | 31.4 | 39.2 | 18.6 | 32.4 | 33.6 |
| 60 | 0.5 | 1 | 0 | 57.2 | 63.8 | 30.2 | 44.4 | 49.8 |
| 60 | 1.0 | 1 | 1 | 89.2 | 92.0 | 82.0 | 86.2 | 87.2 |
| 60 | 1.0 | 0 | 1 | 83.4 | 86.2 | 50.4 | 80.2 | 81.6 |
| 60 | 1.0 | 1 | 0 | 74.2 | 75.6 | 45.4 | 62.4 | 66.0 |
| 100 | 0.5 | 1 | 1 | 97.6 | 98.4 | 97.0 | 96.6 | 97.6 |
| 100 | 0.5 | 0 | 1 | 63.2 | 73.4 | 37.6 | 57.4 | 62.2 |
| 100 | 0.5 | 1 | 0 | 84.8 | 85.6 | 42.0 | 63.4 | 67.4 |
| 100 | 1.0 | 1 | 1 | 99.8 | 97.8 | 98.0 | 99.4 | 99.8 |
| 100 | 1.0 | 0 | 1 | 98.2 | 96.0 | 70.0 | 96.2 | 97.8 |
| 100 | 1.0 | 1 | 0 | 95.4 | 89.2 | 65.4 | 85.8 | 88.0 |

### 3.3.3 Misspecified Link Functions

The proposed $EM$ test is quite robust against the mispecification of the logit link used in modeling the subgroup membership. In this section, we consider three cases of misspecification in the logistic component of the model.

- C1: $\pi(x) = \Phi(x/v)$ with $v = 1.95$, where $\Phi$ is the probability distribution function of the standard normal; that is, the true model for $\delta_i$ is probit. All other aspects of the model are the same as that in Section 3.2.

- C2: $\pi(x) = F_5(x/v)$ with $v = 1.50$, where $F_5$ is the probability distribution function of the t-distribution with 5 degrees of freedom. All other aspects of the model are the same as that in Section 3.2.

- C3: The logistic component of the model has $\boldsymbol{X} = (1, x_1, x_2)^T$ where $x_1 \sim N(-1, 1)$, $x_2 \sim N(0, 1)$, and $\boldsymbol{\gamma} = (1, c, 0.5)^T$, but the variable $x_2$ is missing from our working model. All other aspects of the model are the same as that in Section 3.2.

Under these scenarios, the means of $\delta_i$ at $c = 1$ are roughly the same as that under the model we considered in the previous section. The Type I errors of the $EM$ tests are not affected by the specification of $\pi$, but the powers vary. We report the powers of the $EM^{(9)}$ test in Table 3.4, where the case C0 refers to the case with correctly specified model under Section 3.3.2. It is clear from the table that the $EM$ test retains good power under the moderate misspecifications of the logistic component of our model.

## 3.4 Proof of Lemmas

Here we provide the proofs of the lemmas in details from Section 3.2.2.

*Proof of Lemma III.4.* First, we show that for each $\boldsymbol{\gamma} \in \Gamma$, $\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)$ is positive definite.

Table 3.4: Power (%)of $EM^{(9)}$ under the correctly specified model C0 of Section 3.3.2 and three mis-specified models C1 – C3 of Section 3.3.3 at the 5% level.

| $n$ | $a$ | $b$ | $c$ | C0 | C1 | C2 | C3 |
|---|---|---|---|---|---|---|---|
| 60 | 0.5 | 1 | 1 | 77.4 | 73.8 | 77.0 | 80.2 |
| 60 | 0.5 | 0 | 1 | 33.6 | 27.6 | 35.4 | 29.0 |
| 60 | 0.5 | 1 | 0 | 49.8 | 52.6 | 49.8 | 49.4 |
| 60 | 1.0 | 1 | 1 | 87.2 | 87.0 | 87.4 | 90.4 |
| 60 | 1.0 | 0 | 1 | 81.6 | 80.2 | 83.0 | 75.6 |
| 60 | 1.0 | 1 | 0 | 66.0 | 70.6 | 67.0 | 69.4 |
| 100 | 0.5 | 1 | 1 | 97.6 | 96.4 | 97.6 | 95.8 |
| 100 | 0.5 | 0 | 1 | 62.2 | 58.2 | 62.2 | 47.6 |
| 100 | 0.5 | 1 | 0 | 67.4 | 69.2 | 68.2 | 65.8 |
| 100 | 1.0 | 1 | 1 | 99.8 | 98.6 | 99.8 | 99.2 |
| 100 | 1.0 | 0 | 1 | 97.8 | 97.4 | 97.8 | 92.4 |
| 100 | 1.0 | 1 | 0 | 88.0 | 89.8 | 87.6 | 82.0 |

Recall from Equation (3.7) that

$$
\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0, (Y, \boldsymbol{Z}, \boldsymbol{X})) = \frac{1}{\sigma_0^2}
\begin{pmatrix}
\boldsymbol{A} & 0 & \boldsymbol{B}(\boldsymbol{\gamma}) \\
0 & 2 & 0 \\
\boldsymbol{B}(\boldsymbol{\gamma}) & 0 & \boldsymbol{C}(\boldsymbol{\gamma})
\end{pmatrix},
\tag{3.29}
$$

where $\boldsymbol{A} = \mathbb{E}(\boldsymbol{Z}\boldsymbol{Z}^T), \boldsymbol{B}(\boldsymbol{\gamma}) = \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T), \boldsymbol{C}(\boldsymbol{\gamma}) = \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})^2\boldsymbol{Z}\boldsymbol{Z}^T)$.

For any vector $\boldsymbol{a} \in \mathbb{R}^{2 \times q_1 + 1}$. Write $\boldsymbol{a}^T$ as $(\boldsymbol{a}_1^T, a_2, \boldsymbol{a}_3^T)$, in which $\boldsymbol{a}_1, \boldsymbol{a}_3 \in \mathbb{R}^{q_1}$ and $a_2 \in \mathbb{R}$. Then

$$
\begin{aligned}
\boldsymbol{a}^T \boldsymbol{I}_{\boldsymbol{\gamma}}(\theta_0)\boldsymbol{a} &= (\boldsymbol{a}_1^T, a_2^T, \boldsymbol{a}_3^T)
\begin{pmatrix}
\mathbb{E}(\boldsymbol{Z}\boldsymbol{Z}^T) & 0 & \mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T) \\
0 & 2 & 0 \\
\mathbb{E}(\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T) & 0 & \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})^2\boldsymbol{Z}\boldsymbol{Z}^T)
\end{pmatrix}
\begin{pmatrix}
\boldsymbol{a}_1 \\
a_2 \\
\boldsymbol{a}_3
\end{pmatrix} \\
&= \mathbb{E}((\boldsymbol{a}_1^T + \boldsymbol{a}_3\pi(\boldsymbol{X}^T\boldsymbol{\gamma}))Z)^2 + 2a_2^2.
\end{aligned}
$$

Since from Assumption 1, $Z$ is linearly independent, the matrix $I_{\boldsymbol{\gamma}}(\theta_0)$ is not positive definite if and only if there exist $\boldsymbol{a}_1, \boldsymbol{a}_3 \in \mathbb{R}^{q_1}$ such that $||\boldsymbol{a}_1||^2 + ||\boldsymbol{a}_3||^2 \neq 0$, but

$\boldsymbol{a}_1^T + \boldsymbol{a}_3^T \pi(\boldsymbol{X}^T \boldsymbol{\gamma}) = 0$. This is equivalent to the condition that $\pi(\boldsymbol{X}^T \boldsymbol{\gamma})$ is a constant. But according to the definition of $\Gamma$ in (3.11), for any $\boldsymbol{\gamma} \in \Gamma$, $\pi(\boldsymbol{X}^T \boldsymbol{\gamma})$ is a non-constant random variable. Therefore, under Assumption 1, for each $\boldsymbol{\gamma} \in \Gamma$, the fisher information $\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)$ is positive definite.

Let $h_1(\boldsymbol{\gamma}) = \lambda_{\min}(I_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) > 0$, and $h_2(\boldsymbol{\gamma}) = \lambda_{\max}(\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) < \infty$. We can see that $h_1, h_2$ are continuous functions, and along with the compactness of the parameter set $\Gamma$, we complete the proof of the lemma. $\qquad\square$

*Proof of Lemma III.5.* To apply the uniform law of large numbers, we explore some properties for $\dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}_0)$. By direct calculations, we show that for $\boldsymbol{\theta} \in \Theta$, there exist functions $L_1, L_2$ such that $||\dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}) - \dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}_0)|| < L_1(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i)$, $||\dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta})|| < L_2(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i)$, and $\mathbb{E}_{\boldsymbol{\theta}_0} L_1 < \infty, \mathbb{E}_{\boldsymbol{\theta}_0} L_2 < \infty$.

Then by Theorem 2 in *Jennrich* (1969), we have

$$\sup_{\boldsymbol{\theta} \in \Theta, \boldsymbol{\gamma} \in \Gamma} || \sum_{i=1}^{n} (\dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}) - \dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}_0)/n) - \mathbb{E}_{\boldsymbol{\theta}_0} \dot{l}_{\boldsymbol{\gamma}}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\theta})|| = o_p(1),$$

(3.30)

and

$$\sup_{\boldsymbol{\gamma} \in \Gamma} || \frac{1}{n} \sum_{i=1}^{n} (\dot{l}_{\boldsymbol{\gamma}}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}_0))|| = o_p(1).$$

Therefore, taking $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_n$ in Equation (3.30) gives

$$\sup_{\boldsymbol{\theta} \in \Theta, \boldsymbol{\gamma} \in \Gamma} ||\mathbb{E}_{\boldsymbol{\theta}_0} \dot{l}_{\boldsymbol{\gamma}}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\theta})|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_n}|| = o_p(1).$$

(3.31)

By expansion of the expectation of the first derivative of the log-likelihood at $\boldsymbol{\theta}_0$,

$$\mathbb{E}\dot{l}_{\boldsymbol{\gamma}}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\theta}) = \mathbb{E}\dot{l}_{\boldsymbol{\gamma}}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\theta}_0) + \mathbb{E}\ddot{l}_{\boldsymbol{\gamma}}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\theta})|_{\boldsymbol{\theta} = \boldsymbol{\theta}^*}(\boldsymbol{\theta} - \boldsymbol{\theta}_0),$$

(3.32)

for some $\boldsymbol{\theta}^* \in \Theta$, $||\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*| < ||\boldsymbol{\theta}_0 - \boldsymbol{\theta}||$, denoted as $\boldsymbol{\theta}^*(\boldsymbol{\theta})$.

Let $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_n$ in Equation (3.32), we have

$$\mathbb{E}_{\theta_0}(\dot{l}_{\boldsymbol{\gamma}}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n}) = -\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}^*(\hat{\boldsymbol{\theta}}_n))(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0). \tag{3.33}$$

By Equation (3.31) and Equation (3.33), we have

$$\sup_{\boldsymbol{\gamma} \in \Gamma} |\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}^*(\hat{\boldsymbol{\theta}}_n))(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)| = o_p(1).$$

By Lemma III.4, we have

$$\sup_{\boldsymbol{\gamma} \in \Gamma} ||\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0|| = o_p(1). \tag{3.34}$$

By the uniform law of large numbers (*Jennrich*, 1969), we can have

$$\frac{1}{n}\ddot{l}_{\boldsymbol{\gamma}}(\theta) + \boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}) = o_p(1), \tag{3.35}$$

uniformly in $\boldsymbol{\theta} \in \Theta, \boldsymbol{\gamma} \in \Gamma$.

By the central limit theorem and that $\mathbb{E}\dot{l}(Y, \boldsymbol{Z}, \boldsymbol{X}; \boldsymbol{\theta}_0) = 0$, $\dot{l}_{\boldsymbol{\gamma}}(\theta_0)/\sqrt{n} = O_p(1)$ for each $\boldsymbol{\gamma} \in \Gamma$. Recall from Equation (3.6) that

$$\begin{aligned}\dot{l}_{\boldsymbol{\gamma}}(Y, \boldsymbol{Z}, \boldsymbol{X}; \theta_0) &= \left(\tfrac{1}{\sigma_0^2}(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)\boldsymbol{Z}^T, -\tfrac{1}{\sigma_0} + \tfrac{1}{\sigma_0^2}(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)^2, \right. \\ &\quad \left. \tfrac{1}{\sigma_0^2}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)\boldsymbol{Z}^T\right)^T,\end{aligned} \tag{3.36}$$

and note that only the last term in $\dot{l}_{\boldsymbol{\gamma}}(Y, \boldsymbol{Z}, \boldsymbol{X}; \theta_0)$ involves $\boldsymbol{\gamma}$. So to show $\dot{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)/\sqrt{n} = O_p(1)$ uniformly in $\gamma \in \Gamma$ it suffices to show $\mathcal{F} = \{\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(Y - \boldsymbol{Z}^T\boldsymbol{\beta}_0)\boldsymbol{Z}^T : \boldsymbol{\gamma} \in \Gamma\}$ is P-Donsker (*van der Vaart* (1998)) component-wisely.

Suppose that $\Gamma \subset \Gamma_1 \times \cdots \Gamma_{q_2}$, where $\Gamma_m \subset \mathbb{R}$ is an interval, $m \in \{1, \cdots, q_2\}$. For any $\epsilon > 0$, for any $m \in \{1, \cdots, q_2\}$, grid $\Gamma_m$ by $a_m^1 \leq \cdots \leq a_m^{n_m}$ such that $|a_m^j - a_m^{j-1}| < \epsilon$ and $n_m \leq M_m/\epsilon + 1$ where $M_m$ is the length of $\Gamma_m$. We prove the case when $\boldsymbol{Z}$ is a scaler, otherwise we can prove it component-wisely. Write

$\boldsymbol{X} = (\boldsymbol{X}_1, \cdots, \boldsymbol{X}_{q_2})^T$. Then construct the functions of $(Y, \boldsymbol{Z}, \boldsymbol{X})$ with $\varepsilon = Y - \boldsymbol{Z}\boldsymbol{\beta}_0$,

$$
\begin{aligned}
l_{j_1,\cdots,j_{q_2}} &= \pi\big((a_1^{j_1} I_{(X_1 \varepsilon Z > 0)} + a_1^{j_1+1} I_{(X_1 \varepsilon Z < 0)})X_1, \cdots, (a_{q_2}^{j_{q_2}} I_{(X_{q_2} \varepsilon Z > 0)} \\
&\quad + a_{q_2}^{j_{q_2}+1} I_{(X_{q_2} \varepsilon Z < 0)})X_{q_2}\big)\varepsilon Z,
\end{aligned}
$$

and

$$
\begin{aligned}
u_{j_1,\cdots,j_{q_2}} &= \pi\big((a_1^{j_1+1} I_{(X_1 \varepsilon Z > 0)} + a_1^{j_1} I_{(X_1 \varepsilon Z < 0)})X_1, \cdots, (a_{q_2}^{j_{q_2}+1} I_{(X_{q_2} \varepsilon Z > 0)} \\
&\quad + a_{q_2}^{j_{q_2}} I_{(X_{q_2} \varepsilon Z < 0)})X_{q_2}\big)\varepsilon Z.
\end{aligned}
$$

Then $l_{j_1,\cdots,j_{q_2}} \leq u_{j_1,\cdots,j_{q_2}}$. We have brackets (*van der Vaart (1998); van der Vaart and Wellner (2000)*) $\{(l_{j_1,\cdots,j_{q_2}}, u_{j_1,\cdots,j_{q_2}}) : j_m = 1, \cdots, n_m, m = 1, \cdots, q_2\}$. In total the number of such brackets are bounded by $C_0/\epsilon^{q_2}$ with $C_0 = M_1 \cdots M_{q_2}$.

Since that $|exp(x_1)/(1 + exp(x_1)) - exp(x_2)/(1 + exp(x_2))| = |1/(1 + exp(x_1)) - 1/(1 + exp(x_2))| \leq |x_1 - x_2|$, with direct algebra we have

$$
||u_{j_1,\cdots,j_{q_2}} - l_{j_1,\cdots,j_{q_2}}||_{L_2}^2 \leq \mathbb{E}\varepsilon^2 Z^2 \mathbb{E}||X||^2 \epsilon^2 = C_1^2 \epsilon^2,
$$

where $C_1^2 = \mathbb{E}\varepsilon^2 Z^2 \mathbb{E}||X||^2 < \infty$ by Assumption 1.

Let $\epsilon_0^2 = C_1 \epsilon^2$, so $\epsilon = \epsilon_0/C_1$. Therefore the bracketing numbers

$$
N_{[]}(\epsilon_0, \mathcal{F}, L_2) \leq C_0/\epsilon^{q_2} \leq C_0(C_1)^{q_2}/\epsilon_0^{q_2}.
$$

Moreover, the $L_2$ bracketing integral

$$
\begin{aligned}
J_{[]}(1, \mathcal{F}, L_2) &= \int_0^1 \sqrt{\log N_{[]}(\epsilon_0, \mathcal{F}, L_2)} d\epsilon_0 \\
&< \int_0^1 \sqrt{\log C_0(C_1)^{q_2}/\epsilon_0^{q_2}} d\epsilon_0 < \infty,
\end{aligned}
$$

because that $\int_0^1 \log \epsilon_0 d\epsilon_0 < \infty$. Finally by Theorem 19.5 in *van der Vaart (1998)*, $\mathcal{F} = \{\pi(\boldsymbol{X}^T \boldsymbol{\gamma})(Y - \boldsymbol{Z}^T \boldsymbol{\beta}_0)\boldsymbol{Z}^T : \boldsymbol{\gamma} \in \Gamma\}$ is P-Donsker.

Therefore, $\dot{l}_{\gamma}(\theta_0)/\sqrt{n} = O_p(1)$ uniformly in $\gamma \in \Gamma$.

Applying a first-order Taylor expansion to $\dot{l}_{\gamma}(\theta_0)$, we have for some $\theta^* \in \Theta$,

$$\frac{1}{\sqrt{n}}\dot{l}_{\gamma}(\theta_0) = \{\boldsymbol{I}_{\gamma}(\theta_0) - [\frac{1}{n}\sum_{i=1}^{n}\ddot{l}_{\gamma}(\boldsymbol{X}_i;\theta^*) + \boldsymbol{I}_{\gamma}(\theta^*)] + [\boldsymbol{I}_{\gamma}(\theta^*) - \boldsymbol{I}_{\gamma}(\theta_0)]\}$$
$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) + o_p(1). \tag{3.37}$$

Since now we have $\sup_{\gamma \in \Gamma} ||\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0|| = o_p(1)(3.34)$; $\sup_{\gamma \in \Gamma} ||\boldsymbol{I}_{\gamma}(\theta^*) - \boldsymbol{I}_{\gamma}(\theta_0)|| = o_p(1)$ from the reasons $||\theta^* - \theta_0|| = o_p(1)$, $\boldsymbol{I}_{\gamma}(\theta)$ is continuous in $\gamma$ and $\theta$, and $\Gamma$ is a compact set; $\sup_{\gamma \in \Gamma} ||1/n\sum_{i=1}^{n}\ddot{l}_{\gamma}(\boldsymbol{X}_i, \theta^*) + \boldsymbol{I}_{\gamma}(\theta^*)|| = o_p(1)$ (3.35), $\sup_{\gamma \in \Gamma} ||1/\sqrt{n}\dot{l}_{\gamma}(\theta_0)|| = O_p(1)$, and $\boldsymbol{I}_{\gamma}(\theta_0)$ have lower bounds of the eigenvalues in $\gamma \in \Gamma$ from Lemma III.4, therefore,

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_p(1) \tag{3.38}$$

uniformly in $\gamma \in \Gamma$.

Then from Equation (3.37), we get

$$\frac{1}{\sqrt{n}}\dot{l}_{\gamma}(\theta_0) = \boldsymbol{I}_{\gamma}(\theta_0)\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) + o_p(1) \tag{3.39}$$

uniformly in $\gamma \in \Gamma$.

Then, we have, for some $\theta^* \in \Theta$,

$$\begin{aligned}
l_{\gamma}(\hat{\boldsymbol{\theta}}_n) - l_{\gamma}(\theta_0) &= \dot{l}_{\gamma}(\theta_0)^T(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) + \frac{1}{2}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)^T\ddot{l}_{\gamma}(\theta^*)(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \\
&= \dot{l}_{\gamma}(\theta_0)^T(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) + \frac{1}{2}(\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0))^T(-\boldsymbol{I}_{\gamma}(\theta_0) + (-\boldsymbol{I}_{\gamma}(\theta^*) \\
&\quad + \boldsymbol{I}_{\gamma}(\theta_0)) + (\frac{1}{n}\ddot{l}_{\gamma}(\theta^*) + \boldsymbol{I}_{\gamma}(\theta^*)))(\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)) \\
&= \dot{l}_{\gamma}(\theta_0)^T(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) - \frac{1}{2}(\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0))^T\boldsymbol{I}_{\gamma}(\theta_0)(\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)) \\
&\quad + o_p(1).
\end{aligned} \tag{3.40}$$

uniformly in $\gamma \in \Gamma$.

From Equation (3.39) and (3.40), we have

$$
\begin{aligned}
2(l_{\boldsymbol{\gamma}}(\hat{\boldsymbol{\theta}}_n) - l_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) &= (\tfrac{1}{\sqrt{n}}\dot{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0))^T \boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)^{-1}(\tfrac{1}{\sqrt{n}}\dot{l}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) \\
&= (\boldsymbol{M}_1(\boldsymbol{\gamma}), \boldsymbol{M}_2(\boldsymbol{\gamma}))^T \boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)^{-1}(\boldsymbol{M}_1(\boldsymbol{\gamma}), \boldsymbol{M}_2(\boldsymbol{\gamma})) + o_p(1).
\end{aligned}
$$
$$(3.41)$$

and similarly,

$$
2(l_{\boldsymbol{\gamma}}(\hat{\boldsymbol{\theta}}_0) - l_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)) = \boldsymbol{M}_1(\boldsymbol{\gamma})^T I_{\boldsymbol{\gamma}11}(\boldsymbol{\theta}_0)^{-1} \boldsymbol{M}_1(\boldsymbol{\gamma}) + o_p(1), \tag{3.42}
$$

uniformly in $\boldsymbol{\gamma} \in \Gamma$, where $\boldsymbol{I}_{\boldsymbol{\gamma}11}(\boldsymbol{\theta}_0)$ is the top left $(q_1 + 1) \times (q_1 + 1)$ submatrix of $I_{\boldsymbol{\gamma}}(\theta_0)$, and $\boldsymbol{M}_1(\boldsymbol{\gamma})$ and $\boldsymbol{M}_2(\boldsymbol{\gamma})$ are the same as defined in Equation (3.8) and (3.9), that is, $1/\sqrt{n}\dot{\boldsymbol{l}}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0) = (M_1(\boldsymbol{\gamma}), M_2(\boldsymbol{\gamma}))^T$.

Finally from the same argument about matrix manipulation as Section 3.1.1, we complete the proof of the lemma. □

*Proof of Lemma III.6.* We present the proof in two steps.

*Step 1:* we show that $\boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}_0 + o_p(1)$ for $k = 1, \ldots, K$, a fact that is also stated at the beginning of the Appendix.

With the initial parameter $\boldsymbol{\gamma}^{(0)}$, the $EM$ algorithm finds

$$
\boldsymbol{\theta}^{(0)} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^{n} \log f(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}, \boldsymbol{\gamma}^{(0)}),
$$

where $f$ here denotes the joint density of $(Y, \boldsymbol{Z}, \boldsymbol{X})$. In fact, $\boldsymbol{\theta}^{(0)} = ((\boldsymbol{\beta}_1^{(0)})^T, \sigma^{(0)}, (\boldsymbol{\beta}_2^{(0)})^T)^T$ is the MLE computed with fixed $\boldsymbol{\gamma}^{(0)}$. It follows from the positive definiteness of $\boldsymbol{I}_{\boldsymbol{\gamma}}(\boldsymbol{\theta}_0)$ that the MLE is consistent when the null hypothesis is true, that is, $\boldsymbol{\theta}^{(0)} = \boldsymbol{\theta}_0 + o_p(1)$, when the true parameter is $\boldsymbol{\theta}_0 = ((\boldsymbol{\beta}_0^T), \sigma_0, \boldsymbol{0}_{1 \times q_1})^T$ under the null hypothesis.

Then for each EM iteration as described in Section 3.2, we obtain $\{(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}) : k = 1, \ldots, K\}$, a sequence determined by $\boldsymbol{\gamma}^{(0)}$. Since each $EM$ iteration increases the

log-likelihood, resulting in

$$\sum_{i=1}^n \log f(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}) \geq \sum_{i=1}^n \log f(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}^{(0)}, \boldsymbol{\gamma}^{(0)})$$
$$\geq \sum_{i=1}^n \log f(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}_0, \boldsymbol{\gamma}^{(0)})$$
$$= \sum_{i=1}^n \log f(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\theta}_0, \boldsymbol{\gamma}^{(k)}),$$

then from Theorem 5.14 of *van der Vaart* (1998) it follows that $\boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}_0 + o_p(1)$ for $k = 1, \ldots, K$.

*Step 2:* we show that $\boldsymbol{\gamma}^{(k)} = \boldsymbol{\gamma}_0 + o_p(1)$ for $k = 1, \ldots, K$. Let

$$b(\boldsymbol{\theta}, \boldsymbol{\gamma}, \tilde{\boldsymbol{\gamma}}; Y, \boldsymbol{Z}, \boldsymbol{X}) = (P(\delta = 1 | Y, \boldsymbol{Z}, \boldsymbol{X}; \boldsymbol{\theta}, \boldsymbol{\gamma}) - \pi(\boldsymbol{X}^T \tilde{\boldsymbol{\gamma}})) \boldsymbol{X},$$

and $g_n(\boldsymbol{\theta}, \boldsymbol{\gamma}, \tilde{\boldsymbol{\gamma}}) = \mathbb{P}_n b(\boldsymbol{\theta}, \boldsymbol{\gamma}, \tilde{\boldsymbol{\gamma}}; Y, \boldsymbol{Z}, \boldsymbol{X})$. Note that $b(\boldsymbol{\theta}_0, \boldsymbol{\gamma}, \boldsymbol{\gamma}; Y, \boldsymbol{Z}, \boldsymbol{X}) = 0$ for any $\boldsymbol{\gamma} \in \tilde{\Gamma}$.

The value of $\boldsymbol{\gamma}$ after the $k$th iteration, which satisfies $g_n(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}, \boldsymbol{\gamma}) = 0$, is denoted as $\boldsymbol{\gamma}_{temp}^{(k+1)}$, as given by (3.15). By the empirical process theory used in the proof of Lemma III.5, we have $g_n(\boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\gamma}) = o_p(1)$ uniformly in $(\boldsymbol{\theta}, \boldsymbol{\gamma}) \in \Theta \times \tilde{\Gamma}$. We also know $\mathbb{E}b(\boldsymbol{\theta}, \boldsymbol{\gamma}, \boldsymbol{\gamma}; Y, \boldsymbol{Z}, \boldsymbol{X})$ is uniformly continuous on the compact set $\Theta \times \tilde{\Gamma}$. Therefore,

$$g_n(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}, \boldsymbol{\gamma}^{(k)}) = \mathbb{E}b(\boldsymbol{\theta}_0, \boldsymbol{\gamma}^{(k)}, \boldsymbol{\gamma}^{(k)}; Y, \boldsymbol{Z}, \boldsymbol{X}) + o_p(1) = o_p(1),$$

for $k = 1, \ldots, K$.

Now considering $g_n$ as a function of its last argument, we have

$$
\begin{aligned}
o_p(1) &= g_n(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}, \boldsymbol{\gamma}_{temp}^{(k+1)}) - g_n(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}, \boldsymbol{\gamma}^{(k)}) \\
&= \tfrac{dg_n(\boldsymbol{\theta}^{(k)}, \boldsymbol{\gamma}^{(k)}, \tilde{\boldsymbol{\gamma}})}{d\tilde{\boldsymbol{\gamma}}}|_{\tilde{\boldsymbol{\gamma}} = \boldsymbol{\gamma}^*} (\boldsymbol{\gamma}_{temp}^{(k+1)} - \boldsymbol{\gamma}^{(k)}) \qquad (3.43) \\
&= \mathbb{P}_n \pi(\boldsymbol{X}^T \boldsymbol{\gamma}^*)(1 - \pi(\boldsymbol{X}^T \boldsymbol{\gamma}^*)) \boldsymbol{X} \boldsymbol{X}^T (\boldsymbol{\gamma}_{temp}^{(k+1)} - \boldsymbol{\gamma}^{(k)}),
\end{aligned}
$$

where $\boldsymbol{\gamma}^*$ satisfies $||\boldsymbol{\gamma}^* - \boldsymbol{\gamma}^{(k)}|| \le ||\boldsymbol{\gamma}_{temp}^{(k+1)} - \boldsymbol{\gamma}^{(k)}||$.

We consider the case of $\boldsymbol{\gamma}_{temp}^{(k+1)} \in \tilde{\Gamma}$. In this case, $\boldsymbol{\gamma}^* \in \tilde{\Gamma}$. By the uniform law of large numbers, $\mathbb{P}_n \pi(\boldsymbol{X}^T\boldsymbol{\gamma})(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma}))\boldsymbol{X}\boldsymbol{X}^T$ converges to $P\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma}))\boldsymbol{X}\boldsymbol{X}^T$ uniformly in $\boldsymbol{\gamma}$. In addition, by the same argument for Lemma III.4, we know that the eigenvalues of the matrix $P\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma}))\boldsymbol{X}\boldsymbol{X}^T$ have a positive lower bound uniformly for $\boldsymbol{\gamma} \in \tilde{\Gamma}$. Then it follows from (3.43) that $\boldsymbol{\gamma}^{(k+1)} - \boldsymbol{\gamma}^{(k)} = \boldsymbol{\gamma}_{temp}^{(k+1)} - \boldsymbol{\gamma}^{(k)} = o_p(1)$.

If $\boldsymbol{\gamma}_{temp}^{(k+1)} \notin \tilde{\Gamma}$, then we have $\boldsymbol{\gamma}^{(k+1)} = \boldsymbol{\gamma}^{(k)}$. By induction, we have $\boldsymbol{\gamma}^{(k)} - \boldsymbol{\gamma}^{(0)} = o_p(1)$ for $k = 1, \ldots, K$. $\qquad\square$

*Proof of Theorem V.4.* From Section 3.1.1, we have that under $H_0$, for any $\boldsymbol{\gamma} \in \tilde{\Gamma}$,

$$T_K(\boldsymbol{\gamma}) = ||\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})||^2 + o_p(1), \qquad (3.44)$$

where $\boldsymbol{\psi}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\gamma})$ is given in (3.10).

We also have that, by Theorem 7.2 of van der Vaart (1998, p. 94 ), the log-likelihood ratio at parameters $\boldsymbol{\eta}_a$ and $\boldsymbol{\eta}_0$

$$\log\frac{dP_{\boldsymbol{\eta}_a}^n}{dP_{\boldsymbol{\eta}_0}^n} = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\boldsymbol{h}^T\dot{l}_{\boldsymbol{\beta}_2}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i)|_{\boldsymbol{\eta}_0} - \frac{1}{2}\boldsymbol{h}^T I_{\boldsymbol{\beta}_2}\boldsymbol{h} + o_p(1), \qquad (3.45)$$

under $H_0$, in which

$$\dot{l}_{\boldsymbol{\beta}_2}(y, \boldsymbol{z}, \boldsymbol{x})|_{\boldsymbol{\eta}_0} = \frac{\partial l(y, \boldsymbol{z}, \boldsymbol{x})}{\partial \boldsymbol{\beta}_2}|_{\boldsymbol{\eta}_0} = \sigma_0^{-2}\pi(\boldsymbol{x}^T\boldsymbol{\gamma}_0)(y - \boldsymbol{z}^T\boldsymbol{\beta}_0)\boldsymbol{z}, \qquad (3.46)$$

and

$$I_{\boldsymbol{\beta}_2} = -\mathbb{E}\frac{\partial^2 l(y, \boldsymbol{z}, \boldsymbol{x})}{\partial \boldsymbol{\beta}_2 \partial \boldsymbol{\beta}_2^T}|_{\boldsymbol{\eta}_0} = C(\boldsymbol{\gamma}_0). \qquad (3.47)$$

From (3.45), we have, by the central limit theorem and the fact that $\mathbb{E}\boldsymbol{\psi} = \boldsymbol{0}$ and

$\text{Var}(\boldsymbol{\psi}) = \boldsymbol{I}_{q_1}$, the identity matrix of dimension $q_1$,

$$\begin{pmatrix} \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) \\ \log \frac{dP_{\boldsymbol{\eta}_a}^n}{dP_{\boldsymbol{\eta}_0}^n} \end{pmatrix}$$

converges in distribution to

$$N \left( \begin{pmatrix} \boldsymbol{0} \\ -\frac{1}{2} \boldsymbol{h}^T C(\boldsymbol{\gamma}_0) \boldsymbol{h} \end{pmatrix}, \begin{pmatrix} I_{q_1} & \boldsymbol{s}(\boldsymbol{\eta}_0, \boldsymbol{h}) \\ \boldsymbol{s}(\boldsymbol{\eta}_0, \boldsymbol{h})^T & \boldsymbol{h}^T C(\boldsymbol{\gamma}_0) \boldsymbol{h} \end{pmatrix} \right) \tag{3.48}$$

under $H_0$, where

$$\begin{aligned} \boldsymbol{s}(\boldsymbol{\eta}_0, \boldsymbol{h}) &= \mathbb{E}_{\boldsymbol{\eta}_0} \left( \boldsymbol{h}^T \dot{l}_{\boldsymbol{\beta}_2}(y, \boldsymbol{z}, \boldsymbol{x})|_{\boldsymbol{\eta}_0} \boldsymbol{\psi}(y, \boldsymbol{z}, \boldsymbol{x}; \boldsymbol{\gamma}) \right) \\ &= \sigma_0^{-2} \boldsymbol{I}_{\gamma 22 \cdot 1}^{-1/2} (\mathbb{E}(\pi(\boldsymbol{X}^T \boldsymbol{\gamma}) \pi(\boldsymbol{X}^T \boldsymbol{\gamma}_0) \boldsymbol{Z} \boldsymbol{Z}^T) - B(\boldsymbol{\gamma}) A^{-1} B(\boldsymbol{\gamma}_0)) \boldsymbol{h}. \end{aligned} \tag{3.49}$$

By Le Cam's third lemma, of van der Vaart (1998, p. 90 ), we have that,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \boldsymbol{\psi}(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) \to N \left( \boldsymbol{s}(\boldsymbol{\eta}_0, \boldsymbol{h}), I_{q_1} \right) \tag{3.50}$$

in distribution under $H_a$.

Therefore the test statistic has a noncentral $\chi^2$ distribution with the noncentral parameter $\lambda(\boldsymbol{\gamma}) = ||\boldsymbol{s}(\boldsymbol{\eta}_0, \boldsymbol{h})||^2$, as in (3.22). In particular, if $\boldsymbol{\gamma} = \boldsymbol{\gamma}_0$, we have (3.23). □

# CHAPTER IV

# Application to the AIDS Data

In this chapter, we apply the aforementioned methods to our motivation example in Section 1.1.

In this randomized trial, ACTG 320 Study, AIDS patients are randomly assigned to one of two daily regimens: one is the treatment with the protease inhibitor indinavir in addition to zidovudine and lamivudine, and the second is the treatment with the two nucleosides zidovudine and lamivudine alone. Following the analysis of *Hammer et al.* (1997) and *Zhao et al.* (2013), we analyze the $CD4$ count change at the 24th week as the response with three baseline variables: age, baseline $CD4$ counts, and $RNA$ concentration on the logarithm scale. We ask whether a subgroup of patients has greater benefits from the treatment of adding a protease inhibitor to two nucleoside analogues, and how the baseline variables can be used to predict the subgroup membership. This of course is not meant to be a full investigation of the ACTG 320 trial, but is used to demonstrate how our proposed subgroup analysis can add value to the existing methods.

In the following analysis, we use $cd4.0$, $rna.0$, and $cd4.24$ to denote the baseline $CD4$ counts, baseline $RNA$ concentration, and the $CD4$ change at the 24th week, respectively. We use $trt$ as the treatment indicator, with 1 denoting the treatment of adding a protease inhibitor. We work with the subjects with no missing

values and without extreme $CD4$ counts, i.e., $cd4.0 \in (0, 200]$ and $cd4.24 \leq 400$, giving rise to a sample size of $n = 800$ subjects. We give a summary of the data in Table 4.1. By identifying the covariates $\boldsymbol{Z} = (1, trt, \log(cd4.0), \log_{10}(rna.0), Age)$

Table 4.1: Summary statistics for our ACTG study ($n = 800$)

| | cd4.24 | log(cd4.0) | $log_{10}(rna.0)$ | trt | Age |
|---|---|---|---|---|---|
| Min. | -132.00 | -0.69 | 1.70 | 0 | 15.97 |
| 0.25 Quantile | 0.00 | 3.07 | 4.66 | 0 | 33.33 |
| Median | 34.00 | 4.14 | 5.09 | 0 | 38.33 |
| Mean | 52.69 | 3.84 | 4.98 | 0.47 | 39.43 |
| 0.75 Quantile | 89.75 | 4.82 | 5.45 | 1 | 44.65 |
| Max. | 395.00 | 5.29 | 5.88 | 1 | 73.93 |
| s.d. | 74.51 | 1.15 | 0.68 | 0.50 | 9.01 |

for the normal component and $\boldsymbol{X} = (1, \log(cd4.0), \log_{10}(rna.0), Age)$ for the logistic component in Model (2.1), we use the proposed $EM^{(K)}$ test for the existence of subgroups. The choice of the specific bases in $\log(cd4.0)$ and $\log_{10}(rna.0)$ carries no significance; we simply follow some of the earlier work when including those variables. We take three randomly generated values from $\tilde{\Gamma}$ to form the set of initial values $\boldsymbol{\Gamma} = \{(2.44, -3.35, -2.33, 1.24)^T, (0.95, -4.49, 0.47, 4.46)^T, (1.00, -2.72, 4.64, -2.84)^T\}$, with the resulting $p$-value $< 0.001$ for $K = 0, 3$ and $9$. In fact the $p$ values are insensitive to the choice of $K$. It is clear that we reject the null hypothesis of no subgroups in this study. The estimates of the parameters and their standard errors are given in Table 4.2. As we see from the table, the differential treatment effects were evident, but $Age$ is statistically insignificant for subgroup membership. If a subject has higher baseline $CD4$ counts and higher baseline $RNA$ concentration, he/she is more likely to be in a subgroup where the treatment effect on the 24-week $CD4$ change is much greater, as demonstrated by the estimate of $\boldsymbol{\beta}_2$ for $trt$.

A by-product of our analysis is that we can use $\pi(X^T \boldsymbol{\gamma})$ to score any prospective patient and suggest that the patients with a higher score receive the treatment. A different scoring system was developed in *Zhao et al.* (2013), and a quick comparison

is in order.

Suppose that a subject with a high score of

$$S_1(\boldsymbol{X}) = \pi(-7.89 + 0.44\log(cd4.0) + 1.10\log_{10}(rna.0) - 0.02\text{Age}) \tag{4.1}$$

receives the treatment based on the estimated $\boldsymbol{\gamma}$ in our model. Applying the method of *Zhao et al.* (2013) to the same data set, we obtained the score of

$$S_2(\boldsymbol{X}) = -64.76 + 8.69\log(cd4.0) + 25.87\log_{10}(rna.0) - 0.76\text{Age} \tag{4.2}$$

for the same purpose. We plot the two scores in Figure 4.2, and they have a high rank correlation. Since the two scores are not on the same scale, we use the quantiles

Table 4.2: Parameter estimates and their standard errors when Model (2.1) is used to fit the data in the ACTG study. The variable names such as *trt* and $\log(cd4.0)$ are attached to the coefficients $\boldsymbol{\beta}_1, \boldsymbol{\beta}_2$ and $\boldsymbol{\gamma}$, whereas $\boldsymbol{\beta}_1(1), \boldsymbol{\beta}_2(1)$ and $\boldsymbol{\gamma}(1)$ refer to the intercepts.

|  | $\boldsymbol{\beta}_1(1)$ | $\boldsymbol{\beta}_1(trt)$ | $\boldsymbol{\beta}_1(\log(cd4.0))$ | $\boldsymbol{\beta}_1(\log_{10}(rna.0))$ | $\boldsymbol{\beta}_1(\text{Age})$ |
|---|---|---|---|---|---|
| Parameter | -35.96 | 41.24 | -1.97 | 6.10 | 0.69 |
| Standard error | 23.12 | 5.23 | 2.02 | 3.39 | 0.25 |
|  | $\boldsymbol{\beta}_2(1)$ | $\boldsymbol{\beta}_2(trt)$ | $\boldsymbol{\beta}_2(\log(cd4.0))$ | $\boldsymbol{\beta}_2(\log_{10}(rna.0))$ | $\boldsymbol{\beta}_2(\text{Age})$ |
| Parameter | 97.00 | 112.98 | 22.16 | -24.66 | -0.92 |
| Standard error | 89.40 | 15.79 | 7.56 | 13.14 | 0.90 |
|  | $\boldsymbol{\gamma}(1)$ | $\boldsymbol{\gamma}(\log(cd4.0))$ | $\boldsymbol{\gamma}(\log_{10}(rna.0))$ | $\boldsymbol{\gamma}(\text{Age})$ | $\sigma$ |
| Parameter | -7.89 | 0.44 | 1.10 | -0.02 | 49.78 |
| Standard error | 2.20 | 0.18 | 0.35 | 0.02 | 11.77 |

of the scores in determining subgroups. For any $q \in (0,1)$, we assign any patient whose score is above the $q$ quantile within a scoring system to subgroup 1 and the rest to subgroup 2. To see which scoring system is better, we use a 5-fold cross-validation. We use the training data to estimate the coefficients in the scores and assign subgroup membership to those subjects in the testing data. Then we take the average of the treatment effect of selected Subgroup 1 for each $q$ from the five training

sets, as well as the treatment effect difference of the target subgroup 1 from the rest of the training sets, subgroup 2. Note that the subgroup 1 consists of roughly 100 (1-$q$) percentage of the subjects in the testing data. As $q$ varies, the treatment effect differentials under the two scoring systems are shown in Figure 4.3. (use seed=0 to split the data in R.) The subgroup 1 identified by our proposed method generally enjoys a slightly greater treatment effect from the addition of the protease inhibitor than the subgroup 1 identified by the scoring system used in *Zhao et al.* (2013). In addition, the treatment effect difference is in Figure 4.4.

We hope to have balanced covariates in the treatment and control groups such that the treatment effect is not due to a particular unbalanced covarite. We split the data into half (use seed=0 in R) to obtain a training set and a testing set. On the testing set with $q = 0.5$, we have a subgroup from our method. For this subgroup, we have exactly 50% of the subjects receiving the treatment. We show the Q-Q plot of the covarites and the scores in the treatment group as well as in the control group in this chosen subgroup in Figure 4.1, from which we see the quantiles are relatively similar.

Now we repeat the experiment 100 times. We show the treatment effect and the difference at $q = 0.75$ for the first 20 experiments in Figure 4.7 and Figure 4.8, respectively. We plot the mean of the treatment effect in the chosen subgroups for each $q$ as well as the treatment effect difference of the chosen subgroup from the rest in Figure 4.5 and Figure 4.6, respectively. Overall, the subgroups selected by our method (the circles) have very similar treatment effects and treatment effect differences as those selected by the method in *Zhao et al.* (2013) (the crossings). Since the score of $S_2(\boldsymbol{X})$ was derived under a different set of model assumptions, we take the high agreement between the two scoring systems as another piece of confirmation that our structured normal-logistic mixture model designed primarily for a model-based test on the existence of subgroups captures the subgroup characteristics well.

Figure 4.1: Q-Q plot of the covariates in treatment and control groups in a chosen subgroup $q = 0.5$.



Figure 4.2: Scatter plot of $S_2(X)$ from (4.2) versus $S_1(X)$ from (4.1).

Figure 4.3: One splitting: the treatment effects for Subgroup 1 in the testing data identified by the top $100(1 - q)\%$ scores of two competing scoring systems. The open circles correspond to the scores based on $S_1(X)$ from our proposed method. The crosses correspond to the scores $S_2(X)$ used in *Zhao et al.* (2013).

Figure 4.4: One splitting: the treatment effects the difference of Subgroup 1 and Subgroup 2 in the testing data separated by the top $100(1-q)\%$ scores of two competing scoring systems. The open circles correspond to the scores based on $S_1(X)$ from our proposed method. The crosses correspond to the scores $S_2(X)$ used in *Zhao et al.* (2013).



Figure 4.5: In 20 experiments: the treatment effects for Subgroup 1 in the testing data identified by the top $100(1-0.75)\%$ scores of two competing scoring systems in repeated experiments. The open circles correspond to the scores based on $S_1(X)$ from our proposed method. The crosses correspond to the scores $S_2(X)$ used in *Zhao et al.* (2013).

Figure 4.6: In 20 experiments: the treatment effects difference of Subgroup 1 and Subgroup 2 in the testing data separated by the top $100(1 - 0.75)\%$ scores of two competing scoring systems in repeated experiments. The open circles correspond to the scores based on $S_1(X)$ from our proposed method. The crosses correspond to the scores $S_2(X)$ used in *Zhao et al.* (2013).



Figure 4.7: Mean of the treatment effects for Subgroup 1 in the testing data identified by the top $100(1-q)\%$ scores of two competing scoring systems in repeated experiments. The open circles correspond to the scores based on $S_1(X)$ from our proposed method. The crosses correspond to the scores $S_2(X)$ used in *Zhao et al.* (2013).

53

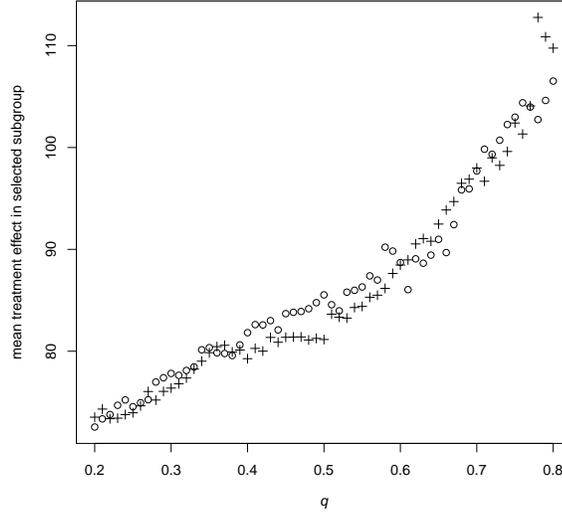Figure 4.8: Mean of the treatment effects for Subgroup 1 in the testing data identified by the top $100(1-q)\%$ scores of two competing scoring systems in repeated experiments. The open circles correspond to the scores based on $S_1(X)$ from our proposed method. The crosses correspond to the scores $S_2(X)$ used in *Zhao et al.* (2013).
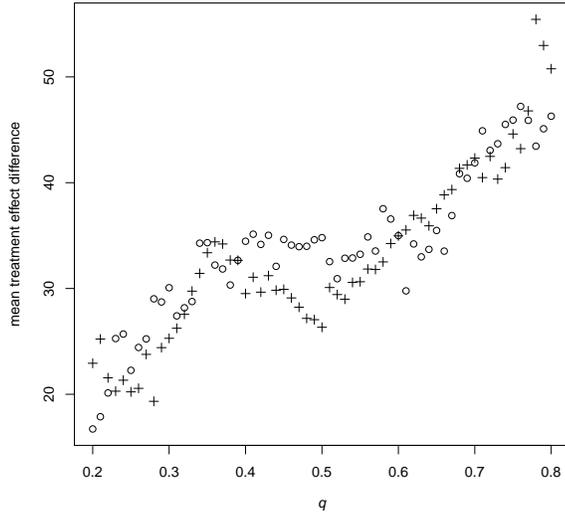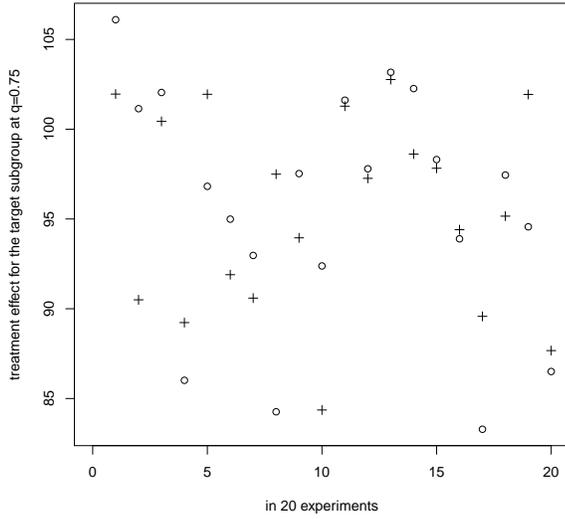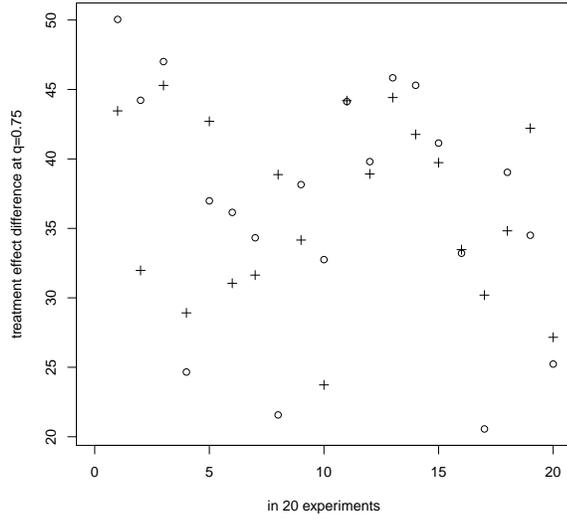
# CHAPTER V

# Subgroups of Heterogeneous Variances

For two subgroups, a difference in the means is often associated with a difference in the variances. The structured logistic-normal mixture model considered in the earlier chapters assumes equal variances in the subgroups. If we apply the equal variance model to the data that are generated from a mixture model with unequal variance, the estimators could be biased and the test might lose power. In this chapter, we consider the cases where we have heterogeneous variances in the two normal components.

## 5.1  Model

Suppose that we have a logistic-normal mixture model with unequal variances in the two normal components. For $i = 1, \ldots, n$,

$$
\begin{aligned}
Y_i &= \boldsymbol{Z}_i^T(\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2\delta_i) + \varepsilon_i\big(\sigma_1\delta_i + \sigma_2(1 - \delta_i)\big), \\
P(\delta_i = 1|\boldsymbol{X}_i, \boldsymbol{Z}_i) &= \pi(\boldsymbol{X}_i^T\boldsymbol{\gamma}) \equiv \exp(\boldsymbol{X}_i^T\boldsymbol{\gamma})/(1 + \exp(\boldsymbol{X}_i^T\boldsymbol{\gamma})), \\
P(\delta_i = 0|\boldsymbol{X}_i, \boldsymbol{Z}_i) &= 1 - P(\delta_i = 1|\boldsymbol{X}_i),
\end{aligned}
\tag{5.1}
$$

where $n$ is the sample size, $Y_i \in \mathbb{R}$ is the outcome, $\delta_i \in \{0, 1\}$ is the subgroup indicator, $\boldsymbol{Z}_i \in \mathbb{R}^{q_1}$ is the covariate associated with the subgroup mean, $\boldsymbol{X}_i \in \mathbb{R}^{q_2}$ is the covariate associated with the group membership, $\boldsymbol{\beta}_1 \in \mathbb{R}^{q_1}, \boldsymbol{\beta}_2 \in \mathbb{R}^{q_1}, \boldsymbol{\gamma} \in \mathbb{R}^{q_2}$ are the corresponding coefficients, and $\varepsilon_i \sim N(0, 1)$ are white noises. The first elements of

$\boldsymbol{X}_i$ and $\boldsymbol{Z}_i$ are 1, and the second element of $\boldsymbol{Z}_i$ is the treatment indicator. We can have overlapping variables in the random vectors of $\boldsymbol{X}_i$ and $\boldsymbol{Z}_i$. The overall parameters are $\boldsymbol{\eta}^T = (\boldsymbol{\gamma}^T, \boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T, \sigma_1, \sigma_2)$. Write $\boldsymbol{\theta}^T = (\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T, \sigma_1, \sigma_2)$ as the parameters without $\boldsymbol{\gamma}$. We observe the data $\{\boldsymbol{W}_i = (Y_i, \boldsymbol{Z}_i^T, \boldsymbol{X}_i^T), i = 1, \ldots n\}$, and $\delta_i$'s are viewed as latent variables. The observations $\boldsymbol{W}_i$'s are independent. In the case of $\boldsymbol{\beta}_2 = 0$ and $\sigma_1 = \sigma_2$, the model reduces to one normal component. Therefore, we consider $\boldsymbol{\beta}_2 = 0$ and $\sigma_1 = \sigma_2$ as our null hypothesis that no subgroup exists.

## 5.2 Penalized Likelihood

For a mixture normal model, with unequal variances, the likelihood is unbounded and the MLE does not exist (*McLachlan and Peel*, 2000). To appreciate this, let $Y_1, \cdots, Y_n$ be i.i.d. from

$$\pi N(\theta_1, \sigma_1^2) + (1 - \pi) N(\theta_2, \sigma_2^2),$$

then the likelihood

$$\Pi_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\{-\frac{-(Y_i - \theta_1)^2}{2\sigma_1^2}\} + \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\{-\frac{-(Y_i - \theta_2)^2}{2\sigma_2^2}\} \right\}$$

goes to infinity by taking $\theta_1 = Y_1$ and letting $\sigma_1$ go to zero. So the maximum likelihood estimator does not exist.

In order to restrict the two variances away from zero, an easy way is to impose a reasonable bound as follows.

In the first step, let $\hat{\sigma}_0$ be the maximum likelihood estimator of $\sigma$ under the equal variance model. Then in the following EM steps, constrain $\sigma_1$ and $\sigma_2$ in $[a\hat{\sigma}_0, b\hat{\sigma}_0]$ for some $a, b \in \mathbb{R}$ such that $0 \leq a < 1 < b \leq \infty$. With this constraint, and if the true values of $\sigma_1$ and $\sigma_2$ do fall into this range, then the estimation and testing problems

work the same way as in the common variance case. Because it is not easy to find $a$ and $b$ to ensure that the true variances are covered, we will consider the alternative approach of penalized likelihood.

We consider a penalty $p_n(\sigma)$, following *Chen and Li* (2009), with certain conditions to be specified later. In particular, we take

$$p_n(\sigma) = -\lambda\left(\frac{S_n^2}{\sigma^2} + \log\left(\frac{\sigma^2}{S_n^2}\right)\right), \tag{5.2}$$

where $S_n^2$ is a reasonable estimator of $\sigma^2$, and $\lambda$ is a tuning parameter. Since under the null hypothesis, the parameters are not identifiable from the equal variance model, but from the proof of Lemma III.6 in Section 3.4, we have consistent estimator of $\sigma^2$ given a $\boldsymbol{\gamma}$ with nonzero slope, so we suggest to use the maximum likelihood estimator of the variance given some $\boldsymbol{\gamma}$ as $S_n^2$. For the tuning parameter $\lambda$, we show how we choose it adaptively in Section 5.6, but a general data adaptive choice of $\lambda$ needs further study.

The penalized log-likelihood function is

$$
\begin{aligned}
pl(\boldsymbol{\eta}) &= \sum_{i=1}^n \log\left[\sum_{j=0}^1 f(Y_i|\boldsymbol{Z}_i, \boldsymbol{X}_i, \delta_i = j; \boldsymbol{\beta}_j, \sigma_j)P(\delta_i = j|\boldsymbol{X}_i; \boldsymbol{\gamma})\right] \\
&\quad + p_n(\sigma_1) + p_n(\sigma_2).
\end{aligned} \tag{5.3}
$$

To maximize the penalized likelihood, the slightly modified EM algorithm goes as follows: at the $(k+1)th$ step,

$$
\begin{aligned}
Q(\boldsymbol{\eta}^{(k+1)}|\boldsymbol{\eta}^{(k)}) &= \sum_{i=1}^n \mathbb{E}_{\delta_i|w_i, \boldsymbol{\eta}^{(k)}}\Big\{ I_{(\delta_i=1)} \log\Big(\frac{\pi(\boldsymbol{X}_i^T\boldsymbol{\gamma})}{\sqrt{2\pi}\sigma_1} \exp\big(-\frac{(Y_i - \boldsymbol{Z}_i^T(\boldsymbol{\beta}_1+\boldsymbol{\beta}_2))^2}{2\sigma_1^2}\big)\Big) \\
&\quad + I_{(\delta_i=0)} \log\Big(\frac{1-\pi(\boldsymbol{X}_i^T\boldsymbol{\gamma})}{\sqrt{2\pi}\sigma_2} \exp\big(-\frac{(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)^2}{2\sigma_2^2}\big)\Big)\Big\} + p_n(\sigma_1) + p_n(\sigma_2),
\end{aligned}
$$

which gives the $E$ step,

$$
\begin{aligned}
a_i^{(k)} &= P(\delta_i = 1 | Y_i, \mathbf{Z}_i, \mathbf{X}_i; \eta^{(k)}) \\
&= f(Y_i | \delta_i = 1, \mathbf{Z}_i; \theta^{(k)}) P(\delta_i = 1 | \mathbf{X}_i; \boldsymbol{\gamma}^{(k)}) / \{ f(Y_i | \delta_i = 1, \mathbf{Z}_i; \theta^{(k)}) P(\delta_i = 1 | \mathbf{X}_i; \boldsymbol{\gamma}^{(k)}) \\
&\quad + f(Y_i | \delta_i = 0, \mathbf{Z}_i; \theta^{(k)}) P(\delta_i = 0 | \mathbf{X}_i; \boldsymbol{\gamma}^{(k)}) \},
\end{aligned}
$$

$$(5.4)$$

$b_i^{(k)} = 1 - a_i^{(k)}$, $\boldsymbol{a}^{(k)} = (a_1^{(k)}, \ldots, a_n^{(k)})$, and $\boldsymbol{b}^{(k)} = (b_1^{(k)}, \ldots, b_n^{(k)})$; and the $M$ step:

$$
\begin{aligned}
\boldsymbol{\gamma}^{(k+1)} &= \operatorname{argmax}_{\boldsymbol{\gamma}} \sum_i \{ a_i^{(k)} \log \pi(\mathbf{X}_i^T \boldsymbol{\gamma}) + b_i^{(k)} \log(1 - \pi(\mathbf{X}_i^T \boldsymbol{\gamma})) \}; \\
(\boldsymbol{\beta}_{temp}^{(k+1)}, \sigma_1^{(k+1)}) &= \operatorname{argmax}_{\boldsymbol{\beta}, \sigma} \sum_i \{ a_i^{(k)} \log(\tfrac{1}{\sigma} \exp(-(Y_i - \mathbf{Z}_i^T \boldsymbol{\beta})^2 / (2\sigma^2))) \} + p_n(\sigma); \\
(\boldsymbol{\beta}_1^{(k+1)}, \sigma_2^{(k+1)}) &= \operatorname{argmax}_{\boldsymbol{\beta}, \sigma} \sum_i \{ b_i^{(k)} \log(\tfrac{1}{\sigma} \exp(-(Y_i - \mathbf{Z}_i^T \boldsymbol{\beta})^2 / (2\sigma^2))) \} + p_n(\sigma),
\end{aligned}
$$

$$(5.5)$$

and $\boldsymbol{\beta}_2^{(k+1)} = \boldsymbol{\beta}_{temp}^{(k+1)} - \boldsymbol{\beta}_1^{(k+1)}$. In the M step, the estimation of $\boldsymbol{\theta}^{(k+1)}$ is a least squares problem; for the particular penalty in Equation (5.2), the estimators of $\sigma_1^{(k+1)}$ and $\sigma_2^{(k+1)}$ given $\boldsymbol{\beta}_1^{(k+1)}$ and $\boldsymbol{\beta}_2^{(k+1)}$ are

$$
\sigma_1^{(k+1)} = \left( \frac{\sum a_i^{(k)} (Y_i - \mathbf{Z}_i^T (\boldsymbol{\beta}_1^{(k+1)} + \boldsymbol{\beta}_2^{(k+1)}))^2 / 2 + \lambda S_n^2}{\sum a_i^{(k)} / 2 + \lambda} \right)^{1/2},
$$

and

$$
\sigma_2^{(k+1)} = \left( \frac{\sum b_i^{(k)} (Y_i - \mathbf{Z}_i^T \boldsymbol{\beta}_1^{(k+1)})^2 / 2 + \lambda S_n^2}{\sum a_i^{(k)} / 2 + \lambda} \right)^{1/2}.
$$

We can see that the new estimators of $\sigma_1^2$ and $\sigma_2^2$ from penalized likelihood are weighted sums of the estimators without penalty and $S_n^2$.

In general, for the penalty $p_n(\sigma)$ which could be data-dependent. For the variables $\mathbf{X}$ and $\mathbf{Z}$, we further impose the following conditions. If we partition the covariate vector $\mathbf{Z}$ into continuous components $\mathbf{V}$ and discrete components $U$, that is, let $\mathbf{Z}^T = (\mathbf{V}^T, \mathbf{U}^T)$, where $\mathbf{V}$ consists of only continuous variables and $\mathbf{U}$ consists of only discrete variables with a finite sample space. Then we impose the following conditions:

C0. The penalty $p_n(\sigma) < 0$ for all $\sigma$, and it goes to negative infinity as $\sigma$ goes to zero almost surely.

C1. For some integer $n_0$ and all $n \geq n_0$, $\inf_{0<\sigma\leq(1/n)} \frac{p_n(\sigma)}{(\log n)^2 \log \sigma} \geq 8$, almost surely.

C2. Under the null hypothesis of $\boldsymbol{\beta_2} = 0$ and $\sigma_1 = \sigma_2 = \sigma_0$, we have $p_n(\sigma_0) = o_p(n)$ almost surely; under the alternative hypothesis, $p_n(\sigma) = o(n)$ almost surely at $\sigma_1$ and $\sigma_2$.

C3. For any unit vector $\alpha$ with the same dimension as the vector $\boldsymbol{V}$, the conditional distribution function of $\boldsymbol{V}^T \alpha | \boldsymbol{U} = \boldsymbol{u}$ is continuous for any $\boldsymbol{u}$, and the density is bounded from above.

C4. The expectation $\mathbb{E}(||\boldsymbol{V}|||\boldsymbol{U} = \boldsymbol{u}) < \infty$ uniformly in $\boldsymbol{u}$.

In particular, it is easy to see that the penalty function (5.2) satisfies conditions C0-C2. Note that for any positive $\lambda$, $p_n(\sigma)$ of (5.2) achieves its maximum at $\sigma^2 = S_n^2$, and goes to negative infinity as $\sigma$ approaches zero or infinity.

From an upcoming paper *Shen et al.* (2014), we study the consistency of the parameter estimators under the alternative hypothesis to get

**Proposition V.1.** *Under the alternative model with* $\boldsymbol{\beta_2} \neq 0$ *or* $\sigma_1 \neq \sigma_2$, *assume Conditions C0-C4, then the estimators from maximizing the penalized likelihood of Equation (5.3) are strongly consistent.*

## 5.3 Penalized $EM$ Test

In this Section, we discuss the hypothesis testing about the existence of subgroups. We test the null hypothesis of $\boldsymbol{\beta_2} = 0$ and $\sigma_1 = \sigma_2$, where the two normal components have the same parameters, and the model is degenerate.

### 5.3.1 Penalized *EM* Test Process

We have a similar EM test process to that of Section 3.2.1, but with the penalty in the objective function in Equation (5.3). For some given nonnegative integer $K$, in the end of the $K$th iteration, assume that we have the estimator $\boldsymbol{\eta}_j^{(K)}$, then, for each $j = 1, 2, \cdots, J$, let

$$pEM_j^{(K)} = 2(pl(\boldsymbol{\eta}_j^{(K)}) - pl(\hat{\boldsymbol{\theta}}_0, \boldsymbol{\gamma}_j)), \tag{5.6}$$

in which $pl(\cdot)$ is defined in Equation (5.3). The test statistic is

$$pEM^{(K)} = \max\{pEM_j^{(K)} : j = 1, 2, \ldots, J\}. \tag{5.7}$$

### 5.3.2 Properties

Now we evaluate the limiting distribution of the proposed test statistic (5.7). From *Shen et al.* (2014), we have the consistency of the parameter estimators in the EM process in the following result.

**Proposition V.2.** *Under the null model such that $\boldsymbol{\beta}_2 = 0$ and $\sigma_1 = \sigma_2 = \sigma_0$, and assume Conditions C0-C4, then for any finite $K \in \mathbb{Z}$, the estimator from the EM process $\boldsymbol{\theta}^{(K)} = (\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T, \sigma_1, \sigma_2)^T$ at the $K$th iteration is strongly consistent.*

By direct calculations, we can see for a given $\boldsymbol{\gamma}$, under the null of $\boldsymbol{\beta}_2 = 0, \sigma_1 = \sigma_2 = \sigma_0$, the "Fisher information matrix" from the penalized likelihood $I_{\boldsymbol{\gamma}}^*(\boldsymbol{\theta})$ for $\boldsymbol{\theta}$ is

$$\frac{1}{\sigma_0^2} \begin{pmatrix} I_1 & 0_{2q_1 \times 2} \\ 0_{2 \times 2q_1} & I_2 \end{pmatrix}. \tag{5.8}$$

where

$$I_1 = \begin{pmatrix} \mathbb{E}\boldsymbol{Z}\boldsymbol{Z}^T & \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T \\ \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T & \mathbb{E}\pi^2(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T \end{pmatrix} \tag{5.9}$$

and

$$I_2 = \begin{pmatrix} 2\mathbb{E}\pi^2(\boldsymbol{X}^T\boldsymbol{\gamma}) - \sigma_0^2 \frac{\mathbb{E}p_n''(\sigma_0)}{n} & 2\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma})) \\ 2\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma})) & 2\mathbb{E}(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma}))^2 - \sigma_0^2 \frac{\mathbb{E}p_n''(\sigma_0)}{n} \end{pmatrix}. \qquad (5.10)$$

If the variables $\boldsymbol{X}$ and $\boldsymbol{Z}$ are not degenerate, $\boldsymbol{\gamma}_X \neq 0$, and $\mathbb{E}p_n''(\sigma_0^2) < 0$, then the information matrix is positive definite.

In order to eliminate the asymptotic effect of the penalty from the final approximation, we further impose the following condition:

C5. Under the null hypothesis, $\mathbb{E}p_n''(\sigma) < 0$, $\mathbb{E}p_n''(\sigma) = o_p(n)$, and $p_n'(\sigma) = o_p(\sqrt{n})$.

Therefore, if C0-C5 hold, the penalized likelihood ratio test has the asymptotic chi-square distribution under the null hypothesis by a quadratic approximation. That is, for a fixed $\boldsymbol{\gamma}$, by the same derivations in Section 3.1 and Theorem V.2, we have the quadratic approximation of the penalized likelihood ratio statistic $T^*(\boldsymbol{\gamma})$:

$$T^*(\boldsymbol{\gamma}) = 2(pl(\hat{\boldsymbol{\theta}}_n, \boldsymbol{\gamma}) - pl(\hat{\boldsymbol{\theta}}_0, \boldsymbol{\gamma})) = ||h(\boldsymbol{\gamma})||^2 + o_p(1), \qquad (5.11)$$

where $\hat{\boldsymbol{\theta}}_n = \text{argmax}_{\boldsymbol{\theta}} \, pl(\boldsymbol{\theta}, \boldsymbol{\gamma})$, and $\hat{\boldsymbol{\theta}}_0 = \text{argmax}_{\boldsymbol{\theta} \in H_0} \, pl(\boldsymbol{\theta}, \boldsymbol{\gamma})$, and

$$h(\boldsymbol{\gamma}) = (h_1(\boldsymbol{\gamma}), h_2(\boldsymbol{\gamma})), \qquad (5.12)$$

in which

$$
\begin{aligned}
h_1(\boldsymbol{\gamma}) &= D(\boldsymbol{\gamma})^{-1/2}(N_2(\boldsymbol{\gamma}) - B(\boldsymbol{\gamma})A^{-1}N_1(\boldsymbol{\gamma})), \\
h_2(\boldsymbol{\gamma}) &= \left\{-2q_n(\sigma_0) + 2\mathbb{E}\left(1 - \pi(\boldsymbol{X}^T\boldsymbol{\gamma})\right)^2 - \frac{\left(2(1 - \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})) - q_n(\sigma_0)\right)^2}{2 - 2q_n(\sigma_0)}\right\}^{-1/2} \\
&\quad \left\{N_4(\boldsymbol{\gamma}) - \frac{-q_n(\sigma_0) + 2(1 - \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma}))}{2 - 2q_n(\sigma_0)}\left(N_3(\boldsymbol{\gamma}) + N_4(\boldsymbol{\gamma})\right)\right\}, \\
q_n(\sigma_0) &= \frac{\mathbb{E}p_n''(\sigma_0)}{n},
\end{aligned}
$$

$$A \quad = \quad \mathbb{E}\boldsymbol{Z}\boldsymbol{Z}^T,$$

$$B(\boldsymbol{\gamma}) \quad = \quad \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T,$$

$$C(\boldsymbol{\gamma}) \quad = \quad \mathbb{E}\pi^2(\boldsymbol{X}^T\boldsymbol{\gamma})\boldsymbol{Z}\boldsymbol{Z}^T,$$

$$D(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sigma_0^2}(C(\boldsymbol{\gamma}) - B(\boldsymbol{\gamma})A^{-1}B(\boldsymbol{\gamma})),$$

$$N_1(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sqrt{n}\sigma_0^2}\sum_{i=1}^{n}(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)\boldsymbol{Z}_i^T,$$

$$N_2(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sqrt{n}\sigma_0^2}\sum_{i=1}^{n}\pi(\boldsymbol{X}_i\boldsymbol{\gamma})(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)\boldsymbol{Z}_i^T,$$

$$N_3(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sqrt{n}\sigma_0^3}\sum_{i=1}^{n}\{\pi(\boldsymbol{X}_i\boldsymbol{\gamma})((Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)^2 - \sigma^2) + \sigma^3\tfrac{p_n'(\sigma_0)}{n}\},$$

$$N_4(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sqrt{n}\sigma_0^3}\sum_{i=1}^{n}\{(1 - \pi(\boldsymbol{X}_i\boldsymbol{\gamma}))((Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)^2 - \sigma_0^2) + \sigma_0^3\tfrac{p_n'(\sigma_0)}{n}\}.$$

In particular, for the penalty (5.2), under the null hypothesis we have $\mathbb{E}p_n''(\sigma_0) = -4\lambda/\sigma_0^2$, and $p_n'(\sigma_0) = 2\lambda(S_n^2 - \sigma_0^2)/\sigma_0^3$. By Condition C5 and by eliminating the penalty-related terms we have

$$T^*(\boldsymbol{\gamma}) = ||h^*(\boldsymbol{\gamma})||^2 + o_p(1), \tag{5.13}$$

$$h^*(\boldsymbol{\gamma}) = (h_1(\boldsymbol{\gamma}), h_2^*(\boldsymbol{\gamma})), \tag{5.14}$$

$$h_2^*(\boldsymbol{\gamma}) \quad = \quad \big(2(\mathbb{E}\pi^2(\boldsymbol{X}\boldsymbol{\gamma}) - (\mathbb{E}\pi(\boldsymbol{X}\boldsymbol{\gamma}))^2)\big)^{-1/2}$$
$$\big(N_4^*(\boldsymbol{\gamma}) - (1 - \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma}))(N_3^*(\boldsymbol{\gamma}) + N_4^*(\boldsymbol{\gamma}))\big),$$

$$N_3^*(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sqrt{n}\sigma_0^3}\sum_{i=1}^{n}\{\pi(\boldsymbol{X}_i\boldsymbol{\gamma})((Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)^2 - \sigma_0^2)\},$$

$$N_4^*(\boldsymbol{\gamma}) \quad = \quad \tfrac{1}{\sqrt{n}\sigma_0^3}\sum_{i=1}^{n}\{(1 - \pi(\boldsymbol{X}_i\boldsymbol{\gamma}))((Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)^2 - \sigma_0^2)\}.$$

We can write

$$T^*(\boldsymbol{\gamma}) = ||\frac{1}{\sqrt{n}}\psi^*(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})||^2 + o_p(1), \tag{5.15}$$

where $\psi^*(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) = (\psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})^T, \psi_0(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}))$, $\psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})$ is the same as defined for the equal variance model in Equation (3.10) of Section 3.1.1,

and the additional term

$$\psi_0(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma}) = \left(2(\mathbb{E}\pi^2(\boldsymbol{X}^T\boldsymbol{\gamma}) - (\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma}))^2)\right)^{-1/2}(\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma}) - \pi(\boldsymbol{X}^T\boldsymbol{\gamma}))$$
$$\left(\frac{(Y_i - \boldsymbol{Z}_i^T\boldsymbol{\beta}_1)^2}{\sigma_0^2} - 1\right).$$

(5.16)

Direct calculations show that both $\psi(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})$ and $\psi_0(Y_i, \boldsymbol{Z}_i, \boldsymbol{X}_i; \boldsymbol{\gamma})$ have mean zero, and the covariance matrix of $\psi^*$ is $I_{q_1+1}$. Therefore, $T^*(\boldsymbol{\gamma})$ has a $\chi^2$ limiting distribution with the degrees of freedom $q_1 + 1$. In addition, from the quadratic representation in Equation (5.15) and Proposition V.2, by the same proof of Theorem III.3, we will have (5.15) holds uniformly in $\boldsymbol{\gamma} \in \Gamma$ where $\Gamma$ is defined in (3.11), and we have the following proposition.

**Proposition V.3.** *Under the null hypothesis and assumptions C0-C5, for any finite integers $J > 0$ and $K \geq 0$, the penalized EM test statistic $pEM^{(K)}$ of the unequal variance model from Equation (5.7) converges to a fixed distribution as $n \to \infty$.*

### 5.3.3  Local Power

We calculate the local power of the $pEM$ test for the heterogeneous cases. Consider the parameters and the local alternative as $\boldsymbol{\eta}_0 = (\boldsymbol{\beta}_0, \boldsymbol{0}, \sigma_0, \sigma_0, \boldsymbol{\gamma}_0)^T$ and $\boldsymbol{\eta}_a^* = (\boldsymbol{\beta}_0, n^{-1/2}\boldsymbol{h}^T, \sigma_0 + n^{-1/2}h_1, \sigma_0, \boldsymbol{\gamma}_0)^T$, respectively. That is,

$$H_0: \quad \boldsymbol{\beta}_2 = \boldsymbol{0}, \sigma_1 = \sigma_2 = \sigma_0, \ v.s.$$
$$H_a^*: \quad \boldsymbol{\beta}_2 = n^{-1/2}\boldsymbol{h}, \sigma_1 = \sigma_2 + n^{-1/2}h_1 = \sigma_0 + n^{-1/2}h_1,$$

(5.17)

where $h \in \mathbb{R}^{q_1}$ and $h_1$ is a constant. By the same proof for Theorem V.4, we obtain the following result:

**Proposition V.4.** *Under $H_a^*$ and assumptions C0-C5, the test statistic $pEM^{(K)}$, with any value $\boldsymbol{\gamma} \in \tilde{\Gamma}$ and for any positive integer $K$, converges to a noncentral chi-square*

*distribution with the degree of freedom $q_1 + 1$ and the noncentrality parameter*

$$\lambda^*(\boldsymbol{\gamma}) = \lambda(\boldsymbol{\gamma}) + \lambda_1(\boldsymbol{\gamma}), \tag{5.18}$$

*where $\lambda(\boldsymbol{\gamma})$ is the same as defined in Equation (3.22), and*

$$\begin{aligned}
\lambda_1(\boldsymbol{\gamma}) &= ||\sigma_0^{-1}2^{1/2}h_1\big(\mathbb{E}\pi^2(\boldsymbol{X}^T\boldsymbol{\gamma}) - (\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma}))^2\big)^{-1/2} \\
&\quad \big(\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0) - \mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma})\mathbb{E}\pi(\boldsymbol{X}^T\boldsymbol{\gamma}_0))||^2.
\end{aligned} \tag{5.19}$$

In Section 3.2.4, we have seen that $\lambda(\boldsymbol{\gamma})$ is maximized at $\boldsymbol{\gamma}_0$. Direct calculations show that $\lambda(\boldsymbol{\gamma}_0)$ is also maximized at $\boldsymbol{\gamma}_0$. Therefore, $\lambda^*(\boldsymbol{\gamma})$ achieves its maximum at the true one $\boldsymbol{\gamma}_0$. The power for the penalized $EM$ test is then $P(\chi^2_{q_1+1;\lambda^*(\boldsymbol{\gamma})} > \chi^2_{q_1+1}(1-\alpha))$. In particular, if $h_1 = 0$, then the noncentral parameters are the same as those for the EM tests developed under the equal variance model, but have different degrees of freedom, $q_1+1$, not $q_1$. We shall compare the two $EM$ tests in the following simulation studies.

## 5.4   Simulations

In this section, we study the finite sample performances of the proposed methods through simulation studies. We will first show the effect of $\lambda$ on the parameter estimations under an alternative model. Second, we evaluate the type I errors and the powers of the penalized $EM$ test. We will compare the powers from the $pEM$ test and the $EM$ test of Section 3.2, for data generated from both the equal and unequal variance models.

### 5.4.1   Estimations

We start with evaluating the parameter estimations under the alternative models. We do a simulation study using similar settings to those of Section 2.6 to show the

performance of the estimators from the penalized likelihood under different penalty parameters and different signal ratios. Data are generated from

$$Y_i = \mu_1 + \nu_1 T_i + \alpha_1 Z_i + (\mu_2 + \nu_2 T_i + \alpha_2 Z_i)\delta_i + \varepsilon_i(\sigma_1 \delta_i + \sigma_2(1 - \delta_i)),$$

$$P(\delta_i = 1|X_i) = \pi(\gamma_0 + \gamma_1 X_i),$$

for $i = 1, \ldots, n$, where $\varepsilon_i \sim N(0, 1)$, independent of $(X_i, Z_i)$. The $n$ observations are independent. We use $X_i = Z_i$ from $N(1, 1)$, and $n = 100$. We collect the means and sample standard deviations of the maximum (penalized) likelihood estimates in 1000 repeated experiments. In the simulation, we use the constriction of $\nu_2 > 0$ to guarantee the identifiability of the parameters. We examine the estimates from both the unequal variance model in Table 5.1 and Table 5.2 and the equal variance model in Table 5.3. Larger $\lambda$ in the penalty term gives slightly larger bias. Overall the biases do not change much with varying $\lambda$, so we will fix $\lambda = 1$ in the power calculations later unless otherwise specified. The estimates from the equal variance model for all the parameters except the $\sigma$'s are quite close to the true ones in Table 5.3 when $\sigma_1/\sigma_2$ is close to 1, but the bias gets larger as the ratio of two $\sigma$'s increases to 2 or 3.

Table 5.1: The biases and the sample standard deviations of the estimates under the unequal variance model with different choices of $\lambda$ .

| parameters | | bias | sd | bias | sd | bias | sd |
|---|---|---|---|---|---|---|---|
| | | $\lambda = 0.05$ | | $\lambda = 1$ | | $\lambda = 50$ | |
| $\mu_1$ | 2.0 | 0.010 | 0.130 | 0.010 | 0.130 | 0.018 | 0.132 |
| $\nu_1$ | 0.0 | 0.001 | 0.121 | 0.001 | 0.121 | 0.003 | 0.121 |
| $\alpha_1$ | 2.0 | 0.004 | 0.070 | 0.004 | 0.070 | 0.008 | 0.071 |
| $\mu_2$ | 3.0 | 0.024 | 0.200 | 0.024 | 0.200 | 0.029 | 0.201 |
| $\nu_2$ | 3.0 | 0.012 | 0.216 | 0.012 | 0.216 | 0.013 | 0.216 |
| $\alpha_2$ | 5.0 | 0.006 | 0.128 | 0.006 | 0.128 | 0.008 | 0.128 |
| $\gamma_0$ | 1.0 | 0.057 | 0.384 | 0.056 | 0.384 | 0.040 | 0.382 |
| $\gamma_1$ | -1.0 | 0.051 | 0.296 | 0.050 | 0.296 | 0.039 | 0.294 |
| $\sigma_1$ | 0.4 | 0.015 | 0.041 | 0.010 | 0.040 | 0.061 | 0.036 |
| $\sigma_2$ | 0.6 | 0.019 | 0.062 | 0.022 | 0.061 | 0.077 | 0.047 |

Table 5.2: The biases and the sample standard deviations of the estimates under the unequal variance model with different $\sigma$ ratios.

| $(\sigma_1, \sigma_2)$ | | $(0.4, 0.6)$ | | $(0.5, 1.0)$ | | $(0.5, 1.5)$ | |
|---|---|---|---|---|---|---|---|
| parameters | | bias | sd | bias | sd | bias | sd |
| $\mu_1$ | 2.0 | 0.010 | 0.130 | 0.025 | 0.171 | 0.034 | 0.181 |
| $\nu_1$ | 0.0 | 0.001 | 0.121 | 0.004 | 0.153 | 0.006 | 0.153 |
| $\alpha_1$ | 2.0 | 0.004 | 0.070 | 0.011 | 0.091 | 0.015 | 0.095 |
| $\mu_2$ | 3.0 | 0.024 | 0.200 | 0.071 | 0.305 | 0.150 | 0.422 |
| $\nu_2$ | 3.0 | 0.012 | 0.216 | 0.031 | 0.338 | 0.070 | 0.487 |
| $\alpha_2$ | 5.0 | 0.006 | 0.128 | 0.023 | 0.202 | 0.046 | 0.285 |
| $\gamma_0$ | 1.0 | 0.057 | 0.384 | 0.060 | 0.407 | 0.059 | 0.431 |
| $\gamma_1$ | -1.0 | 0.051 | 0.296 | 0.053 | 0.311 | 0.053 | 0.326 |
| $\sigma_1$ | | 0.015 | 0.041 | 0.096 | 0.049 | 0.019 | 0.049 |
| $\sigma_2$ | | 0.019 | 0.062 | 0.361 | 0.103 | 0.061 | 0.158 |

Table 5.3: The biases and the sample standard deviations of the estimates under the equal variance model.

| $(\sigma_1, \sigma_2)$ | | $(0.4, 0.6)$ | | $(0.5, 1.0)$ | | $(0.5, 1.5)$ | |
|---|---|---|---|---|---|---|---|
| parameters | | bias | sd | bias | sd | bias | sd |
| $\mu_1$ | 2.0 | 0.007 | 0.129 | 0.032 | 0.179 | 0.103 | 0.229 |
| $\nu_1$ | 0.0 | 0.001 | 0.121 | 0.006 | 0.153 | 0.017 | 0.169 |
| $\alpha_1$ | 2.0 | 0.002 | 0.069 | 0.013 | 0.094 | 0.044 | 0.118 |
| $\mu_2$ | 3.0 | 0.007 | 0.200 | 0.010 | 0.315 | 0.003 | 0.470 |
| $\nu_2$ | 3.0 | 0.006 | 0.216 | 0.004 | 0.338 | 0.007 | 0.488 |
| $\alpha_2$ | 5.0 | 0.003 | 0.128 | 0.005 | 0.207 | 0.007 | 0.317 |
| $\gamma_0$ | 1.0 | 0.003 | 0.366 | 0.050 | 0.391 | 0.151 | 0.419 |
| $\gamma_1$ | -1.0 | 0.008 | 0.282 | 0.022 | 0.295 | 0.082 | 0.305 |
| $\sigma_1$ | | 0.093 | 0.041 | 0.258 | 0.071 | 0.554 | 0.115 |
| $\sigma_2$ | | 0.107 | 0.041 | 0.158 | 0.071 | 0.495 | 0.115 |

## 5.4.2 Type I Errors

We evaluate the proposed EM test by examining the accuracy of the type I errors. We use the same setting as that for Table 3.2 in Section 3.3.1. The resulting type I errors are summarized in Table 5.4 and 5.5 for $\lambda = 1$ and $\lambda = 50$, respectively. We see the type I errors are quite close to the nominal levels for $K = 0, 3$, and 9.

Table 5.4: Type I errors of the $EM$ tests with bootstrap approximations based on 1000 data sets, unequal variance case, with $\lambda = 1$.

| $n$ | Nominal level $\alpha$ | $pEM^{(0)}$ | $pEM^{(3)}$ | $pEM^{(9)}$ |
|---|---|---|---|---|
| $n$=60 | 0.01 | 0.013 | 0.010 | 0.010 |
| | 0.05 | 0.049 | 0.045 | 0.047 |
| | 0.10 | 0.086 | 0.099 | 0.094 |
| $n$=100 | 0.01 | 0.012 | 0.012 | 0.012 |
| | 0.05 | 0.053 | 0.054 | 0.053 |
| | 0.10 | 0.106 | 0.107 | 0.110 |

Table 5.5: Type I errors of the $EM$ tests with bootstrap approximations based on 1000 data sets, unequal variance case, with $\lambda = 50$.

| $n$ | Nominal level $\alpha$ | $pEM^{(0)}$ | $pEM^{(3)}$ | $pEM^{(9)}$ |
|---|---|---|---|---|
| $n$=60 | 0.01 | 0.013 | 0.013 | 0.014 |
| | 0.05 | 0.044 | 0.050 | 0.051 |
| | 0.10 | 0.089 | 0.088 | 0.094 |
| $n$=100 | 0.01 | 0.010 | 0.010 | 0.008 |
| | 0.05 | 0.049 | 0.049 | 0.048 |
| | 0.10 | 0.103 | 0.116 | 0.113 |

## 5.4.3 Power Comparison

Power is calculated using the same setting as in the equal variance case in Section 3.3.2. The power is obtained from the EM test from both the true unequal variance model and also the equal variance model. We fix $\lambda = 1$. When we have equal or close variances from Table 5.6 and Table 5.7, the penalized $EM$ test give comparable power compared to the $EM$ test for all the settings. As we increase the ratio of $\sigma_2/\sigma_1$

to 2 and 3, as in Table 5.8 where $\sigma_1 = 0.5$ and $\sigma_2 = 1.0$, and in Table 5.9 where $\sigma_1 = 0.5$ and $\sigma_2 = 1.5$, the penalized test is significantly more powerful.

Table 5.6: Power (%) of the $EM$ and $pEM$ tests at the 5% level. Both test uses $\mathbf{\Gamma} = \{(1,2)^T, (1,-2)^T\}$. The parameters of Model (1) are $\boldsymbol{\beta}_1 = (1,0,2)^T$, $\boldsymbol{\beta}_2 = (1,a,b)^T$, $\boldsymbol{\gamma} = (1,c)^T$, $\sigma_1 = 0.5$ and $\sigma_2 = 0.5$.

| $n$ | $a$ | $b$ | $c$ | $pEM^{(0)}$ | $pEM^{(3)}$ | $pEM^{(9)}$ |
|-----|-----|-----|-----|-------------|-------------|-------------|
| 60  | 0.5 | 1   | 1   | 75.2        | 72.2        | 74.0        |
| 60  | 0.5 | 0   | 1   | 33.0        | 34.6        | 36.8        |
| 60  | 1.0 | 1   | 1   | 89.2        | 88.2        | 88.0        |
| 60  | 1.0 | 0   | 1   | 79.0        | 83.4        | 84.6        |
| 100 | 0.5 | 1   | 1   | 94.2        | 92.0        | 93.8        |
| 100 | 0.5 | 0   | 1   | 49.0        | 57.4        | 57.6        |
| 100 | 1.0 | 1   | 1   | 99.0        | 99.2        | 99.0        |
| 100 | 1.0 | 0   | 1   | 95.0        | 97.4        | 98.0        |
| $n$ | $a$ | $b$ | $c$ | $EM^{(0)}$  | $EM^{(3)}$  | $EM^{(9)}$  |
| 60  | 0.5 | 1   | 1   | 74.0        | 76.2        | 77.8        |
| 60  | 0.5 | 0   | 1   | 25.6        | 36.8        | 36.0        |
| 60  | 1.0 | 1   | 1   | 86.2        | 86.8        | 87.8        |
| 60  | 1.0 | 0   | 1   | 69.0        | 80.8        | 84.8        |
| 100 | 0.5 | 1   | 1   | 96.2        | 95.8        | 96.8        |
| 100 | 0.5 | 0   | 1   | 35.4        | 49.2        | 54.8        |
| 100 | 1.0 | 1   | 1   | 99.0        | 99.2        | 99.4        |
| 100 | 1.0 | 0   | 1   | 86.2        | 95.6        | 97.6        |

## 5.5 Discussion

For the choice of $J$, $K$, and the initial $\boldsymbol{\gamma}$'s, we suggest the same principles as in the equal variance case of Section 3.2.3. For the penalty in Equation (5.2), we can use any fixed positive $S_n^2$ and $\lambda$ in theory, and since $p_n(\cdot)$ is maximized as $\sigma^2 = S_n^2$, we prefer $S_n^2$ to be the estimator under the equal variance model from the $EM$ algorithm without iterating $\boldsymbol{\gamma}$ as a reasonable estimator of the variance term. We do not perform a full $EM$ algorithm due to the identifiability issue of the parameters under the null model. We can also use the variance estimator from the null model, but when the data is generated under the alternative model, the variance estimator under the null

Table 5.7: Power (%) of the $EM$ and $pEM$ tests at the 5% level. Both test uses $\mathbf{\Gamma} = \{(1,2)^T, (1,-2)^T\}$. The parameters of Model (1) are $\boldsymbol{\beta}_1 = (1,0,2)^T$, $\boldsymbol{\beta}_2 = (1,a,b)^T$, $\boldsymbol{\gamma} = (1,c)^T$, $\sigma_1 = 0.4$ and $\sigma_2 = 0.6$.

| $n$ | $a$ | $b$ | $c$ | $pEM^{(0)}$ | $pEM^{(3)}$ | $pEM^{(9)}$ |
|-----|-----|-----|-----|-------------|-------------|-------------|
| 60  | 0.5 | 1   | 1   | 75.8        | 72.8        | 75.2        |
| 60  | 0.5 | 0   | 1   | 43.2        | 44.6        | 43.4        |
| 60  | 1.0 | 1   | 1   | 92.4        | 88.2        | 89.8        |
| 60  | 1.0 | 0   | 1   | 84.6        | 87.6        | 88.4        |
| 100 | 0.5 | 1   | 1   | 95.6        | 92.0        | 94.2        |
| 100 | 0.5 | 0   | 1   | 70.2        | 72.8        | 74.0        |
| 100 | 1.0 | 1   | 1   | 99.4        | 98.2        | 98.4        |
| 100 | 1.0 | 0   | 1   | 98.0        | 98.4        | 99.4        |
| $n$ | $a$ | $b$ | $c$ | $EM^{(0)}$  | $EM^{(3)}$  | $EM^{(9)}$  |
| 60  | 0.5 | 1   | 1   | 69.0        | 70.6        | 73.6        |
| 60  | 0.5 | 0   | 1   | 18.8        | 31.8        | 37.6        |
| 60  | 1.0 | 1   | 1   | 81.0        | 83.8        | 87.8        |
| 60  | 1.0 | 0   | 1   | 57.4        | 77.8        | 82.0        |
| 100 | 0.5 | 1   | 1   | 92.4        | 94.4        | 94.8        |
| 100 | 0.5 | 0   | 1   | 23.2        | 36.8        | 49.6        |
| 100 | 1.0 | 1   | 1   | 98.2        | 98.4        | 98.8        |
| 100 | 1.0 | 0   | 1   | 77.0        | 93.8        | 96.0        |

Table 5.8: Power (%) of the $EM$ and $pEM$ tests at the 5% level. Both test uses $\mathbf{\Gamma} = \{(1,2)^T, (1,-2)^T\}$. The parameters of Model (1) are $\boldsymbol{\beta}_1 = (1,0,2)^T$, $\boldsymbol{\beta}_2 = (1,a,b)^T$, $\boldsymbol{\gamma} = (1,c)^T$, $\sigma_1 = 0.5$ and $\sigma_2 = 1.0$.

| $n$ | $a$ | $b$ | $c$ | $pEM^{(0)}$ | $pEM^{(3)}$ | $pEM^{(9)}$ |
|---|---|---|---|---|---|---|
| 60 | 0.5 | 1 | 1 | 51.8 | 49.4 | 50.4 |
| 60 | 0.5 | 0 | 1 | 33.2 | 31.6 | 38.4 |
| 60 | 1.0 | 1 | 1 | 68.8 | 63.6 | 65.4 |
| 60 | 1.0 | 0 | 1 | 62.4 | 63.4 | 65.2 |
| 100 | 0.5 | 1 | 1 | 81.4 | 78.2 | 80.8 |
| 100 | 0.5 | 0 | 1 | 65.0 | 66.2 | 68.0 |
| 100 | 1.0 | 1 | 1 | 92.0 | 89.6 | 91.4 |
| 100 | 1.0 | 0 | 1 | 88.4 | 90.0 | 90.4 |
| $n$ | $a$ | $b$ | $c$ | $EM^{(0)}$ | $EM^{(3)}$ | $EM^{(9)}$ |
| 60 | 0.5 | 1 | 1 | 31.8 | 35.4 | 38.4 |
| 60 | 0.5 | 0 | 1 | 12.2 | 20.6 | 27.2 |
| 60 | 1.0 | 1 | 1 | 43.2 | 44.2 | 47.2 |
| 60 | 1.0 | 0 | 1 | 19.8 | 34.2 | 44.8 |
| 100 | 0.5 | 1 | 1 | 58.2 | 57.8 | 60.6 |
| 100 | 0.5 | 0 | 1 | 11.8 | 23.0 | 34.2 |
| 100 | 1.0 | 1 | 1 | 70.4 | 73.0 | 75.4 |
| 100 | 1.0 | 0 | 1 | 26.0 | 43.2 | 58.6 |

Table 5.9: Power (%) of the $EM$ and $pEM$ tests at the 5% level. Both test uses $\mathbf{\Gamma} = \{(1,2)^T, (1,-2)^T\}$. The parameters of Model (1) are $\boldsymbol{\beta}_1 = (1,0,2)^T$, $\boldsymbol{\beta}_2 = (1,a,b)^T$, $\boldsymbol{\gamma} = (1,c)^T$, $\sigma_1 = 0.5$ and $\sigma_2 = 1.5$.

| $n$ | $a$ | $b$ | $c$ | $pEM^{(0)}$ | $pEM^{(3)}$ | $pEM^{(9)}$ |
|-----|-----|-----|-----|-------------|-------------|-------------|
| 60  | 0.5 | 1   | 1   | 55.2 | 54.2 | 58.0 |
| 60  | 0.5 | 0   | 1   | 54.2 | 54.2 | 56.8 |
| 60  | 1.0 | 1   | 1   | 68.6 | 66.8 | 65.6 |
| 60  | 1.0 | 0   | 1   | 70.8 | 71.6 | 74.2 |
| 100 | 0.5 | 1   | 1   | 83.0 | 82.4 | 83.2 |
| 100 | 0.5 | 0   | 1   | 82.4 | 83.8 | 85.2 |
| 100 | 1.0 | 1   | 1   | 91.8 | 89.0 | 90.8 |
| 100 | 1.0 | 0   | 1   | 92.2 | 92.0 | 92.8 |

| $n$ | $a$ | $b$ | $c$ | $EM^{(0)}$ | $EM^{(3)}$ | $EM^{(9)}$ |
|-----|-----|-----|-----|------------|------------|------------|
| 60  | 0.5 | 1   | 1   | 20.4 | 26.4 | 31.0 |
| 60  | 0.5 | 0   | 1   | 19.2 | 33.2 | 40.8 |
| 60  | 1.0 | 1   | 1   | 28.2 | 33.4 | 39.6 |
| 60  | 1.0 | 0   | 1   | 23.4 | 38.0 | 47.8 |
| 100 | 0.5 | 1   | 1   | 32.0 | 36.4 | 42.0 |
| 100 | 0.5 | 0   | 1   | 18.8 | 38.6 | 51.2 |
| 100 | 1.0 | 1   | 1   | 40.8 | 47.4 | 51.8 |
| 100 | 1.0 | 0   | 1   | 21.8 | 44.2 | 58.6 |

model often tends to be larger than that estimated under the equal variance model, and therefore the test might be less powerful.

Under the unequal variance model, we are testing the null of $\boldsymbol{\beta}_2 = 0$ and $\sigma_1 = \sigma_2$ simultaneously. A standard second step is necessary to include the confidence interval of $\boldsymbol{\beta}_2$ for understanding if the null is rejected due to the differences only in the variances.

Between the $EM$ test developed under the equal variance model, and the $pEM$ test developed under the unequal variance model, we have seen that $pEM$ is more powerful than the $EM$ test in general when the variance ratios are away from 1, and does not lose much under equal variance model. In reality, the differences in the mean and variances between the two subgroups are often present simultaneously, therefore, if heterogeneity is suspected, the penalized $EM$ test is recommended.

## 5.6   AIDS Data

In this section, we revisit the ACTG 320 study described in Section 1.1 and analyzed by the equal variance model in Chapter IV. Using the same starting values of the $EM$ test in Chapter IV, that is, $\boldsymbol{\Gamma} = \{(2.44, -3.35, -2.33, 1.24)^T, (0.95, -4.49, 0.47, 4.46)^T,$ $(1.00, -2.72, 4.64, -2.84)^T\}$, the $p$-values of $pEM$ is less than 0.001 for $K = 3$ and 9 at $\lambda = 1$ using bootstrap sample size 5000. The $p$-values remain less than 0.001 for $\lambda = 200, 400$ and 800. So the null hypothesis of no subgroups is rejected.

For parameter estimation, we consider a range of $\lambda$'s and select one adaptively as follows. For each $\lambda$, define the score function of falling into the subgroup of better treatment effect to be $S_3(\boldsymbol{X}; \lambda) = \pi(\boldsymbol{X}^T\boldsymbol{\gamma})$, where $\boldsymbol{\gamma}$ is estimated using the penalized likelihood at $\lambda$. Let

$$TT(\lambda; q) = \mathbb{E}\{Y|(trt = 1, S_3(\boldsymbol{X}; \lambda) > Q_q(S_3))\} - \mathbb{E}\{Y|(trt = 0, S_3(\boldsymbol{X}; \lambda) > Q_q(S_3))\}$$
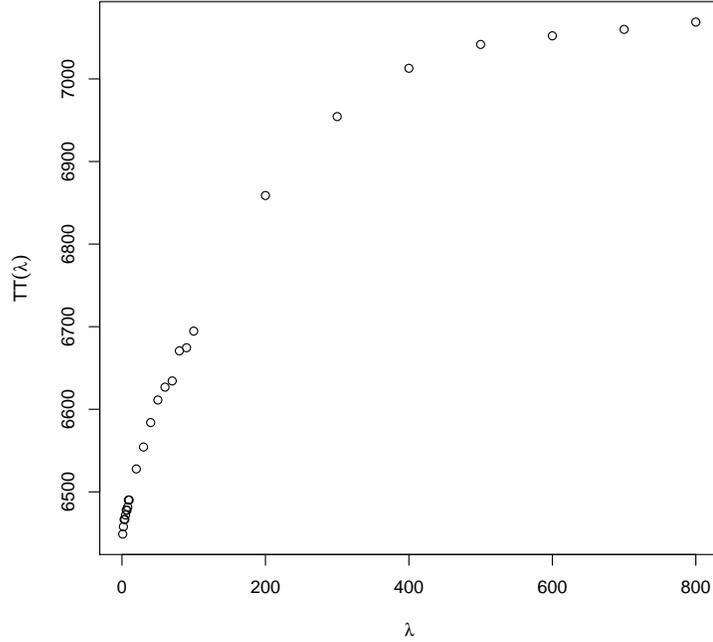
$$(5.20)$$

Figure 5.1: The overall treatment effect in (5.21) with different $\lambda$'s.

be the treatment effect in the selected subgroup where $S_3(\boldsymbol{X}; \lambda) > Q_q(S_3)$, and $Q_q(S_3)$ is the qth quantile of the scores $S_3(\boldsymbol{X}; \lambda)$. We use an overall treatment effect

$$TT(\lambda) = \int\limits_{0.1}^{0.9} TT(\lambda; q)dq. \tag{5.21}$$

The expectation and integral for (5.20) and (5.21) are estimated by sample means and sum over $q = 0.10, 0.11, \cdots 0.90$, respectively. We show $TT(\lambda)$ for $\lambda \in [1, 800]$ in Figure 5.1. We see that $TT(\lambda)$ increases as $\lambda$ increases, but after $\lambda$ reaches 400, $TT(\lambda)$ is stable. For a given $q = 0.8$, the treatment effect in the targeted subgroup is shown in Figure 5.2, from which the subgroup with $\lambda = 400$ gives the highest treatment effect. Given the information, we choose $\lambda = 400$ and use it in the penalized maximum likelihood method. The parameter estimates are given in Table 5.1.

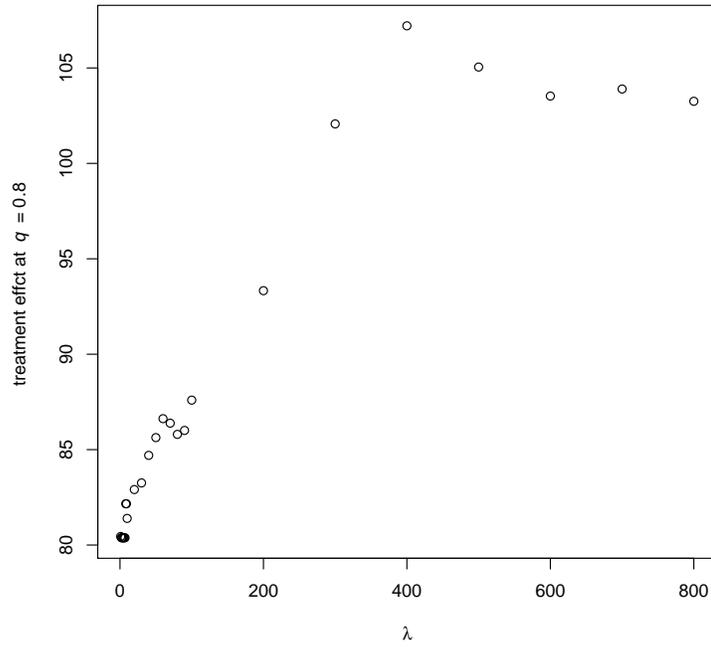The treatment effects in the two subgroups are 41.74 and 93.50, respectively. The

Figure 5.2: The treatment effect in (5.20) in selected subgroup with $q = 0.8$ and different $\lambda$'s.

Table 5.10: Parameter estimates and their standard errors when Model (5.1) is used to fit the data in the ACTG study with $\lambda = 400$.

| $\boldsymbol{\beta}_1$ | 1 | $trt$ | $\log(cd4.0)$ | $\log_{10}(rna.0)$ | Age |
|---|---|---|---|---|---|
| est | -45.97 | 41.74 | -0.68 | 7.40 | 0.72 |
| se | 44.48 | 6.34 | 3.60 | 6.73 | 0.38 |
| $\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2$ | 1 | $trt$ | $\log(cd4.0)$ | $\log_{10}(rna.0)$ | Age |
| est | -42.60 | 93.50 | 8.05 | 6.14 | 0.01 |
| se | 49.38 | 7.79 | 4.06 | 7.44 | 0.45 |
| $\boldsymbol{\gamma}$ | 1 | | $\log(cd4.0)$ | $\log_{10}(rna.0)$ | Age |
| Parameter | -9.18 | | 0.68 | 1.41 | -0.02 |
| se | 1.03 | | 0.08 | 0.16 | 0.01 |
| | $\sigma_1$ | $\sigma_2$ | | | |
| est | 57.65 | 48.28 | | | |
| se | 0.97 | 1.24 | | | |

difference of the treatment effects is smaller than that under the equal variance model. The mean probability of falling into the subgroup of higher treatment is around 0.42, and the estimated $\pi(\boldsymbol{X\gamma})$ values are more spread out on both side of 0.5, than what we obtained under the equal variance model. The ratio of the two $\sigma$'s is around 1.2. The BIC values of equal variance model and unequal variance models are 8902 and 8985, respectively, based on which the equal variance model is preferred.

If we use the unequal variance model, from the coefficients in Table 5.10, we obtain the scores describing the probability of getting a higher treatment effect to be

$$S_3(\boldsymbol{X}) = \pi(-9.18 + 0.68 \log(cd4.0) + 1.41 \log_{10}(rna.0) - 0.02\text{Age}). \quad (5.22)$$

We have the scores $S_1(X)$ in (4.1) obtained from the equal variance model . The scatter plot of the two scores from equal and unequal variance models are given in Figure 5.3. They have linear and rank correlations around 0.98.

In this example, the equal variance model and the unequal variance model lead to very similar subgroup scores and the two models might not be highly distinguishable. The later model however allows a more relaxed condition on the variances of the subgroups, and resulted in a more interpretable $\pi$ values for subgroup assignments.
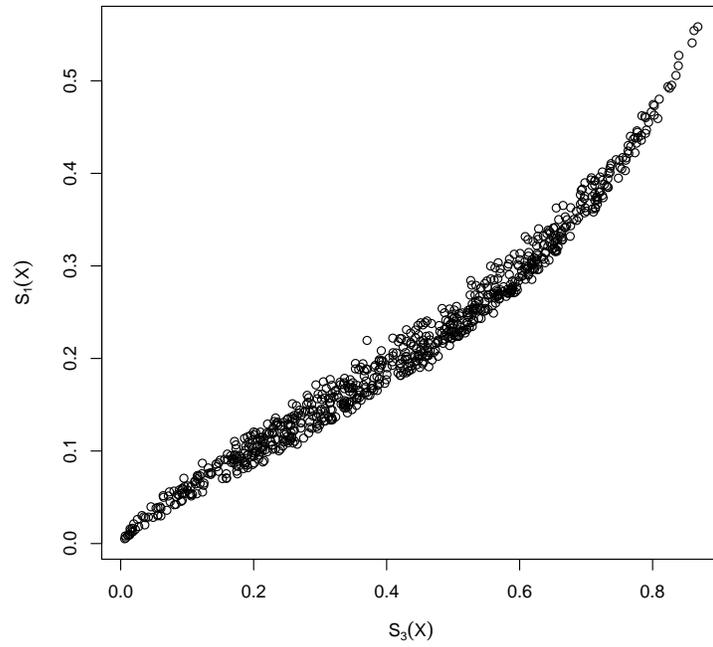
Figure 5.3: Scatter plot of $S_1(X)$ from (4.1) versus$S_3(X)$ from (5.22).

# CHAPTER VI

# Summary

We propose a model-based framework for the dual purposes of a confirmatory subgroup analysis and a predictive modeling of subgroup membership. In contrast to the existing work on subgroup identification in clinical and medical statistics, our disciplined approach aims to reduce false positive subgroup identification. In the meantime, our model can generate a scoring system that can be used to predict subgroup membership, as demonstrated in the analysis of an AIDS study. We propose a (penalized) $EM$ test based on a small number of $EM$ iterations towards the likelihood of the logistic-normal mixture model and obtain the asymptotic representation of the test statistic. The proposed test avoids some of the challenges and complications, both computational and theoretical, associated with the likelihood ratio tests for mixture models. Through simulation studies and a real data example, we demonstrate that the proposed methodology is a valuable addition to subgroup analysis.

Like any model-based inference, the proposed $EM$ test needs to be understood in conjunction with an appropriate sensitivity analysis against model mis-specifications. For the structured logistic-normal model, our empirical work shows that the test is quite robust against moderate deviations from the logistic component of the model, but the normal component of the model is rather critical. We hope that future research will address a broader set of questions on robustness in subgroup analysis.

Finally, we discuss briefly the case of non-normal components. Consider

$$
\begin{aligned}
P(Y_i|(\boldsymbol{Z}_i, \delta_i)) &= \phi(Y_i|\boldsymbol{Z}_i; \boldsymbol{\beta}_1 + \boldsymbol{\beta}_2\delta_i, \sigma), \\
P(\delta_i = 1|\boldsymbol{X}_i, \boldsymbol{Z}_i) &= \pi(\boldsymbol{X}_i^T\boldsymbol{\gamma}) \equiv \frac{\exp(\boldsymbol{X}_i^T\boldsymbol{\gamma})}{1+\exp(\boldsymbol{X}_i^T\boldsymbol{\gamma})}, \\
P(\delta_i = 0|\boldsymbol{X}_i, \boldsymbol{Z}_i) &= 1 - P(\delta_i = 1|\boldsymbol{X}_i),
\end{aligned}
\tag{6.1}
$$

where the notations are similar to those in Equation (2.1), $\sigma$ is a common effect parameter, and $\phi(\cdot)$ is a probability density function.

We can derive similar EM algorithm for estimation of the parameters. In *Wang* (1994); *Wang et al.* (1996), and *Wang and Puterman* (1998), mixed Poisson models and mixed binomial models with mixing proportions depending on the covariates through a logistic link are discussed for a fixed number of groups. If $\phi$ is a Poisson density, any finite group models have identifiable parameters. Note that for the non-normal cases, with some conditions, without the undesired degenerate Fisher information matrix, the constraint that the slope of $\boldsymbol{\gamma}$ has to be bounded away from zero along the EM process can be relaxed. In the work of *Zhu and Zhang* (2004, 2006), the likelihood ratio test of this problem is developed, but involves intensive computation for p-values. The identifiability of the parameters for binary response needs further research.

# BIBLIOGRAPHY

# BIBLIOGRAPHY

Altstein, L. L., G. Li, and R. M. Elashoff (2011), A Method to Estimate Treatment Efficacy among Latent Subgroups of a Randomized Clinical Trial, *Statistics in Medicine*, *30*(7), 709–717.

Bonetti, M., and R. D. Gelber (2004), Patterns of Treatment Effects in Subsets of Patients in Clinical Trials, *Biostatistics*, *5*(3), 465–481.

Cai, T., L. Tian, P. H. Wong, and L. Wei (2011), Analysis of Randomized Comparative Clinical Trial Data for Personalized Treatment Selections, *Biostatistics*, *12*(2), 270–282.

Chen, H., and J. Chen (2003), Tests for Homogeneity in Normal Mixtures in the Presence of a Structural Parameter, *Statistica Sinica*, *13*, 351–365.

Chen, H., J. Chen, and J. D. Kalbfleisch (2001), A Modified Likelihood Ratio Test for Homogeneity in Finite Mixture Models, *Journal of the Royal Statistical Society, Series B: Statistical Methodology*, *63*(1), 19–29.

Chen, J. (1995), Optimal Rate of Convergence for Finite Mixture Models, *The Annals of Statistics*, *23*(1), 221–233.

Chen, J., and P. Li (2009), Hypothesis Test for Normal Mixture Models: the EM Approach, *The Annals of Statistics*, *37*(5A), 2523–2542.

Dempster, A., N. Laird, and D. Rubin (1977), Maximum Likelihood from Incomplete Data via the EM Algorithm, *Journal of the Royal Statistical Society, Series B: Statistical Methodology*, *39*(1), 1–38.

Dixon, D. O., and R. Simon (1991), Bayesian Subset Analysis, *Biometrics*, *47*(3), 871–881.

Foster, J. C., J. M. Taylor, and S. J. Ruberg (2011), Subgroup Identification from Randomized Clinical Trial Data, *Statistics in Medicine*, *30*(24), 2867–2880.

Frasure-Smith, N., F. Lesperance, R. H. Prince, et al. (1997), Randomised Trial of Home-based Psychosocial Nursing Intervention for Patients Recovering from Myocardial Infarction, *The Lancet*, *350*(9076), 473–479.

Fritsch, J., M. Finke, and A. Waibel (1997), Adaptively growing hierarchical mixtures of experts, in *In Advances in Neural Information Processing Systems 9*, pp. 459–465, MIT Press.

Goeffinet, B., P. Loisel, and B. Laurent (1992), Testing in Normal Mixture Models when the Proportions are Known, *Biometrika*, *79*(4), 842–846.

Hammer, S., K. Squires, M. Hughes, J. Grimes, et al. (1997), A controlled trial of two nucleoside analogues plus indinavir in persons with human immunodeficiency virus infection and cd4 cell counts of 200 per cubic millimeter or less, *The New England Journal of Medicine*, *337*(11), 725–733.

Jennrich, R. I. (1969), Asymptotic Properties of Non-Linear Least Squares Estimator, *The Annals of Mathematical Statistics*, *40*(2), 633–643.

Jiang, W., and M. A. Tanner (1999a), Hierarchical Mixtures-of-Experts for Exponential Family Regression Models: Approximation and Maximum Likelihood Estimation, *The Annals of Statistics*, *27*(3), 987–1011.

Jiang, W., and M. A. Tanner (1999b), On the Identifiability of Mixtures-of-Experts, *Neural Networks*, *12*(9), 1253–1258.

Jiang, W., and M. A. Tanner (1999c), On the Approximation Rate of Hierarchical Mixtures-of-Experts for Generalized Linear Models, *Neural Computation*, *11*(5), 1183–1198.

Jordan, M. I., and R. A. Jacobs (1994), Hierarchical Mixtures of Experts and the EM Algorithm, *Neural Computation*, *6*(2), 181–214.

Li, P., and J. Chen (2010), Testing the Order of a Finite Mixture, *Journal of the American Statistical Association*, *105*(491), 1084–1092.

Lipkovich, I., A. Dmitrienko, J. Denne, and G. Enas (2011), Subgroup identification based on differential effect search–a recursive partitioning method for establishing response to treatment in patient subpopulations, *Statistics in Medicine*, *30*(21), 2601–2621.

Lipkovicha, I., and A. Dmitrienkoa (2014), Strategies for Identifying Predictive Biomarkers and Subgroups with Enhanced Treatment Effect in Clinical Trials Using SIDES, *Journal of Biopharmaceutical Statistics*, *24*(1), 130–153.

Lo, Y., N. R. Mendell, and D. B. Rubin (2001), Testing the Number of Components in a Normal Mixture, *Biometrika*, *88*(3), 767–778.

Louis, T. A. (1982), Finding the Observed Information Matrix when Using the EM Algorithm, *Journal of the Royal Statistical Society, Series B: Statistical Methodology*, *44*(2), 226–233.

Luo, R., C.-L. Tsai, and H. Wang (2008), On Mixture Regression Shrinkage and Selection via the Mr. Lasso, *International Journal of Pure and Applied Mathematics*, *46*, 403–414.

McLachlan, G., and D. Peel (2000), *Finite Mixture Models*, 1 ed., Wiley-Interscience.

Muthén, B., and T. Asparouhov (2009), Multilevel Regression Mixture Analysis, *Journal of the Royal Statistical Society, Series A: Statistics in Society*, *172*(3), 639–657.

Muthén, B., and K. Shedden (1999), Finite Mixture Modeling with Mixture Outcomes Using the EM Algorithm, *Biometrics*, *55*, 463–469.

Naik, P. A., P. Shi, and C.-L. Tsai (2007), Extending the akaike information criterion to mixture regression models, *Journal of the American Statistical Association*, *102*(477), 244–254.

Peng, F., R. A. Jacobs, and M. A. Tanner (1996), Bayesian inference in mixtures-of-experts and hierarchical mixtures-of-experts models with an application to speech recognition, *Journal of the American Statistical Association*, *91*(435), 953–960.

Shen, J., Y. Wang, and X. He (2014), *Inference for Logistic-Normal Mixtures with Heterogeneous Components*, manuscript.

Simon, R. (2002), Bayesian Subset Analysis: Application to Studying Treatment-by-Gender Interactions, *Statistics in Medicine*, *21*(19), 2909–2916.

Song, X., and M. S. Pepe (2004), Evaluating Markers for Selecting a Patient's Treatment, *Biometrics*, *60*(4), 874–883.

Song, Y., and G. Y. Chi (2007), A Method for Testing a Pre-specified Subgroup in Clinical Trials, *Statistics in Medicine*, *26*(19), 3535–3549.

Su, X., C.-L. Tsai, H. Wang, D. M. Nickerson, and B. Li (2009), Subgroup Analysis via Recursive Partitioning, *Journal of Machine Learning Research*, *10*, 141–158.

Teicher, H. (1961), Identifiability of Mixtures, *Annals of Mathematical Statistics*, *32*(1), 244–248.

Teicher, H. (1963), Identifiability of Finite Mixtures, *Annals of Mathematical Statistics*, *34*(4), 1265–1269.

Tibshirani, R. (1996), Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society, Series B: Statistical Methodology*, *58*(1), 267–288.

Ueda, N., and Z. Ghahramani (2002), Bayesian model search for mixture models based on optimizing variational bounds, *Neural Networks*, *15*, 1223–1241.

van der Vaart, A. W. (1998), *Asymptotic Statistics*, Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press.

van der Vaart, A. W., and J. A. Wellner (2000), *Weak Convergence and Empirical Processes: With Applications to Statistics*, Springer, New York.

Wang, P. (1994), Mixed Regression Models for Discrete Data, Ph.D. thesis, The University of British Columbia.

Wang, P., and M. L. Puterman (1998), Mixed Logistic Regression Models, *Journal of Agricultural, Biological, and Environmental Statistics*, *3*(2), 175–200.

Wang, P., M. L. Puterman, I. Cockburn, and N. Le (1996), Mixed Poisson Regression Models with Covariate Dependent Rates, *Biometrics*, *52*(2), 381–400.

Wong, C. S., and W. K. Li (2001), On a Logistic Mixture Autoregressive Model, *Biometrika*, *88*(3), 833–846.

Yuksel, S. E., J. N. Wilson, and P. D. Gader (2012), Twenty Years of Mixture of Experts, *IEEE Transactions on Neural Networks and Learning Systems*, *23*(8), 1177–1193.

Zhao, L., L. Tian, T. Cai, B. Claggett, and L. Wei (2013), Effectively Selecting a Target Population for a Future Comparative Study, *Journal of the American Statistical Association*, *108*(502), 527–539.

Zhu, H., and H. Zhang (2004), Hypothesis Testing in Mixture Regression Models, *Journal of the Royal Statistical Society, Series B: Statistical Methodology*, *66*(1), 3–16.

Zhu, H., and H. Zhang (2006), Asymptotics for Estimation and Testing Procedures under Loss of Identifiability, *Journal of Multivariate Analysis*, *97*, 19–45.