# MACHINE LEARNING METHODS FOR PRECISION MEDICINE USING PATIENT ELECTRONIC HEALTH RECORDS AND MOBILE SENSOR DATA

Jitong Lou

A dissertation submitted to the faculty at the University of North Carolina at Chapel Hill in partial fulfillment of the requirements for the degree of Doctor of Philosophy in the Department of Biostatistics in the Gillings School of Global Public Health.

Chapel Hill
2021

Approved by:

Donglin Zeng

Xinming An

Gary Koch

Xianming Tan

Yuanjia Wang

## ABSTRACT

Jitong Lou: Machine Learning Methods for Precision Medicine using
Patient Electronic Health Records and Mobile Sensor Data
(Under the direction of Donglin Zeng)

In the field of precision medicine, researchers adopt machine learning techniques to solve health-related problems, while applying such methods needs substantial health data. Electronic health records (EHRs) and mobile sensor data have become two important and abundant sources of health data. However, the modeling techniques for applying such data are still under development. The objective of this dissertation is to develop innovative frameworks of machine learning methods to use EHRs and/or mobile sensor data for disease prediction and precision medicine.

The first problem we address is using retrospectively collected EHRs data to learn latent patterns that can inform patient's health status. To handle data challenges in EHRs, we propose an approach that is based on multivariate generalized linear models in which latent Gaussian processes are introduced to model between-marker correlations over time. Using the inferred latent processes, we integrate irregularly measured health markers of mixed types into composite scores and apply hierarchical clustering to learn latent subgroup structures among patients. We demonstrate the utility of the proposed model through simulation studies and an EHRs dataset for type 2 diabetes (T2D) patients.

The next topic we investigate is recommending optimal individualized treatments to patients in EHRs data. To handle the multicategory comparison of treatments and confounding effects among patients, we incorporate the latent subgroups and use the one-versus-one approach to extend a matched learning model. Each matched learning for binary treatments is implemented by a weighted support vector machine with matched pairs of patients. Using the proposed method, we select the optimal treatments from four classes of T2D treatments and achieve a better control of glycated hemoglobin than one-size-fits-all rules for an EHRs dataset.

The last problem we explore is using mobile sensor data to predict outcomes and identify

objective biomarkers related to adverse posttraumatic neuropsychiatric sequelae. To overcome the difficulties in utilizing mobile sensor data, we develop a two-stage model that considers the measurement resolution and temporal pattern of features collected from mobile sensors. Finally, we apply our method to predict the pain experience of participants who experienced traumatic events, using the data collected from a large-scale cohort study.

To my parents.

## ACKNOWLEDGEMENTS

I would first like to express my deepest gratitude to my dissertation advisor, Dr. Donglin Zeng. Your talent and expertise were invaluable in formulating the research questions and methodology of this dissertation. From every meeting with you, I could learn new insights into the projects and statistics. More importantly, your dedication and passion in teaching, research, and supervision set a perfect example for me and made me a better researcher.

I would also like to sincerely thank my committee members, Dr. Xinming An, Dr. Gary Koch, Dr. Xianming Tan, and Dr. Yuanjia Wang. Dr. An made a significant contribution to the data source and data analysis in Chapter 4. I also want to thank you for supervising me in my final year and guiding me to the right directions in industry. Dr. Koch proposed some novel extensions that I had never thought about for the dissertation. Other than that, I am deeply grateful to you for your generous help on the GRA position at UCB inc. during my second year. This experience helped me a lot in job hunting. Dr. Tan was kind to be my committee member and provided helpful feedback to sharpen my thinking. Dr. Wang was intensively engaged in the research in Chapters 2 and 3. Besides your insightful comments on electronic health records and precision medicine, I really appreciate your time for revising our manuscripts and pushing them to the publications. Without any of you, I could not have completed this dissertation, and all of your help brought my work to a higher level.

In addition, I would like to thank Dr. Di Wu, Dr. Cyrus Vaziri, and Dr. Yang Yang for your generous financial support in my final year as well as for the 4-year collaboration. It was you that led me to the field of bioinformatics and cancer genomics, and these memorable research experience made me become versatile. Also, I am very thankful to Dr. Richard Bilsborrow for being my academic advisor during my master's degree, and to all other excellent faculty and staff members at UNC for their various assistance.

Next, many thanks to my friends for sharing valuable personal experience with me and for inviting me to various events outside of academic life.

Last but not least, I would like to give special thanks to my parents and my girlfriend, Juliana Yuan. Dad and Mom, thank you for always encouraging me to learn what I like and for supporting me unconditionally to pursue my dreams. I finally made it. Thank you, Juliana, for always being there for me. Your love and companionship are the best support for me during the graduate program, and I cannot wait to enter the next stage of life with you.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| Acrophase | acrophase of cosinor rhythmometry. |
| Amplitude | amplitude of cosinor rhythmometry. |
| ApEn | approximate entropy. |
| APNS | Adverse posttraumatic neuropsychiatric sequelae. |
| AURORA | Advancing Understanding of RecOvery afteR traumA. |
| avgsqi | average Signal-Quality-Index. |
| BMI | body mass index. |
| dc | heart rate deceleration capacity. |
| DD | whether diabetic drugs were prescribed at a clinical encounter. |
| ECG | electrocardiogram. |
| EHR | electronic health record. |
| GLMs | generalized linear models. |
| GPS | global positioning systems. |
| HbA1c | glycated hemoglobin. |
| HBP | hypertension/high blood pressure. |
| HDL | high-density lipoprotein. |
| HRV | heart rate variability. |
| ITR | individualized treatment rule. |
| L5 | least active five hours. |
| LASSO | least absolute shrinkage and selection operator. |
| LF/HF | ratio of low and high frequency spectral contents. |
| logMed | logarithm of the number of medications prescribed at each encounter. |
| meanACC | mean of activity counts. |
| NN intervals | interbeat intervals from which artifacts have been removed. |
| NNiqr | interquartile range of NN Intervals. |
| NNkurt | kurtosis of NN Intervals. |
| NNmean | mean of NN Intervals. |
| NNskew | skewness of NN Intervals. |

| | |
|---|---|
| OSU-WMCIW | Ohio State University Wexner Medical Center Information Warehouse. |
| PPG | photoplethysmography. |
| PTSD | posttraumatic stress disorder. |
| RA | relative amplitude. |
| RBF | radial basis function. |
| RCT | randomized clinical trial. |
| RDoC | Research Domain Criteria. |
| RMSSD | root-mean square differences of successive RR intervals. |
| RR intervals | intervals between all successive heartbeats. |
| SBP | systolic blood pressure. |
| SD1/SD2 | ratio of two standard deviation measures for a Poincare plot. |
| SDNN | standard deviation of NN Intervals. |
| stdACC | standard deviation of activity counts. |
| SVM | support vector machine. |
| SVR | support vector regression. |
| SWfragmentation | sleep-wake fragmentation. |
| T2D | type 2 diabetes. |
| TC | total cholesterol. |
| wakePercentage | wake percentage. |
| WMSE | weighted mean squared error. |

# CHAPTER 1: INTRODUCTION

In the modern era of healthcare, precision medicine has attracted great attentions from researchers in fields of statistical and biomedical science. Precision medicine is a medical paradigm that evolves individual characteristics of each patient such as demographics, lab test results, and genetic information, to optimize treatments (Ginsburg and Phillips, 2018).

As the emergence of large-scaled electronic systems, one important source of patient's health data is the electronic health record (EHR) which automatically captures patients health information through normal medical practice (Gunter and Terry, 2005; Cebul et al., 2011; Herrin et al., 2012). EHRs store patient health information in a digit format and the data can be shared across different institutions, so doctors and researchers have more opportunities to utilize the vast amounts of longitudinal data recorded in every moment. In general, EHRs record the information including patient demographics, vital signs, laboratory test results, medications, disease diagnosis codes, and medical insurances, on a large population over long time frames. Compared to experimental research like the randomized clinical trial (RCT), such observational studies provide stronger real world evidences because of larger patient populations, more flexible patient eligibility criteria, and longer duration for observations. Furthermore, studies using large-scaled EHRs may reflect real-world patterns of treatment pathway which can not be conducted and observed in RCTs (Hripcsak et al., 2016). Therefore, integrative analyses of these information across time provide great opportunities to understand individual patient's disease progression and susceptibility in real world settings, so as to predict disease prognosis and optimize personalized treatments adapted to evolving patient-specific features.

Although EHR data have advantages to the research in precision medicine, there are still some major challenges in applying such data. First of all, health markers in EHRs are often correlated and irregularly measured. Data sparsity and missing data can also be serious issues. In addition, EHR data have a mixture of health marker types such as continuous, binary, and count data. Moreover, there are underlying homogeneities and heterogeneities among thousands of patients in the data

source. Thus, existing approaches for longitudinal data have limitations and cannot handle all the data challenges.

In Chapter 2, we propose an innovative framework to take advantage of the rich health information in retrospectively collected EHRs and identify latent patient subgroups. The framework is built on multivariate generalized linear models (GLMs) which jointly analyze correlated and mixed type of health markers over time. In multivariate GLMs, covariate effects are time dependent and latent Gaussian processes are introduced to characterize between-marker correlations over time. Using inferred latent processes, we integrate the irregularly measured health markers of mixed types into composite scores and apply hierarchical clustering to learn latent subgroup structures among patients. We use the method of moments and kernel-weighted local estimating equations to estimate parameters that represent the covariate effects and between-marker correlations. Also, we adopt an inverse weighted method to standardize the intensity of health marker measurements. In this way, the bias, which is caused by heterogeneous temporal patterns of different health markers, in parameter estimation is reduced. We prove theoretical properties of the proposed estimators such as Fisher consistencies and asymptotic distributions. To demonstrate the performance of the proposed framework on finite samples, we apply our method to type 2 diabetes (T2D) patients in an EHR dataset collected from the Ohio State University Wexner Medical Center Information Warehouse (OSU-WMCIW). The analysis shows different trends of age, sex, and race effects on hypertension/high blood pressure (HBP), total cholesterol (TC), glycated hemoglobin (HbA1c), high-density lipoprotein (HDL), and medications. The associations among these markers vary over time during the study window. The hierarchical clustering of patients reveals four subgroups, and each patient subgroup is summarized by a unique profile based on patient's health status. The latent patterns are further confirmed by another split of the EHRs for the same cohort, suggesting that an effective healthcare management for these patients should be performed separately for each subgroup.

In pace with the rapid development of computational power, an increasing number of recent literature have adopted machine learning techniques to precision medicine (Mesko, 2017). Broadly speaking, these machine learning research estimate the individualized treatment rule (ITR) through maximizing some clinical outcomes. Most of the existing methods are designed for the comparison between binary treatments and/or for RCTs. Nevertheless, in real world data such as the EHR

dataset we handle with, there are multicategory treatments (86 types of T2D drugs) and large samples of patients (over 50,000 patients) with diverse backgrounds.

To fill up the gap between existing methods and real world evidence, in Chapter 3, we extend a matched learning method (Wu et al., 2020) to recommend the most effective treatment for each individual in the EHR dataset. We handle multicategory treatments by the one-versus-one approach and majority voting strategy (Bishop, 2006). Each matched learning classifier for two treatments is implemented by a weighted support vector machine (SVM) (Cortes and Vapnik, 1995) with matched sets of patients. The matched set of a target patient is defined as a group of patients who have the same group membership and have similar characteristics as him/her, but they receive an alternative treatment. By comparing a target patient only with patients in his/her matched set, the confounding effects are reduced, so, in this case, the difference in the clinical outcome reflects the treatment effect. If the average clinical outcome over the matched set is more beneficial, then the target patient should switch from the assigned treatment to the alternative treatment. Otherwise, the assigned treatment is already the optimal treatment for the target patient. Lastly, the treatment with the highest vote across all binary comparisons is set to be the optimal treatment of all classes. In a real data application, we estimate the optimal individualized treatment among four types of T2D treatments for over five thousands patients (a different cohort from Chapter 2) in the EHR dataset of the OSU-WMCIW. We compare our method with four one-size-fits-all rules and two treatment rules estimated by Q-learning (Watkins and Dayan, 1992; Murphy, 2005; Qian and Murphy, 2011), and the comparison shows the proposed method has a better management of HbA1c level than other models by at least 5%-13%.

Besides EHRs, personal sensing data collected from mobile sensor also have received increasing interests from various areas, especially in health care. Mobile sensors embedded in smartphones (short messaging service, microphone, global positioning systems (GPS)), and smartwatches (electrocardiogram (ECG), accelerometer, time stamp) continuously monitor the health-related indices and activities of subjects. Thus, doctors and patients can precisely understand the treatment response, disease progression, and health status through the continuous measurement of health markers. However, modeling techniques for mobile sensor data are still under development. The challenges of handling mobile sensor data come from three major aspects. First of all, mobile sensor data are high dimensional. Secondly, mobile sensor data are often correlated and have lagged effects on

outcome variables. Lastly, data collected from different sensors can have different time scales. For example, heart rate data have thousands of measurements for each day and the data have cyclic patterns for the sample person. In the meantime, actigraphy data just have daily measurements. Most of the existing literature simply treat all the features as independent predictors and fit them using complicated machine learning or deep learning models, but these data challenges will bring in biases to the prediction indeed.

In Chapter 4, we create a two-stage semi-parametric method for modeling clinical outcomes and identifying objective biomarkers related to psychiatric disorder using mobile sensor data. The first stage adopts a linear regression model to describe the relationship between different domains/sensors of features and health markers. The model handles the effect of time scales and retrospective measurements on the features by assigning weights to features measured at different time points. We apply the least absolute shrinkage and selection operator (LASSO) method (Tibshirani, 1996; Tibshirani et al., 2004) to select the most informative features and improve the model interpretability. In the second stage, we implement the selected features and estimated weights to a support vector regression (SVR) model (Drucker et al., 1997) to improve the prediction accuracy. This step accounts for non-linear interactions and between-domain comorbidities of features that are not captured in the first stage. We verify the proposed method on a sample of participants in the Advancing Understanding of RecOvery afteR traumA (AURORA) study (McLean et al., 2020). For participants who experienced traumatic events, this application studies the relationship between construct scores, which quantify the pain experience of participants, and two domains of mobile sensor features at six follow-up time points after the events. The two domains of mobile sensor features include eight activity features and eleven heart rate variability (HRV) features collected from smartwatches. Compared to a SVR model that does not adjust for feature selection, measurement resolutions, and temporal patterns, the proposed method achieves a better prediction performance by using only 12% of the total features.

Each of Chapters 2 to 4 defines the research problem, introduces the background and existing literature, explains the proposed method, demonstrates the method through numeric examples, and summarizes the contributions and conclusions. Chapter 5 discusses the limitations of this dissertation and provides potential directions for future research.

## CHAPTER 2: LEARNING LATENT HETEROGENEITY FOR TYPE 2 DIABETES PATIENTS USING LONGITUDINAL HEALTH MARKERS IN ELECTRONIC HEALTH RECORDS

## 2.1 Introduction

In the modern era of precision medicine, one important source of patient's health data is EHRs. EHRs data consist of longitudinal medical records from a large number of patients in one or more electronic healthcare systems that digitally capture measurements of patients health status through normal medical practices (Gunter and Terry, 2005; Cebul et al., 2011; Herrin et al., 2012), including patient's vital signs, laboratory measurements, disease diagnosis codes, procedure codes, and medications. Benefits of EHRs include cost effectiveness, real time updates, and reflections on patients disease course and healthcare managements in realistic settings. Therefore, integrative analyses of this information over time provide great opportunities to understand the heterogeneity of patient's disease progression and susceptibility in real world settings, which is useful for monitoring disease prognosis and optimizing personalized healthcare management.

Due to the retrospective nature of EHRs, the analysis of EHRs is complicated by the following challenges: first, the health markers measured over time are multivariate and the measurements can be either continuous (e.g., lab measures), binary (e.g., disease diagnoses), or counts (e.g., number of medications); second, for each patient, the health marker data are collected at each clinical encounter so the measurement times can be irregular, sparse, and heterogeneous across patients; third, the measurement times are often informative to patients health status or health care processes.

This work is motivated by the analyses of EHRs of T2D patients obtained from the OSU-WMCIW. The data collection spanned a time period of 8 years (between 2011 and 2018) from a total of 58,490 patients. The data contained patients medical records of glycated hemoglobin, high-density lipoprotein, total cholesterol, hypertension, and all medications prescribed at each clinical encounter. Because these markers were of different types and were not measured at the same time across and within patients, directly combining the values from these markers is neither meaningful nor feasible. For example, Figure 2.1 gives a snapshot of the measurement time of

several health markers from 20 randomly selected patients. Clearly, each marker was measured sparsely at irregular times for each patient, and the measurement time patterns vary significantly from patient to patient.



Figure 2.1: Observation time patterns of 5 health markers for 20 randomly selected T2D patients in the EHRs at the OSU-WMCIW.

Joint models based on linear or generalized mixed effects models have been commonly used for analyzing multivariate longitudinal data (Verbeke et al., 2014). In the joint models, various distribution families are used (Verbeke and Molenberghs, 2000; Davidian and Giltinan, 2003; Molenberghs and Verbeke, 2005), and subject-specific random effects are shared across all health markers to explain their dependence due to a finite number of latent variables. For example, Lambert and Vandenhende (2002) jointly analyzed three repeatedly measured longitudinal outcomes using copula models in a dose titration safety study; Gueorguieva and Sanacora (2006) proposed correlated probit models for joint analysis of repeated measurements with ordinal and continuous health markers. Some extensions allowed time-dependent effects (Huang et al., 2002; Fan and Zhang, 2008), but assumed constant between-marker dependence over time. However, assuming parametric patterns or attributing the dependence to a few time-invariant random effects is rather restrictive

6

especially for modeling EHRs over a long period of time, since in EHRs, the trajectories of the health markers and their dependence may vary over time depending on the disease progression and medication usage for each patient. Moreover, it is computationally challenging to maximize a joint likelihood in the presence of a large number of patients and many health markers.

Machine learning approaches have been also proposed to perform EHR analysis, such as deep Poisson factor models (Henao et al., 2016), tensor factorization and non-negative matrix factorization (Ho et al., 2014), and deep exponential families (Miscouridou et al., 2018). These approaches, although more flexible than aforementioned statistical models, are less interpretable and are highly computationally intensive, requiring substantial work for data engineering and model tuning. More importantly, none of these approaches can account for irregular but informative measurement patterns as seen in EHRs.

In this chapter, we seek to strike a balance between the complex statistical modelling and flexible machine learning methods, while accounting for the unique challenges in EHRs. To conduct an integrative analysis of EHRs, we extend the multivariate GLMs by assuming appropriate distribution and link functions depending on the marker type. We allow the effects of covariates on health markers to be time-varying. Moreover, to account for the time-varying dependence among health markers, we introduce latent Gaussian processes into the models, where the covariance matrix is assumed to vary over time. For estimation, we adopt kernel smoothing method to pool information across time points and patients and apply weights to account for the heterogeneous patterns of measurement times. The inferred latent processes represent patients underlying health status, so in order to integrate these mixed-type health markers, we use the inferred latent processes to calculate the distances between any two patients using the Mahalanobis distance (De Maesschalck et al., 2000). Finally, we apply hierarchical clustering to identify patients health patterns and characterize between-group heterogeneities.

The remaining parts of this chapter are organized as follows. In Section 2.2, we propose our models and describe main ideas. We then provide inferences on estimating model parameters and procedures to perform numerical computations. In Section 2.3, we derive the asymptotic distributions of the estimators. We conduct simulation studies in Section 2.4. In Section 2.5, we apply our method to an integrative analysis on health markers for T2D patients using EHRs from the OSU-WMCIW.

## 2.2 Methodologies

### 2.2.1 Statistical Models for Integrative Analysis

Suppose EHR data are obtained from $n$ patients. For the $i$th patient, let $\boldsymbol{X}_i$ be $m$-dimensional baseline covariates. Among $p$ health markers, let $Y_{ik}(t)$ denote the measurement of the $k$th health marker at time $t$. We suppose $Y_{ik}(t)$ is measured at time points $t_{ik1}$, $t_{ik2}$, ..., $t_{ikn_{ik}}$, where $n_{ik}$ is the total count of observations on the $k$th health marker for the $i$th patient. The total number of observations up to time $t$ can be represented by a counting process $N_{ik}(t) \equiv \sum_{j=1}^{n_{ik}} I(t_{ikj} \leq t)$, where $I(\cdot)$ is the indicator function. Since the documentation times are patient's clinical encounters in the EHR system, patterns of these documentation/measurement time points may carry information on patients health status. Thus, we model the intensity of $N_{ik}(t)$ as

$$\mathbb{E}\left[dN_{ik}(t)|\boldsymbol{X}_i\right] = \lambda_k(t)\exp\left\{\boldsymbol{X}_i^T\boldsymbol{\gamma}_k\right\}dt, \tag{2.1}$$

where $\lambda_k(t)$ is a baseline intensity function, and $\gamma_k$ is a vector of intensity parameters. By modeling the intensity of EHR measurement rates, one can adjust for the bias of informative measurement patterns and account for between patient heterogeneity.

We further assume $Y_{ik}(t)$ follows a distribution in an exponential family model as follows:

$$f_{ik}(y; \theta_{ik}, \phi_{ik}, t) = \exp\left\{\frac{y\theta_{ik}(t) - b_k\left(\theta_{ik}(t)\right)}{a_k\left(\phi_{ik}(t)\right)} + c_k\left(y, \phi_{ik}(t)\right)\right\}, \tag{2.2}$$

where $\theta_{ik}(t)$ and $\phi_{ik}(t)$ are the canonical parameter and the dispersion parameter, respectively, specific to each patient and each health marker. $a_k(\cdot)$, $b_k(\cdot)$, and $c_k(\cdot)$ are known functions. Let $\theta_{ik}(t) = g_k(\mu_{ik}(t))$, where $g_k(\cdot)$ is the canonical link function, and $\mu_{ik}(t)$ is the mean of $Y_{ik}(t)$. To capture the patient heterogeneity and dependence, we assume, at time $t$,

$$g_k(\mu_{ik}(t)) = \boldsymbol{X}_i^T\boldsymbol{\beta}_k(t) + \epsilon_{ik}(t),$$

$$\boldsymbol{\epsilon}_i(t) \sim \mathcal{N}_p\left(\boldsymbol{0}, \boldsymbol{\Omega}(t)\right), \tag{2.3}$$

where $\boldsymbol{\beta}_k(t)$ is a vector of regression coefficients for covariates $\boldsymbol{X}_i$. $\epsilon_{ik}(t)$ is the $k$th element of the latent Gaussian process $\boldsymbol{\epsilon}_i(t) = \{\epsilon_{i1}(t), \epsilon_{i2}(t), \ldots, \epsilon_{ip}(t)\}^T$. $\boldsymbol{\epsilon}_i(t)$ is independent of $\boldsymbol{X}_i$, and it follows

a mean-zero multivariate Gaussian distribution with a covariance matrix $\boldsymbol{\Omega}(t)$. Estimating variances locally will requires dense measurements from the same health marker, which is not the case for the EHRs. Moreover, in our empirical application the estimated variances do not vary much across time (Section 2.6). Thus, to ensure numerical stability in subsequent analysis, we assume each latent process to have a constant variance and the constant is estimated using historical records. Hence, in $\boldsymbol{\Omega}(t)$, only the correlations among health markers, that is, the off-diagonal elements need to be estimated.

Under the proposed models (2.2) and (2.3), each measurement $Y_{ik}(t)$ can be uniquely represented by the latent process $\epsilon_{ik}(t)$. Since $\epsilon_{ik}(t)$ has the same scale for different $k$, one can integrate the latent processes $\{\epsilon_{ik}(t) : k = 1, 2, \ldots, p\}$ as an alternative way to integrate the mixed-type health markers. The integration can use the Mahalanobis distance as follows,

$$D_{ij} = \left\{ \int_t [\boldsymbol{\epsilon}_i(t) - \boldsymbol{\epsilon}_j(t)]^T \boldsymbol{\Omega}^{-1}(t) [\boldsymbol{\epsilon}_i(t) - \boldsymbol{\epsilon}_j(t)] \, dt \right\}^{1/2}. \tag{2.4}$$

Thus, there are several important advantages of using the proposed models to perform an integrative analysis of mixed-type health markers. First of all, despite the health markers are irregularly measured and mixed-type, we can map them onto the same scale to align patients and characterize the between-patients heterogeneity. In addition, the dimension of latent processes can be further reduced to some lower dimensional subspaces than the number of health markers. Therefore, through the representation of latent processes, we achieve a dimension reduction.

### 2.2.2 Model Parameter Estimation

First, we use marker-specific Anderson-Gill intensity models (Andersen and Gill, 1982) to estimate $\boldsymbol{\gamma}_k$ in (2.1). With the estimator $\widehat{\boldsymbol{\gamma}}_k$, we normalize the counting process $N_{ik}(t)$ by letting $\widetilde{N}_{ik}(t) = N_{ik}(t) \exp\left\{-\boldsymbol{X}_i^T \widehat{\boldsymbol{\gamma}}_k\right\}$. Thus, the normalized counting process is homogeneous across different patients and different health markers.

Next, to estimate $\boldsymbol{\beta}_k(t)$ for any fixed time point $t$, we solve the following kernel-weighted local estimating equation

$$U_{n,k}(\boldsymbol{\beta}_k(t)) \equiv \frac{1}{n} \sum_{i=1}^n \int \boldsymbol{X}_i \left[Y_{ik}(s) - \mathbb{E}\left[Y_{ik}(t)|\boldsymbol{X}_i\right]\right] K_{h_{1n}}(s-t) d\widetilde{N}_{ik}(s) = \boldsymbol{0}, \tag{2.5}$$

where $K_h(z) = h^{-1}K(z/h)$ with $K(z)$ being a symmetric kernel function, and $h_{1n}$ is the bandwidth of $K_h(z)$. Essentially, we assign weights to the observed measurements $Y_{ik}(s)$ near $t$, and we pool them together across all patients to estimate the mean (first moment) of $Y_{ik}(t)$. This pooling process relies on the kernel smoothing. Also, pooling information across observations nearby and across patients overcomes the difficulty in parameter estimations that some sparsely measured health markers do not have sufficient samples at some time points. Moreover, using $d\widetilde{N}_{ik}(s)$ instead of $dN_{ik}(s)$, we remove the heterogeneity of informative measurement time points among patients in a similar spirit as inverse probability weighting.

Similarly, to estimate the correlation between two latent processes, $\sigma_{kl}(t) = \mathrm{Cov}(\epsilon_{ik}(t), \epsilon_{il}(t))$, we propose to solve the following kernel-weighted local estimating equation, for $k \neq l$,

$$U_{n,k,l}(\sigma_{kl}(t)) \equiv \frac{1}{n^2} \sum_{i=1}^{n} \iint \left[Y_{ik}(s)Y_{il}(s') - \mathbb{E}\left[Y_{ik}(t)Y_{il}(t)|\boldsymbol{X}_i\right]\right] \widetilde{K}_{h_{2n}}(s-t, s'-t)d\widetilde{N}_{ik}(s)d\widetilde{N}_{il}(s') = 0,$$
(2.6)

where $\widetilde{K}_h(z_1, z_2)$ is a bivariate kernel function with bandwidth $h_{2n}$.

### 2.2.3 Numerical Computation

When the link functions in (2.3) take some simple forms, $\mathbb{E}\left[Y_{ik}(t)|\boldsymbol{X}_i\right]$ in (2.5) and $\mathbb{E}\left[Y_{ik}(t)Y_{il}(t)|\boldsymbol{X}_i\right]$ in (2.6) can be explicitly computed. Specifically, for $g_k(z) = g_l(z) = z$,

$$\mathbb{E}\left[Y_{ik}(t)|\boldsymbol{X}_i\right] = \boldsymbol{X}_i^T\boldsymbol{\beta}_k(t),$$

and

$$\mathbb{E}[Y_{ik}(t)Y_{il}(t)|\boldsymbol{X}_i] = \boldsymbol{X}_i^T\boldsymbol{\beta}_k(t)\boldsymbol{X}_i^T\boldsymbol{\beta}_l(t) + \sigma_{kl}(t).$$

When $g_k(z)$ takes a general form, we can compute the above expectations using the Gauss-Hermite quadrature method (Abramowitz and Stegun, 1965).

Since $U_{n,k}(\boldsymbol{\beta}_k(t))$ is only related to the parameter $\boldsymbol{\beta}_k(t)$, we can solve (2.5) and obtain $\widehat{\boldsymbol{\beta}}_k(t)$ for each health marker $k$, separately. Similarly, plugging $\widehat{\boldsymbol{\beta}}_k(t)$ and $\widehat{\boldsymbol{\beta}}_l(t)$ to (2.6), we can solve the equation and obtain $\widehat{\sigma}_{kl}(t)$ for each pair of health markers, separately. Therefore, even with many health markers, that is, $p$ is moderate or large, our algorithm can efficiently handle the computation burden by solving the estimating equations separately. Finally, we apply the above procedures for time grids $t_1, t_2, \ldots, t_N$ to obtain the parameter estimators over the whole range of the follow-up.

A distance matrix $D$ can be obtained by computing the Mahalanobis distance in (2.4) between each pair of patients. In particular, with the estimated latent processes, the distance is approximated by

$$D_{ij} = \left\{ \sum_{t=t_1}^{t_N} [\widehat{\boldsymbol{\epsilon}}_i(t) - \widehat{\boldsymbol{\epsilon}}_j(t)]^T \widetilde{\boldsymbol{\Omega}}^{-1}(t) [\widehat{\boldsymbol{\epsilon}}_i(t) - \widehat{\boldsymbol{\epsilon}}_j(t)] \right\}^{1/2}, \tag{2.7}$$

and

$$\widehat{\boldsymbol{\epsilon}}_i(t) = \mathbb{E}\left[ \boldsymbol{\epsilon}_i(t) \Big| \boldsymbol{Y}_i(t), \widehat{\boldsymbol{\beta}}_k(t), \widehat{\sigma}_{kl}(t) \right], \tag{2.8}$$

where $\widetilde{\boldsymbol{\Omega}}(t)$ is the covariance matrix of $\widehat{\boldsymbol{\epsilon}}_i(t)$. In particular,

$$\mathbb{E}\left[ \boldsymbol{\epsilon}_i(t) \Big| \boldsymbol{Y}_i(t), \widehat{\boldsymbol{\beta}}_k(t), \widehat{\sigma}_{kl}(t) \right] = \frac{\int P\left( \boldsymbol{Y}_i(t) \Big| \boldsymbol{\epsilon}_i(t), \widehat{\boldsymbol{\beta}}_k(t), \widehat{\sigma}_{kl}(t) \right) P\left( \boldsymbol{\epsilon}_i(t) | \widehat{\sigma}_{kl}(t) \right) \boldsymbol{\epsilon}_i(t) d\boldsymbol{\epsilon}_i(t)}{\int P\left( \boldsymbol{Y}_i(t) \Big| \boldsymbol{\epsilon}_i(t), \widehat{\boldsymbol{\beta}}_k(t), \widehat{\sigma}_{kl}(t) \right) P\left( \boldsymbol{\epsilon}_i(t) | \widehat{\sigma}_{kl}(t) \right) d\boldsymbol{\epsilon}_i(t)}.$$

The subsequent steps can be calculated using the Gauss-Hermite quadrature method as well, and the details are given in Appendix A.1.

### 2.2.4    Data-adaptive Selection of Bandwidths

Our asymptotic results in Appendix A suggest the bandwidths $h_{1n}$ and $h_{2n}$ can be chosen, respectively, on the order of $n^{-1/3}$ and $n^{-1/4}$. However, for practical applications, we consider a data-adaptive method for selecting the bandwidths (Cao et al., 2015). The key idea is using observed data to obtain the empirical bias and variability of the estimators in terms of the bandwidths. Consequently, we search for the bandwidths that minimize the empirical mean squared error of selecting them.

Specifically, to choose the optimal bandwidth $h_{1n}$ for estimating $\widehat{\boldsymbol{\beta}}_k(t)$, we first consider a reasonable range of bandwidths. For a fixed bandwidth $h$ and a fixed time point $t$, we denote $\widehat{\boldsymbol{\beta}}_{kh}(t)$ to the estimator for $\boldsymbol{\beta}_k(t)$. To estimate the bias of $\widehat{\boldsymbol{\beta}}_{kh}(t)$, we fit a least squares regression by regressing $\widehat{\boldsymbol{\beta}}_{kh}(t)$ on $h^2$. We denote the regression coefficient of $h^2$ as $\widehat{\boldsymbol{C}}_k(t)$. Since the bias of $\widehat{\boldsymbol{\beta}}_{kh}(t)$ is on the order of $h^2$, as shown in the asymptotic result, $\left\| \widehat{\boldsymbol{C}}_k(t) \right\| h^2$ is an estimator for the bias of $\widehat{\boldsymbol{\beta}}_{kh}(t)$. Next we investigate the variability of $\widehat{\boldsymbol{\beta}}_{kh}(t)$. We randomly split the data into two equal parts. Using either one of the split data, we obtain $\widehat{\boldsymbol{\beta}}_{kh}^{*1}(t)$ as the estimator for $\boldsymbol{\beta}_{kh}(t)$ in this case. Similarly, using the other half, we obtain $\widehat{\boldsymbol{\beta}}_{kh}^{*2}(t)$. Thus, $\frac{1}{4} \left\| \widehat{\boldsymbol{\beta}}_{kh}^{*1}(t) - \widehat{\boldsymbol{\beta}}_{kh}^{*2}(t) \right\|^2$ can be used as an unbiased estimator of the variance of $\widehat{\boldsymbol{\beta}}_{kh}(t)$. Finally, given all the time points, we select the

optimal bandwidth as $\arg\min_h \sum_t \text{MSE}_\beta^h(t)$, where

$$\text{MSE}_\beta^h(t) = \sum_{k=1}^p \left\{ \widehat{\text{Var}} \left[\widehat{\boldsymbol{\beta}}_{kh}(t)\right] + \left(\widehat{\text{Bias}} \left[\widehat{\boldsymbol{\beta}}_{kh}(t)\right]\right)^2 \right\} = \sum_{k=1}^p \left\{ \frac{1}{4} \left\|\widehat{\boldsymbol{\beta}}_{kh}^{*1}(t) - \widehat{\boldsymbol{\beta}}_{kh}^{*2}(t)\right\|^2 + \left\|\widehat{\boldsymbol{C}}_k(t)\right\|^2 h^4 \right\}.$$
(2.9)

We denote the optimal $h_{1n}$ as $H_1$ and denote the corresponding estimators for $\boldsymbol{\beta}_k(t)$ as $\widehat{\boldsymbol{\beta}}_{kH_1}(t)$. Next, given $h_{1n} = H_1$ and $\boldsymbol{\beta}_k(t) = \widehat{\boldsymbol{\beta}}_{kH_1}(t)$, we select the optimal $h_{2n}$, the bandwidth for estimating $\sigma_{kl}(t)$'s, by minimizing the empirical mean squared error of the corresponding estimators, which is numerically calculated in the similar way to above.

## 2.3 Theoretical Results

We first state the following required conditions.

**Condition 1.** True parameters $\lambda_k^0(t)$, $\boldsymbol{\beta}_k^0(t)$, and $\sigma_{kl}^0(t)$ are continuously twice differentiable for any $t \in [0, \tau]$, where $k, l = 1, 2, \ldots, p$ and $k \neq l$. In addition, $\lambda_k^0(t)$ is strictly positive. Furthermore, the second moments of $\text{Cov}(dN_{ik}(t), dN_{ik}(s)|\boldsymbol{X}_i)/dtds$ and temporal covariances $\text{Cov}(\epsilon_{ik}(t), \epsilon_{il}(s))$ are continuously twice-differentiable.

**Condition 2.** The vector of baseline covariate $\boldsymbol{X}$ is bounded. If there exists a vector $\boldsymbol{b}$ such that $\boldsymbol{X}^T\boldsymbol{b} = 0$, then $\boldsymbol{b} = 0$.

**Condition 3.** $h_{1n}, h_{2n} \to 0$ and $nh_{1n}, nh_{2n}^2 \to \infty$. Furthermore, $nh_{1n}^5, nh_{2n}^6 \to 0$.

**Condition 4.** The kernel function $K(z)$ is a symmetric density function satisfying $\int z^2 K(z)dz < \infty$. Similarly, $\widetilde{K}(z_1, z_2)$ is a symmetric bivariate density function with bounded fourth moments.

Condition 1 is used to give the asymptotic distribution for the parameter estimators in (2.1), and it assumes some smoothness properties of the time-varying coefficients and covariance matrices. From condition 3, the choice of $h_{1n}$ and $h_{2n}$ can be $n^{-1/3}$ and $n^{-1/4}$, respectively. A potential choice of the kernel satisfying condition 4 can be the Gaussian kernel or the Epanechnikov kernel. Theorem 2.3.1 states the asymptotic distribution of parameters $\widehat{\boldsymbol{\beta}}_k(t)$, $k = 1, 2, \ldots, p$. Theorem 2.3.2 establishes the asymptotic distribution of parameters $\widehat{\sigma}_{kl}(t)$, $k, l = 1, 2, \ldots, p$, and $k \neq l$.

**Theorem 2.3.1** (asymptotic distribution of $\widehat{\boldsymbol{\beta}}_k(t)$). *Under conditions 1 to 4, for any fixed $t$,*

$$(nh_{1n})^{1/2} A_k(t) \left[\widehat{\boldsymbol{\beta}}_k(t) - \boldsymbol{\beta}_k^0(t)\right] \to_d \mathcal{N}_m \left(\boldsymbol{0}, \boldsymbol{\Sigma}_k(t)\right),$$
(2.10)

12

*where*

$$A_k(t) = \lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\boldsymbol{X}^T \int \left[g_k^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon_k(t))\right]' f(\epsilon_k(t))d\epsilon_k(t)\right],$$

*and the asymptotic variance*

$$\boldsymbol{\Sigma}_k(t) = \lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\boldsymbol{X}^T\sigma^2(t, \boldsymbol{X}, \epsilon_k(t)) \exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_k^0\right\}\right] \int_z K^2(z)dz,$$

*where $\sigma^2(t, \boldsymbol{X}, \epsilon_k(t))$ is a function of $\epsilon_k(t)$. Its definition and the proof of theorem 2.3.1 are given in Appendix A.2.*

**Theorem 2.3.2** (asymptotic distribution of $\widehat{\sigma}_{kl}(t)$)**.** *Under conditions 1 to 4, for any fixed $t$,*

$$(nh_{2n}^2)^{1/2}B_{kl}(t)\left[\widehat{\sigma}_{kl}(t) - \sigma_{kl}^0(t)\right] \to_d \mathcal{N}(0, \Sigma_{kl}(t)), \tag{2.11}$$

*where*

$$
\begin{aligned}
B_{kl}(t) &= \lambda_k^0(t)\lambda_l^0(t)\mathbb{E}\left[\iint g_k^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon_k(t))g_l^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_l^0(t) + \epsilon_l(t))\right. \\
&\quad \times \left. \frac{\partial f(\epsilon_k(t), \epsilon_l(t); \sigma_{kl}(t))}{\partial\sigma_{kl}(t)}\bigg|_{\sigma_{kl}(t)=\sigma_{kl}^0(t)} d\epsilon_k(t)d\epsilon_l(t)\right],
\end{aligned}
$$

*is assumed to be nonsingular, and the asymptotic variance*

$$\Sigma_{kl}(t) = \lambda_k^0(t)\lambda_l^0(t)\mathbb{E}\left[\psi^2(t, t, \boldsymbol{X}, \epsilon_k(t), \epsilon_l(t)) \exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_k^0\right\} \exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_l^0\right\}\right] \iint \widetilde{K}^2(z_1, z_2)dz_1dz_2,$$

*where $\psi^2(t, t, \boldsymbol{X}, \epsilon_k(t), \epsilon_l(t))$ is a function of $\epsilon_k(t)$ and $\epsilon_l(t)$. Its definition and the proof of theorem 2.3.2 are given in Appendix A.3.*

Since the asymptotic variances in theorem 2.3.1 and theorem 2.3.2 do not have simple expressions, we use the bootstrap method to estimate the asymptotic variances in practice.

## 2.4 Simulation Studies

In the simulation studies, we simulated data of six health markers for 5,000 subjects. For the $i$th subject, we generated two covariates $X_{i1} \sim \text{Uniform}(-1, 1)$ and $X_{i2} \sim \text{Bernoulli}(0.5) - 0.5$. Thus, $\boldsymbol{X}_i = (1, X_{i1}, X_{i2})^T$ was a three-dimensional vector of baseline variables. The maximum observation time $T_i$ for each subject was set to 12. The measured time points for simulated markers were generated

from a Poisson process whose intensity function was $\mathbb{E}\left[dN_{ik}(t)|\boldsymbol{X}_i\right] = 0.5 \exp\left\{0.5X_{i1} + 0.25X_{i2}\right\}dt$. For the variances of latent processes, we assumed $c_k = 1$, $k = 1, 2, \ldots, 6$. Suppose there were $N_i$ unique measured time points $t_{i1}, t_{i2}, \ldots, t_{iN_i}$ for all latent processes of the subject $i$, we sampled $\boldsymbol{\epsilon}_i(t_{i1})$, $\boldsymbol{\epsilon}_i(t_{i2})$, $\ldots$, $\boldsymbol{\epsilon}_i(t_{iN_i})$ from a mean-zero multivariate Gaussian distribution with a covariance matrix $\boldsymbol{\Omega}(\boldsymbol{t}_i) = \boldsymbol{\Sigma}_2(\boldsymbol{t}_i) \otimes \boldsymbol{\Sigma}_1$, where $\boldsymbol{t}_i = (t_{i1}, t_{i2}, \ldots, t_{iN_i})$,

$$
\boldsymbol{\Sigma}_1 = \begin{pmatrix}
1 & 0.34 & 0.48 & 0.58 & 0.03 & 0.05 \\
0.34 & 1 & 0.80 & -0.49 & -0.78 & 0.80 \\
0.48 & 0.80 & 1 & -0.16 & -0.36 & 0.53 \\
0.58 & -0.49 & -0.16 & 1 & 0.80 & -0.69 \\
0.03 & -0.78 & -0.36 & 0.80 & 1 & -0.85 \\
0.05 & 0.80 & 0.53 & -0.69 & -0.85 & 1
\end{pmatrix},
$$

and

$$
\boldsymbol{\Sigma}_2(\boldsymbol{t}_i) = \begin{pmatrix}
1 & e_{12} & \ldots & e_{1N_i} \\
e_{21} & 1 & \ldots & e_{2N_i} \\
\vdots & \vdots & \ddots & \vdots \\
e_{N_i 1} & e_{N_i 2} & \ldots & 1
\end{pmatrix},
$$

where $e_{kl} = \exp\left\{-(t_{ik} - t_{il})^2\right\}$, $k, l = 1, 2, \ldots, N_i$. Thus, at each measured time point, $\boldsymbol{\Omega}(t)$ is constant and equals to $\boldsymbol{\Sigma}_1$, but there exist underlying dependencies in the time intervals between these time points.

The values of simulated markers were generated according to (2.2) and (2.3). To assess the ability of our models in Section 2.2.1 to handle mixed-type markers, we assumed $Y_{i1}(t)$ and $Y_{i4}(t)$ were Gaussian distributed. $Y_{i2}(t)$ was Poisson distributed. $Y_{i3}(t)$, $Y_{i5}(t)$, and $Y_{i6}(t)$ were Bernoulli distributed. Thus, $g_1^{-1}(z) = g_4^{-1}(z) = z$, $g_2^{-1}(z) = e^z$, and $g_3^{-1}(z) = g_5^{-1}(z) = g_6^{-1}(z) = e^z/(1 + e^z)$. Furthermore, since the distributions of $Y_{i1}(t)$ and $Y_{i4}(t)$ had dispersion parameters, we set $\phi_{i1}(t) = \phi_{i4}(t) = 0.5$. The true values of $(\boldsymbol{\beta}_1(t), \boldsymbol{\beta}_2(t), \boldsymbol{\beta}_3(t), \boldsymbol{\beta}_4(t), \boldsymbol{\beta}_5(t), \boldsymbol{\beta}_6(t))$ were assumed to be

$$
\begin{pmatrix}
-0.44 - \frac{t}{8} & -0.93 + \frac{t}{9} & 0.35 + \frac{t}{10} & -1.36 + \frac{t}{10} & \cos(-0.25 + t) & 0.91 + \frac{(t-6)^3}{216} \\
0.6 + \frac{\sqrt{t}}{3} & -0.53 - \frac{\sqrt{t}}{2} & -2 + \sqrt{t} & \sin(0.76 + t) & 0.37 + \frac{t}{10} & \frac{t}{10} \\
-0.5 + \frac{\sqrt[3]{t}}{2} & 0.4 + \frac{\sqrt[3]{t}}{2} & 1.9 - \sqrt[3]{t} & \cos(-0.3 + t) & \sin(-0.68 + t) & 1.23 + \frac{(t-6)^2}{36}
\end{pmatrix}.
$$

The scaled Epanechnikov kernel was chosen as the kernel function in (2.5), that is,

$$K_{h_{1n}}(z) = \frac{3}{4h_{1n}} \left[ 1 - \left( \frac{z}{h_{1n}} \right)^2 \right]_+ . \tag{2.12}$$

Furthermore, the kernel function in (2.6) was set to the product of two scaled univariate Epanechnikov kernels, that is,

$$\widetilde{K}_{h_{2n}}(z_1, z_2) = \frac{9}{16h_{2n}^2} \left[ 1 - \left( \frac{z_1}{h_{2n}} \right)^2 \right]_+ \left[ 1 - \left( \frac{z_2}{h_{2n}} \right)^2 \right]_+ . \tag{2.13}$$

Since the data-adaptive method for selecting bandwidths was computationally intensive, we first conducted a preliminary study on the simulated data. We used the method in Section 2.2.4 and selected the optimal bandwidths among $h = cn^{-1/z}$, where $n = 5000$, $c = \{5, 10, 20, 30\}$, and $z = 1, 2, \ldots, 10$. Hence, the potential bandwidths ranged from 0.001 to 12.800. We found $h_{1n} = 5n^{-1/3} = 0.292$ and $h_{2n} = 10n^{-1/3} = 0.585$ were close to the optimal. This set of $h_{1n}$ and $h_{2n}$ was used in all subsequent simulations.

For time points $t = 0, 1, \ldots, 12$, we solved (2.5) and (2.6), and we obtained $\widehat{\boldsymbol{\beta}}_k(t)$ and $\widehat{\sigma}_{kl}(t)$. We evaluated the accuracies of the asymptotic approximations by calculating the average bias and the sample standard deviation of $\widehat{\boldsymbol{\beta}}_k(t)$ and $\widehat{\sigma}_{kl}(t)$, respectively. In addition, using the bootstrap method, we calculated the bootstrap estimators for standard errors of $\widehat{\boldsymbol{\beta}}_k(t)$ and $\widehat{\sigma}_{kl}(t)$. Specifically, for each dataset, we resampled 5000 observations with replacement from $\boldsymbol{X}$ to produce a bootstrap dataset $\boldsymbol{X}^{*1}$. We could use $\boldsymbol{X}^{*1}$ to produce a new bootstrap estimator for $\boldsymbol{\beta}_k(t)$, which we called $\widehat{\boldsymbol{\beta}}_k^{*1}(t)$. This procedure was repeated $B$ times in order to produce $B$ different bootstrap datasets, $\boldsymbol{X}^{*1}$, $\boldsymbol{X}^{*2}$, $\ldots$, $\boldsymbol{X}^{*B}$, and $B$ corresponding $\boldsymbol{\beta}_k(t)$ estimators, $\widehat{\boldsymbol{\beta}}_k^{*1}(t)$, $\widehat{\boldsymbol{\beta}}_k^{*2}(t)$, $\ldots$, $\widehat{\boldsymbol{\beta}}_k^{*B}(t)$. Next we computed the sample variance of these bootstrap estimators and treated it as the estimated variance. Similar procedures were also applicable to $\widehat{\sigma}_{kl}(t)$. Afterwards, 95% confidence intervals of each parameter were constructed. Finally, we counted how many times true parameters $\boldsymbol{\beta}_k(t)$ and $\sigma_{kl}(t)$ fell in their confidence intervals to obtain coverage probabilities.

Table 2.1 and Table 2.2 summarize the main results over 100 simulations at $t = 1$. From Tables 2.1 and 2.2, we can conclude that, at $t = 1$, our method yields estimators $\widehat{\boldsymbol{\beta}}_k(t)$ which are close to the true parameters. All the estimators deviate from true parameters by less than 0.03. On the other hand, the absolute values of biases between estimators $\widehat{\sigma}_{kl}(t)$ and true parameters

15

become a little greater, but most of them are still less than 0.1. In addition, the bootstrap based standard errors are reasonable estimators for the standard deviations of $\widehat{\boldsymbol{\beta}}_k(t)$ and $\widehat{\sigma}_{kl}(t)$. Almost all the differences between SD and SE are smaller than 0.03, except for $\widehat{\sigma}_{34}(t)$. Also, excluding $\widehat{\sigma}_{13}(t)$, all the coverage probabilities are greater than or equal to 0.9, and the majority of them are around 0.95.

Table 2.1: Summary statistics for $\boldsymbol{\beta}_k(t)$ at $t = 1$ based on 100 simulations.

| Marker | Parameter | True value | Bias | SD | SE | CP |
|---|---|---|---|---|---|---|
| $Y_1$ | $\beta_{10}$ | -0.565 | 0.002 | 0.035 | 0.039 | 0.98 |
| Continuous | $\beta_{11}$ | 0.933 | 0.001 | 0.059 | 0.067 | 0.98 |
| | $\beta_{12}$ | 0.000 | -0.002 | 0.085 | 0.078 | 0.94 |
| $Y_2$ | $\beta_{20}$ | -0.819 | 0.007 | 0.050 | 0.058 | 0.98 |
| Count | $\beta_{21}$ | -1.030 | 0.026 | 0.112 | 0.112 | 0.95 |
| | $\beta_{22}$ | 0.900 | -0.010 | 0.117 | 0.132 | 0.97 |
| $Y_3$ | $\beta_{30}$ | 0.450 | -0.006 | 0.074 | 0.077 | 0.94 |
| Binary | $\beta_{31}$ | -1.000 | 0.013 | 0.112 | 0.136 | 0.99 |
| | $\beta_{32}$ | 0.900 | 0.011 | 0.157 | 0.151 | 0.93 |
| $Y_4$ | $\beta_{40}$ | -1.260 | -0.006 | 0.038 | 0.039 | 0.93 |
| Continuous | $\beta_{41}$ | 0.982 | -0.010 | 0.063 | 0.068 | 0.97 |
| | $\beta_{42}$ | 0.765 | -0.005 | 0.078 | 0.077 | 0.93 |
| $Y_5$ | $\beta_{50}$ | 0.732 | 0.001 | 0.077 | 0.074 | 0.95 |
| Binary | $\beta_{51}$ | 0.470 | 0.001 | 0.149 | 0.134 | 0.92 |
| | $\beta_{52}$ | 0.315 | -0.014 | 0.163 | 0.150 | 0.92 |
| $Y_6$ | $\beta_{60}$ | 0.331 | -0.018 | 0.085 | 0.077 | 0.90 |
| Binary | $\beta_{61}$ | 0.100 | 0.004 | 0.144 | 0.136 | 0.95 |
| | $\beta_{62}$ | 1.924 | -0.021 | 0.163 | 0.156 | 0.94 |

*Note*: "Bias" is the bias of the average estimates; "SD" is the sample standard deviation of the estimates; "SE" is the average of the estimated standard errors based on 100 bootstrap samples; "CP" is the coverage probability of the 95% confidence intervals.

After examining the estimators at a fixed time point, we also investigated the estimation performance as time changes. For instance, Figure 2.2 presents true parameters vs estimators across the 13 time points for $\beta_{52}(t)$ and $\sigma_{34}(t)$, respectively. From Figure 2.2, we can conclude $\widehat{\beta}_{52}(t)$ is very close to the true parameter at each time point, and it well captures the underlying smooth function of $\beta_{52}(t)$ across time. Although the bias between $\sigma_{34}(t)$ and $\widehat{\sigma}_{34}(t)$ is greater than that between $\beta_{52}(t)$ and $\widehat{\beta}_{52}(t)$, all of $\sigma_{34}(t)$ are in the interquartile range of $\widehat{\sigma}_{34}(t)$. Thus, the estimators perform consistently and the deviations are reasonable.

Figure 2.2: Top panel: true $\beta_{52}(t)$ versus $\widehat{\beta}_{52}(t)$ across 13 time points based on 100 simulations. Bottom panel: true $\sigma_{34}(t)$ versus $\widehat{\sigma}_{34}(t)$ across 13 time points based on 100 simulations. Red triangles: true values of the parameter. Blue triangles: average estimators of the parameter. Red curve: the true function of the parameter.

Table 2.2: Summary statistics for $\sigma_{kl}(t)$ at $t = 1$ based on 100 simulations.

| Parameter | True value | Bias | SD | SE | CP |
|---|---|---|---|---|---|
| $\sigma_{12}$ | 0.342 | -0.061 | 0.149 | 0.136 | 0.91 |
| $\sigma_{13}$ | 0.484 | -0.058 | 0.202 | 0.213 | 0.98 |
| $\sigma_{14}$ | 0.578 | -0.086 | 0.127 | 0.121 | 0.87 |
| $\sigma_{15}$ | 0.034 | 0.030 | 0.216 | 0.218 | 0.95 |
| $\sigma_{16}$ | 0.047 | -0.009 | 0.210 | 0.207 | 0.96 |
| $\sigma_{23}$ | 0.799 | -0.150 | 0.388 | 0.382 | 0.90 |
| $\sigma_{24}$ | -0.493 | 0.065 | 0.232 | 0.233 | 0.95 |
| $\sigma_{25}$ | -0.779 | 0.078 | 0.241 | 0.242 | 0.94 |
| $\sigma_{26}$ | 0.796 | -0.143 | 0.371 | 0.366 | 0.91 |
| $\sigma_{34}$ | -0.163 | 0.048 | 0.216 | 0.252 | 0.97 |
| $\sigma_{35}$ | -0.363 | -0.024 | 0.252 | 0.261 | 0.97 |
| $\sigma_{36}$ | 0.530 | -0.024 | 0.257 | 0.249 | 0.95 |
| $\sigma_{45}$ | 0.802 | -0.076 | 0.212 | 0.219 | 0.94 |
| $\sigma_{46}$ | -0.686 | 0.089 | 0.228 | 0.244 | 0.94 |
| $\sigma_{56}$ | -0.846 | -0.019 | 0.160 | 0.181 | 0.97 |

*Note*: "Bias" is the bias of the average estimates; "SD" is the sample standard deviation of the estimates; "SE" is the average of the estimated standard errors based on 100 bootstrap samples; "CP" is the coverage probability of the 95% confidence intervals.

## 2.5 Real Data Application

### 2.5.1 Data Prepocessing

We applied the proposed method to analyze EHRs of T2D patients from the OSU-WMCIW. In our application, we included three baseline variables $\boldsymbol{X}_i$: baseline age, race (1: white; 0: non-white), and sex (1: male; 0: female). Besides, there were five health markers $Y_{ik}(t)$ related to T2D: HBP, TC, HbA1c, HDL, and medications prescribed at each clinical encounter. Here, we dichotomized HBP as HBP=1 if a patient's systolic blood pressure is higher than 140 mmHg and 0, otherwise. The medications served as one strong indicator of patient's comorbidity and they could be T2D related or not. Thus, the health markers in the analysis consisted of three continuous markers (TC, HbA1c, HDL), one binary marker (HBP) and one count marker (number of medications).

For analysis, we split the data into three parts for different purposes. The first data consisted of the records collected between 2011 and 2012 and was used to estimate the variances of individual latent processes by fitting univariate generalized linear mixed models. The second part included the records from 24,975 patients between 2013 and 2017 who had at least one marker measurement. This part of the data was used for training our models and learning latent groups among the

patients. The third part was the data collected in 2018 and would be used for validation purpose. The flow-chart for this application is illustrated in Figure 2.3.



Figure 2.3: Flow-chart of the proposed analysis framework of EHRs to dissect patient heterogeneity using a diverse set of health markers.

In our model fitting using the second part of the data, after checking normal ranges for the health markers (Stone et al., 2014; Whelton et al., 2018; American Diabetes Association, 2018), we removed extreme records such as TC $\leq 0$ or $\geq 500$ mg/dL, HbA1c $\leq 3$ or $\geq 20\%$, and HDL $\leq 0$ or $\geq 120$ mg/dL. This led to a deletion of 1% of the data and a total number of 24,655 patients for analysis. Among these patients, 52.08% were female, 63.42% were white, and their ages in years ranged from 18.30 to 97.67 with a mean of 56.06. All of them had at least one observation for at least one health marker in the 5 years, but not necessarily for other health markers. Specifically, the average numbers of records for HBP, TC, HbA1c, HDL, and the number of medications per patient during these 5 years were 17.50, 4.01, 5.95, 3.64, and 53.21, respectively. In order to minimize the influence of different scales on the numeric stability, we normalized all continuous variables before identifying patient subgroups. Each of them has zero mean and unit variance.

### 2.5.2 Results

Table 2.3 shows the effect of each demographic variable on the pattern of the measurement times for each marker. From Table 2.3, we conclude that elder patients tend to have more observations for all health markers and females appeared to have more observations for HBP, HbA1c, and the

19

number of medications, while males tend to have more TC measurements. Finally, whites have significantly less observations for HBP, HbA1c, and the number of medications than non-whites.

Table 2.3: Effects of demographic variables on the frequency of health marker measurements

| Marker | Demographic | Est | HR | SE | Z | P-value |
|---|---|---|---|---|---|---|
| HBP | age | 0.065 | 1.067 | 0.006 | 10.552 | $< 0.001$ |
| | sex | 0.064 | 1.066 | 0.013 | 4.813 | $< 0.001$ |
| | race | -0.129 | 0.879 | 0.014 | -9.425 | $< 0.001$ |
| TC | age | 0.035 | 1.035 | 0.006 | 6.238 | $< 0.001$ |
| | sex | -0.035 | 0.965 | 0.012 | -3.000 | 0.003 |
| | race | -0.012 | 0.988 | 0.013 | -0.968 | 0.333 |
| HbA1c | age | 0.008 | 1.008 | 0.005 | 1.721 | 0.085 |
| | sex | 0.034 | 1.034 | 0.009 | 3.678 | $< 0.001$ |
| | race | -0.044 | 0.957 | 0.009 | -4.650 | $< 0.001$ |
| HDL | age | 0.047 | 1.048 | 0.005 | 10.090 | $< 0.001$ |
| | sex | -0.010 | 0.990 | 0.010 | -1.007 | 0.314 |
| | race | -0.007 | 0.993 | 0.010 | -0.715 | 0.475 |
| Medications | age | 0.042 | 1.043 | 0.006 | 7.262 | $< 0.001$ |
| | sex | 0.086 | 1.090 | 0.012 | 7.069 | $< 0.001$ |
| | race | -0.113 | 0.893 | 0.013 | -8.988 | $< 0.001$ |

*Note*: "Est" is the regression coefficient estimator; "HR" is the hazard ratio; "SE" is the standard error of the coefficient estimator; "Z" is the statistic for a z-test; "P-value" is the p-value for the z-test.

To estimate the parameters in the joint models, we first implemented the adaptive method of bandwidth selection as stated in Section 2.2.4, and results are shown in Figure 2.4. We ended up to choose $h_{1n} = 564.112$ days and $h_{2n} = 494.687$ days as the optimal bandwidths. Using the optimal bandwidths, we estimated $\beta_k(t)$ and $\sigma_{kl}(t)$ at 61 time points. The results are presented in Figure 2.5 and Figure 2.6, respectively. The salmon-colored ribbons in these two figures are 95% confidence intervals for the parameters based on 100 bootstrap datasets.

Figure 2.5 presents the relationships between each pair of health markers and covariates. In general, all health markers exhibit changes over time. Mean HbA1c $(\widehat{\beta}_{30}(t))$ decreases during the first 1.5 years and has an increasing trend afterward, which may suggest the difficulty to achieve long-term control of glycemic levels in a chronically ill patient population. Mean HDL $(\widehat{\beta}_{40}(t))$ shows a similar quadratic pattern over time, suggesting difficulty of long-term TC control. The estimated regression coefficients for covariates, i.e., the estimated effects of covariates on health markers, do not show any pattern of drastic changes over time. Instead, the estimated values across time fluctuate

Figure 2.4: Bandwidth selection results for the real data application. (a): the plot for $\sum_t \mathrm{MSE}_\beta^{h_1}(t)$ vs. $h_1$. The optimal $h_{1n} = 564.112$ days. (b): the plot for $\sum_t \mathrm{MSE}_\sigma^{h_2}(t)$ vs. $h_2$. The optimal $h_{2n} = 494.687$ days. Red triangles: optimal bandwidths.

around mean values. However, we can observe decreasing trends for $\widehat{\beta}_{20}(t)$ and $\widehat{\beta}_{50}(t)$, suggesting that as time increases, the expected means of TC and the number of medications decrease. $\widehat{\beta}_{11}(t)$ and $\widehat{\beta}_{41}(t)$ are positive across time, while $\widehat{\beta}_{21}(t)$ and $\widehat{\beta}_{31}(t)$ are negative. $\widehat{\beta}_{51}(t)$ is negative but close to 0. Hence, estimators $\widehat{\beta}_{\cdot1}(t)$ suggest that elder subjects on average have higher HBP and HDL, but they have lower TC and HbA1c. There is no apparent difference in the average number of medications between elder subjects and younger subjects. Similarly, estimators of sex effect, $\widehat{\beta}_{\cdot2}(t)$, suggest that compared with men, women tend to have higher expected means of TC and HDL, but they have lower values of HBP and the number of medications. Although women have slightly lower expected means of HbA1c than men, the difference is inapparent. For race, the estimators of $\widehat{\beta}_{\cdot3}(t)$ indicate that white people have lower or equal expected means than non-white people in almost all five health markers.

Figure 2.6 presents the correlations between each pair of health markers. The results suggest the concurrent correlations between HBP and TC, HBP and medications, TC and HbA1c, TC and HDL are positive and moderate. Moreover, there exist negative and observable concurrent correlations between HbA1c and HDL, HDL and medications. The correlation between HbA1c

Figure 2.5: Estimated regression coefficients $\widehat{\boldsymbol{\beta}}_k(t)$ across 61 time points using $h_{1n} = 564.112$ days and $h_{2n} = 494.687$ days from EHRs at the OSU-WMCIW. Salmon-colored ribbons: 95% confidence intervals for the estimators.

Figure 2.6: Estimated correlations $\hat{\sigma}_{kl}(t)$ across 61 time points using $h_{1n} = 564.112$ days and $h_{2n} = 494.687$ days from EHRs at the OSU-WMCIW. Salmon-colored ribbons: 95% confidence intervals for the estimators.

and HDL decreases as time increases. On the opposite, the positive correlation between TC and HDL decreases at the beginning, but increases after about 1 year. The positive correlation between TC and HbA1c has a similar pattern as it decreases at first and increases after 1000 days. The correlations of HBP and TC, HbA1c and number of medications increase in first 500 days, but they start to decrease during 500 to 1000 days, and bounce back afterward. The correlations of HBP and HbA1c, HBP and HDL, HDL and number of medications decrease in first 500 days, and then increase, but decrease again after 1000 days.

One interesting observation from Figure 2.6 is that the estimated correlation between the number of medications and HBP is as high as 0.6 but its correlations with TC and HDL are both negative, fluctuating around -0.30. However, there does not appear to be a strong association between the number of medications and HbAc1 over time. This may suggest that the patients in this cohort were most likely to take medications that aimed to control the levels of TC and HDL, but not necessarily for controlling the level of HbA1c. The latter is consistent with the fact that over 90% of drugs recorded in this database are non-diabetic drugs. One possible interpretation of the observed time-dependent correlation pattern is that there might exists another unobserved disease health marker that influences the two observed markers temporally. Thus, the estimated correlation pattern could be potentially useful to identify such "common cause" health markers so as to better understand the mechanism of disease progression.

Finally, we computed the similarity between each pair of patients using the distance defined in (2.7). To compute $\widehat{\epsilon}_i(t)$ as (2.8), we substituted $\widehat{Y}_i(t)$ with the nearest neighbor observation of time $t$ for patient $i$. Using the between-patient similarity matrix, we performed a cluster analysis on the 24,655 patients, and the results are given in Figure 2.7. We observed 4 clusters within which patients had similar health marker profiles.

To better understand the health patterns of patients in each subgroup, we calculated the average of normalized measurements for each health marker in each group, as shown in Figure 2.8. In the top panel of Figure 2.8, the value in each cell is averaged over all patients and all clinical encounters between January 1, 2013 and December 31, 2017. We compared these values to the average of each health marker in the entire study sample. A higher value of HDL and a lower value of HBP, TC, and HbA1c represent healthier T2D status. The number of medications prescribed at each clinical encounter does not directly reflect the disease status, but a lower count usually indicates a less

Figure 2.7: Dendrogram of Mahalanobis distances for 24,655 patients at the OSU-WMCIW. Group index numbers are assigned according to group sizes.

severe state. Group 4 contains 2,163 patients, whose TC was slightly higher than the overall average. Their HBP, HDL, and the number of medications were lower than the overall averages. In addition, they had the highest HDL and it was substantially higher than the overall average. Thus, group 4 is the relatively healthy group in which patients did not take many medications. Group 1 contains 10,705 patients who were less healthy since they had lower-than-average HDL, but other health markers were favorable or roughly neutral. The TC of 6,930 patients in group 2 was higher than the overall average, while other health markers were lower or around the averages. We conclude that group 2 is a moderately ill group. For the 4,857 patients in group 3, their TC levels were slightly lower than the overall average, however, they had the highest HbA1c. Also, other markers indicated bad health status. Therefore, group 3 patients were in the most severe state of T2D.

To examine whether the subgroups inferred by the clustering truly represent patients health profiles, we validated the detected patterns using the third fold of the split data that consisted of the EHR data collected after January 1, 2018. These data were not used in any other analyses of this application. The average values of normalized measurements for each health marker in each group are shown in the bottom panel of Figure 2.8. We conclude that the patients health patterns identified prior to year 2018 are consistent with those patterns afterward. Therefore, the patient groups are not only meaningful, but also represent some true underlying patient patterns over time. This robustness is particularly important to the long-term health management of T2D patients.

| | group 1 | group 2 | group 3 | group 4 |
|---|---|---|---|---|
| Number of medications | 0.111 | −0.196 | 0.185 | −0.336 |
| High−density lipoprotein | −0.274 | 0.021 | −0.238 | 1.823 |
| Glycated hemoglobin | −0.402 | −0.176 | 1.221 | −0.19 |
| Total cholesterol | −0.49 | 0.832 | −0.18 | 0.161 |
| High blood pressure | −0.026 | 0.036 | 0.147 | −0.317 |

Severe

Healthy

(b)

| | group 1 | group 2 | group 3 | group 4 |
|---|---|---|---|---|
| Number of medications | 0.005 | −0.03 | 0.114 | −0.149 |
| High−density lipoprotein | −0.201 | −0.078 | −0.242 | 1.514 |
| Glycated hemoglobin | −0.23 | −0.146 | 0.794 | −0.139 |
| Total cholesterol | −0.318 | 0.479 | −0.139 | 0.165 |
| High blood pressure | −0.021 | 0.043 | 0.069 | −0.17 |

Severe

Healthy

Figure 2.8: Averages of normalized measurements by health markers and patient subgroups. (a): using data from 1/1/2013 to 12/31/2017. (b): using data after 1/1/2018. Red: more severe status than the overall sample average in terms of a health marker; blue: healthier status than the overall sample average in terms of a health marker; white: overall sample average status in terms of a health marker.

## 2.6 Discussion

In this chapter, we proposed a latent temporal process model to integrate health markers in EHRs and characterize patient heterogeneities. The proposed method is capable of handling unbalanced records and informative visits, that is, patients can have missing health markers at some encounters or with visit times depending on their health status. Additionally, our model can both fit different types of health marker, capture the dependence structures among health markers, and takes into account informative patterns of visit times, via the intensity function of health markers. The real data application shows the capability of the proposed method on addressing the data challenges of EHRs, integrating different types of health markers, and identifying meaningful and robust patient subgroups. Therefore, the proposed method may shed lights on the detection of patient homogeneities and heterogeneities, and serve as a step towards applications of personalized medicine.

In the parameter estimation process, we assumed that variances of the latent variables $\epsilon_i(t)$ were fixed and they were estimated using the EHR data of 2011 and 2012. To study whether the constant variance was reasonable, we estimated the changes in variances from six different time periods in windows of 2 years as well as using the whole 5-year data, and the results, as shown in Tables A.1 and A.2, indicate that the estimates varied little. Thus, the constant variance assumption seems to be reasonable for our application. In addition, we re-estimated $\beta_k(t)$ and $\sigma_{kl}(t)$ using the same proposed parameter estimation approach but with the 5-year variance estimates in Table A.2. Figures A.1 and A.2 reveal slight changes in the estimated coefficients. In fact, the absolute percentage changes between the two sets of coefficients are less than 1%, except for $\widetilde{\beta}_{1.}(t)$ and $\widetilde{\sigma}_{1.}(t)$ which have changes of up to 3%. Therefore, we could conclude that the estimation results are robust to the constant variance estimates.

Moreover, to investigate the effect of bandwidth selection on parameter estimation, we report $\beta_k(t)$ using two suboptimal bandwidths that are close to the optimal bandwidth in Section 2.5.2. Figure A.3 shows that the suboptimal estimators preserve the similar pattern to $\widehat{\beta}_k(t)$. The Canberra distances (Lance and Williams, 1966) between the optimal estimators and suboptimal estimators of

$\beta_{kj}(t)$ across time, $k = 1, \ldots, p$, $j = 0, \ldots, m$, are calculated as

$$d(\boldsymbol{\beta}_{kj,H_1}, \boldsymbol{\beta}_{kj,H_1'}) = \frac{1}{N} \sum_{t=t_1}^{t_N} \frac{\left| \boldsymbol{\beta}_{kj,H_1}(t) - \boldsymbol{\beta}_{kj,H_1'}(t) \right|}{|\boldsymbol{\beta}_{kj,H_1}(t)| + \left| \boldsymbol{\beta}_{kj,H_1'}(t) \right|}, \tag{2.14}$$

where $\boldsymbol{\beta}_{kj,H_1}$ is the vector of estimated $\{\beta_{kj}(t) : t = t_1, \ldots, t_N\}$ using the optimal bandwidth $H_1$ and $\boldsymbol{\beta}_{kj,H_1'}$ is the vector of these estimators using a suboptimal bandwidth $H_1'$. Most of the distances are as smaller than 0.05 as given in Table A.3, confirming the estimates using the optimal and suboptimal bandwidths are close. The conclusions could also be drawn for estimating $\sigma_{kl}(t)$ (cf. Figure A.4 and Table A.4).

In our models, we assumed that the intensity function of the counting process only depended on the baseline covariates. This assumption can be violated if the intensity also depends on the historical marker values. However, directly incorporating time-dependent marker values, which are missing for most of time points, is challenging. To examine how this assumption may affect our results, we included an ad hoc marker value, defined as the mean value of HbA1c in the past 12 months, in the intensity model (2.1). From Table A.5, the effects of the historical HbA1c level on frequencies of HBP, TC, and HbA1c are significant, while the historical HbA1c level has lower impacts on frequencies of HDL and the number of medications. Figures A.5 and A.6 also reflect this phenomenon that there are slight differences between two versions of estimators for HDL and the number of medications. Although differences between two versions of estimators for HBP and TC are moderate, the new estimators still locate within or around the 95% bootstrapped confidence intervals for the original estimators. However, for HbA1c, the differences could not be ignored since the estimated curves present some unusual shapes. Therefore, further investigation is needed regarding what time-dependent marker values should be used and how missing data issues should be addressed.

As stated in Section 2.1, the latent processes can be also viewed as projections of the health markers onto a lower dimensional space. Therefore, our method can be used for identifying latent clusters among patients as illustrated in our application, and at the same time can also play a role in learning personalized disease prognosis and personalized disease management. For example, the summary of latent processes can be used to improve the understanding of treatment propensity scores in EHRs when learning individualized treatment rules. Lastly, the latent processes can be

28

included in disease outcome models as prognostic or predictive health markers, and we will show this extension in Chapter 3.

# CHAPTER 3: ESTIMATING INDIVIDUALIZED TREATMENT RULES FOR MULTICATEGORY TYPE 2 DIABETES TREATMENTS USING ELECTRONIC HEALTH RECORDS

## 3.1 Introduction

T2D is the most common type of diabetes which causes millions of people to suffer from severe diabetes-related complications such as heart attacks, stroke, blindness, and kidney failure (Roglic, 2016). To treat T2D, American Diabetes Association (2018) recommended to use metformin monotherapy as the initial treatment and select additional therapies based on patient-centered considerations. A treatment guideline from the United Kingdom also suggested metformin as the first-line drug, unless it is contraindicated or not tolerated (McGuire et al., 2016). Palmer et al. (2016) summarized 301 clinical trials (1.4 million patient-month) in which metformin and other 8 available classes of glucose-lowering drugs were compared. This meta analysis reported that when compared with other drugs given as monotherapy, metformin only had better or similar effects on managing HbA1c levels among adults with T2D. Nevertheless, there was no significant difference in all-cause mortality or other complications between any glucose-lowering drugs alone or combined. There is a lack of conclusive evidence for the best T2D management strategy from clinical trials.

With the emergence of large-scale electronic systems such as EHRs, which usually contain patient demographics, vital signs, laboratory test results, medications, diagnosis, and medical insurances documented at the point of care, there has been an increasing trend of using EHRs as an observational database to study T2D treatment patterns in real world practices. For example, Montvida et al. (2018) selected 1.02 million adults with T2D from the U.S. Centricity Electronic Medical Records and concluded that, from 2005 to 2016, first-line use increased for metformin (60% to 77%) and decreased for sulfonylureas (20% to 8%). Canivell et al. (2019) used a 5-year-EHR for 15,205 patients with T2D from the SIDIAP database and assessed glycemic controls after treatment intensification. Compared to experiments such as RCTs, observational studies use real-world information from larger patient populations and contain a longer duration of observations, and thus may offer valuable complements to RCTs. Studies using large-scale EHRs may reflect real-world patterns of treatment

pathways which can neither be conducted nor observed in RCTs (Hripcsak et al., 2016). More importantly, EHRs provide a great opportunity to study the heterogeneity of treatment responses in a large population so that we can learn optimal ITRs for T2D patients to fulfill the goal of precision medicine, a medical paradigm that utilizes individual patient's characteristics such as demographics, lab test results, and genetic information, to optimize treatments (Ginsburg and Phillips, 2018).

There has been intensive methods development for precision medicine in the fields of statistics and machine learning over the last decade (Mesko, 2017). These methods include regression model-based methods such as Q-learning (Watkins and Dayan, 1992; Murphy, 2005; Qian and Murphy, 2011), A-learning (Murphy, 2003; Robins, 2004), regret-regression (Henderson et al., 2010), and subgroup analysis (Foster et al., 2011; Lipkovich et al., 2011; Fu et al., 2016). Through directly optimizing ITR-related value functions, Zhao et al. (2012) proposed an outcome weighted learning approach that converted the estimation of ITRs to a weighted classification problem. Similar methods were later developed in contrast weighted learning (Tao and Wang, 2017) and augmented outcome weighted learning (Liu et al., 2018). More recently, Wu et al. (2020) proposed a matched learning approach, called M-learning, to learn ITRs based on pairs of patients who shared similar pre-treatment health profiles. This approach was demonstrated to be more robust than weighting methods.

However, the above methods are confronted by the following challenges when applied to EHRs. First, characterizing individual patient's pretreatment condition is difficult since their health markers measured over time are multivariate and the measurements can be continuous (e.g., lab measures), binary (e.g., disease diagnoses) or counts (e.g., number of medications). Moreover, these measurements are taken at patient's clinical encounters which potentially depend on their underlying health status. Thus, not accounting for informative measurement patterns of health markers may cause selection bias (Haneuse, 2016; Haneuse and Daniels, 2016). Second, there are often many observed treatment options and patient's propensity to receive one specific treatment is complex and heterogeneous, which may not be captured by parametric models. Furthermore, there presents substantial heterogeneity among patients in terms of treatments and outcomes that need to be accounted for when learning ITRs. Standard weighting methods suffer from numerical instability due to low representation of patients with some treatments.

In this chapter, to address the challenges in EHRs, we propose a general framework for learning

ITRs for T2D patients and use one concrete dataset as an example to demonstrate the framework. Specifically, we propose a multivariate longitudinal model to model the time-trajectory of different types of health markers through a generalized exponential family of distributions, while accounting for their dependence through a latent multivariate Gaussian temporal process. We also adopt inverse intensity weighting to adjust for potential informative times of measurements. Through the joint models, we can identify several T2D patient subgroups using clustering algorithm to summarize patients' health profiles based on their pre-treatment EHRs. To learn ITRs within each subgroup, we create a few classes of treatments and apply nonparametric methods to estimate treatment propensity scores. Finally, to handle the challenge of multiple treatments, we extend matched learning method in Wu et al. (2020) to multicategory treatments. Particularly, we develop an one-versus-one matched learning method to estimate ITRs. The derived rules are further validated through cross-validation.

The remaining part of this chapter is organized as follows. In Section 3.2, we provide the details of the proposed models and learning methods. In Section 3.3, we demonstrate an implementation of our methods to EHRs from the OSU-WMCIW. In Section 3.4, we describe the estimated ITRs for T2D patients and compare with the observed treatments in EHRs. Concluding remarks are given in Section 3.5.

## 3.2 A General Framework to Learn Optimal ITRs Using EHRs

We use $A$, $\boldsymbol{Z}$, and $R$ to denote a T2D patient's treatment at a decision time (referred to as time zero), pre-treatment features, and reward outcome, respectively. We assume no unobserved confounding and stable unit treatment value assumption, which are two crucial assumptions to allow using the EHRs for learning the optimal treatment rules. The first assumption implies that the treatment assignment is independent of potential outcomes given $\boldsymbol{Z}$, so there will not be any hidden bias due to unobserved confounding; while the second assumption implies that there is no treatment interference between the patients. The assumptions are not testable due to the observational nature of the EHRs but may be plausible if $\boldsymbol{Z}$ contains sufficient information about why each patient received one particular treatment and one patient's response does not depend the other patients' treatments or responses. Under these assumptions, it is known that the optimal ITR is a function mapping $\boldsymbol{Z}$ to $A$'s domain and it is given as the treatment that yields the maximum value of $\mathbb{E}\left(R|\boldsymbol{Z}, A = a\right)$. Many methods have been developed to estimate such optimal ITR using RCTs,

but our goal is to instead use EHRs to estimate the optimal ITR.

Data from EHRs consist of patient's health marker measurements, for example, body mass index (BMI), cholesterol level, and HbA1c for T2D patients, as well as received medications, at clinical encounters over a span of calendar time windows. Time zero is usually set to be the index date when a patient received treatment $A$, and the reward outcome, $R$, is a pre-defined measure indicating disease improvement since time zero (for example, HbA1c reduction within 6 months after taking the treatment). However, obtaining a reasonable set of feature variables for $\boldsymbol{Z}$ is challenging, since they not only include patient's demographics (age, gender, race), but more importantly, should reflect patient's preconditions that are useful for the treatment decision. The latter must be extracted from patient's longitudinal health markers before time zero.

In the following sections, we first extend the method in Chapter 2 to extract patient's pre-treatment health profiles using EHRs that will be included as feature variables for learning ITRs. We then propose a matching-based learning algorithm to estimate optimal treatment rules that will maximize patient's outcomes.

### 3.2.1 Characterizing Patient's Pre-treatment Health Conditions

Given the heterogeneity among patients in EHRs, it is important to characterize patient's pre-treatment conditions based on longitudinal marker measurements in EHRs. We present methods to handle two challenges: the first challenge is that patterns of measurement time points may depend on patients' underlying health status and thus are informative; the second challenge is that health markers are of mixed types and collected at different time points.

We use the same notations of $n$, $p$, $\boldsymbol{X}_i$, $Y_{ik}(t)$, and $N_{ik}(t)$ as those in Section 2.2.1. However, in this framework, we extend the equation (2.1) to

$$\mathbb{E}\left[dN_{ik}(t)|\text{observed data up to }t\right] = \lambda_k(t)\exp\left\{\boldsymbol{X}_i^T\boldsymbol{\gamma}_k + \boldsymbol{L}_{ik}^T(t)\boldsymbol{\eta}_k\right\}dt, \quad (3.1)$$

where $\lambda_k(t)$ is a baseline intensity function, $\boldsymbol{L}_{ik}(t)$ is a vector of observed health history up to time $t$, and $\boldsymbol{\gamma}_k$ and $\boldsymbol{\eta}_k$ are intensity parameters. For example, in real data, $\boldsymbol{L}_{ik}(t)$ can take the value of average blood pressure measurements in past 3 months, and a patient with higher blood pressures is likely to have a revisit in a shorter time than a patient with normal blood pressures. This effect is measured by the intensity parameter $\boldsymbol{\eta}_k$. We still assume $Y_{ik}(t)$ follows a generalized exponential

family model as (2.2). Also, the regression coefficients $\boldsymbol{\beta}_k(t)$ and the covariance matrix $\boldsymbol{\Omega}(t)$ are defined and modeled as (2.3).

The parameters in the intensity model for $N_{ik}(t)$ can be estimated by fitting a standard Andersen-Gill proportional intensity model. We denote the coefficient estimators as $\widehat{\boldsymbol{\gamma}}_k$, and $\widehat{\boldsymbol{\eta}}_k$. To estimate the parameters in model (2.3) and further account for irregular time intervals, we define $d\widetilde{N}_{ik}(s) = dN_{ik}(s)/\exp\left\{\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k + \boldsymbol{L}_{ik}^T(s)\widehat{\boldsymbol{\eta}}_k\right\}$, and then we solve the following kernel-weighted local estimating equations for each health marker to pool information across times and patients:

$$U_{n,k}(\boldsymbol{\beta}_k(t)) = \frac{1}{n}\sum_{i=1}^n \int K_{h_{1n}}(s-t)\boldsymbol{X}_i[Y_{ik}(s) - E_{ik}(t)]d\widetilde{N}_{ik}(s) = \boldsymbol{0}, \tag{3.2}$$

where $E_{ik}(t) = \mathbb{E}\left[g_k^{-1}\left(\boldsymbol{X}_i^T\boldsymbol{\beta}_k(t) + \epsilon_{ik}(t)\right)\Big|\boldsymbol{X}_i\right]$, $K_h(u) = K(u/h)/h$ with $K(u)$ being a symmetric kernel function, usually taken to be the Epanechnikov kernel or Gaussian kernel, and $h$ is its bandwidth.

Denote the estimators obtained by solving (3.2) as $\widehat{\boldsymbol{\beta}}_k(t)$. Similarly, for each element in $\boldsymbol{\Omega}(t) = (\sigma_{kl}(t))$, we solve the following kernel-weighted local estimating equations

$$U_{n,k,l}(\sigma_{kl}(t)) = \frac{1}{n^2}\iint \widetilde{K}_{h_{2n}}(s-t,s'-t)\left[Y_{ik}(s)Y_{il}(s') - E_{ikl}(t)\right]d\widetilde{N}_{ik}(s)d\widetilde{N}_{il}(s') = \boldsymbol{0}, \tag{3.3}$$

where the double integration excludes $s = s'$ if $k = l$. In (3.3),

$$E_{ikl}(t) = \mathbb{E}\left[g_k^{-1}\left(\boldsymbol{X}_i^T\widehat{\boldsymbol{\beta}}_k(t) + \epsilon_{ik}(t)\right)g_l^{-1}\left(\boldsymbol{X}_i^T\widehat{\boldsymbol{\beta}}_l(t) + \epsilon_{il}(t)\right)\Big|\boldsymbol{X}_i\right],$$

and $\widetilde{K}_h(u_1, u_2) = \widetilde{K}(u_1/h, u_2/h)/h^2$. Here, $\widetilde{K}(u_1, u_2)$ is a bivariate kernel function, usually taken to be the product of univariate Epanechnikov or Gaussian kernel, and $h$ is its bandwidth. Denote the estimators as $\widehat{\sigma}_{kl}(t)$. Unlike (2.6), in this case, we estimate the diagonal elements of $\boldsymbol{\Omega}(t)$ directly from the data instead of historical records or external resources. The bandwidths in the estimating equations (3.2) and (3.3) are determined using the same approach in (2.9). In Appendix B.1, we provide numerical evidence through a simulation study to demonstrate the good performance of the proposed estimation method.

Finally, to characterize patient's pre-treatment health status into clusters where patients within

the same cluster share similar health profiles, we compute the similarity distance between each pair of patients as (2.7), and then we perform a hierarchical clustering based on the distance matrix. Because the Mahalanobis distance naturally accounts for the between-marker correlation, the use of such distance can effectively remove the redundant information regarding the patient's health status. We choose the hierarchical clustering because it is a powerful approach to identify homogeneous, interpretable groups of the patients. Thus, even though the original heath markers are measured irregularly and are of very different data types, our joint models enable one to combine them using the latent processes on the same scales and account for dependence over time. With the estimated subgroups from clustering, the subsequent ITRs will be estimated for each subgroup separately.

### 3.2.2 Matched Learning for Multicategory Treatments

When estimating ITRs for binary treatments, Wu et al. (2020) showed that M-learning could outperform other commonly used methods for observational databases. Thus, we generalize M-learning to handle multicategory treatments which are commonly seen in EHRs. First, in each patient subgroup $s \in \{1, \ldots, S\}$ that was identified before, the comparison among a total of $K$ treatments can be converted to $K(K-1)/2$ comparisons between two treatments, which can be integrated using the one-versus-one method to yield an optimal ITR for all treatments.

Specifically, for each patient $i$ in each subgroup $s$, let $\boldsymbol{Z}_i$ denote the baseline covariates $\boldsymbol{X}_i$ and some additional pre-treatment health marker information, for example, the average BMI in the past year. We let $A_i$ and $R_i$ be the treatment at time zero and the reward outcome post-treatment. For each treatment pair $(u, v)$, let $T_i = 1$ if $A_i = u$ and $T_i = -1$ if $A_i = v$. Assume there are $N_{u,v}^s$ patients who received treatment $u$ or $v$ in group $s$. Antonelli et al. (2018) and Wu et al. (2020) proposed a doubly robust matching method to improve the efficiency of matching methods. This method uses not only covariates but also propensity scores and prognostic scores, denoted by $\pi(\boldsymbol{Z}_i) \equiv P(T_i = 1 | \boldsymbol{Z}_i)$ and $\psi(\boldsymbol{Z}_i) \equiv \mathbb{E}[R_i | \boldsymbol{Z}_i]$, respectively, to create matched sets. Thus, for the $i$th patient, the improved matched set, denoted by $M_{is}$, has an expression as follows:

$$M_{is} = \{j : A_j \neq A_i, d(\boldsymbol{H}_j, \boldsymbol{H}_i) \leq \delta\},$$

where $d(\cdot, \cdot)$ is a distance function, $\boldsymbol{H}_i = \left(\boldsymbol{Z}_i, \widehat{\pi}(\boldsymbol{Z}_i), \widehat{\psi}(\boldsymbol{Z}_i)\right)$, and $\delta$ is a threshold which may vary with $i$. In our implementation, we use random forests (Ho, 1995) to perform a multicategory

classification for $T_i$ given $\boldsymbol{Z}_i$ to obtain $\widehat{\pi}(\boldsymbol{Z}_i)$ and use gradient boosting machines (Friedman, 2001) to estimate $R_i$ given $\boldsymbol{Z}_i$ to obtain $\widehat{\psi}(\boldsymbol{Z}_i)$. As suggested by Antonelli et al. (2018) and Wu et al. (2020), the doubly robust method may lead to the optimal treatment rules even if the model for the propensity score, or the model for the prognostic score is misspecified, but not both, and including prognostic scores was empirically shown to perform better than the methods without using them. Following Wu et al. (2020), we adopt a weighted SVM (Cortes and Vapnik, 1995) with weights for estimating ITR when comparing treatments $u$ and $v$. Specifically, we minimize the following objective function

$$
\begin{aligned}
V_{u,v}^s(f;g) \;=\; & (N_{u,v}^s)^{-1} \sum_{i \in \text{subgroup } s, A \in \{u,v\}} |M_{is}|^{-1} \sum_{j \in M_{is}} |R_j - R_i| \\
& \times \phi\left(-f(\boldsymbol{Z}_i) T_i \text{sign}(R_j - R_i)\right) + \lambda_{u,v}^s \|f\|_{\mathcal{H}_k},
\end{aligned}
\tag{3.4}
$$

where $\phi(x)$ is the hinge loss given by $\max(1 - x, 0)$, $f(\cdot)$ is a function such that the decision rule $D(\boldsymbol{Z}) = \text{sign}(f(\boldsymbol{Z}))$, $|M_{is}|$ is the size of matched set $M_{is}$, $\lambda_{u,v}^s$ is a tuning parameter, and $\mathcal{H}_k$ is a reproducing kernel Hilbert space with a kernel function $k(\cdot, \cdot)$. Using the weight $|R_j - R_i|$ ensures that the estimated treatment rule is driven by comparing the pairs of patients who have large outcome differences. Using a weighted SVM to minimize $V_{u,v}^s(f;g)$, we obtain the optimal ITR, $D_{u,v}^{*s}$, for comparing treatment $u$ to $v$ in group $s$. Similarly, for the remaining treatment pairs, we estimate the corresponding decision rules. Therefore, for the $i$th patient in group $s$, we derive the optimal ITR, $D^{*s}(\boldsymbol{Z}_i)$, as the majority vote recommended by $\left\{D_{u,v}^{*s}(\boldsymbol{Z}_i) : u, v \in \{1, \ldots, K\}, u \neq v\right\}$.

In the above learning algorithm, the tuning parameters are chosen using cross-validation. The ITRs estimated from a training sample (i.e., $\widehat{D}^{*s}(\cdot)$) are evaluated using an independent testing sample by calculating an empirical value function defined as

$$
\frac{\sum_{u=1}^K \sum_{i \in \text{test sample}} I(A_i = \widehat{D}^{*s}(\boldsymbol{Z}_i) = u) R_i / \widehat{P}(A_i = u | \boldsymbol{Z}_i)}{\sum_{u=1}^K \sum_{i \in \text{test sample}} I(A_i = \widehat{D}^{*s}(\boldsymbol{Z}_i) = u) / \widehat{P}(A_i = u | \boldsymbol{Z}_i)},
\tag{3.5}
$$

where $\widehat{P}(A_i = u | \boldsymbol{Z}_i)$ is the estimated probability of $A_i = u$ from the propensity score estimation.

## 3.3 Implementation Using OSU-WMCIW EHRs

### 3.3.1 Data Preparation

The EHRs from OSU-WMCIW contain demographics, laboratory test measures, vital signs, and diagnosis codes for 58,490 patients diagnosed with T2D between 2011 and 2018. We set time for the treatment decision for each patient (time zero) as the last clinical encounter in year 2016 when he or she received T2D medications. This choice was based on two facts in order to learn ITRs from data with limited time periods: we needed a sufficient time window of the longitudinal history before time zero to precisely characterize patient's health status and subgroups; and we needed a reasonable follow-up period after time zero to precisely calculate the outcome variable, the HbA1c level at 6 months after time zero.

In the EHRs, there were four health markers $Y_{ik}(t)$ associated with T2D: systolic blood pressure (SBP), HbA1c, HDL, and BMI. After checking normal ranges for the health markers (Stone et al., 2014; Whelton et al., 2018; American Diabetes Association, 2018), we removed missing, duplicated, and extreme entities such as SBP$\geq$250 mmHg, HbA1c$\leq$3 or $\geq$20 %, HDL$\leq$0 or $\geq$120 mg/dL, and BMI$\leq$10 or $\geq$60 kg/m$^2$. In addition, we created a binary variable to denote whether diabetic drugs were prescribed at a clinical encounter (DD) and a continuous variable as the logarithm of the number of medications prescribed at each encounter (logMed). Both variables were considered as an important indication of T2D patient's comorbidity status. Therefore, our analysis included one binary longitudinal marker and other 5 continuous health markers over time (SBP, HbA1c, HDL, BMI, logMed). We require that at least one measurement for at least one marker is available before time zero. With this restriction, there were a total of 8,456 subjects with 497,763 longitudinal records before time zero date and they were used to learn patient's pre-treatment subgroups using the method in Section 3.2.1. Among these patients, 53.43% were female, 62.03% were white, and their ages in years ranged from 17.89 to 100.84 with a mean of 59.67. The number of records for SBP, HbA1c, HDL, BMI, and DD was 9.8, 2.3, 4.0, 2.2, 14.6, and 29.3, respectively, when averaged all the patients.

The medications at time zero were considered to be treatments, $A_i$, for learning ITRs. There were 3,978 types of medications observed in the EHRs, and we classified the medications to either the 163 diabetic drugs for T2D (Drugs.com, 2019) or the remaining non-T2D medications before

further grouping them. As stated in Section 3.1, metformin is commonly considered as the first-line drug for T2D and it may have a better control of HbA1c levels than other T2D drugs given as monotherapy. On the other hand, there is some evidence that basal insulin also serves as an important T2D treatment. Thus, we compared four classes of treatments: metformin monotherapy, insulin monotherapy, other T2D monotherapy or combinations of other T2D drugs, and combinations of at least two classes of treatment among the aforementioned three classes. We referred the third class as "other T2D drugs" and the fourth class as "multiple T2D drugs". Thus, $A_i$ was one of four treatments including metformin, insulin, other T2D drugs and multiple T2D drugs.

We were interested in the treatment effects on reducing HbA1c level after time zero. For the $i$th patient, we constructed the outcome variable $R_i$ as the expected HbA1c level 6 months after the date of time zero. In particular, we first collected all available HbA1c measures from lab tests, which were conducted from time zero to up to one year after. For each patient, we performed a linear interpolation model as $\alpha_i + \beta_i(t_{ij} - t_{i,\text{baseline}})$, where $t_{ij}$ was the date (in years) for the $j$th measurement for patient $i$, $t_{i,\text{baseline}}$ was the date for time zero, and $\alpha_i$ and $\beta_i$ were respectively the intercept and slope for the patient's trajectory. Finally, based on the least-square estimates, we defined the outcome as the expected HbA1c level at 6 month after treatment for each patient, which was given as $R_i = \widehat{\alpha}_i + 0.5\widehat{\beta}_i$. In OSU-WMCIW data, only 5,458 patients had at least two HbA1c measurements during the year after time zero so their outcomes could be calculated. Furthermore, we excluded 333 patients whose estimated slope coefficient was either greater than 5 or less than $-5$, which were not sensible clinically.

To construct the feature variables, $\boldsymbol{Z}_i$, we first used the proposed method in Section 3.2.1 to obtain subgroups of all patients using all available heath markers before time zero. More specifically, when fitting the joint models, the covariates entering the intensity models (3.1) for the measurement times included the demographics and time-dependent covariates indicating whether there was any measurement of longitudinal marker $k$ in the past 6 months and if there was, what was the average marker value. Using about 2-year data prior to time zero, we estimated the intensity parameters $\boldsymbol{\gamma}_k$ and $\boldsymbol{\eta}_k$, along with $\boldsymbol{\beta}_k(t)$ and $\boldsymbol{\Omega}(t)$ at 25 equally spaced time points (in days), where the time length between two consecutive time points were 30 days. Given $\widehat{\boldsymbol{\beta}}_k(t)$ and $\widehat{\boldsymbol{\Omega}}(t)$, in this particular application, we integrated health markers over time by calculating $\widehat{\boldsymbol{\epsilon}}_i(t)$ through the 20-point Gaussian quadrature. Using the between-patient similarity matrix, we performed a hierarchical

clustering analysis on all 8,456 patients. Appendix B.2 presents parameter estimation and clustering analysis results on identifying patient subgroups.

The derived groups using our models represented patient's chronic preconditions so might not capture patient's most recent health status before the treatment. Therefore, we also included in $\boldsymbol{Z}_i$ the average values of the heath markers during the most recent year before time zero. Consequently, $\boldsymbol{Z}_i$ consisted of the derived group membership, the average values of SBP, HbA1c, HDL, BMI, DD, and logMed in the past one year before time zero, as well as age, gender, and race variables. Finally, we had data of $(A_i, \boldsymbol{Z}_i, R_i)$ from 5,125 patients for learning optimal ITRs.

### 3.3.2 Learning ITRs Using Proposed Method

Figure 3.1 describes the flow-chart of our proposed framework, along with methods used at each step. As illustrated in this figure, we used EHRs before time zero (last clinical encounter in year 2016) to fit joint models and learn subgroups for these patients. After using the 12-month-data of HbA1c after time zero to define the outcome, we applied the multicategory matched learning method in Section 3.2.2 to estimate optimal ITRs in each subgroup.



Figure 3.1: Flow-chart of learning optimal ITRs using the EHRs at the OSU-WMCIW.

Before estimating propensity scores, prognostic scores, and ITRs, we normalized all the continuous variables to alleviate the bias introduced by scaling. Propensity scores, $\pi(\boldsymbol{Z}_i)$, were estimated by a 10-fold cross-validation random forest with 3 repeats. The misclassification rates on the whole

training data were 3.1%, < 0.001%, 16.4%, 11.0% and 15.2% for each of groups 1 to 5, respectively. The 5th and 95th percentiles of estimated propensities are around 0.01 and 0.80. To avoid extreme weights in the calculation of value functions, we truncated probabilities less than 1% or greater than 80%. Similarly, prognostic scores, $\psi(\boldsymbol{Z}_i)$, were estimated by a gradient boosting model with 5,000 trees of which maximum depth was 4 in each patient subgroup. The model provided a good fit to the clinical outcome with a mean squared error less than $10^{-4}$. The most important covariate in estimating both propensity and prognostic scores is the recent one-year HbA1c.

Finally, in the doubly robust matching step in Section 3.2.2, we used the Euclidean metric as the distance function $d(\cdot, \cdot)$ and the 1-nearest neighbor to create matched pairs. In other words, for the $i$th patient, among other patients with the same gender and race but different baseline treatments, we searched for the patient who has the closest Euclidean distance to him/her in terms of demographics variables, recent measurements of health markers, and two estimated scores. We applied the radial basis function (RBF) kernel to optimize the objective function (3.4) in the matched learning model and 2-fold cross-validation with 100 repeats were used to learn optimal ITRs. The tuning parameter was selected from $\left\{2^k : k = 0, \pm 1, \ldots, \pm 15\right\}$ using 2-fold cross-validation. We calculated the bandwidth parameter of each RBF kernel according to a data-driven method which could be implemented using the kernlab (Karatzoglou et al., 2004) R package.

Given ITRs, we used (3.5) in Section 3.2.2 to evaluate the optimal ITR and compare with the universal rules of the four classes of treatments. In addition to the universal rules, we also compared the model performance of our method with Q-learning. In Q-learning, we used all the feature variables and the treatments in the model. Random forests and SVMs with RBF kernels were applied to obtain parameter estimates through 2-fold cross-validation with 100 repeats via the caret (Kuhn, 2020) R package. The cost parameter in SVMs was also selected from $\left\{2^k : k = 0, \pm 1, \ldots, \pm 15\right\}$ using 2-fold cross-validation. The parameter that controls the number of features being randomly selected at each node in random forests was chosen from a pool of 10 potential values which were automatically determined by the caret package, using 2-fold cross-validation.

### 3.4 Analysis Results of OSU-WMCIW EHRs

#### 3.4.1 Identified Latent Subgroups

We identified 5 subgroups in the EHR datasets and patients in each cluster have similar health profiles. Figure 3.2 shows the averages of normalized health marker measurements for patients in each cluster. The value in each cell is averaged across all patients and clinical visits for the corresponding health marker and patient subgroup. A higher value of HDL and a lower value of HbA1c, HDL, and BMI indicate a healthier status. The value of DD and logMed do not directly reflect the disease state. However, a lower value indicates that physicians tend to prescribe less medications to this group of patients, and thus a less severe state. We compared these values to the sample average of each health marker. A healthier T2D status is indicated in blue and a severe condition is indicated in red.

| | group 1 | group 2 | group 3 | group 4 | group 5 |
|---|---|---|---|---|---|
| Log number of medications (logMed) | 0.15 | 0.049 | −0.1 | 0.342 | −0.213 |
| Recieve diabetic drugs (DD) | −0.234 | 0.111 | −0.233 | 0.51 | 0.018 |
| Body mass index (BMI) | −0.16 | −0.508 | 1.024 | −0.112 | −0.49 |
| High−density lipoprotein (HDL) | −0.098 | −0.337 | −0.269 | −0.455 | 1.519 |
| Glycated hemoglobin (HbA1c) | −0.262 | −0.242 | −0.06 | 1.904 | −0.163 |
| Systolic blood pressure (SBP) | 1.415 | −0.255 | −0.003 | 0.129 | −0.238 |

Figure 3.2: Averages of normalized measurements by health markers and patient subgroups using 24-month-data before baseline treatment dates. Red: more severe status than the overall sample average in terms of a health marker; Blue: healthier status than the overall sample average in terms of a health marker; White: overall sample average status in terms of a health marker.

Group 5 is comprised of 1,231 patients. All the health markers convey the information that this

group of people is relative healthier. Compared to other groups, patients in this group did not take excessive number of medications. Group 2 contains 3,360 patients whose HDL is slightly lower than the average. This suggests that they might have some difficulties in controlling the cholesterol level. However, other health markers reflect a relatively healthy status of patients in this group. The SBP of 737 patients in group 1 is apparently higher than the average, while other health markers are below or around the averages. This pattern indicate that these patients might have hypertension and a moderate status of T2D. Group 3 have 2,446 patients and their BMIs are above the average. Besides BMI, the value of HDL in this group represent a bad signal as well. Thus, patients in group were at a moderately severe state. For the 682 patients in group 4, almost all the health markers show the most severe severe state of T2D. In particular, their HbA1c level is much higher. Another interesting fact is physicians had prescribed more-than-average diabetic drugs and ancillary medications to this group of patients, but their diseases were not controlled well. This result implies that this group might not have received the optimal treatments, and, therefore, we would focus on this group in the following analysis.

### 3.4.2 ITRs for Multicategory Treatments

As described in Section 3.3.1, we used 5,125 patients out of 8,456 patients in the finalized dataset for learning optimal ITRs. Before estimating optimal ITRs on the whole dataset, we compared the performance of the proposed method with the universal rules of four treatment classes and Q-learning. The results of 100 cross-validation repetitions are displayed in Figure 3.3, and the summary statistics are listed in Table B.1.

In general, the empirical HbA1c value of estimated ITRs is lower than any universal rules (i.e., "one-size-fits-all" rules) in any of the 100 repetitions. Compared to Q-learning using random forests and SVMs, the proposed method has lower average empirical values in the majority of cases as well. For example, in group 4, the weighted mean outcomes in (3.5) are 8.479, 9.199, 8.866, and 9.121 for patients prescribed metformin only, insulin only, other T2D drugs, and multiple treatments, respectively. ITRs estimated by the proposed approach achieve a mean empirical value function of 7.752, which is lower than that for Q-learning using random forests (9.052) and Q-learning using SVMs with RBF kernels (7.794). Furthermore, the standard deviations of our ITRs are relatively small across all patient subgroups. Thus, we conclude the proposed M-learning model outperforms universal rules and Q-learning on estimating ITRs for HbA1c control with a higher value and a

Figure 3.3: The empirical value function for the expected HbA1c level using 2-fold cross-validation with 100 repeats (a lower value means more beneficial).

lower variance.

After evaluating the model performance, we estimated the ITRs using the whole dataset. The distributions of four treatment classes in baseline assignments and ITRs are displayed in Figure 3.4. Taking group 4 as an example, there are 53 (14.2%) patients who received metformin only; 152 (40.8%) patients received insulin only; 88 (23.6%) received other T2D drugs; and 80 (21.4%) received at least two treatments. However, in other patient subgroups, 25% to 35% of patients were assigned metformin monotherapy, insulin monotherapy, and other T2D drugs, respectively. The proportion of multiple treatments is around 8% to 15%. Thus, the proportion of patients in group 4 receiving metformin monotherapy is much lower than that in other groups, while the normalized HbA1c level of group 4, as shown in Figure 3.2, is the highest among all groups.

Compared to the observed baseline treatments, the proportion of assigning either metformin monotherapy or other T2D drugs increases in almost every patient group. In contrast, the estimated ITRs suggest to prescribe insulin to less patients than observed. The proportions of metformin monotherapy and other T2D drugs recommended by the ITRs are fairly close; however, in group 4, metformin has a drastic increment in the assignment proportion from 14.2% to 40.2%. Similarly, in group 5, the proportion of insulin monotherapy decreases from 38.1% to 20.4%, together with

Figure 3.4: The distribution of observed treatments vs treatments recommended by estimated ITRs within each subgroup.

observable increases in the assignment of metformin monotherapy and other T2D drugs.

Table 3.1 displays the contingency table of observed treatments and the estimated ITRs recommendations. For metformin monotherapy, observed treatments and ITRs are matched by about 50% to 60% times. The proportion of other T2D drugs ranges from 40% to 60% across patient groups. Nevertheless, the observed insulin monotherapy merely has about a 30% to 40% match rate with ITRs. Particularly, in group 3, 4 and 5 of which patients have the highest HbA1c measurements, ITRs tend to assign metformin and other T2D drugs to more than 65% of patients who originally received insulin. We also investigated the agreement of the observed prescriptions and estimated ITRs using the Cohen's kappa coefficient $\kappa$ (Cohen, 1960). As shown in Table 3.2, all of the $\kappa$ coefficients for the $4 \times 4$ contingency tables of group 1 to 5 (rows "overall") are between 0.179 and 0.247. McHugh (2012) suggested any $\kappa$ value below 0.4 indicates at least moderate disagreement between the two categorical variables. Therefore, this result reveals the moderate disagreement between the observed prescriptions and estimated ITRs in all the identified patient subgroups. Moreover, we decomposed the "overall" disagreement to individual treatment class level by comparing each of the four treatment classes versus the rest. The $\kappa$ values for all comparisons are less than 0.4, suggesting the consistent disagreement of all treatment classes. Particularly, for

44

Table 3.1: Contingency tables of four treatment classes by observed prescriptions (rows) and estimated ITRs (columns).

| Group | | Other drugs | Insulin only | Metformin only | Multiple drugs | Total in prescriptions |
|---|---|---|---|---|---|---|
| Group 1 | Other drugs | 42 | 41 | 5 | 12 | 100 |
| | Insulin only | 43 | 55 | 38 | 9 | 145 |
| | Metformin only | 21 | 38 | 68 | 4 | 131 |
| | Multiple drugs | 11 | 28 | 15 | 15 | 69 |
| | Total in ITRs | 117 | 162 | 126 | 40 | 445 |
| | | | | | | |
| Group 2 | Other drugs | 242 | 200 | 55 | 40 | 557 |
| | Insulin only | 160 | 199 | 124 | 27 | 510 |
| | Metformin only | 319 | 17 | 342 | 7 | 685 |
| | Multiple drugs | 44 | 44 | 109 | 40 | 237 |
| | Total in ITRs | 765 | 480 | 630 | 114 | 1989 |
| | | | | | | |
| Group 3 | Other drugs | 168 | 82 | 116 | 19 | 385 |
| | Insulin only | 183 | 152 | 95 | 13 | 443 |
| | Metformin only | 69 | 127 | 237 | 14 | 447 |
| | Multiple drugs | 71 | 46 | 40 | 46 | 203 |
| | Total in ITRs | 491 | 407 | 488 | 92 | 1478 |
| | | | | | | |
| Group 4 | Other drugs | 52 | 20 | 10 | 6 | 88 |
| | Insulin only | 13 | 50 | 87 | 2 | 152 |
| | Metformin only | 20 | 2 | 28 | 3 | 53 |
| | Multiple drugs | 23 | 5 | 25 | 27 | 80 |
| | Total in ITRs | 108 | 77 | 150 | 38 | 373 |
| | | | | | | |
| Group 5 | Other drugs | 93 | 30 | 38 | 34 | 195 |
| | Insulin only | 115 | 94 | 88 | 23 | 320 |
| | Metformin only | 37 | 36 | 156 | 25 | 254 |
| | Multiple drugs | 17 | 11 | 25 | 18 | 71 |
| | Total in ITRs | 262 | 171 | 307 | 100 | 840 |

Table 3.2: Cohen's kappa coefficients for the agreement of observed prescriptions and estimated ITRs.

| Group | Treatment | Value | SE | Z | P-value |
|---|---|---|---|---|---|
| Group 1 | Overall | 0.179 | 0.032 | 5.576 | <0.001 |
| | Other drugs | 0.191 | 0.051 | 3.757 | <0.001 |
| | **Insulin only** | **0.022** | 0.047 | 0.463 | **0.643** |
| | Metformin only | 0.338 | 0.048 | 6.973 | <0.001 |
| | Multiple drugs | 0.182 | 0.060 | 3.052 | 0.002 |
| | | | | | |
| Group 2 | Overall | 0.180 | 0.015 | 11.877 | <0.001 |
| | Other drugs | 0.062 | 0.022 | 2.818 | 0.005 |
| | Insulin only | 0.204 | 0.024 | 8.477 | <0.001 |
| | Metformin only | 0.284 | 0.022 | 12.622 | <0.001 |
| | Multiple drugs | 0.163 | 0.031 | 5.318 | <0.001 |
| | | | | | |
| Group 3 | Overall | 0.181 | 0.018 | 10.314 | <0.001 |
| | Other drugs | 0.129 | 0.026 | 4.886 | <0.001 |
| | Insulin only | 0.099 | 0.027 | 3.696 | <0.001 |
| | Metformin only | 0.280 | 0.026 | 10.610 | <0.001 |
| | Multiple drugs | 0.247 | 0.036 | 6.872 | <0.001 |
| | | | | | |
| Group 4 | Overall | 0.247 | 0.031 | 7.851 | <0.001 |
| | Other drugs | 0.366 | 0.053 | 6.839 | <0.001 |
| | Insulin only | 0.224 | 0.047 | 4.744 | <0.001 |
| | **Metformin only** | **0.083** | 0.042 | 1.957 | **0.050** |
| | Multiple drugs | 0.371 | 0.060 | 6.163 | <0.001 |
| | | | | | |
| Group 5 | Overall | 0.218 | 0.022 | 9.835 | <0.001 |
| | Other drugs | 0.192 | 0.036 | 5.396 | <0.001 |
| | Insulin only | 0.160 | 0.033 | 4.917 | <0.001 |
| | Metformin only | 0.337 | 0.034 | 9.970 | <0.001 |
| | Multiple drugs | 0.124 | 0.044 | 2.794 | 0.005 |

*Note*: In column "Treatment", for each group, "Overall" represents the $4 \times 4$ contingency table in Table 3.1. The other four labels represent the $2 \times 2$ contingency tables for each of the four treatment classes versus the rest, respectively; "Value" is the value of Cohen's kappa coefficient $\kappa$ for the corresponding contingency table; "SE" is the asymptotic standard error of $\kappa$; "Z" is the test statistic for a z-test with the null hypothesis $\kappa = 0$; "P-value" is the p-value for the z-test.

the assignments of insulin monotherapy in group 1 and metformin monotherapy in group 4, the $\kappa$ coefficients are not significantly different from 0 (p-values are 0.643 and 0.050, respectively), so we could not reject the null hypothesis that the agreement is the same as chance agreement in these two cases.

Bringing all the comparisons together, we can conclude that metformin monotherapy has the optimal effect on HbA1c control, i.e., with the smallest HbA1c level at 6 months, especially for patients with low or moderate HbA1c levels. Whereas, the insulin monotherapy does not have noticeable advantages over other T2D drugs, and may even have worse HbA1c management when the baseline HbA1c level is high. Also, the insulin monotherapy may not be optimal to people who are relatively healthy and do not suffer from T2D complications. In this case, physicians may consider to prescribe either metformin monotherapy or other T2D drugs. These conclusions agree with the findings in Palmer et al. (2016) and may provide new directions to test treatment effects of T2D drugs.

### 3.4.3 Interpretable ITRs

To improve the interpretability of optimal ITRs, we employed the SHAP value approach (Lundberg and Lee, 2017) to examine the importance of features by each patient group. For each of the six binary classifiers in the weighted SVM in the M-learning model, we used the fastshap (Greenwell, 2020) R package to compute the approximate Shapley values using 10 Monte Carlo simulations for each row in the dataset and each feature. Taking the absolute value of all Shapley values, we computed the average across all incidences for each value, and treated these averages as the importance of features. The results are presented in Figure 3.5.

The most informative features in estimating ITRs are the prognostic scores and three propensity scores. The propensity scores summarize the observed treatment patterns that patients received. This result suggests that incorporating the propensity and prognostic scores is crucial to not only the doubly robust matching estimator framework but also the determination of ITRs. Gender and race have less important effects on optimal treatments, while age and the remaining six features related to health markers have moderate and non-negligible impacts on the estimated ITRs.

## 3.5 Discussion

In this chapter, we propose a general framework to estimate optimal ITRs for multicategory treatments using EHRs. Our first contribution is to create a novel latent process model, which

Figure 3.5: Importances for feature variables based on the absolute values of Shapley values. Prog: prognostic score. PropI: propensity score for insulin monotherapy. PropO: propensity score for other T2D drugs. PropM: propensity score for metformin monotherapy.

jointly analyzes different types of health markers and accounts for informative measurement patterns. Using patient similarities estimated from the latent process model, we identify subgroups of patients with homogeneous health profiles. The identified patient subgroups show different health patterns, and the cluster membership of patients is an important feature in M-learning to match patients among a more homogeneous pool.

The second contribution is the generalization of M-learning to multicategory treatments. We reduce the problem of selecting the optimal option among multicategory treatments into multiple binary classification problems using the one-versus-one strategy, and use the majority voting to integrate the results of the binary classification problems. Using this doubly robust multicategory M-learning, which incorporates propensity scores and prognostic scores, we reduce the confounding due to both covariates and recent patterns of health markers. Thus, our approach tackles the challenges in EHRs, and takes full advantage of information available from the health markers. Compared to any universal rules, our ITRs lead to a better control of the HbA1c level up to 13%. This result suggests our method is practical, and assists in the prescription of T2D treatments for HbA1c management.

We proposed a general pipeline (Figure 3.1) to learn the optimal treatment rules using the EHR data. However, there are certainly alternative ways to use in each component of this pipeline. For example, instead of using the hierarchical clustering, other clustering methods such as $k$-means (Forgy, 1965) or Gaussian mixture models may be used to identify the latent patterns in the EHRs. Also, the optimal number of patient subgroups may be selected by objective statistics (Rousseeuw, 1987; Tibshirani et al., 2002) using automated algorithms (Kassambara and Mundt, 2020). Another potential extension is that the outcome, HbA1c level at 6 months, was interpolated based on a linear function using one-year data since the HbA1c value changes slowly and smoothly. However, more sophisticated interpolation models such as splines may be useful if the measurements are taken intensively or over multiple years. In our application, there is a significant difference in the patients that are recommended metformin monotherapy between the observed prescriptions and estimated ITRs, especially for the patient group with the highest baseline HbA1c level. One possible reason for the discrepancy is we merely focus on glycaemic control, while, in real world, clinicians have to concern about other clinical outcomes or dose-related adverse events such as hypoglycaemia, gastrointestinal disorders, and renal failure. Thus, we can extend our work by considering the balance control of clinical outcomes and adverse events. Finally, our proposed framework focuses on finding the optimal treatment in a short period of time at a single decision point (e.g., at a patient visit). The method is not designed to optimize the long-term outcome over years and after multiple treatment phases. To adjust for the delayed effect in sequential decision making or long-term health management, our methods can be extended to learn the optimal dynamic treatment regimes to maximize the long term reward.

## CHAPTER 4: IDENTIFICATION OF OBJECTIVE BIOMARKERS FOR ADVERSE POSTTRAUMATIC NEUROPSYCHIATRIC SEQUELAE USING MOBILE SENSOR DATA

## 4.1   Introduction

Adverse posttraumatic neuropsychiatric sequelae (APNS) are common among trauma survivors after experiencing traumatic events such as car crash, physical and sexual assault, and natural disasters (Kessler, 2000). Common APNS disorders include posttraumatic stress disorder (PTSD), depression, post-concussion syndrome, and regional or widespread pain (McLean et al., 2020). Similar to many other mental health disorders, APNS classification and diagnosis are mainly based on subjective self-report measures that are not well mapped to underlying neurobiological mechanisms. Consequently, patients diagnosed with the same APNS disorder often experience very heterogeneous symptoms. Another issue of the self-report measures stems from the fact that more than half of people with any mental illness do not receive treatments or avoid treatments (Mental Health America, 2020). Even people who seek help often have a delay between the occurrence of the illness and clinical visits. In this case, the self-report measures, laboratory tests, and clinical diagnoses at a single point of time are probably inaccurate due to the delay. As a result, lack of objective measures has greatly impacted the research for APNS disorders, and the identification of objective biomarkers is critical to advancing APNS and mental health research.

Jain et al. (2015) firstly attempted to define the concept of digital phenotyping and described the opportunities in incorporating mobile sensor data into healthcare. They argued that if a person suffered from any disease, the mobile sensors they wear could trace the disease expressions, so the patterns in digital data could reflect disease symptoms. Later, Insel (2018) discussed the future of digital phenotyping and its potential functions in psychiatry and mental health. He stated the primary merit of digital phenotyping was that the digital phenotyping continuously recorded the objective signals of a person in their daily life, which could be different from the information they reported to clinicians later. Thus, besides the traditional clinical data, additional sources of continuous measurements can provide psychiatrists new insights into the assessment, treatment,

and prevention of mental illnesses. On the other hand, the development in sensor technology and information science makes digital data become an objective and ecological source of such measurements. In the past decade, a large variety of portable and interconnected mobile devices, such as smartphones, smartwatches, and wristband activity trackers, have been developed and widely used in daily activities by individuals of all ages and ethnicities around the world (Reinertsen and Clifford, 2018). Various types of sensors are built into these mobile devices, and massive amount of objective psychophysical signals are collected continuously by passive sensing from mobile device users over time. For example, smartphones have seven physical sensors which monitor the physical activity, sleep, and heart rate signals of users (Cornet and Holden, 2018).

Mobile sensor data provide invaluable information about the daily behaviors of users and have attracted an increasing attention from researchers in the field of mental health since 2010. For instance, Moshe et al. (2021) recruited 60 adult participants to explore the use of daily-life behaviour markers in the prediction of symptoms of depression and anxiety, using smartphones and wearable devices. Participants were continuously monitored over a 2-week period during which measures related to GPS, phone usage, activity, sleep, and HRV were recorded. The findings based on these data demonstrated a number of features had significant correlations with symptoms of depression and anxiety. Furthermore, smartphone features, wearable device features, and patient self-reported mood scores together provided the strongest prediction of depression. Depp et al. (2019) conducted a research to check the association between GPS data and symptom clusters in schizophrenia. Depp and colleagues collected GPS locations, which were tracked every 5 minutes by smartphones, and ecological momentary assessment reports of locations and behaviors from a total of 142 participants with schizophrenia (n=86) or healthy comparison subjects (n=56). They found the less GPS mobility was related to negative symptoms of schizophrenia. Thus, they concluded passive GPS sensing could be used for interventions for schizophrenia. Haines-Delmont et al. (2020) recently published an implementation of digital phenotyping in the suicide risk prediction. 66 qualified patients consented to participate the experiment, and their health information such as sleep behavior, mobility, and phone usage were collected up to a week through a smartphone application ("app"), which was linked to commercial wristbands and social networks. As a result, the research team proposed an algorithm that revealed the potential for discriminant suicide risk predictions, utilizing smartphone-generated and passive sensing data. Other representative research using mobile sensor data, classified by

mental illnesses, can be found in the fields of bipolar disorder (Faurholt-Jepsen et al., 2015; Abdullah et al., 2016; Beiwinkel et al., 2016; Palmius et al., 2017; Faurholt-Jepsen et al., 2019), schizophrenia (Wang et al., 2016; Staples et al., 2017; Barnett et al., 2018), depression (Saeb et al., 2015; Canzian and Musolesi, 2015; Ben-Zeev et al., 2015; Jacobson et al., 2019), PTSD (Karstoft et al., 2015; Minassian et al., 2015; Pyne et al., 2016; Place et al., 2017; Bourla et al., 2018), suicidal thoughts (Husky et al., 2014; Hallensleben et al., 2017; Kleiman et al., 2017, 2018), and stress (Muaremi et al., 2013; Sano and Picard, 2013; Sano et al., 2018; Goodday and Friend, 2019; DaSilva et al., 2019).

However, Linnstaedt et al. (2020) pointed out that there had been few progresses on the classification and ontology of APNS, and this was partially due to the lack of large-scale longitudinal studies that tracked APNS in large populations to achieve a sufficient statistical power and replicable findings. To help address these challenges, the National Institutes of Mental Health, collaborating with the US Army Medical Research and Material Command and several foundations, institutions, and companies, developed the AURORA study (McLean et al., 2020). The AURORA study is an ongoing large-scale (n=5000 target sample) cohort study that enrolls trauma survivors at the emergency department and follows them for one year. During the 1-year period, self-report indicator variables, genomic, neuroimaging, psychophysical, physiological, neurocognitive, and mobile sensor data are collected. These multi-layered longitudinal data from the AURORA study provide the field with unparalleled opportunities to advance the research on the classification and ontology of APNS. According to McLean et al. (2020), one of the primary goals of the AURORA study was to develop clinical decision tools for multidimensional APNS outcomes using the range of biobehavioral data collected. Motivated by this, we aimed to build a framework to predict construct scores, which were based on self-report indicator variables, and identify objective biomarkers related to APNS.

Although mobile sensor data provide great opportunities to advance mental health research, they also bring some unique analytical challenges. First of all, mobile sensor data are often in the form of high-frequency high-dimensional time series data. For example, the heart rate measurements traditionally are monitored every minute, so the number of features is 1440 for each day without any feature engineering. Additionally, mobile sensor data can be highly correlated. For instance, the heart rate measurements of different days have cyclic patterns for the same person. Finally, multimodal sensing may cause different time scales in mobile sensor data across domains. Again, heart rate measurements are monitored every minute, while measurements of sleep quality are

collected daily or every 12 hours. As a result, the features with finer resolutions will overwhelm others. These challenges make the modeling techniques for mobile sensor data be different from those for the data used by clinicians, which are subjective, infrequently sampled, and small-scale (Reinertsen and Clifford, 2018).

To overcome aforementioned data challenges, existing literature in mobile sensor data applications have adopted various data-mining techniques to solve formulated supervised learning tasks as ours. Among these literature, the most commonly used analytical methods are support vector network models (Ferdous et al., 2015; Abdullah et al., 2016; Wahle et al., 2016; Kelly et al., 2017; Sano et al., 2018), decision tree based models (Stütz et al., 2015; Garcia-Ceja et al., 2016; Jacobson et al., 2019; Faurholt-Jepsen et al., 2019), and variations of linear models (Saeb et al., 2015; Place et al., 2017; DaSilva et al., 2019; Moshe et al., 2021). Support vector network models consist of SVMs and SVRs, and they have been proved to be powerful methods for handling high dimensional data and solving classification/regression problems. Common passive sensing features considered by support vector network models are extracted from ECG and photoplethysmography (PPG). However, support vector network models require careful data preprocessing and parameter tuning. Furthermore, an increase in the sample size or the number of features can drastically increase the computing source of running support vector network models. Lastly, it can be difficult to explain the support vector network models to non-statisticians and identify the informative biomarkers if the linear kernel is not chosen. Decision tree based models, including random forests and gradient boosting methods, have been employed for handling classification/regression tasks on a mixture of categorical and continuous features. Thus, they are suitable for data collected from multiple digital sensors across domains. Additionally, compared to support vector network models, decision tree models are simpler and easier to interpret. Nevertheless, the performance of decision tree models can be limited under the case of overwhelming noise features. Therefore, decision tree models are not usually applied to high-frequent and/or high-dimensional sparse data. Linear models have been well studied and they can serve as the baseline models for most tasks. The generalized forms of linear models enable them to handle different types of outcomes. It is fast to train a linear model and make predictions. Also, it is easy to understand and explain the prediction results of linear models. Moreover, in high-dimensional spaces, linear models often have a powerful performance. The use of regularization parameters (L1, L2, elastic net) gives linear models the capability of performing feature selection.

In particular, LASSO (Tibshirani, 1996) is appropriate for using when the interpretability of the model is given the greatest consideration. The primary downside of linear models is the complexity of specifying interaction terms between features. Hence, the nonlinear relationship between the outcome and features can hardly be captured, and the predictions probably are biased.

Accounting for the challenges of using mobile sensor data and the findings in existing literature, it is unreasonable to simply input raw features from all types of mobile sensors into any single model or just use the features filtered by marginal screening methods (Lo et al., 2015). Instead, we propose a two-stage model to predict the continuous construct scores and identify objective biomarkers. In the first stage, we use a linear regression model with LASSO penalties to perform variable selection so that we improve the model interpretability. Specifically, in the linear model, we assign weights to features measured on different days and in different hours, and then we penalize the objective function for these weights and estimate the weights. Purposes of the first stage are to reduce the dimension of feature sets and select those features carrying the majority of the model effect. Also, the use of weights unifies the heterogeneous time scales of different domains of features, while keeping the most representative time-lag signals as much as possible. However, it is possible that the linear structure between construct scores and features is misspecified. Hence, to capture the nonlinear interactions between different domains of features, we apply a SVR model in the second stage. In this step, we plug the weights estimated from the first stage into a customized kernel function to further reduce the heterogeneity of different scales. Using the two-stage model, we not only select the most informative features but also combine the effects of various resolutions and temporal patterns. Therefore, we strike a balance between the prediction accuracy and model explainability.

The remaining part of this chapter are organized as follows. In Section 4.2, we introduce the data source of this chapter. In specific, we explain, in the AURORA study, the procedures of collecting raw data from mobile sensors, creating constructed scores, and selecting mobile sensor feature sets. In Section 4.3, we define the proposed model and describe our step-by-step algorithm for estimating model parameters. In Section 4.4, we illustrate an implementation of our method to a mobile sensor dataset collected by the AURORA study. We explore the relationship between three sets of features – demographic, activity, HRV features – and construct scores related to the pain experience of participants. Finally, we compare the performance of our method with three

alternative models and report the result.

## 4.2 Materials

### 4.2.1 Data Collection

The AURORA study has collected a large amount of digital sensing data, including HRV, activity, and sleep data from Verily's smartwatches, as well as keystroke, GPS, text, and voice data from the Mindstrong Discovery™ app on smartphones. The current research focuses on the identification of objective biomarker of APNS based on HRV data and activity data collected from smartwatches. The unprocessed HRV data and constructed HRV features were derived based on PPG signals on 5-minute sliding windows, which overlapped continuously for four and a half minutes. Meanwhile, the raw data and features of physical activities were obtained daily from accelerometer signals. In this chapter, we used the Research Domain Criteria (RDoC) framework (National Institute of Mental Health, 2021) to classify APNS, and we defined 10 RDoC related constructs - pain, loss, sleep discontinuity, nightmare, anxiety, hyperarousal, avoidance, re-experience, somatic and mental fatigue - by self-report indicator variables of participants. Self-report indicator variables were selected by domain experts from a rotating battery of smartphone-based survey. Each survey was administered at 6 different time points within the first 8 weeks after the traumatic event of each participant. To create factor scores for the 10 constructs, a joint measurement model (factor analysis model) across all 6 time points were developed for each construct, and then the performances of model fit were evaluated by different indices such as the comparative fit index (Bentler, 1990), standardized root mean square residual (Hooper et al., 2008), and Tucker-Lewis index (Tucker and Lewis, 1973). Finally, construct (factor) scores at all 6 time points were calculated for each construct using the joint measurement model.

### 4.2.2 Construct Scores

Preliminary results of the AURORA study revealed, among the 10 RDoC constructs, pain had a relatively strong association with the collected mobile sensor data, especially those from smartwatches. Thus, for this chapter, we selected the construct scores of pain as the outcome variables. The scale of the self-report indicator variable for pain ranged from 0 to 10, representing no pain to severe pain. The details of self-report indicator variable for pain in smartphone-based follow-up surveys can be found in Table 4.1.

Table 4.1: Questions for the self-report indicator variable for pain in smartphone-based follow-up surveys.

| Outcome Variables | Timepoints | Questions |
|---|---|---|
| Pain | Days: 1,9,21,31,43,53 | a. How would you rate your pain in the past 24 hours at its worst? |
| | | b. How would you rate your pain in the past 24 hours on average? |

### 4.2.3 Feature Sets

To predict the outcome scores, we used two fields of feature variables collected from mobile sensors: activity features and HRV features. 8 activity features in Table 4.2 were extracted from the preprocessed accelerometer signals, and these features were extracted daily. In particular, the activity feature set consisted of descriptive statistics of activity counts, cosinor-based rhythmometry (Cornelissen, 2014) metrics, and movement measurements. The mean of activity counts (meanACC) and standard deviation of activity counts (stdACC) were calculated for each 24-hour epoch of the accelerometer data, and they belonged to descriptive statistics of activity counts. After fitting a cosine model to the accelerometer data, the acrophase of cosinor rhythmometry (Acrophase) and amplitude of cosinor rhythmometry (Amplitude) were estimated. Acrophase represented the difference between the time point when daytime/nighttime switched and the starting time point of the 24-hour epoch. Amplitude indicated the difference between activity counts during daytime and nighttime. Finally, movement measurements included the wake percentage (wakePercentage), sleep-wake fragmentation (SWfragmentation), least active five hours (L5), and relative amplitude (RA). wakePercentage and SWfragmentation revealed the restfulness or restlessness due to sleep disturbances, while L5 and RA quantified the average activity in the waking and sleeping periods.

11 HRV features in Table 4.3 were derived from PPG signals, and these features were computed over 5-minute time windows and updated every 30 seconds. Specifically, the HRV feature set consisted of time-domain measures, frequency-domain measures, and non-linear measures (Shaffer and Ginsberg, 2017). Most HRV time-domain measurements were related to the descriptive statistics of interbeat intervals from which artifacts have been removed (NN intervals) and intervals between all successive heartbeats (RR intervals). For instance, we used the mean of NN Intervals (NNmean), interquartile range of NN Intervals (NNiqr), skewness of NN Intervals (NNskew),

Table 4.2: Activity feature set

| Feature Name | Description |
| --- | --- |
| meanACC | Mean of actigraphy in a 24-hour time window. |
| stdACC | Standard deviation of actigraphy in a 24-hour time window. |
| Acrophase* | Parameter in the cosine model which fits to actigraphy data. |
| Amplitude* | Parameter in the cosine model which fits to actigraphy data. |
| wakePercentage** | Number of waking epochs divided by data length. |
| SWfragmentation** | Number of transitions between waking and sleeping epochs divided by data length. |
| L5 | Least active 5 hour period in an the average 24 hour pattern. |
| RA | Normalized difference between the most active 10h period and least active 5h period in an average 24h pattern. |

*: The cosine model was fitted following Cornelissen (2014).
**: The waking and sleeping epochs were classified by the Cole-Kripke algorithm (Cole et al., 1992).

kurtosis of NN Intervals (NNkurt), average Signal-Quality-Index (avgsqi), standard deviation of NN Intervals (SDNN), root-mean square differences of successive RR intervals (RMSSD), and heart rate deceleration capacity (dc) as time-domain measurements. These features quantified the short-term variability of heartbeats. Meanwhile, using the Fast Fourier Transformation or autoregressive model, the variability of heart rates were separated into different rhythms that operated in ultra-low-frequency, very-low-frequency, low-frequency, and high-frequency bands (Task Force Report, 1996). Consequently, frequency-domain measures were created to estimate the distribution of signal energy in these four frequency bands. We included one such measurement into our model, and it was the ratio of low and high frequency spectral contents (LF/HF). Lastly, two non-linear measures - the ratio of two standard deviation measures for a Poincare plot (SD1/SD2) and approximate entropy (ApEn) - helped us evaluate the unpredictability, regularity, and complexity of RR time series.

Figures 4.1 and 4.2 display the scatter plots and correlations of the activity and HRV feature variables of 1000 randomly selected entries. Most of the feature variables have weak or moderate correlations between each other, but none of the absolute values of the correlations exceeds 0.85.

## 4.3   Methods and Algorithms

In this section, we build a regression model to predict the construct scores, using HRV and activity data collected before/on each of the 6 time points. By looking for informative features in the finalized model, we could identify the objective biomarkers for APNS related constructs. Figure

Table 4.3: HRV feature set

| Feature Name | Description |
| --- | --- |
| NNmean | Mean of NN intervals calculated in five-minute windows. |
| NNiqr | Interquartile range of NN intervals calculated in five-minute windows. |
| NNskew | Skewness of NN intervals calculated in five-minute windows. |
| NNkurt | Kurtosis of NN intervals calculated in five-minute windows. |
| avgsqi | Average Signal-Quality-Index in five-minute windows. |
| SDNN | Standard deviation of NN intervals calculated in five-minute windows. |
| RMSSD | Root-mean square differences of successive RR intervals. |
| LF/HF | Ratio of frequency activities in the low frequency (0.04 - 0.15Hz) range to the high frequency (0.15 - 0.40Hz) range. |
| dc | Average capacity of an autonomic nervous system to decelerate heartbeats (Jänig, 1989). |
| SD1/SD2 | The ratio of SD1 to SD2 (Tulppo et al., 1996). |
| ApEn | Entropic measurement to quantify the regularity of medical data (Pincus et al., 1991). |

*Note*: for each feature, we calculated four hourly summary statistics (mean, max, min, standard deviation). We used labels such as $NNmean_{mean}$, $NNmean_{max}$, $NNmean_{min}$, and $NNmean_{std}$ to denote these metrics.



Figure 4.1: Scatter plots and correlations of activity feature variables.

Figure 4.2: Scatter plots and correlations of HRV feature variables.

4.3 describes the flow-chart of our proposed framework and the method used at each step.

### 4.3.1 Penalized Least Squares Estimation

Suppose there are $N$ subjects in the mobile sensor dataset, and the finest time points for outcomes are $T_1, T_2, \ldots, T_n$. Let $\boldsymbol{Y}_t = (Y_{1t}, \ldots, Y_{Nt})^T$ be the outcome for all subjects at time $t$, and $\boldsymbol{X}_t = (\boldsymbol{X}_{t1}, \ldots, \boldsymbol{X}_{tp})^T$, where $\boldsymbol{X}_{tj}$ contains $s_j$ feature variables from the $j$th domains at time $t$. In this specific example, $\boldsymbol{Y}_t$ is the vector of construct scores of pain (see Section 4.2.2), and $\boldsymbol{X}_{tj}$ represents the activity features and HRV features as described in Section 4.2.3.

To quantify the relationship between outcome $\boldsymbol{Y}_t$ and feature $\boldsymbol{X}_{tj}$, we propose the following working model: for the $i$th subject,

$$Y_{it} = \boldsymbol{Z}_i^T \boldsymbol{\gamma} + \sum_{j=1}^{p} \sum_{l=1}^{k_j} \omega_{jl} \widetilde{\boldsymbol{X}}_{itjl}^T \boldsymbol{\beta}_j + \epsilon_{it}, \tag{4.1}$$

where $\boldsymbol{Z}_i$ is the vector of baseline covariates, which are set to demographic variables in this example. $\epsilon_{it} \sim \mathcal{N}(0, \sigma_t^2)$, $i = 1, \ldots, N$, are i.i.d. random errors. $\omega_{jl}$ indicates the time-lag weight for the $j$th domain at $l$ days before time $t$. $\widetilde{\boldsymbol{X}}_{itjl}$ is the vector of feature variables measured $l$ days before time

| Data resource | Methods/Models | Details/Explanations |
|---|---|---|
| HRV 5-min-window data, measured every 30 seconds | Calculate HRV hourly summary statistics $X_{HRV}$ | For each hour, if there are more than 20 HRV 5-min-window data, calculate hourly mean/std/min/max. |
| | Create HRV daily features $\tilde{X}_{HRV}$ | For each day, if all $X_{HRV}^b$ of the first 6 hours are non-missing, let $\tilde{X}_{HRV} = \sum \lambda_b X_{HRV}^b$. Time scale weights $\lambda_b \geq 0$ and $\sum \lambda_b = 1$. |
| ACT daily features, $\tilde{X}_{ACT}$ Demographics, $Z$ Pain scores, $Y$ | Determine sample size, and split data to training (80%) and test sets (20%). | $j = HRV, ACT$. $k = 1,2$. For 2 consecutive days prior to an observed pain score $Y$, if $\tilde{X}_{jk}$ are non-missing, include this $Y$ into the sample. Denote the dimension of $\tilde{X}_j$ to $s_j$. |

**Stage 1**
Assume $Y = \gamma^T Z + \sum_j \sum_k \omega_{jk}\beta_j^T \tilde{X}_{jk} + \epsilon$, and estimate $\gamma, \omega, \beta, \lambda$ — Time lag weights $\omega_{jk} \geq 0$ and $\sum_k \omega_{jk} = 1$. Apply lasso penalty to $\beta$ and fused lasso to $\omega$.

Feature selection — Exclude features (demographics, ACT, HRV) with regression coefficients $|\beta| < 0.01$.

**Stage 2**
Use the selected features to train SVM — Exploit a custom kernel on two samples $H$ and $H'$,
$K(H, H') = \exp\left\{-\frac{1}{2\sigma^2}\left[\|Z - Z'\|^2 + \sum_j \frac{\sum_k \hat{\omega}_{jk}\|\tilde{X}_{jk} - \tilde{X}'_{jk}\|^2}{s_j}\right]\right\}$

Predict pain scores, $\hat{Y}$, on test sets — Use 5-fold cross-validations to tune hyperparameters.

Repeat the above procedures for 10 times and take average results. Compare the results with sophisticated ML methods — Evaluation metric: $WMSE = \left(\frac{1}{N}\sum_{i=1}^N \frac{1}{|t_i|}\left(Y_{it} - \hat{Y}_{it}\right)^2\right)^{1/2}$

Figure 4.3: Flow-chart of predicting construct scores of pain using mobile sensor data in the AURORA study.

$t$. Here, $\widetilde{\boldsymbol{X}}_{itjl} = \boldsymbol{X}_{itjl}$ if feature variables from the $j$th domain are measured on the same time scale as the outcomes. Otherwise, if the measurement time scale of feature variables are finer, $\widetilde{\boldsymbol{X}}_{itjl}$ is an aggregated value of $\boldsymbol{X}_{itjl}$ given as

$$\widetilde{\boldsymbol{X}}_{itjl} = \sum_{b=1}^{m_j} \lambda_{jb}\boldsymbol{X}_{itjlb}, \tag{4.2}$$

where $\boldsymbol{X}_{itjlb}$ is the value of $\boldsymbol{X}_{itjl}$ at grid $b$, and $m_j$ is the ratio of the time unit of outcome measurements to the time unit of feature measurements.

We assume $k_j$, $\boldsymbol{\gamma}$, $\boldsymbol{\beta}_j$, $\boldsymbol{\omega}_j$, and $\boldsymbol{\lambda}_j$ are unknown and need to be estimated. $k_j = 1, \ldots, K_j$ (that is, for daily outcomes, retrieve the feature variables of previous $K_j$ days at most). $\boldsymbol{\gamma}$ is a vector of length $m + 1$. $\boldsymbol{\beta}_j$ is a vector of length $s_j$. $\boldsymbol{\omega}_j$ is a vector of length $k_j$. $\boldsymbol{\lambda}_j$ is a vector of length $m_j$. To estimate these parameters, we apply the following penalized least squares method with the LASSO and fused LASSO (Tibshirani et al., 2004) penalty terms:

$$
\begin{aligned}
\underset{k,\gamma,\beta,\omega,\lambda}{\arg\min} \quad & \frac{1}{N}\sum_{i=1}^N \frac{1}{|t_i|}\sum_t \Bigg[ Y_{it} - \boldsymbol{Z}_i^T\boldsymbol{\gamma} - \sum_{j=1}^p I(m_j = 1)\sum_{l=1}^{k_j}\omega_{jl}\boldsymbol{X}_{itjl}^T\boldsymbol{\beta}_j \\
& - \sum_{j=1}^p I(m_j > 1)\sum_{l=1}^{k_j}\omega_{jl}\left(\sum_{b=1}^{m_j}\lambda_{jb}\boldsymbol{X}_{itjlb}^T\right)\boldsymbol{\beta}_j \Bigg]^2 \\
& + \alpha_1\sum_{j=1}^p\sum_{l=1}^{s_j}|\beta_{jl}| + \alpha_2\sum_{j=1}^p\sum_{l=1}^{k_j}|\omega_{jl} - \omega_{j,l-1}|,
\end{aligned}
\tag{4.3}
$$

60

subject to

$$\omega_{jl} \geq 0, \sum_{l=1}^{k_j} \omega_{jl} = 1, \lambda_{jb} \geq 0, \sum_{b=1}^{m_j} \lambda_{jb} = 1.$$

The detailed algorithm for estimating the parameters in (4.3) is described in Algorithm 1.

The proposed model has important advantages on handling the challenges of mobile sensor data. Firstly, for each of the 6 time points at which surveys from participants were collected, we estimate the relationship between construct scores and temporal features. Thus, it is unnecessary to assume the construct scores of the same subject were identically distributed. Instead, the proposed model takes into account the potential heterogeneities of instances due to the change in time points. Additionally, we use the time-scale weights $\boldsymbol{\lambda}_j$ to rescale the effects of different domains of temporal features on construct scores, so the effects become comparable and the bias introduced by different domains can be alleviated. Moreover, we use the time-lag weights $\boldsymbol{\omega}_j$ to account for the auto-correlation among temporal features observed at nearby time stamps and quantify the lagged effects of temporal features on the construct scores. Another advantage of the proposed method is that the use of the LASSO penalty on $\boldsymbol{\beta}_j$ can select the most informative features and overcome the curse of dimensionality. For example, $\boldsymbol{\beta}_j = 0$ implies the features in the $j$th domain do not have observable effects on the outcome. Thus, the model has the function of group-typed feature selection. Lastly, the use of fused LASSO penalty on $\boldsymbol{\omega}_j$ guarantees the smoothness in lagged observations of temporal features.

### 4.3.2 Prediction through SVR with Customized Kernel

The parameter estimation in Section 4.3.1 assigns weights to each modality but does not consider between-modality interactions. Thus, we next use kernel machines and apply the SVR to incorporate potential nonlinear interactions and improve predictions. In this method, the key step is to define the between-subject distance based on the vector of features $\boldsymbol{H}_{it}^T = \left( \boldsymbol{Z}_i^T, \widetilde{\boldsymbol{X}}_{itjl}^T \right)$. For our purpose, we employ the time-lag weights to define the kernel distance

$$K(\boldsymbol{H}_{it_i}, \boldsymbol{H}_{i't_{i'}}) = \exp \left\{ -\frac{1}{2\sigma^2} \left( \|\boldsymbol{Z}_i - \boldsymbol{Z}_{i'}\|^2 + \sum_{j=1}^{p} \frac{1}{\sqrt{s_j}} \sum_{l=1}^{k_j} \hat{\omega}_{jl} \left\| \widetilde{\boldsymbol{X}}_{i,t_i,jl} - \widetilde{\boldsymbol{X}}_{i',t_{i'},jl} \right\|^2 \right) \right\}, \quad (4.4)$$

where $\hat{\omega}_{jl}$ are the estimators of time-lag weights obtained from (4.3). In (4.4), we firstly compute the lagged differences between rescaled temporal features, which are on the same scale across domains,

**Algorithm 1:** Algorithm for estimating parameters in (4.3).

---

**Result:** Estimates of $k_j$, $\boldsymbol{\gamma}$, $\boldsymbol{\beta}_j$, $\boldsymbol{\omega}_j$, and $\boldsymbol{\lambda}_j$, $j = 1, \ldots, p$.

Denote set $\mathcal{K} = \{(k_1, \ldots, k_p)\}$ and set $\mathcal{A} = \{(\alpha_1, \alpha_2)\}$.

**for** $(\alpha_1, \alpha_2) \in \mathcal{A}$ **do**

  **for** $(k_1, \ldots, k_p) \in \mathcal{K}$ **do**

    1. Initialize $\boldsymbol{\theta}$. Fit a linear regression of $\boldsymbol{Y}$ on $\boldsymbol{Z}$ to obtain $\boldsymbol{\gamma}^{(0)}$. Set $\boldsymbol{\omega}_j^{(0)} = \boldsymbol{I}_{k_j}/k_j$, $\boldsymbol{\lambda}_j^{(0)} = \boldsymbol{I}_{m_j}/m_j$, $j = 1, \ldots, p$. Minimize (4.3) w.r.t $\boldsymbol{\beta}$ to obtain $\boldsymbol{\beta}_j^{(0)}$. Denote $\boldsymbol{\theta}^{(s)} = \left(\boldsymbol{\gamma}^{(s)}, \boldsymbol{\beta}^{(s)}, \boldsymbol{\omega}^{(s)}, \boldsymbol{\lambda}^{(s)}\right)$. In this step, $s = 0$.

    2. Update $\boldsymbol{\theta}$. **while** $\left\|\boldsymbol{\theta}^{(s+1)} - \boldsymbol{\theta}^{(s)}\right\|_2 < 10^{-4}$ *or* $s \leq 100$ **do**

      (a) Given $\boldsymbol{\beta}_j^{(s)}$, $\boldsymbol{\omega}_j^{(s)}$, and $\boldsymbol{\lambda}_j$, minimizing (4.3) w.r.t $\boldsymbol{\gamma}$ is equivalent to solving a least square problem. Denote the estimates to $\boldsymbol{\gamma}^{(s+1)}$.

      (b) Given $\boldsymbol{\beta}_j^{(s)}$, $\boldsymbol{\omega}_j^{(s)}$, and $\boldsymbol{\gamma}^{(s+1)}$, minimizing (4.3) w.r.t $\boldsymbol{\lambda}_j$ is equivalent to solving a least square problem under nonnegativity and linear constraints. Denote the estimates to $\boldsymbol{\lambda}_j^{(s+1)}$.

      (c) Let $\delta_{j1} = \omega_{j1}$ and $\delta_{jl} = \omega_{jl} - \omega_{j,l-1}$. Then $\sum_{l=1}^{k_j} \omega_{jl} = 1$ yields $\sum_{l=1}^{k_j} \sum_{c=1}^{l} \delta_{jc} = 1$. Thus, $\delta_{j1} = [1 - \sum_{l=2}^{k_j}(k_j + 1 - l)\delta_{jl}]/(k_j + 1)$. Substitute $\omega_{j1}, \ldots, \omega_{jk_j}$ with $\delta_{j2}, \ldots, \delta_{jk_j}$ in (4.3). Given $\boldsymbol{\beta}_j^{(0)}$, $\boldsymbol{\gamma}^{(1)}$, and $\boldsymbol{\lambda}_j^{(1)}$, minimizing (4.3) w.r.t $\boldsymbol{\delta}_j$ is equivalent to solving a least square problem with LASSO penalties and linear constraints. According to Gainesa et al. (2018), this problem can be converted to a quadratic programming. Denote the estimates to $\widehat{\boldsymbol{\delta}}_j$, and then transform $\widehat{\boldsymbol{\delta}}_j$ to $\boldsymbol{\omega}_j^{(s+1)}$.

      (d) Given $\boldsymbol{\gamma}^{(s+1)}$, $\boldsymbol{\omega}_j^{(s+1)}$, and $\boldsymbol{\lambda}_j^{(s+1)}$, minimizing (4.3) w.r.t $\boldsymbol{\beta}_j$ is equivalent to solving a least square problem with LASSO penalties. Denote the estimates to $\boldsymbol{\beta}_j^{(s+1)}$.

    **end**

    3. Denote the estimates at the end of step 2 to $\widehat{\boldsymbol{\gamma}}$, $\widehat{\boldsymbol{\beta}}_j$, $\widehat{\boldsymbol{\omega}}_j$, and $\widehat{\boldsymbol{\lambda}}_j$, $j = 1, \ldots, p$.

      (a) Obtain $\hat{Y}_{it} = \widehat{\boldsymbol{\gamma}}^T \boldsymbol{Z}_i + \sum_{j=1}^{p} \sum_{l=1}^{k_j} \hat{\omega}_{jl} \widehat{\boldsymbol{\beta}}_j^T \widetilde{\boldsymbol{X}}_{itjl}$.

      (b) Estimate the variance of $\epsilon_t$ as the sample variance of $Y_{it} - \hat{Y}_{it}$. Denote it to $\hat{\sigma}_t^2$.

      (c) Calculate the joint likelihood of model (4.1) as

$$L = \prod_{i=1}^{N} \prod_{t} \left( \frac{1}{\sqrt{2\pi\hat{\sigma}_t^2}} \exp\left\{ -\frac{\left(Y_{it} - \hat{Y}_{it}\right)^2}{2\hat{\sigma}_t^2} \right\} \right).$$

  **end**

  4. Choose $(k_1, \ldots, k_p)$ which minimizes $\text{BIC} = [P \ln(\sum_i \sum_t 1) - 2\ln(L)]$, where $P = 1 + m + \sum_{j=1}^{p}(s_j + k_j + m_j)$, the number of parameters estimated by the model. Denote the optimal values to $\hat{k}_j$, $j = 1, \ldots, p$.

**end**

Select the optimal $\alpha_1$ and $\alpha_2$ from $\mathcal{A}$ by cross-validation (for each fold of data, conduct step 1 to 4). The corresponding estimates of $k_j$, $\boldsymbol{\gamma}$, $\boldsymbol{\beta}_j$, $\boldsymbol{\omega}_j$, and $\boldsymbol{\lambda}_j$, $j = 1, \ldots, p$, are the final estimates.

---

of two subjects. Next, within each domain, we use time-lag weights $\boldsymbol{\omega}_j$ to integrate all the lagged differences, and we normalize the differences by dividing them by the squared root of the dimension of features in the domain, $\sqrt{s_j}$. A critical advantage of (4.4) is that we take into account the delayed effects of features on the outcome by implementing the time-lag weights $\boldsymbol{\omega}_j$. Without the time-lag weights, then the feature measurements taken at different days have equal weights, so the model cannot utilize the informative time patterns. Also, compared to common kernel functions such as the linear kernel and RBF kernel, the proposed kernel distance can alleviate the bias introduced by different time scales of features through the normalization procedure. Otherwise, features with finer resolutions tend to dominate the computation of between-subject distances such that the outcome may heavily rely on the corresponding domains. Therefore, considering the effects of time-lag and time scales, the kernel distance can be regarded as a composite score which comprehensively reflects the difference between two samples.

## 4.4 Applications

### 4.4.1 Data Preparation

The database of the AURORA study contained 3139 participants. For the $i$th participant, at time point $t$, the outcome variable, $Y_{it}$, in (4.1) was the construct score of pain. The two domains of feature variables in this application, $\boldsymbol{X}_{t1}$ and $\boldsymbol{X}_{t2}$, were activity features in Table 4.2 and HRV features in Table 4.3, respectively. We set the maximum number of days of retrieving previous feature variables to 2 days for both activity and HRV feature variables, so we investigated the time-lag effect in a short term.

As the pattern in Figure 4.4 suggests, HRV features were often available during the first a few hours of each day. During the remaining hours, the HRV features had a great proportion of missing values. Considering the pattern of missing data, we created 4 summary statistics that calculated the hourly mean, max, min, and standard deviation of each HRV feature. We used labels such as $\mathrm{NNmean_{mean}}$, $\mathrm{NNmean_{max}}$, $\mathrm{NNmean_{min}}$, and $\mathrm{NNmean_{std}}$ to denote these metrics. We used the 44 hourly HRV features instead of the 11 original HRV features in (4.3) and (4.4).

We only treated a summary statistic as non-missing if it was calculated from one of the first 6 hours in an epoch and, in any of the 6 hours, at least 10 minutes of the original HRV features were available. Also, we excluded construct scores whose previous records of activity or HRV features
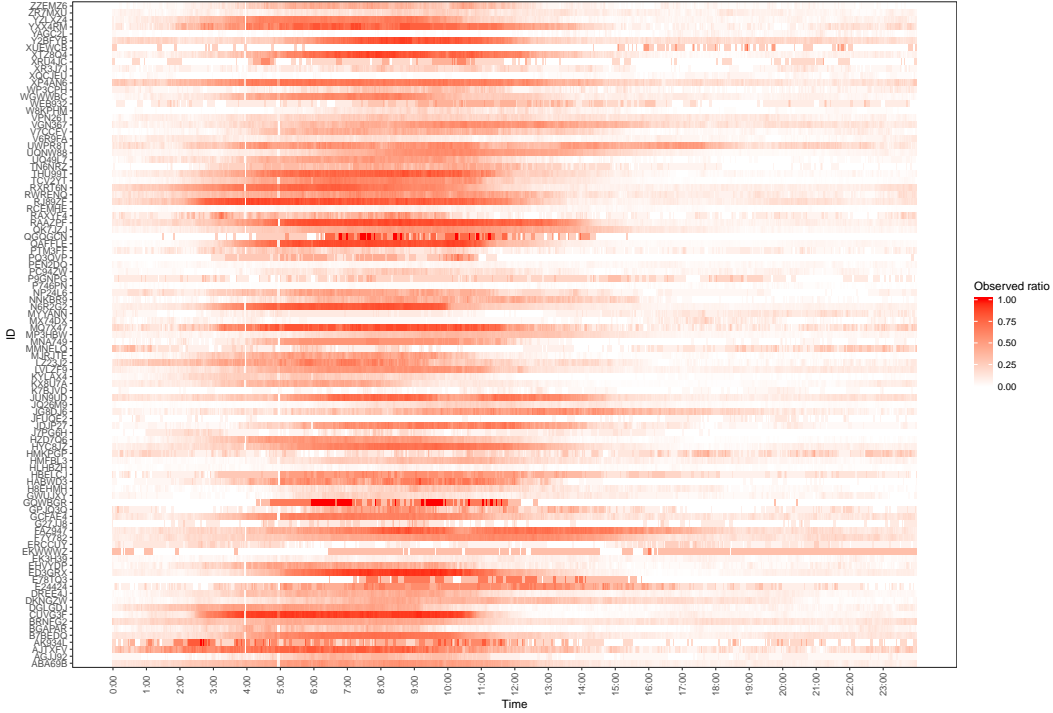
Figure 4.4: Missing data pattern of HRV features for 100 randomly selected participants. Mobile devices were refreshed at 4am or 5am UTC of each day.

were missing. Therefore, the finalized dataset included 1285 construct scores of pain collected from 632 participants. Besides activity and HRV features, we used 5 demographic variables $\boldsymbol{Z}_i$ to help predict construct scores of pain. The description and distribution of demographic variables are listed in Table 4.4.

Finally, to overcome the problem of over-fitting, we randomly divided 80% of the participants into training data and 20% of the participants into validation data. We adopted training data and the technique of cross-validation to derive optimal estimates, and then we used the validation data to evaluate our method and compare it with selected machine learning models in the existing literature.

### 4.4.2 Implementation of Proposed Methods

Before implementing the proposed method, we transformed the original construct scores of pain using the ordered quantile normalization in the bestNormalize (Peterson and Cavanaugh, 2019) R package to fulfill the normality assumption of (4.1). In addition, we normalized all the continuous feature variables to alleviate the bias introduced by scaling. To estimate time-lag weights, $\boldsymbol{\omega}_j$, and regression parameters, $\boldsymbol{\beta}_j$ and $\boldsymbol{\gamma}$, in (4.3), we minimized the objective function using the pracma

Table 4.4: Demographic information of participants for predicting construct scores.

| Demographic Variable | Statistic | Value |
|---|---|---|
| Highest Grade | mean(sd) | 15.7(2.4) |
| | median(min, max) | 15.0(9.0, 21.0) |
| Age | mean(sd) | 38.6(13.9) |
| | median(min, max) | 36.0(18.0, 74.0) |
| Marital Status | sample size(proportion) | Never married: 327(51.7%) |
| | | Other: 305(48.3%) |
| Gender | sample size(proportion) | Female: 462(73.1%) |
| | | Male: 170(26.9%) |
| Race | sample size(proportion) | Black: 235(37.2%) |
| | | Other: 397(62.8%) |

(Borchers, 2019) and ADMM (You and Zhu, 2018) R packages. The tuning parameters $\alpha_1$ and $\alpha_2$ were selected from $\{2^k : k = 0, \pm 1, \ldots, \pm 15\}$ using a 5-fold cross-validation. Informative feature variables were selected if the absolute values of their coefficients exceeded 0.01. Given the selected features $\boldsymbol{H}_{it}$, we refit (4.3) and updated parameter estimates $\hat{k}_j$, $\widehat{\boldsymbol{\gamma}}$, $\widehat{\boldsymbol{\beta}}_j$, $\widehat{\boldsymbol{\omega}}_j$, and $\widehat{\boldsymbol{\lambda}}_j$. Next, we plugged the time-lag weights $\widehat{\boldsymbol{\omega}}_j$ to (4.4) and employed the SVR model with this kernel function. In the SVR model, tuning parameters $C$ and $\epsilon$ were selected from $\{2^k : k = 0, \pm 1, \ldots, \pm 15\}$ using a 5-fold cross-validation. The model fitting was conducted on training data by kernlab (Karatzoglou et al., 2004) R package, and we evaluated the model through a weighted mean squared error (WMSE) metric as follows:

$$\text{WMSE} = \left( \frac{1}{N} \sum_{i=1}^{N} \frac{1}{|t_i|} \sum_{t} \left( Y_{it} - \hat{Y}_{it} \right)^2 \right)^{1/2} . \tag{4.5}$$

We chose the SVR model with the minimum WMSE as the final model. Subsequently, we applied the finalized model to the validation data and obtained predicted construct scores of pain of the validation data. Finally, we conducted the inverse transformation on the predicted outcome scores and calculated the WMSE for the proposed approach.

As a comparison of the proposed method, we used five alternative models. The first two models were linear regression models with no predictor variables and with only demographic variables. These two models served as baseline models. The third approach was the "naive" SVR model which treated all of the activity and HRV features at different time as independent predictors. This approach represented a typical way to handle mobile sensor data in the existing literature. The

fourth model was the linear regression model in (4.1), but we used the informative features remained after feature selection instead of the whole feature set. Through this model, we tested whether the proposed approach captured the nonlinear between-domain interactions between features. The last candidate was the same as the proposed method, however, with the RBF kernel in the second stage. By comparing with this model, we verified the necessity of using the proposed kernel (4.4) to integrate different domains and improve the prediction accuracy.

Table 4.5 presents the parameter estimation and feature selection results for construct scores of pain. We selected $\gamma$ and $\beta$ parameters whose absolute values were greater than 0.01. After the filtering, we reconstructed the linear model in (4.1) with the left features and estimated their parameters again using (4.3). As shown in Table 4.5, all the demographic variables remain in the model. On average, less educated, elder, currently married, female, and black people have greater construct scores of pain, and this result indicates such participants have experienced more severe pain. Among all the activity and HRV features, Amplitude, $\text{RMSSD}_{\text{mean}}$, $\text{SDNN}_{\text{max}}$, $\text{NNskew}_{\text{min}}$, $\text{SDNN}_{\text{std}}$, and $\text{ApEn}_{\text{std}}$ have relatively greater effects on construct scores of pain, and the effects are all negative. From another aspect, we could learn the time-lag effects of mobile device features on construct scores of pain from estimates of $\omega_{jl}$ in Table 4.5. Since all the $\omega_{jl}$ estimates are close to 0.5, we conclude, within each set of mobile sensor features, the features monitored on the same day as the outcome observation and those monitored one day before have roughly equal daily effects on the construct scores of pain. In addition, $\lambda_{25} = 1$ indicates the summary statistics of HRV features of the fifth hour in an epoch fully determine the daily effect.

Finally, we used the selected features to fit the second stage model and predict the construct scores of pain on the validation data. The prediction results of the proposed method as well as the five alternative models are listed in Table 4.6. Compared with the intercept-only regression model, the use of demographic features reduces the WMSE of predicted construct scores of pain from 3.100 to 2.803. Meanwhile, the models employing mobile sensor features further improve the WMSE by 0.067 to 0.283. Thus, in our data, not only demographic features but also mobile sensors features are informative to the prediction of construct scores. The WMSE for the "naive" SVR model is 2.575 and the model includes 549 feature variables. The proposed method achieves a WMSE of 2.520, but it only uses 67 feature variables.

Also, Figure 4.5 reveals a moderate linear correlation between predicted and observed construct

66

Table 4.5: Estimated parameters of features left in (4.1) for predicting construct scores of pain

| Domain | Feature | Parameter | Before selection | After selection |
|---|---|---|---|---|
| Demographics | Intercept | $\gamma_0$ | -0.1803 | -0.1994 |
| | Highest Grade | $\gamma_1$ | -0.1608 | -0.1662 |
| | Age | $\gamma_2$ | 0.2172 | 0.2210 |
| | Marital Status | $\gamma_3$ | 0.1453 | 0.1470 |
| | Gender | $\gamma_4$ | 0.1510 | 0.1615 |
| | Race | $\gamma_5$ | 0.2577 | 0.2915 |
| | | | | |
| Activity | Amplitude | $\beta_{1,4}$ | -0.1358 | -0.1445 |
| | Lagged days | $K_1$ | 2 | 2 |
| | 1 day ago | $\omega_{11}$ | 0.4648 | 0.4712 |
| | Today | $\omega_{12}$ | 0.5352 | 0.5288 |
| | Daily | $\lambda_{11}$ | 1 | 1 |
| | | | | |
| HRV | $\text{RMSSD}_{\text{mean}}$ | $\beta_{2,7}$ | -0.0288 | -0.0183 |
| | $\text{SDNN}_{\text{max}}$ | $\beta_{2,17}$ | -0.0106 | -0.0855 |
| | $\text{NNskew}_{\text{min}}$ | $\beta_{2,25}$ | -0.0375 | -0.0264 |
| | $\text{SDNN}_{\text{std}}$ | $\beta_{2,39}$ | -0.0596 | -0.0509 |
| | $\text{ApEn}_{\text{std}}$ | $\beta_{2,44}$ | -0.0102 | -0.0075 |
| | Lagged days | $K_2$ | 2 | 2 |
| | Today | $\omega_{21}$ | 0.5000 | 0.5000 |
| | 1 day ago | $\omega_{22}$ | 0.5000 | 0.5000 |
| | 1st hour | $\lambda_{21}$ | 0 | 0 |
| | 2nd hour | $\lambda_{22}$ | 0 | 0 |
| | 3rd hour | $\lambda_{23}$ | 0 | 0 |
| | 4th hour | $\lambda_{24}$ | 0 | 0 |
| | 5th hour | $\lambda_{25}$ | 1 | 1 |
| | 6th hour | $\lambda_{26}$ | 0 | 0 |

*Note*: estimators in the full model with $|\gamma| < 0.01$ and $|\beta| < 0.01$ are not listed in this table.

Table 4.6: Predicted results of construct scores on validation data for four models.

| Model | WMSE | Number of features |
|---|---|---|
| Proposed model | 2.520 | 67 |
| Linear regression model (intercept only) | 3.100 | 0 |
| Linear regression model (demographic features) | 2.803 | 5 |
| "Naive" SVR model | 2.575 | 549 |
| Linear regression model (selected features) | 2.736 | 67 |
| 2-stage model with RBF kernel | 2.626 | 67 |

scores of pain (the Pearson correlation coefficient $\rho = 0.4219$). If we ignore the nonlinear interactions among features and use the linear regression model (4.1) with features left in Table 4.5 to predict construct scores of pain, the WMSE increases from 2.520 to 2.736. Also, compared to the RBF kernel, the proposed kernel in (4.4) improves the prediction and reduce the WMSE by 0.106. Therefore, we draw the conclusion that the proposed method has a close performance to the sophisticated SVR model, but has much fewer feature variables. Lastly, the use of the second stage model handles the potential issues of the linear model, and the proposed kernel has advantages over the commonly used RBF kernel.
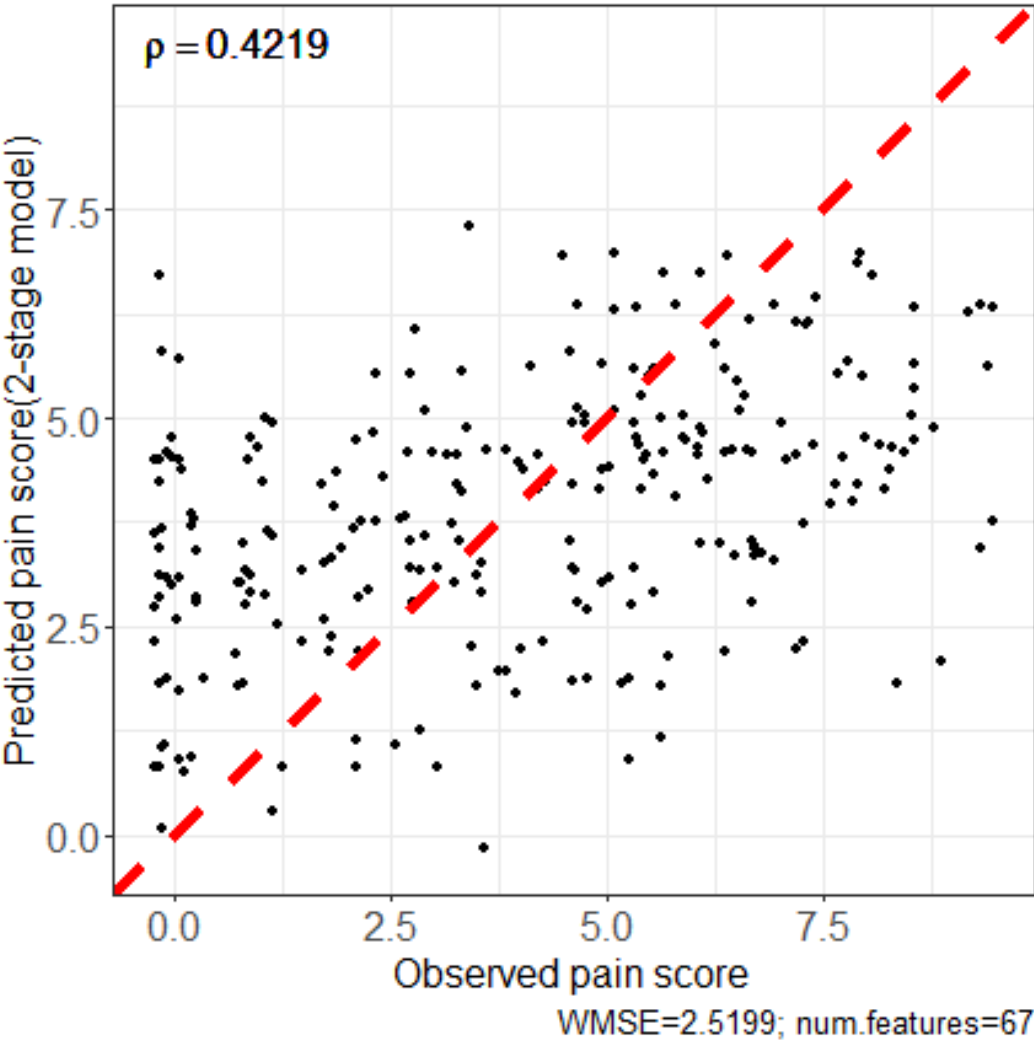


Figure 4.5: Scatter plots of predicted construct scores versus observed construct scores in validation data. Red dashed lines: lines y=x. The Pearson correlation coefficient between predicted and observed values, $\rho$, is 0.4219.

### 4.5 Discussion

In this chapter, we proposed a framework to estimate RDoC construct scores and identify objective biomarkers related to APNS, using mobile sensor data from the AURORA study. The major contribution of this framework is to create a novel two-stage model, which jointly implements different domains of features and accounts for complicated interactions between them. In the first stage, we select the most informative features to the estimation of outcomes through a linear regression model with LASSO and fused LASSO penalties. To adjust for the heterogeneity between different domains of features on different time scales, we aggregate the measurements and unify them to the same level. We also incorporate the effect of historical feature information on the current observations of outcomes by estimating the time-lag weights. In the second stage, we employ the SVR model to account for between-domain interactions and estimate the possible non-linear relationship between outcomes and features. Especially, we apply a customized kernel function as (4.4) in the SVR model. Compared to common kernel functions such as the RBF kernel, our kernel function reduces the heterogeneities between different domains of features by normalizing time-lag weights and time scales. Thus, our approach tackles the challenges in the high-dimensional, correlated, and multimodal measured mobile sensor data and takes full advantage of the available information. Compared to popular nonparametric methods which treat features observed at different time as independent predictors, our model leads to a better prediction accuracy but only requires 12% of the feature variables. This result suggests our method is a practical assist in the prediction of RDoC constructs and identification of objective biomarkers.

## CHAPTER 5: EXTENSIONS AND FUTURE WORK

### 5.1 Identify Patient Subgroups

In Chapter 2, the estimation of both regression coefficients and correlations among latent processes only relies on one or two health markers, so our method can be easily extended to handle a large number of health markers. In this case, one can perform the computation by parallel computing to save computing time and cost. Inferences on the estimators can be made based on subsampling subsets of the data. Another extensions to the proposed method is to estimate all regression coefficients $\boldsymbol{\beta}$'s simultaneously by incorporating the entire covariance matrix $\boldsymbol{\Omega}(t)$ to the estimating equations for $\boldsymbol{\beta}$'s, or to allow the marker-specific and time-sensitive bandwidth selection during the parameter estimation (especially when the smoothness of health marker trajectories are expected to be substantially different). Lastly, instead of the nonparametric estimation, explicitly modeling the temporal dependence within the same health marker as well as across health markers probably will improve the estimation. Despite the increased computational burden, an advantage of this extension is the potential to obtain a more precise assessment of the latent process given the entire history of health markers.

### 5.2 Matched Learning Model for Multicategory Treatments

As discussed in Section 3.5, the matched learning model still has room to extend and handle with more complicated observatory data. Firstly, to better interpret the obtained results, we will consult with physicians about the clinical meanings of identified patient subgroups and ITRs. Also, we plan to establish theoretical properties for the proposed method. For example, proving the Fisher consistency of the estimated decision rules. Moreover, we will further validate the proposed method by comparing with newly published methods (Lou et al., 2018; Huang et al., 2019; Zhang et al., 2020; Qi et al., 2020). In addition, we consider to test and incorporate the proposed method using other data resources such as the high-dimensional measurements of biomarkers collected by mobile sensors to recommend individual treatments. Another extension is to take other aspects of T2D control into consideration. For instance, besides lowering the level of HbA1c, we will take into account the

balance control of adverse events such as hypoglycemia. Lastly, considering the potential switch of treatments after time zero and the research interest in the long-term health management, we plan to extend the proposed method to learn the optimal dynamic treatment regimes. This extension can be realized by substituting the value function in (3.4) with a matching-based value function for multiple stages and using a backward induction (Liu et al., 2018).

## 5.3  Analysis of Mobile Sensor Data

For this topic, we applied the SVR model to capture between-modality interactions as the SVR model performs stable for high-dimensional data. However, other nonparametric models such as random forests, gradient boosting methods, and neural networks can be similarly implemented. In future research, we will also verify our approach on simulated datasets, additional RDoC constructs, and different types of outcome variables such as binary or categorical outcomes. Another extension is to take other domains of feature variables into consideration such as GPS and phone usage data. Finally, to further improve the explainability of the model, we will explore more meaningful summary statistics for HRV features and use the SHAP value approach (Lundberg and Lee, 2017) to evaluate the importance of features in the final model.

# APPENDIX A: APPENDIX FOR CHAPTER 2

## A.1 Gauss-Hermite Quadrature Method for Parameter Estimation

When $g_k^{-1}(z)$ takes a general form, we can compute $\mathbb{E}\left[Y_{ik}(t)|\boldsymbol{X}_i\right]$ and $\mathbb{E}\left[Y_{ik}(t)Y_{il}(t)|\boldsymbol{X}_i\right]$ using the Gauss-Hermite quadrature method (Abramowitz and Stegun, 1965). Suppose $Q$ is the number of mass points, $p_q$ are the mass points, and $w_q$ are weights. Then

$$
\mathbb{E}\left[Y_{ik}(t)|\boldsymbol{X}_i\right] \approx \sum_{q=1}^{Q} \frac{w_q}{\sqrt{\pi}} g_k^{-1}(\boldsymbol{X}_i^T \boldsymbol{\beta}_k(t) + \sqrt{2c_k} p_q),
$$

$$
\begin{aligned}
\mathbb{E}\left[Y_{ik}(t)Y_{il}(t)|\boldsymbol{X}_i\right] \approx & \sum_{q_1=1}^{Q} \sum_{q_2=1}^{Q} \frac{w_{q_1}}{\sqrt{\pi}} \frac{w_{q_2}}{\sqrt{\pi}} \\
& \times g_k^{-1}\left(\boldsymbol{X}_i^T \widehat{\boldsymbol{\beta}}_k(t) + \sqrt{2c_k} p_{q_1}[\boldsymbol{R}_{kl}(t)]_{1,1} + \sqrt{2c_k} p_{q_2}[\boldsymbol{R}_{kl}(t)]_{2,1}\right) \\
& \times g_l^{-1}\left(\boldsymbol{X}_i^T \widehat{\boldsymbol{\beta}}_l(t) + \sqrt{2c_l} p_{q_1}[\boldsymbol{R}_{kl}(t)]_{1,2} + \sqrt{2c_l} p_{q_2}[\boldsymbol{R}_{kl}(t)]_{2,2}\right),
\end{aligned}
$$

where $\boldsymbol{R}_{kl}(t)$ is chosen as the $2 \times 2$ square-root matrix of $\boldsymbol{\Sigma}_{kl}(t)$, and $[\boldsymbol{R}_{kl}(t)]_{i,j}$ denotes the entry in the $i$th row and $j$th column of $\boldsymbol{R}_{kl}(t)$.

On the other hand, we can compute $\mathbb{E}\left[\boldsymbol{\epsilon}_i(t)\Big|\boldsymbol{Y}_i(t), \widehat{\boldsymbol{\beta}}_k(t), \widehat{\sigma}_{kl}(t)\right]$ as follows,

$$
\begin{aligned}
& \mathbb{E}\left[\boldsymbol{\epsilon}_i(t)\Big|\boldsymbol{Y}_i(t), \widehat{\boldsymbol{\beta}}_k(t), \widehat{\sigma}_{kl}(t)\right] \\
= & \frac{\int \mathbb{P}\left(\boldsymbol{Y}_i(t)|\boldsymbol{\epsilon}_i(t)\right) \mathbb{P}\left(\boldsymbol{\epsilon}_i(t)\right) \boldsymbol{\epsilon}_i(t) d\boldsymbol{\epsilon}_i(t)}{\int \mathbb{P}\left(\boldsymbol{Y}_i(t)|\boldsymbol{\epsilon}_i(t)\right) \mathbb{P}\left(\boldsymbol{\epsilon}_i(t)\right) d\boldsymbol{\epsilon}_i(t)} \\
= & \frac{\int_D \prod_{k=1}^{p} \left[f_{ik}\left(y_{ik}(t); \boldsymbol{X}_i^T \widehat{\boldsymbol{\beta}}_k(t) + \epsilon_{ik}(t), \hat{\phi}_{ik}\right) \epsilon_{ik}(t)\right] f\left(\boldsymbol{\epsilon}_i(t); \widehat{\boldsymbol{\Omega}}(t)\right) d\boldsymbol{\epsilon}_i(t)}{\int_D \prod_{k=1}^{p} f_{ik}\left(y_{ik}(t); \boldsymbol{X}_i^T \widehat{\boldsymbol{\beta}}_k(t) + \epsilon_{ik}(t), \hat{\phi}_{ik}\right) f\left(\boldsymbol{\epsilon}_i(t); \widehat{\boldsymbol{\Omega}}(t)\right) d\boldsymbol{\epsilon}_i(t)} \\
\approx & \frac{\sum_{q.=1}^{Q} \prod_{k=1}^{p} \left[f_{ik}\left(y_{ik}(t); \boldsymbol{X}_i^T \widehat{\boldsymbol{\beta}}_k(t) + h_k\left(p_q, \widehat{\boldsymbol{\Omega}}(t)\right), \hat{\phi}_{ik}\right) h_k\left(p_q, \widehat{\boldsymbol{\Omega}}(t)\right) w_{q_k}\right]}{\sum_{q.=1}^{Q} \prod_{k=1}^{p} \left[f_{ik}\left(y_{ik}(t); \boldsymbol{X}_i^T \widehat{\boldsymbol{\beta}}_k(t) + h_k\left(p_q, \widehat{\boldsymbol{\Omega}}(t)\right), \hat{\phi}_{ik}\right) w_{q_k}\right]},
\end{aligned}
$$

where $\sum_{q.=1}^{Q} \equiv \sum_{q_1=1}^{Q} \cdots \sum_{q_p=1}^{Q}$. Furthermore, $f_{ik}(\cdot)$ is the density function of $Y_{ik}(t)$, and $f\left(\boldsymbol{\epsilon}_i(t); \widehat{\boldsymbol{\Omega}}(t)\right)$ is the density function of a multivariate normal distribution with mean $\boldsymbol{0}$ and covariance matrix $\widehat{\boldsymbol{\Omega}}(t)$. Lastly, $h_k\left(p_q, \widehat{\boldsymbol{\Omega}}(t)\right) = \sqrt{\frac{2c_k}{\pi}} \sum_{l=1}^{p} p_{q_l}\left[\widetilde{\boldsymbol{R}}(t)\right]_{k,l}$ with $\widetilde{\boldsymbol{R}}(t)$ being the square-root matrix of $\widehat{\boldsymbol{\Omega}}(t)$, and $[\widetilde{\boldsymbol{R}}(t)]_{k,l}$ is the entry in the $k$th row and $l$th column of $\widetilde{\boldsymbol{R}}(t)$.

## A.2 Proof of Theorem 2.3.1

Without loss of generality, we prove the theorem 2.3.1 for one health marker $k$, and the results can be generalized to other health markers. The key idea is to establish the following relationship, for $\|\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\| < M(nh_{1n})^{-1/2}$,

$$
\begin{aligned}
&\sup_{\boldsymbol{\beta}_k(t)} (nh_{1n})^{1/2}\left| U_{n,k}(\boldsymbol{\beta}_k(t)) - \left\{ U_{n,k}(\boldsymbol{\beta}_k^0(t)) - \mathbb{E}\left[U_{n,k}(\boldsymbol{\beta}_k^0(t))\right]\right\} - A_k(t)\left[\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right]\right| \\
&= O_p\left(n^{1/2}h_{1n}^{5/2}\right) + o_p\left(1 + h_{1n}^{1/2} + (nh_{1n})^{1/2}\left\|\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right\|\right),
\end{aligned}
\tag{A.1}
$$

where $A_k(t)$ is defined in (2.10).

We introduce notations $\boldsymbol{P}_n$ and $\boldsymbol{P}$ to denote the empirical and true probability measure respectively. Then we obtain

$$
\begin{aligned}
U_{n,k}(\boldsymbol{\beta}_k(t)) &= (\boldsymbol{P}_n - \boldsymbol{P})\left[\int K_{h_{1n}}(s-t)\boldsymbol{X}[Y_k(s) - E_k(\boldsymbol{\beta}_k(t),t)]d\tilde{N}_k(s)\right] \\
&\quad + \mathbb{E}\left[\int K_{h_{1n}}(s-t)\boldsymbol{X}[Y_k(s) - E_k(\boldsymbol{\beta}_k(t),t)]d\tilde{N}_k(s)\right] \\
&= I + II,
\end{aligned}
\tag{A.2}
$$

where $E_k(\boldsymbol{\beta}_k(t),t) = \mathbb{E}\left[Y_k(t)|\boldsymbol{X}\right] = \mathbb{E}_\epsilon\left[g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k(t) + \epsilon_k(t)\right)\right]$.

For the second term on the right-hand side of (A.2), we have

$$
\begin{aligned}
&\mathbb{E}\left[\int K_{h_{1n}}(s-t)\boldsymbol{X}[Y_k(s) - E_k(\boldsymbol{\beta}_k(t),t)]d\tilde{N}_k(s)\right] \\
&= \int K_{h_{1n}}(s-t)\mathbb{E}\left[\boldsymbol{X}\left[Y_k(s) - E_k(\boldsymbol{\beta}_k(t),t)\right]\exp\left\{-\boldsymbol{X}^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_k^T(s)\widehat{\boldsymbol{\eta}}_k\right\}dN_k(s)\right] \\
&= \int_s K_{h_{1n}}(s-t)\mathbb{E}\left[\boldsymbol{X}\left[Y_k(s) - E_k(\boldsymbol{\beta}_k(t),t)\right]\exp\left\{-\boldsymbol{X}^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_k^T(s)\widehat{\boldsymbol{\eta}}_k\right\}\lambda_k^0(s) \right. \\
&\quad \left. \times \exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_k^0 + \boldsymbol{L}_k^T(s)\boldsymbol{\eta}_k^0\right\}\right]ds \\
&= \int_s K_{h_{1n}}(s-t)\mathbb{E}\left[\boldsymbol{X}\left\{\int g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(s) + \epsilon_k(s)\right)\exp\left\{\boldsymbol{L}_k^T(s)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}f(\epsilon_k(s))d\epsilon_k(s) \right.\right. \\
&\quad \left.\left. - \int g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k(t) + \epsilon_k(t)\right)f(\epsilon_k(t))d\epsilon_k(t)\int \exp\left\{\boldsymbol{L}_k^T(s)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}f(\epsilon_k(s))d\epsilon_k(s)\right\}\right. \\
&\quad \left. \times \lambda_k^0(s)\exp\left\{\boldsymbol{X}^T(\boldsymbol{\gamma}_k^0 - \widehat{\boldsymbol{\gamma}}_k)\right\}\right]ds,
\end{aligned}
\tag{A.3}
$$

where $f(\epsilon_k(s))$ is the density function of a normal distribution with mean 0 and variance $c_k$. To simplify the expressions in the following proofs, we denote $\epsilon_k(s)$, $f(\epsilon_k(s))$, and $h_{1n}$ to $\epsilon$, $f_{k,s}(\epsilon)$,

and $h$, respectively. Let $s = t + hz$ and perform the Taylor expansion of (A.3) at $t$. Noticing $\int K(z)dz = 1$ and $\int zK(z)dz = 0$, we have

$$
\mathbb{E}\left[\int K_h(s-t)\boldsymbol{X}[Y_k(s) - E_k(\boldsymbol{\beta}_k(t), t)]d\widetilde{N}_k(s)\right]
$$

$$
= \int_z \frac{K(z)}{h}\lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\left\{\int_\epsilon g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)\exp\left\{\boldsymbol{L}_k^T(t)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}f_{k,t}(\epsilon)d\epsilon\right.\right.
$$

$$
\left.- \int_\epsilon g_k^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_k(t) + \epsilon)f_{k,t}(\epsilon)d\epsilon\int_\epsilon \exp\left\{\boldsymbol{L}_k^T(t)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}f_{k,t}(\epsilon)d\epsilon\right\}
$$

$$
\left.\times\exp\left\{\boldsymbol{X}^T(\boldsymbol{\gamma}_k^0 - \widehat{\boldsymbol{\gamma}}_k)\right\}\right]hdz + O_p\left(h^2\right)
$$

$$
= \lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\exp\left\{\boldsymbol{X}^T(\boldsymbol{\gamma}_k^0 - \widehat{\boldsymbol{\gamma}}_k)\right\}\left\{\int_\epsilon g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)\exp\left\{\boldsymbol{L}_k^T(t)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}f_{k,t}(\epsilon)d\epsilon\right.\right.
$$

$$
\left.\left.- \int_\epsilon g_k^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_k(t) + \epsilon)f_{k,t}(\epsilon)d\epsilon\int_\epsilon \exp\left\{\boldsymbol{L}_k^T(t)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}f_{k,t}(\epsilon)d\epsilon\right\}\right] + O_p\left(h^2\right). \qquad (A.4)
$$

After the Taylor expansion of (A.4) at $\boldsymbol{\gamma}_k^0$, $\boldsymbol{\eta}_k^0$, and $\boldsymbol{\beta}_k^0(t)$, we have

$$
\mathbb{E}\left[\int K_h(s-t)\boldsymbol{X}[Y_k(s) - E_k(\boldsymbol{\beta}_k(t), t)]d\widetilde{N}_k(s)\right]
$$

$$
= \lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\left[1 - \boldsymbol{X}^T\left(\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_k^0\right) + o_p\left(\left\|\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_k^0\right\|\right)\right]\times\left\{\int_\epsilon g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)f_{k,t}(\epsilon)d\epsilon\right.\right.
$$

$$
- \int_\epsilon g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)\boldsymbol{L}_k^T(t)f_{k,t}(\epsilon)d\epsilon\left(\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right) + o_p\left(\left\|\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right\|\right)
$$

$$
- \left[\int_\epsilon g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)f_{k,t}(\epsilon)d\epsilon + \int_\epsilon\left[g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)\right]'\boldsymbol{X}^Tf_{k,t}(\epsilon)d\epsilon\left[\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right]\right.
$$

$$
\left.+o_p\left(\left\|\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right\|\right)\right]\times\left[1 - \int_\epsilon\boldsymbol{L}_k^T(t)f_{k,t}(\epsilon)d\epsilon\left(\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right) + o_p\left(\left\|\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right\|\right)\right]\right\}\right]
$$

$$
+O_p\left(h^2\right)
$$

$$
= -\lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\boldsymbol{X}^T\int_\epsilon\left[g_k^{-1}\left(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon\right)\right]'f_{k,t}(\epsilon)d\epsilon\right]\left[\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right] + A_1(X, t)\left(\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right)
$$

$$
+A_2(X, t)\left[\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right]\left(\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right) + A_3(X, t)\left[\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right]\left(\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_k^0\right)
$$

$$
+A_4(X, t)\left(\boldsymbol{\gamma}_k - \boldsymbol{\gamma}_k^0\right)\left(\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right) + A_5(X, t)\left[\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right]\left(\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_k^0\right)\left(\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right)
$$

$$
+o_p\left(\left\|\boldsymbol{\beta}_k(t) - \boldsymbol{\beta}_k^0(t)\right\|\right) + o_p\left(\left\|\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_k^0\right\|\right) + o_p\left(\left\|\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0\right\|\right) + O_p\left(h^2\right). \qquad (A.5)
$$

From Zeng and Lin (2006), we know

$$
n^{1/2}\left\{(\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{\gamma}_k^0)^T, (\widehat{\boldsymbol{\eta}}_k - \boldsymbol{\eta}_k^0)^T\right\}^T \to_d \mathcal{N}(\boldsymbol{0}, \boldsymbol{V}). \qquad (A.6)
$$

Hence, $n^{1/2} \|\widehat{\gamma}_k - \gamma_k^0\| = O_p(1)$ and $n^{1/2} \|\widehat{\eta}_k - \eta_k^0\| = O_p(1)$. By (A.5) and (A.6), we have

$$
\begin{aligned}
(nh)^{1/2}II &= -(nh)^{1/2}A_k(t)\left[\beta_k(t) - \beta_k^0(t)\right] + o_p\left((nh)^{1/2}\left\|\beta_k(t) - \beta_k^0(t)\right\|\right) \\
&\quad + O_p\left(n^{1/2}h^{5/2}\right) + o_p\left(h^{1/2}\right), \tag{A.7}
\end{aligned}
$$

where $A_k(t) = \lambda_k^0(t)\mathbb{E}\left[\boldsymbol{X}\boldsymbol{X}^T \int_\epsilon \left[g_k^{-1}\left(\boldsymbol{X}^T\beta_k^0(t) + \epsilon\right)\right]' f_{k,t}(\epsilon)d\epsilon\right]$. For any fixed $t$, if there exists a $\boldsymbol{\zeta}$ such that $\boldsymbol{\zeta}^T A_k(t)\boldsymbol{\zeta} = 0$, then $\boldsymbol{\zeta}^T \boldsymbol{X} = 0$, so $\boldsymbol{\zeta} = 0$ from condition 2. Thus, $A_k(t)$ is a non-singular matrix. Furthermore, using the similar arguments in Cao et al. (2015), we could obtain the first term in the right-hand side of (A.2), for a $M > 0$ and $\|\beta_k(t) - \beta_k^0(t)\| < M(nh)^{-1/2}$, is equal to

$$
\begin{aligned}
&(nh)^{1/2}(\boldsymbol{P}_n - \boldsymbol{P})\left[\int K_h(s-t)\boldsymbol{X}\left[Y_k(s) - E_k(\beta_k^0(t), t)\right]d\widetilde{N}_k(s)\right] + o_p(1) \\
&= (nh)^{1/2}\left\{U_{n,k}(\beta_k^0(t)) - \mathbb{E}\left[U_{n,k}(\beta_k^0(t))\right]\right\} + o_p(1). \tag{A.8}
\end{aligned}
$$

Combining (A.7) and (A.8), we obtain (A.1). From (A.1), we conclude that there exists a solution to $U_{n,k}(\beta_k(t)) = 0$, say, $\widehat{\beta}_k(t)$, which is $(nh)^{-1/2}$ consistent. Moreover,

$$
\begin{aligned}
&(nh)^{1/2}\left\{U_{n,k}(\beta_k^0(t) - \mathbb{E}\left[U_{n,k}(\beta_k^0(t))\right]\right\} \\
&= (nh)^{1/2}A_k(t)\left[\widehat{\beta}_k(t) - \beta_k^0(t)\right] + o_p\left((nh)^{1/2}\left\|\widehat{\beta}_k(t) - \beta_k^0(t)\right\|\right) \\
&\quad + O_p\left(n^{1/2}h^{5/2}\right) + o_p\left(1 + h^{1/2}\right). \tag{A.9}
\end{aligned}
$$

It remains to obtain the distribution of $(nh)^{1/2}\left\{U_{n,k}(\beta_k^0(t) - \mathbb{E}\left[U_{n,k}(\beta_k^0(t))\right]\right\}$. For $s, t \in [0, \tau]$, where $\tau$ is the maximum observation time, denote

$$
\begin{aligned}
W_i(t) &= \sqrt{\frac{h}{n}}\left\{\int K_h(s-t)\boldsymbol{X}_i\left[Y_{ik}(s) - E_{ik}(\beta_k^0(t), t)\right]d\widetilde{N}_k(s)\right. \\
&\quad \left. - E\left[\int K_h(s-t)\boldsymbol{X}_i\left[Y_{ik}(s) - E_{ik}(\beta_k^0(t), t\right]d\widetilde{N}_k(s)\right]\right\}. \tag{A.10}
\end{aligned}
$$

Then $(nh)^{1/2}\left\{U_{n,k}(\beta_k^0(t) - \mathbb{E}[U_{n,k}(\beta_k^0(t))]\right\} = \sum_{i=1}^n W_i(t)$. Since $(\boldsymbol{X}_i, Y_{ik})$, $i = 1, \ldots, n$ are i.i.d, we can calculate $\mathrm{Var}(\sum_{i=1}^n W_i(t)) = \sum_{i=1}^n \mathrm{Var}(W_i(t))$ as follows:

$$
\sum_{i=1}^n \mathrm{Var}\left[W_i(t)\right] = \sum_{i=1}^n \left\{\mathbb{E}\left[\mathrm{Var}\left[W_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]\right] + \mathrm{Var}\left[\mathbb{E}\left[W_i(t)|\boldsymbol{X}, \epsilon_{ik}(s), N_{ik}(s)\right]\right]\right\}. \tag{A.11}
$$

Since

$$\text{Var}\left[W_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]$$

$$= \frac{h}{n}\text{Var}\left[\int K_h(s-t)\boldsymbol{X}\left[Y_{ik}(s) - E_{ik}(\boldsymbol{\beta}_k^0(t), t)\right]d\widetilde{N}_{ik}(s)\Big|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]$$

$$= \frac{h}{n}\text{Var}\left[\int K_h(s-t)\boldsymbol{X}_i Y_{ik}(s)d\widetilde{N}_{ik}(s)\Big|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]$$

$$= \frac{h}{n}\left[\iint K_h(s_1-t)K_h(s_2-t)\boldsymbol{X}_i\boldsymbol{X}_i^T\mathbb{E}\left[Y_{ik}(s_1)\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s_1)\widehat{\boldsymbol{\eta}}_k\right\}\right.\right.$$

$$\left.\left.\times Y_{ik}(s_2)\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s_2)\widehat{\boldsymbol{\eta}}_k\right\}|\boldsymbol{X}_i, \epsilon_{ik}(s_1), \epsilon_{ik}(s_2)\right]dN_{ik}(s_1)dN_{ik}(s_2)\right]$$

$$-\frac{h}{n}\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\int K_h(s-t)\mathbb{E}\left[Y_{ik}(s)\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s)\widehat{\boldsymbol{\eta}}_k\right\}|\boldsymbol{X}_i, \epsilon_{ik}(s)\right]dN_{ik}(s)\right]^2 \text{(A.12)}$$

We assume, for $s_1 \neq s_2$, $\boldsymbol{P}\{dN_k(s_1) = 1|N_k(s_2) - N_k(s_2-) = 1\} = p_k(s_1, s_2)ds_1$, where $p_k(s_1, s_2)$ is continuous for $s_1 \neq s_2$, and $p_k(s_1\pm, s_2\pm)$ exists. To simplify the expressions in (A.12), we denote

$$\begin{aligned}
F_k(t) &= Y_{ik}(t)\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - L_{ik}^T(t)\widehat{\boldsymbol{\eta}}_k\right\}, \\
G_k(t) &= g_k^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_k^0(t) + \epsilon_{ik}(t))\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(t)\widehat{\boldsymbol{\eta}}_k\right\}, \\
\Lambda_k(t) &= \lambda_k^0(t)\exp\left\{\boldsymbol{X}_i^T\boldsymbol{\gamma}_k^0 + \boldsymbol{L}_{ik}^T(t)\boldsymbol{\eta}_k^0\right\}.
\end{aligned} \qquad \text{(A.13)}$$

We further denote

$$\begin{aligned}
\text{Var}\left[F_k(t)|\boldsymbol{X}, \epsilon_k(t)\right] &= \sigma^2\left(t, \boldsymbol{X}, \epsilon_k(t)\right), \\
\text{Cov}\left[F_k(s), F_k(t)|\boldsymbol{X}, \epsilon_k(s), \epsilon_k(t)\right] &= r\left(s, t, \boldsymbol{X}, \epsilon_k(s), \epsilon_k(t)\right).
\end{aligned} \qquad \text{(A.14)}$$

Using conditioning arguments, the expectation of (A.12) is

$$
\mathbb{E}\left[\mathrm{Var}\left[W_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]\right]
$$

$$
= \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T \iint K_h(s_1-t)K_h(s_2-t)\left[r(s_1,s_2,\boldsymbol{X}_i,\epsilon_{ik}(s_1),\epsilon_{ik}(s_2)) + G_k(s_1)G_k(s_2)\right]\right.
$$

$$
\left. \times dN_{ik}(s_1)dN_{ik}(s_2)\right] - \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\int K_h(s-t)G_k(s)dN_{ik}(s)\right]^2\right]
$$

$$
= \frac{h}{n}\int_{s_2}\int_{s_1}K_h(s_1-t)K_h(s_2-t)\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[r(s_1,s_2,\boldsymbol{X}_i,\epsilon_{ik}(s_1),\epsilon_{ik}(s_2)) + G_k(s_1)G_k(s_2)\right]\right.
$$

$$
\left. \times p_k(s_1,s_2)\Lambda_k(s_2)\right]ds_1 ds_2
$$

$$
+ \frac{h}{n}\int_{s_1}K_h^2(s_1-t)\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\sigma^2(s_1,\boldsymbol{X}_i,\epsilon_{ik}(s_1)) + G_k^2(s_1)\right]\Lambda_k(s_1)\right]ds_1
$$

$$
- \frac{h}{n}\int_s K_h^2(s-t)\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T G_k^2(s)\Lambda_k(s)\right]ds
$$

$$
= I_{11} + I_{12} - I_2. \tag{A.15}
$$

Let $s_1 = t + hz_1$, $s_2 = t + hz_2$, and $s = t + hz$ , do the Taylor expansion of (A.15) at $t$. Since $\iint K(z_1)K(z_2)dz_1 dz_2 = 1$, $\iint z_1 K(z_1)K(z_2)dz_1 dz_2 = \iint z_2 K(z_1)K(z_2)dz_1 dz_2 = 0$, $\int K(z)dz = 1$, $\int zK(z)dz = 0$, and $p_k(t,t) = 1$, we have

$$
\begin{aligned}
I_{11} &= \frac{h^3}{n}\int_{z_2}\int_{z_1}\frac{K(z_1)K(z_2)}{h^2}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\sigma^2(t,\boldsymbol{X}_i,\epsilon_{ik}(t)) + G_k^2(t)\right]p_k(t,t)\Lambda_k(t)\right]dz_1 dz_2 \\
&\quad + O_p\left(h^2\right) \\
&= \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\sigma^2(t,\boldsymbol{X}_i,\epsilon_{ik}(t)) + G_k^2(t)\right]\Lambda_k(t)\right] + O_p\left(h^2\right) \\
&= C_1\frac{h}{n} + O_p\left(h^2\right), \\
I_{12} &= \frac{h^2}{n}\int_{z_1}\frac{K^2(z_1)}{h^2}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\sigma^2(t,\boldsymbol{X}_i,\epsilon_{ik}(t)) + G_k^2(t)\right]\Lambda_k(t)\right]dz_1 + O_p\left(h\right) \\
&= \frac{1}{n}\boldsymbol{X}_i\boldsymbol{X}_i^T\int_z K^2(z)dz\,\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\sigma^2(t,\boldsymbol{X}_i,\epsilon_{ik}(t)) + G_k^2(t)\right]\Lambda_k(t)\right] + O_p\left(h\right), \\
I_2 &= \frac{h^2}{n}\int_z\frac{K^2(z)}{h^2}\mathbb{E}[\boldsymbol{X}_i\boldsymbol{X}_i^T G_k^2(t)\Lambda_k(t)]dz + O_p\left(h\right) \\
&= \frac{1}{n}\int_z K^2(z)dz\,\mathbb{E}[\boldsymbol{X}_i\boldsymbol{X}_i^T G_k^2(t)\Lambda_k(t)] + O_p\left(h\right). \tag{A.16}
\end{aligned}
$$

By (A.12), (A.16), and condition 3, we have

$$
\mathbb{E}\left[\operatorname{Var}\left[W_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]\right]
$$
$$
= \frac{1}{n}\int_z K^2(z)dz\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\sigma^2(t, \boldsymbol{X}_i, \epsilon_{ik}(t))\Lambda_k(t)\right] + O_p\left(h\right). \tag{A.17}
$$

For the second term on the right-hand side of (A.11), we have

$$
\operatorname{Var}\left[\mathbb{E}\left[W_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]\right]
$$
$$
= \frac{h}{n}\operatorname{Var}\left[\int K_h(s-t)\boldsymbol{X}_i\mathbb{E}\left[[Y_{ik}(s) - E_{ik}(\beta_k^0(t),t)]\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s)\widehat{\boldsymbol{\eta}}_k\right\}\Big|\boldsymbol{X}_i, \epsilon_{ik}(s)\right]\right.
$$
$$
\left. \times dN_{ik}(s)\right]
$$
$$
= \frac{h}{n}\operatorname{Var}\left[\int K_h(s-t)\boldsymbol{X}_i\left[g_k^{-1}(\boldsymbol{X}_i^T\beta_k^0(s) + \epsilon_{ik}(s)) - g_k^{-1}(\boldsymbol{X}_i^T\beta_k^0(t) + \epsilon_{ik}(t))\right]\right.
$$
$$
\left. \times \exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s)\widehat{\boldsymbol{\eta}}_k\right\}dN_{ik}(s)\right]. \tag{A.18}
$$

Denote $D_k(s,t) = \exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s)\widehat{\boldsymbol{\eta}}_k\right\}[g_k^{-1}(\boldsymbol{X}_i^T\beta_k^0(s) + \epsilon_{ik}(s)) - g_k^{-1}(\boldsymbol{X}_i^T\beta_k^0(t) + \epsilon_{ik}(t))]$, and then (A.18) can be calculated as

$$
\frac{h}{n}\operatorname{Var}\left[\int K_h(s-t)\boldsymbol{X}_iD(s,t)dN_{ik}(s)\right]
$$
$$
= \frac{h}{n}\mathbb{E}\left[\iint K_h(s_1-t)K_h(s_2-t)\boldsymbol{X}_i\boldsymbol{X}_i^TD_k(s_1,t)D_k(s_2,t)dN_{ik}(s_1)dN_{ik}(s_2)\right]
$$
$$
- \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\int K_h(s-t)D_k(s,t)dN_{ik}(s)\right]^2\right]
$$
$$
= \frac{h}{n}\int_{s_2}\int_{s_1} K_h(s_1-t)K_h(s_2-t)\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^TD_k(s_1,t)D_k(s_2,t)p_k(s_1,s_2)\Lambda_k(s_2)\right]ds_1ds_2
$$
$$
+ \frac{h}{n}\int_{s_1} K_h^2(s_1-t)\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^TD_k^2(s_1,t)\Lambda_k(s_1)\right]ds_1
$$
$$
- \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\int_s K_h(s-t)D_k(s,t)\Lambda_k(s)ds\right]^2\right]
$$
$$
= I_{31} + I_{32} - I_4. \tag{A.19}
$$

Letting $s_1 = t+hz_1$, $s_2 = t+hz_2$, and $s = t+hz$ , we do the Taylor expansion of (A.19) at $t$. Since $\iint K(z_1)K(z_2)dz_1dz_2 = 1$, $\iint z_1K(z_1)K(z_2)dz_1dz_2 = \iint z_2K(z_1)K(z_2)dz_1dz_2 = 0$, $\int K(z)dz = 1$,

$\int zK(z)dz = 0$, $p_k(t,t) = 1$, and $D_k(t,t) = 0$, we have

$$
\begin{aligned}
I_{31} &= \frac{h}{n}\int_{z_2}\int_{z_1}\frac{K(z_1)K(z_2)}{h^2}\mathbb{E}[\boldsymbol{X}_i\boldsymbol{X}_i^T D_k^2(t,t)p_k(t,t)\Lambda_k(t)]h^2 dz_1 dz_2 + O_p\left(h^2\right)\\
&= \frac{h}{n}\mathbb{E}[\boldsymbol{X}_i\boldsymbol{X}_i^T \cdot 0 \cdot \Lambda_k(t)] + O_p\left(h^2\right)\\
&= O_p\left(h^2\right),\\
I_{32} &= \frac{h}{n}\int_{z_1}\frac{K^2(z_1)}{h^2}\mathbb{E}[\boldsymbol{X}_i\boldsymbol{X}_i^T D_k(t,t)\Lambda_k(t)]^2 h dz_1 + O_p\left(h\right)\\
&= \frac{1}{n}\int_z K^2(z)dz\mathbb{E}[\boldsymbol{X}_i\boldsymbol{X}_i^T \cdot 0 \cdot \Lambda_k(t)]^2 + O_p\left(h\right)\\
&= O_p\left(h\right),\\
I_4 &= \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T\left[\int_z\frac{K(z)}{h}D_k(t,t)\Lambda_k(t)hdz\right]^2\right]\\
&= \frac{h}{n}\mathbb{E}\left[\boldsymbol{X}_i\boldsymbol{X}_i^T \cdot 0^2\right] + O_p\left(h^2\right)\\
&= O_p\left(h^2\right).
\end{aligned}
\tag{A.20}
$$

By (A.18), (A.20), and condition 3, we have

$$
\mathrm{Var}\left[\mathbb{E}\left[W_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s), N_{ik}(s)\right]\right] = O_p\left(h\right).
\tag{A.21}
$$

Therefore, by (A.17) and (A.21), we obtain

$$
\begin{aligned}
&\sum_{i=1}^n \mathrm{Var}\left[W_i(t)\right]\\
&= \lambda_k^0(t)\int_z K^2(z)dz\mathbb{E}[\boldsymbol{X}\boldsymbol{X}^T\sigma^2(t,\boldsymbol{X},\epsilon_k(t))\exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_k^0 + \boldsymbol{L}_k^T(t)\boldsymbol{\eta}_k^0\right\}] + O_p\left(h\right)\\
&= \boldsymbol{\Sigma}_k(t) + O_p\left(h\right).
\end{aligned}
\tag{A.22}
$$

Thus, for each fixed time point $t$, $\sum_{i=1}^n \mathrm{Var}\left[W_i(t)\right]$ converges to a constant $\boldsymbol{\Sigma}_k(t)$ as $h \to 0$. To prove the asymptotic normality, we verify the Lyapunov condition. Through the similar calculations to $\boldsymbol{\Sigma}_k(t)$, we could obtain

$$
\sum_{i=1}^n \mathbb{E}\left\{\|W_i - \mathbb{E}(W_i)\|^3\right\} = nO_p\left\{(nh)^{3/2}n^{-3}h^{-2}\right\} = O_p\left\{(nh)^{-1/2}\right\}.
$$

Therefore, by the Lyaponov central limit theorem,

$$(nh)^{1/2} \left\{ U_{n,k}(\boldsymbol{\beta}_k^0(t)) - \mathbb{E}[U_{n,k}(\boldsymbol{\beta}_k^0(t))] \right\} \to_d \mathcal{N}\left(\mathbf{0}, \boldsymbol{\Sigma}_k(t)\right). \tag{A.23}$$

By proving (A.9) and (A.23), we conclude that for any fixed $t$,

$$(nh)^{1/2} \left[ \widehat{\boldsymbol{\beta}}_k(t) - \boldsymbol{\beta}_k^0(t) \right] \to_d \mathcal{N}\left(\mathbf{0}, \left[ A_k^{-1}(t) \right] \boldsymbol{\Sigma}(t) \left[ A_k^{-1}(t) \right]^T \right). \tag{A.24}$$

## A.3  Proof of Theorem 2.3.2

Without loss of generality, we prove the theorem 2.3.2 for two health markers $k$ and $l$, and the results can be generalized to other health markers. The key idea is to establish the following relationship, for $\left| \sigma_{kl}(t) - \sigma_{kl}^0(t) \right| < M(nh_{2n}^2)^{-1/2}$,

$$
\begin{aligned}
& \sup_{\sigma_{kl}(t)} (nh_{2n}^2)^{1/2} \left| U_{n,k,l}(\sigma_{kl}(t)) - \left\{ U_{n,k,l}(\sigma_{kl}^0(t) - \mathbb{E}[U_{n,k,l}(\sigma_{kl}^0(t))] \right\} - B_{kl}(t) \left| \sigma_{kl}(t) - \sigma_{kl}^0(t) \right| \right| \\
= \ & O_p(n^{1/2}h_{2n}^3) + o_p\left( 1 + h_{2n} \left| \sigma_{kl}(t) - \sigma_{kl}^0(t) \right| \right),
\end{aligned}
\tag{A.25}
$$

where $B_{kl}(t)$ is defined in (2.10).

We introduce notations $\boldsymbol{P}_n$ and $\boldsymbol{P}$ to denote the empirical and true probability measure respectively. Then we obtain

$$
\begin{aligned}
U_{n,k,l}(\sigma_{kl}(t)) &= (\boldsymbol{P}_n - \boldsymbol{P}) \left[ \iint \widetilde{K}_{h_{2n}}(s_k - t, s_l - t) \left[ Y_k(s_k)Y_l(s_l) - E_{kl}(\sigma_{kl}(t), t) \right] \right. \\
& \qquad \left. \times d\widetilde{N}_k(s_k)d\widetilde{N}_l(s_l) \right] \\
& \quad + \mathbb{E}\left[ \iint \widetilde{K}_{h_{2n}}(s_k - t, s_l - t) \left[ Y_k(s_k)Y_l(s_l) - E_{kl}(\sigma_{kl}(t), t) \right] d\widetilde{N}_k(s_k)d\widetilde{N}_l(s_l) \right] \\
&= I + II, \tag{A.26}
\end{aligned}
$$

where $E_{kl}(\sigma_{kl}(t), t) = \mathbb{E}[Y_{ik}(t)Y_{il}(t)|\boldsymbol{X}_i] = \mathbb{E}_\epsilon \left[ g_k^{-1}(\boldsymbol{X}^T\widehat{\boldsymbol{\beta}}_k(t) + \epsilon_k(t))g_l^{-1}(\boldsymbol{X}^T\widehat{\boldsymbol{\beta}}_l(t) + \epsilon_l(t)) \right].$

80

For the second term on the right-hand side of (A.26), we have

$$
\mathbb{E}\left[\iint \widetilde{K}_{h_{2n}}(s_k - t, s_l - t)\left[Y_k(s_k)Y_l(s_l) - E_{kl}(\sigma_{kl}(t), t)\right] d\widetilde{N}_k(s_k)d\widetilde{N}_l(s_l)\right]
$$

$$
= \iint \widetilde{K}_{h_{2n}}(s_k - t, s_l - t)\mathbb{E}\Big[\left[Y_k(s_k)Y_l(s_l) - E_{kl}(\sigma_{kl}(t), t)\right]
$$

$$
\times \exp\left\{-\boldsymbol{X}^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_k^T(s_k)\widehat{\boldsymbol{\eta}}_k\right\}\exp\left\{-\boldsymbol{X}^T\widehat{\boldsymbol{\gamma}}_l - \boldsymbol{L}_l^T(s_l)\widehat{\boldsymbol{\eta}}_l\right\} dN_k(s_k)dN_l(s_l)\Big]
$$

$$
= \int_{s_l}\int_{s_k}\widetilde{K}_{h_{2n}}(s_k - t, s_l - t)\mathbb{E}\Big[\left[g_k^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(s_k) + \epsilon_k(s_k))g_l^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_l^0(s_l) + \epsilon_l(s_l))\right.
$$

$$
-g_k^{-1}(\boldsymbol{X}^T\widehat{\boldsymbol{\beta}}_k(t) + \epsilon_k(t))g_l^{-1}(\boldsymbol{X}^T\widehat{\boldsymbol{\beta}}_l(t) + \epsilon_l(t))]
$$

$$
\times\lambda_k^0(s_k)\exp\left\{\boldsymbol{X}^T(\boldsymbol{\gamma}_k^0 - \widehat{\boldsymbol{\gamma}}_k) + \boldsymbol{L}_k^T(s_k)(\boldsymbol{\eta}_k^0 - \widehat{\boldsymbol{\eta}}_k)\right\}
$$

$$
\times\lambda_l^0(s_l)\exp\left\{\boldsymbol{X}^T(\boldsymbol{\gamma}_l^0 - \widehat{\boldsymbol{\gamma}}_l) + \boldsymbol{L}_l^T(s_l)(\boldsymbol{\eta}_l^0 - \widehat{\boldsymbol{\eta}}_l)\right\}\Big]ds_k ds_l.
$$

$$(A.27)$$

To simplify the expressions in the following proofs, we denote $h_{2n}$ to $h$. Using the similar calculations as those in proving theorem 2.3.1, we apply the Taylor expansion of (A.27) at $\boldsymbol{\gamma}_k^0(t)$, $\boldsymbol{\gamma}_l^0(t)$, $\boldsymbol{\eta}_k^0(t)$, $\boldsymbol{\eta}_l^0(t)$, $\boldsymbol{\beta}_k^0(t)$, $\boldsymbol{\beta}_l^0(t)$, and $\sigma_{kl}^0(t)$. Combining the conclusions in (A.6) and (A.24), we obtain

$$
(nh^2)^{1/2}II
$$

$$
= -(nh^2)^{1/2}\lambda_k^0(t)\lambda_l^0(t)\iint g_k^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_k^0(t) + \epsilon_k(t))g_l^{-1}(\boldsymbol{X}^T\boldsymbol{\beta}_l^0(t) + \epsilon_l(t))
$$

$$
\times\frac{\partial f\left(\epsilon_k(t), \epsilon_l(t); \sigma_{kl}(t)\right)}{\partial\sigma_{kl}(t)}\Big|_{\sigma_{kl}(t)=\sigma_{kl}^0(t)}d\epsilon_k(t)d\epsilon_l(t)\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right|
$$

$$
+O_p(n^{1/2}h^3) + o_p\left(h\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right|\right)
$$

$$
\equiv -(nh^2)^{1/2}B_{kl}(t)\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right| + O_p(n^{1/2}h^3) + o_p\left(h\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right|\right). \quad (A.28)
$$

Furthermore, we could obtain the first term in the right-hand side of (A.26) for a $M > 0$ and $\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right| < M(nh^2)^{-1/2}$ is equal to

$$
(nh^2)^{1/2}\left\{U_{n,k,l}(\sigma_{kl}^0(t)) - \mathbb{E}[U_{n,k,l}(\sigma_{kl}^0(t))]\right\} + o_p(1)
$$

$$
= (nh^2)^{1/2}(\boldsymbol{P}_n - \boldsymbol{P})\left[\iint \widetilde{K}_h(s_k - t, s_l - t)\left[Y_k(s_k)Y_l(s_l) - E_{kl}(\sigma_{kl}^0(t), t)\right]\right.
$$

$$
\times d\widetilde{N}_k(s_k)d\widetilde{N}_l(s_l)\bigg] + o_p(1). \quad (A.29)
$$

Combining (A.28) and (A.29), we obtain (A.25). From (A.25), we conclude that there exists a solution to $U_{n,k,l}(\sigma_{kl}(t)) = 0$, say, $\hat{\sigma}_{kl}(t)$, is $(nh^2)^{-1/2}$ consistent and moreover,

$$(nh^2)^{1/2}\left\{U_{n,k,l}(\sigma_{kl}^0(t)) - \mathbb{E}[U_{n,k,l}(\sigma_{kl}^0(t))]\right\}$$
$$= (nh^2)^{1/2}B_{kl}(t)\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right| + O_p(n^{1/2}h^3) + o_p\left(1 + h\left|\sigma_{kl}(t) - \sigma_{kl}^0(t)\right|\right). \quad \text{(A.30)}$$

It remains to obtain the distribution of $(nh^2)^{1/2}\{U_{n,k,l}(\sigma_{kl}^0(t)) - \mathbb{E}[U_{n,k,l}(\sigma_{kl}^0(t))]\}$. Denote

$$V_i(t) = \sqrt{\frac{h^2}{n}}\left\{\iint \widetilde{K}_h(s_k - t, s_l - t)\left[Y_{ik}(s_k)Y_{il}(s_l) - E_{ikl}(\sigma_{kl}^0(t), t)\right]d\tilde{N}_{ik}(s_k)d\tilde{N}_{il}(s_l)\right.$$
$$\left. - \mathbb{E}\left[\iint \widetilde{K}_h(s_k - t, s_l - t)\left[Y_{ik}(s_k)Y_{il}(s_l) - E_{ikl}(\sigma_{kl}^0(t), t)\right]d\tilde{N}_{ik}(s_k)d\tilde{N}_{il}(s_l)\right]\right\}. $$
$$\text{(A.31)}$$

Then $(nh^2)^{1/2}\{U_{n,k,l}(\sigma_{kl}^0(t)) - \mathbb{E}[U_{n,k,l}(\sigma_{kl}^0(t))]\} = \sum_{i=1}^n V_i(t)$. Since $(\boldsymbol{X}_i, Y_{ik})$, $i = 1, \ldots, n$ are i.i.d, we can calculate $\text{Var}\left[\sum_{i=1}^n V_i(t)\right] = \sum_{i=1}^n \text{Var}\left[V_i(t)\right]$ as follows:

$$\sum_{i=1}^n \text{Var}\left[V_i(t)\right] = \sum_{i=1}^n \mathbb{E}\left[\text{Var}\left[V_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s_k), \epsilon_{ik}(s_l), N_{ik}(s_k), N_{il}(s_l)\right]\right]$$
$$+ \sum_{i=1}^n \text{Var}\left[\mathbb{E}\left[V_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s_k), \epsilon_{ik}(s_l), N_{ik}(s_k), N_{il}(s_l)\right]\right]. \quad \text{(A.32)}$$

For $t_1, t_2, s_1, s_2 \in [0, \tau]$, we denote

$$\text{Var}\left[F_k(t_1)F_l(t_2)|\boldsymbol{X}, \epsilon_k(t_1), \epsilon_l(t_2)\right] = \psi^2(t_1, t_2, \boldsymbol{X}, \epsilon_k(t_1), \epsilon_l(t_2)) \quad \text{(A.33)}$$

and

$$\text{Cov}\left[F_k(t_1)F_l(t_2), F_k(s_1)F_l(s_2)|\boldsymbol{X}, \epsilon_k(t_1), \epsilon_l(t_2), \epsilon_k(s_1), \epsilon_l(s_2)\right]$$
$$= u(t_1, t_2, s_1, s_2, \boldsymbol{X}, \epsilon_k(t_1), \epsilon_l(t_2), \epsilon_k(s_1), \epsilon_l(s_2)). \quad \text{(A.34)}$$

Similarly to the calculation of $\boldsymbol{\Sigma}_k(t)$, we obtain

$$
\mathbb{E}\left[\text{Var}\left[V_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s_k), \epsilon_{ik}(s_l), N_{ik}(s_k), N_{il}(s_l)\right]\right]
$$

$$
= \frac{h^2}{n}\mathbb{E}\left[\text{Var}\left[\iint \widetilde{K}_h(s_k - t, s_l - t)Y_{ik}(s_k)Y_{il}(s_l)\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s_k)\widehat{\boldsymbol{\eta}}_k\right\}\right.\right.
$$

$$
\left.\left.\times\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_l - \boldsymbol{L}_{il}^T(s_l)\widehat{\boldsymbol{\eta}}_l\right\}dN_{ik}(s_k)dN_{il}(s_l)\right]\right]
$$

$$
= \frac{h^2}{n}\mathbb{E}\left[\iiiint \widetilde{K}_h(t_k - t, t_l - t)\widetilde{K}_h(s_k - t, s_l - t)\left[u(t_k, t_l, s_k, s_l, \boldsymbol{X}_i, \epsilon_{ik}(t_k), \epsilon_{il}(t_l),\right.\right.
$$

$$
\epsilon_{ik}(s_k), \epsilon_{il}(s_l)) + \{G_k(t_k)G_l(t_l) + r(t_k, t_l, \boldsymbol{X}_i, \epsilon_{ik}(t_k), \epsilon_{il}(t_l))\}\{G_k(s_k)G_l(s_l)
$$

$$
\left.\left.+r(s_k, s_l, \boldsymbol{X}_i, \epsilon_{ik}(s_k), \epsilon_{il}(s_l))\}\right] \times dN_{ik}(t_k)dN_{il}(t_l)dN_{ik}(s_k)dN_{il}(s_l)\right]
$$

$$
-\frac{h^2}{n}\mathbb{E}\left[\left[\iint \widetilde{K}_h(s_k - t, s_l - t)\{G_k(s_k)G_l(s_l) + r(s_k, s_l, \boldsymbol{X}_i, \epsilon_{ik}(s_k), \epsilon_{il}(s_l))\}\right.\right.
$$

$$
\left.\left.\times dN_{ik}(s_k)dN_{il}(s_l)\right]^2\right]
$$

$$
= \frac{h^2}{n}\int_{z_2}\int_{z_1}\frac{\widetilde{K}^2(z_1, z_2)}{h^4}\mathbb{E}\left[\psi^2(t, t, \boldsymbol{X}_i, \epsilon_{ik}(t), \epsilon_{il}(t))\Lambda_k(t)\Lambda_l(t)\right]h^2dz_1dz_2 + O_p\left(h^2\right)
$$

$$
= \frac{1}{n}\int_{z_2}\int_{z_1}\widetilde{K}^2(z_1, z_2)dz_1dz_2\mathbb{E}\left[\psi^2(t, t, \boldsymbol{X}_i, \epsilon_{ik}(t), \epsilon_{il}(t))\Lambda_k(t)\Lambda_l(t)\right] + O_p\left(h^2\right). \tag{A.35}
$$

Denote

$$
D_{kl}(s_k, s_l, t)
$$

$$
= \exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s_k)\widehat{\boldsymbol{\eta}}_k\right\}\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_l - \boldsymbol{L}_{il}^T(s_l)\widehat{\boldsymbol{\eta}}_l\right\}\left[g_k^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_k^0(s_k) + \epsilon_{ik}(s_k))\right.
$$

$$
\left.\times g_l^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_l^0(s_l) + \epsilon_{il}(s_l)) - g_k^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_k^0(t) + \epsilon_{ik}(t))g_l^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_l^0(t) + \epsilon_{il}(t))\right].
$$

We have

$$
\begin{aligned}
& \mathrm{Var}\left[\mathbb{E}[V_i(t)|\boldsymbol{X}_i, \epsilon_{ik}(s_k), \epsilon_{ik}(s_l), N_{ik}(s_k), N_{il}(s_l)]\right] \\
= \quad & \frac{h^2}{n}\mathrm{Var}\left[\iint \widetilde{K}_h(s_k - t, s_l - t)\exp\left\{-\boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_k - \boldsymbol{L}_{ik}^T(s_k)\widehat{\boldsymbol{\eta}}_k - \boldsymbol{X}_i^T\widehat{\boldsymbol{\gamma}}_l - \boldsymbol{L}_{il}^T(s_l)\widehat{\boldsymbol{\eta}}_l\right\}\right. \\
& \times\left[g_k^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_k^0(s_k) + \epsilon_{ik}(s_k))g_l^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_l^0(s_l) + \epsilon_{il}(s_l)) - g_k^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_k^0(t) + \epsilon_{ik}(t))\right. \\
& \left.\left.\times g_l^{-1}(\boldsymbol{X}_i^T\boldsymbol{\beta}_l^0(t) + \epsilon_{il}(t))\right]dN_{ik}(s_k)dN_{il}(s_l)\right] \\
= \quad & \frac{h^2}{n}\mathbb{E}\left[\iiiint \widetilde{K}_h(t_k - t, t_l - t)\widetilde{K}_h(s_k - t, s_l - t)D_{kl}(t_k, t_l, t)D_{kl}(s_k, s_l, t)\right. \\
& \left.\times dN_{ik}(t_k)dN_{il}(t_l)dN_{ik}(s_k)dN_{il}(s_l)\right] \\
& - \frac{h^2}{n}\left[\int_{s_l}\int_{s_k}\widetilde{K}_h(s_k - t, s_l - t)\mathbb{E}[D_{kl}(s_k, s_l, t)\Lambda_k(s_k)\Lambda_l(s_l)]ds_kds_l\right]^2 \\
= \quad & O_p\left(h^2\right).
\end{aligned}
\tag{A.36}
$$

Therefore, by (A.35) and (A.36) , we obtain

$$
\begin{aligned}
& \sum_{i=1}^n Var(W_i(t)) \\
= \quad & \mathbb{E}[\psi^2(t, t, \boldsymbol{X}, \epsilon_k(t), \epsilon_l(t))\exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_k^0 + \boldsymbol{L}_k^T(t)\boldsymbol{\eta}_k^0\right\}\exp\left\{\boldsymbol{X}^T\boldsymbol{\gamma}_l^0 + \boldsymbol{L}_l^T(t)\boldsymbol{\eta}_l^0\right\}] \\
& \times\lambda_k^0(t)\lambda_l^0(t)\int_{z_2}\int_{z_1}\widetilde{K}^2(z_1, z_2)dz_1dz_2 + O_p(h^2) \\
= \quad & \Sigma_{kl}(t) + O_p(h^2).
\end{aligned}
\tag{A.37}
$$

Thus, for each fixed time point $t$, $\sum_{i=1}^n \mathrm{Var}[W_i(t)]$ converges to a constant $\Sigma_{kl}(t)$ as $h \to 0$. To prove the asymptotic normality, we verify the Lyapunov condition. Through the similar calculations to $\Sigma_{kl}(t)$, we could obtain

$$
\sum_{i=1}^n \mathbb{E}\left\{|V_i - \mathbb{E}(V_i)|^3\right\} = nO_p\left\{(nh^2)^{3/2}n^{-3}(h^2)^{-2}\right\} = O_p\left\{(nh^2)^{-1/2}\right\}.
$$

Therefore, by the Lyaponov central limit theorem,

$$
(nh^2)^{1/2}\left\{U_{n,k,l}(\sigma_{kl}^0(t)) - \mathbb{E}[U_{n,k,l}(\sigma_{kl}^0(t))]\right\} \to_d \mathcal{N}(0, \Sigma_{kl}(t)).
\tag{A.38}
$$

By proving (A.30) and (A.38), we conclude that for any fixed $t$,

$$(nh^2)^{1/2} \left[ \widehat{\sigma}_{kl}(t) - \sigma_{kl}^0(t) \right] \to_d \mathcal{N} \left( 0, B_{kl}^{-2}(t) \Sigma_{kl}(t) \right). \tag{A.39}$$

## A.4 Validation on Assumption of Constant Variance

Table A.1: Estimated variances using different 2-year data from EHRs at the OSU-WMCIW.

| Marker | Start Time | End Time | Estimated Variance |
|---|---|---|---|
| HBP | **01/01/2011** | **12/31/2012** | **1.9177** |
| | 01/01/2012 | 12/31/2013 | 1.8068 |
| | 01/01/2013 | 12/31/2014 | 1.9136 |
| | 01/01/2014 | 12/31/2015 | 1.9112 |
| | 01/01/2015 | 12/31/2016 | 1.8812 |
| | 01/01/2016 | 12/31/2017 | 1.8301 |
| TC | **01/01/2011** | **12/31/2012** | **0.4766** |
| | 01/01/2012 | 12/31/2013 | 0.5346 |
| | 01/01/2013 | 12/31/2014 | 0.5120 |
| | 01/01/2014 | 12/31/2015 | 0.5192 |
| | 01/01/2015 | 12/31/2016 | 0.5502 |
| | 01/01/2016 | 12/31/2017 | 0.5631 |
| HbA1c | **01/01/2011** | **12/31/2012** | **0.6387** |
| | 01/01/2012 | 12/31/2013 | 0.7135 |
| | 01/01/2013 | 12/31/2014 | 0.7175 |
| | 01/01/2014 | 12/31/2015 | 0.6707 |
| | 01/01/2015 | 12/31/2016 | 0.6515 |
| | 01/01/2016 | 12/31/2017 | 0.6860 |
| HDL | **01/01/2011** | **12/31/2012** | **0.6578** |
| | 01/01/2012 | 12/31/2013 | 0.7238 |
| | 01/01/2013 | 12/31/2014 | 0.7057 |
| | 01/01/2014 | 12/31/2015 | 0.7278 |
| | 01/01/2015 | 12/31/2016 | 0.7199 |
| | 01/01/2016 | 12/31/2017 | 0.7520 |
| Number of Medications | **01/01/2011** | **12/31/2012** | **0.3301** |
| | 01/01/2012 | 12/31/2013 | 0.3496 |
| | 01/01/2013 | 12/31/2014 | 0.3640 |
| | 01/01/2014 | 12/31/2015 | 0.3558 |
| | 01/01/2015 | 12/31/2016 | 0.3439 |
| | 01/01/2016 | 12/31/2017 | 0.3447 |

Table A.2: Estimated variances using 5-year data from EHRs at the OSU-WMCIW.

| Marker | Start Time | End Time | Estimated Variance |
|---|---|---|---|
| HBP | 01/01/2013 | 12/31/2017 | 1.681 |
| TC | 01/01/2013 | 12/31/2017 | 0.513 |
| HbA1c | 01/01/2013 | 12/31/2017 | 0.648 |
| HDL | 01/01/2013 | 12/31/2017 | 0.721 |
| Number of medications | 01/01/2013 | 12/31/2017 | 0.308 |



Figure A.1: Blue curves: new estimators for regression coefficients $\boldsymbol{\beta}_k(t)$ across 61 time points using variances estimated from 2013 to 2017 EHRs data at the OSU-WMCIW. Red curves: original estimators $\widehat{\boldsymbol{\beta}}_k(t)$.

Figure A.2: Blue curves: new estimators for correlation coefficients $\sigma_{kl}(t)$ across 61 time points using variances estimated from 2013 to 2017 EHRs data at the OSU-WMCIW. Red curves: original estimators $\widehat{\sigma}_{kl}(t)$.

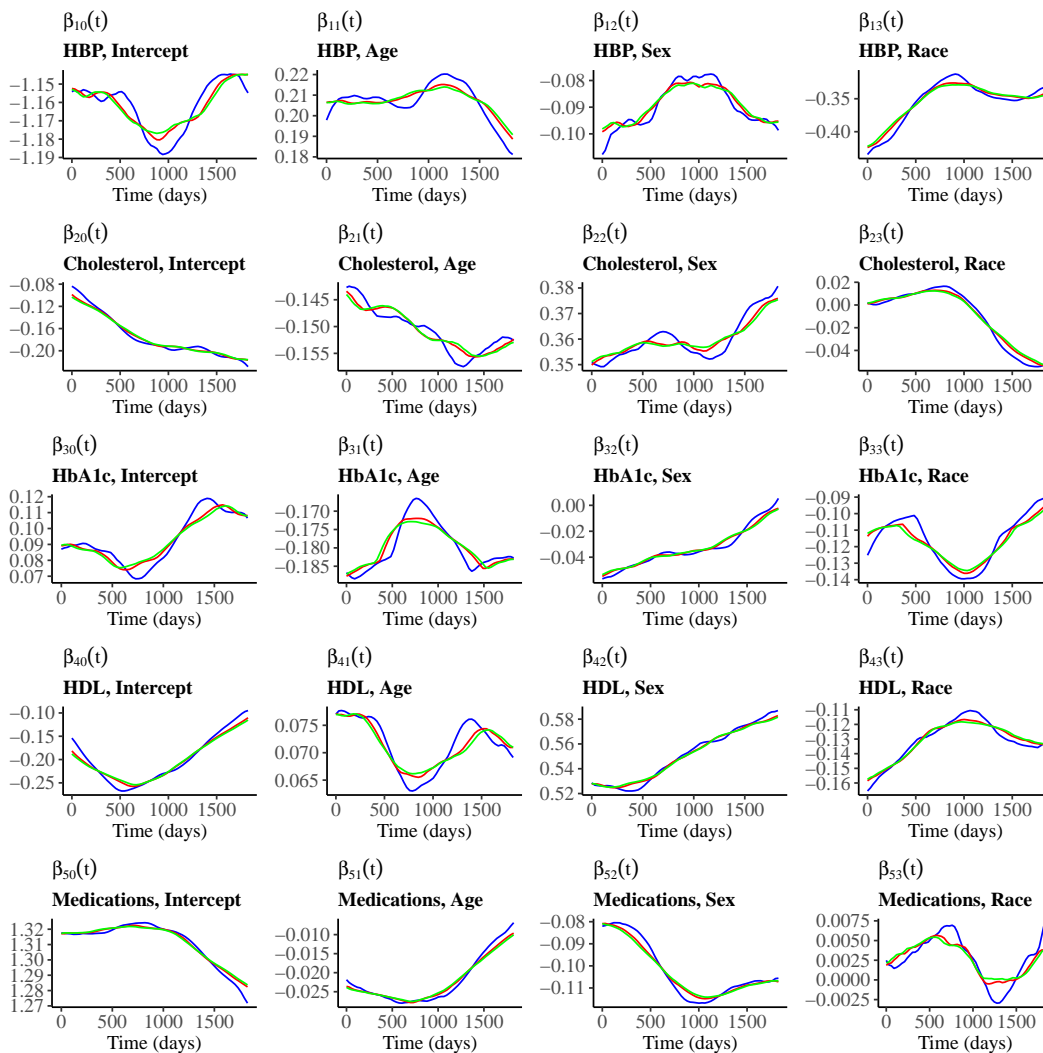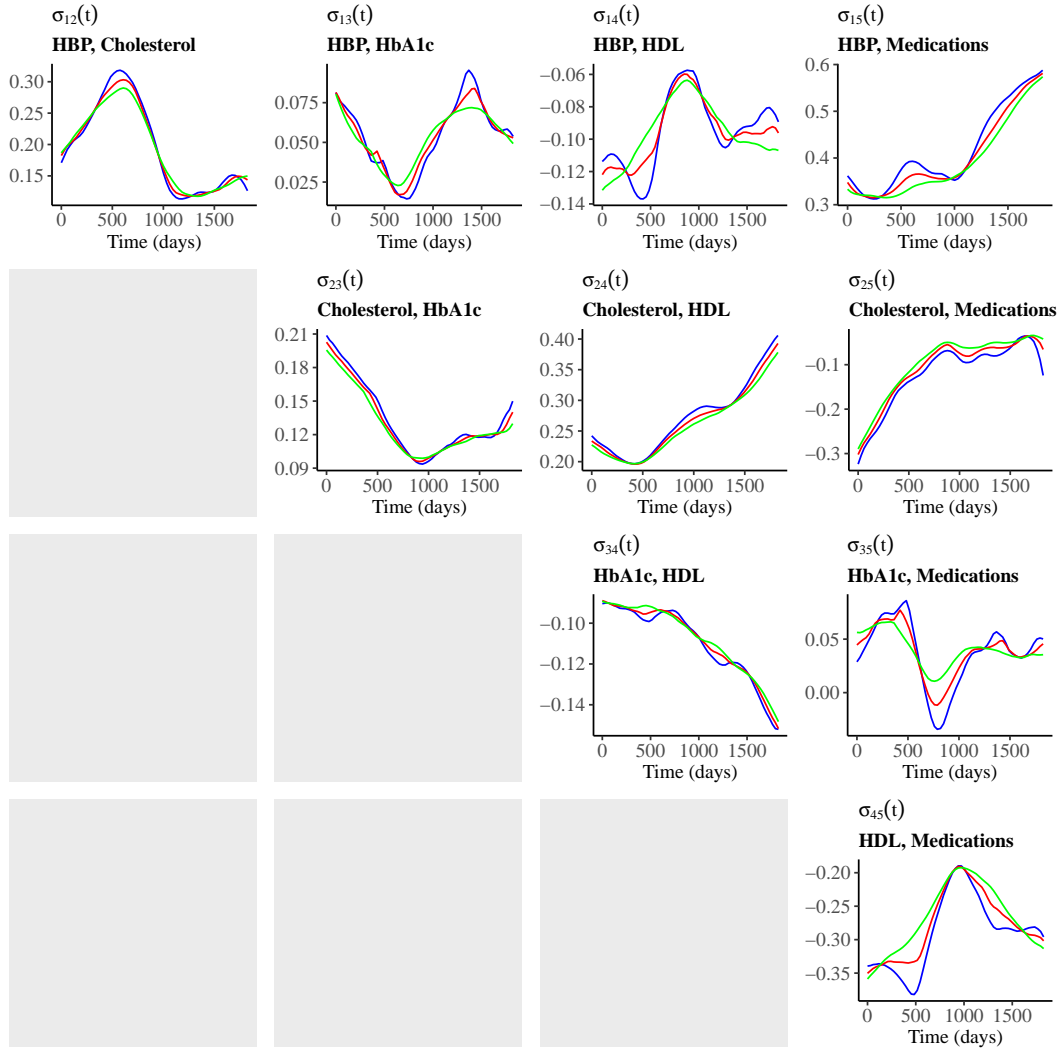## A.5  Validation on Selection of Optimal Bandwidth



Figure A.3: Blue curves: new estimators for regression coefficients $\boldsymbol{\beta}_k(t)$ across 61 time points using a suboptimal bandwidth $h_{1n} = 429.864$ days. Green curves: new estimators for regression coefficients $\boldsymbol{\beta}_k(t)$ across 61 time points using a suboptimal bandwidth $h_{1n} = 604.132$ days. Red curves: original estimators $\widehat{\boldsymbol{\beta}}_k(t)$ from EHRs at the OSU-WMCIW using the optimal bandwidth $h_{1n} = 564.112$ days.

Figure A.4: Blue curves: new estimators for correlation coefficients $\sigma_{kl}(t)$ across 61 time points using a suboptimal bandwidth $h_{2n} = 429.864$ days. Green curves: new estimators for correlation coefficients $\sigma_{kl}(t)$ across 61 time points using a suboptimal bandwidth $h_{2n} = 564.112$ days. Red curves: original estimators $\widehat{\sigma}_{kl}(t)$ from EHRs at the OSU-WMCIW using the optimal bandwidth $h_{2n} = 494.687$ days.

Table A.3: Canberra distances between estimators of regression coefficients $\boldsymbol{\beta}_k(t)$ using optimal and suboptimal bandwidths.

| Marker | Parameter $(H_1 = 564.112 \text{ days})$ | Distance $(H_1' = 429.864 \text{ days})$ | Distance $(H_1' = 604.032 \text{ days})$ |
|---|---|---|---|
| HBP | $\beta_{10}$ | 0.0022 | 0.0005 |
| Binary | $\beta_{11}$ | 0.0094 | 0.0019 |
| | $\beta_{12}$ | 0.0123 | 0.0039 |
| | $\beta_{13}$ | 0.0091 | 0.0025 |
| TC | $\beta_{20}$ | 0.0160 | 0.0042 |
| Continuous | $\beta_{21}$ | 0.0046 | 0.0009 |
| | $\beta_{22}$ | 0.0038 | 0.0009 |
| | $\beta_{23}$ | 0.1722 | 0.0500 |
| HbA1c | $\beta_{30}$ | 0.0207 | 0.0062 |
| Continuous | $\beta_{31}$ | 0.0059 | 0.0016 |
| | $\beta_{32}$ | 0.0923 | 0.0214 |
| | $\beta_{33}$ | 0.0199 | 0.0049 |
| HDL | $\beta_{40}$ | 0.0254 | 0.0061 |
| Continuous | $\beta_{41}$ | 0.0105 | 0.0028 |
| | $\beta_{42}$ | 0.0020 | 0.0006 |
| | $\beta_{43}$ | 0.0115 | 0.0028 |
| Number of Medications | $\beta_{50}$ | 0.0007 | 0.0002 |
| Count | $\beta_{51}$ | 0.0277 | 0.0072 |
| | $\beta_{52}$ | 0.0094 | 0.0029 |
| | $\beta_{53}$ | 0.2615 | 0.2525 |

Table A.4: Canberra distances between estimators of correlation coefficients $\sigma_{kl}(t)$ using optimal and suboptimal bandwidths.

| Parameter $(H_2 = 494.687 \text{ days})$ | Distance $(H_2' = 429.864 \text{ days})$ | Distance $(H_2' = 564.112 \text{ days})$ |
|---|---|---|
| $\sigma_{12}$ | 0.0168 | 0.0139 |
| $\sigma_{13}$ | 0.0458 | 0.0502 |
| $\sigma_{14}$ | 0.0338 | 0.0356 |
| $\sigma_{15}$ | 0.0171 | 0.0156 |
| $\sigma_{23}$ | 0.0131 | 0.0111 |
| $\sigma_{24}$ | 0.0116 | 0.0108 |
| $\sigma_{25}$ | 0.0739 | 0.0664 |
| $\sigma_{34}$ | 0.0080 | 0.0067 |
| $\sigma_{35}$ | 0.1856 | 0.1951 |
| $\sigma_{45}$ | 0.0268 | 0.0242 |

## A.6 Validation on Structure of Intensity Function

Table A.5: Effects of demographic variables and historical HbA1c levels on the frequency of health marker measurements.

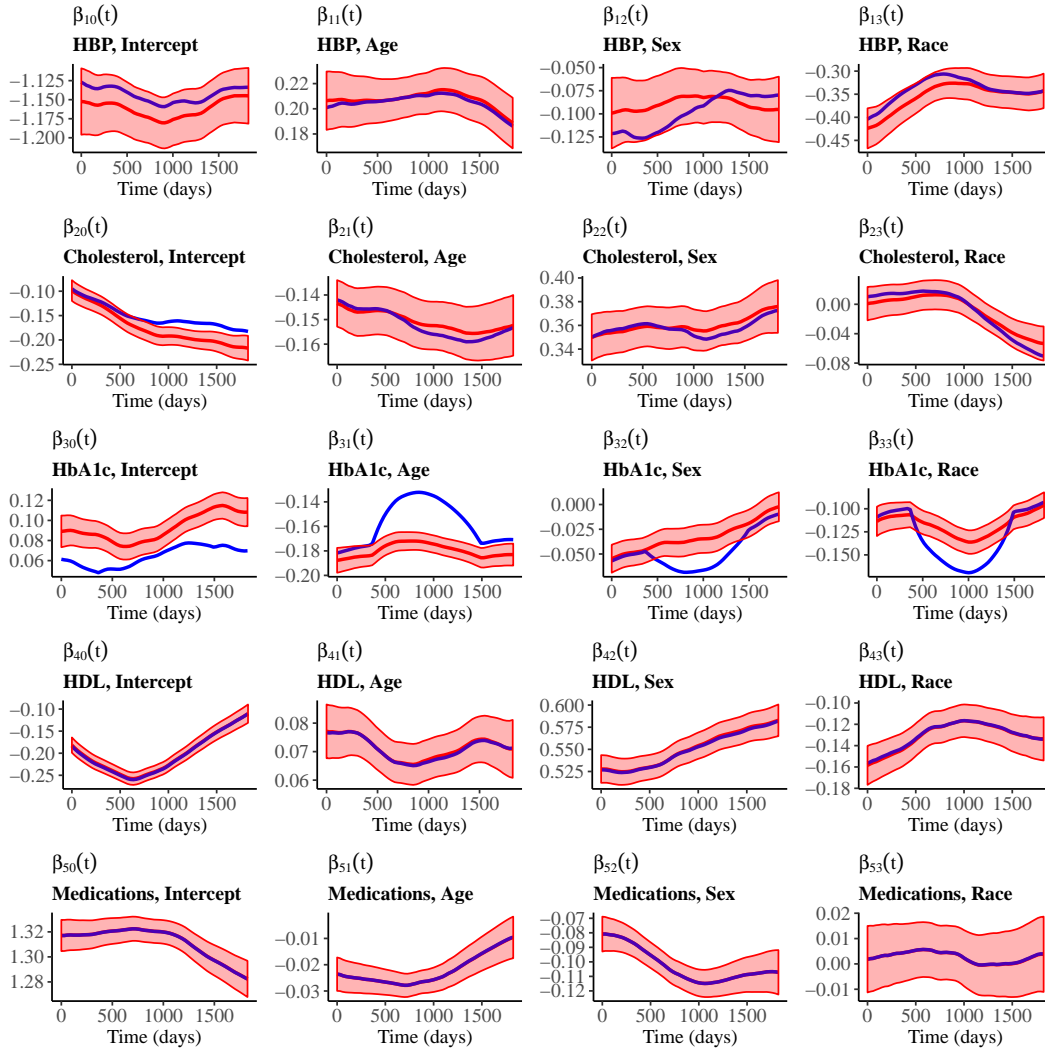| Marker | Demographic | Est | HR | SE | Z | P-value |
|---|---|---|---|---|---|---|
| HBP | age | 0.062 | 1.064 | 0.006 | 9.820 | < 0.001 |
| | sex | 0.060 | 1.062 | 0.014 | 4.290 | < 0.001 |
| | race | -0.132 | 0.876 | 0.015 | -9.011 | < 0.001 |
| | historical HbA1c | -0.029 | 0.971 | 0.009 | -3.110 | 0.002 |
| TC | age | 0.015 | 1.015 | 0.006 | 2.251 | 0.024 |
| | sex | -0.058 | 0.944 | 0.013 | -4.336 | < 0.001 |
| | race | -0.020 | 0.981 | 0.015 | -1.356 | 0.175 |
| | historical HbA1c | -0.077 | 0.926 | 0.010 | -7.606 | < 0.001 |
| HbA1c | age | 0.020 | 1.020 | 0.005 | 4.314 | < 0.001 |
| | sex | 0.026 | 1.026 | 0.010 | 2.704 | 0.007 |
| | race | -0.021 | 0.979 | 0.010 | -2.213 | 0.027 |
| | historical HbA1c | 0.094 | 1.098 | 0.005 | 18.345 | < 0.001 |
| HDL | age | 0.046 | 1.047 | 0.005 | 9.624 | < 0.001 |
| | sex | -0.011 | 0.989 | 0.010 | -1.068 | 0.285 |
| | race | -0.008 | 0.992 | 0.010 | -0.822 | 0.411 |
| | historical HbA1c | -0.012 | 0.988 | 0.007 | -1.807 | 0.071 |
| Medications | age | 0.043 | 1.044 | 0.006 | 7.308 | < 0.001 |
| | sex | 0.086 | 1.090 | 0.012 | 7.074 | < 0.001 |
| | race | -0.113 | 0.893 | 0.013 | -8.960 | < 0.001 |
| | historical HbA1c | 0.006 | 1.006 | 0.008 | 0.703 | 0.482 |

Figure A.5: Blue curves: new estimators for regression coefficients $\boldsymbol{\beta}_k(t)$ across 61 time points using historical HbA1c levels. Red curves: original estimators $\widehat{\boldsymbol{\beta}}_k(t)$ from EHRs at the OSU-WMCIW. Salmon-colored ribbons: 95% confidence intervals for original estimators $\widehat{\boldsymbol{\beta}}_k(t)$.
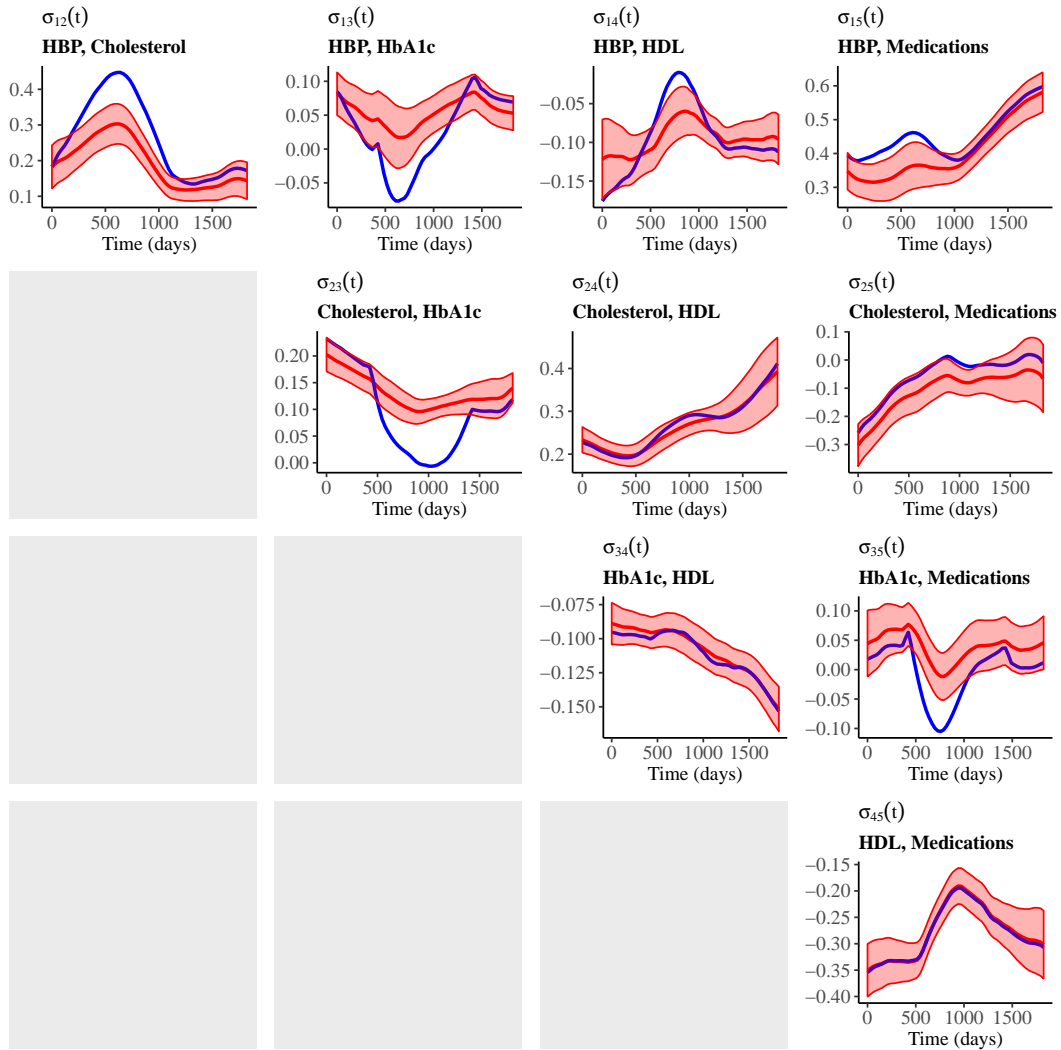
Figure A.6: Blue curves: new estimators for correlation coefficients $\sigma_{kl}(t)$ across 61 time points using historical HbA1c levels. Red curves: original estimators $\widehat{\sigma}_{kl}(t)$ from EHRs at the OSU-WMCIW. Salmon-colored ribbons: 95% confidence intervals for original estimators $\widehat{\sigma}_{kl}(t)$.

# APPENDIX B:  APPENDIX FOR CHAPTER 3

## B.1  Simulation Study for Multivariate Latent Modelling Method

We conducted a simulation study to examine the finite sample performance of the proposed latent modelling approach in Chapter 3 of the main text. In this simulation study, we simulated 100 datasets of two health markers for 10,000 subjects. We assumed $Y_{i1}(t)$ was Gaussian distributed and $Y_{i2}(t)$ was Bernoulli distributed. Thus, $g_1^{-1}(z) = z$ and $g_2^{-1}(z) = e^z/(1+e^z)$. Since the distribution of $Y_{i1}(t)$ had a dispersion parameter, we set $\phi_{i1}(t) = 0.1$. For the $i$th subject, we generated two covariates $X_{i1} \sim \mathcal{N}(0, 1/3)$ and $X_{i2} \sim Bernoulli(0.5) - 0.5$. Thus, $\boldsymbol{X}_i = (1, X_{i1}, X_{i2})^T$ was a 3-dimensional vector of baseline variables. The maximum observation time $T_i$ for each subject was set to 12 days. The measured time points for simulated markers were generated from two Poisson processes, and their intensity functions were $\mathbb{E}\left[dN_{i1}(t)|\boldsymbol{X}_i\right] = \exp\left\{0.5X_{i1} + 0.25X_{i2} + 0.3L_{i11}(t) - 0.1L_{i12}(t)\right\}dt$ and $\mathbb{E}\left[dN_{i2}(t)|\boldsymbol{X}_i\right] = 1.2\exp\left\{0.5X_{i1} + 0.25X_{i2} + 0.3L_{i21}(t) - 0.1L_{i22}(t)\right\}dt$. In the intensity functions, we let $L_{ik1}(t) = 1$ if there existed measurements of $k$th marker in $[t-3, t)$; otherwise, $L_{ik1}(t) = 0$. If $L_{ik1}(t) = 1$, then $L_{ik2}(t)$ was the average value of all $Y_{ik}(t)$ in $[t-3, t)$; otherwise, $L_{ik2}(t) = 0$. The true values of $\boldsymbol{\beta}_k(t)$ were assumed to be

$$\begin{pmatrix} \boldsymbol{\beta}_1^T(t) \\ \boldsymbol{\beta}_2^T(t) \end{pmatrix} = \begin{pmatrix} -1.36 + \frac{t}{10} & \sin(0.76 + t) & \cos(-0.3 + t) \\ \cos(-0.25 + t) & 0.37 + \frac{t}{10} & \sin(-0.68 + t) \end{pmatrix}.$$

Furthermore, we assumed the correlation structure of multivariate latent processes to be

$$\boldsymbol{\Omega}(t) = \begin{pmatrix} 0.5 & -0.25 \\ -0.25 & 0.5 \end{pmatrix} + \frac{1}{20}\begin{pmatrix} \sin(t+2) & \cos(t-0.5) \\ \cos(t-0.5) & \sin(t+3) \end{pmatrix}$$

and

$$\text{Cov}(\boldsymbol{\epsilon}(t), \boldsymbol{\epsilon}(s)) = \exp\left\{-\left(\frac{t-s}{b}\right)^2\right\}\frac{\boldsymbol{\Omega}(t) + \boldsymbol{\Omega}(s)}{2},$$

where $b = 0.5$.

The scaled Epanechnikov kernel was chosen as the kernel function to estimate $\boldsymbol{\beta}_k(t)$, i.e.,

$$K_{h_{1n}}(u) = \frac{3}{4h_{1n}}\left[1 - \left(\frac{u}{h_{1n}}\right)^2\right]_+.$$

Similarly, the kernel function for estimating $\sigma_{kl}(t)$ was set to the product of two scaled univariate Epanechnikov kernels, i.e.,

$$\widetilde{K}_{h_{2n}}(u_1, u_2) = \frac{9}{16h_{2n}^2} \left[1 - \left(\frac{u_1}{h_{2n}}\right)^2\right]_+ \left[1 - \left(\frac{u_2}{h_{2n}}\right)^2\right]_+.$$

We extended the data-adaptive method in Cao et al. (2015) and selected the optimal bandwidths among $0.1, 0.2, \ldots, 0.5$. We found $h = 0.3$ for $\boldsymbol{\beta}_k(t)$ and $h = 0.2$ for $\sigma_{kl}(t)$ were close to the optimal values. This set of bandwidth was used in all subsequent simulations. Since the proposed estimation method was expected to have more stable performance at time points that not on two ends, we estimated $\boldsymbol{\beta}_k(t)$, $\sigma_{kl}(t)$, and correlation coefficients $\rho_{kl}(t)$ at time points $t = 3, 4, \ldots, 11$ days.

Detailed examples on the code scripts of this simulation study can be accessed via `https://github.com/jitonglou/EHR_ITR`. Using one of the 100 simulated dataset, we provide explanatory codes and corresponding R workspaces which follow the workflows shown in Figure B.1.
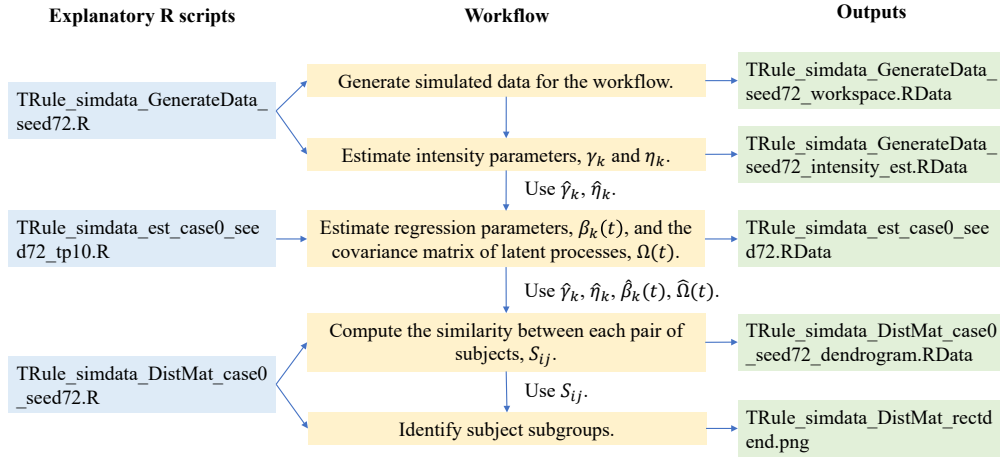


Figure B.1: Flow-chart of identifying latent subject subgroups using a simulated dataset.

Figure B.2 presents the true parameters versus estimators across the nine time points for $\boldsymbol{\beta}_k(t)$, $k = 1, 2$. From Figure B.2, we can conclude $\beta_{1\cdot}(t)$ are close to the true parameters at each time point of interest, and they well capture the underlying smooth function of $\boldsymbol{\beta}_k(t)$ across time. Compared with $\widehat{\beta}_{1\cdot}(t)$, the biases and variations of $\widehat{\beta}_{2\cdot}(t)$ seem greater. However, all of $\beta_{2\cdot}(t)$ are in the interquartile ranges of $\widehat{\beta}_{2\cdot}(t)$, indicating reasonable estimation results.

Figure B.3 displays the true parameters versus estimators across the nine time points for $\sigma_{12}(t)$
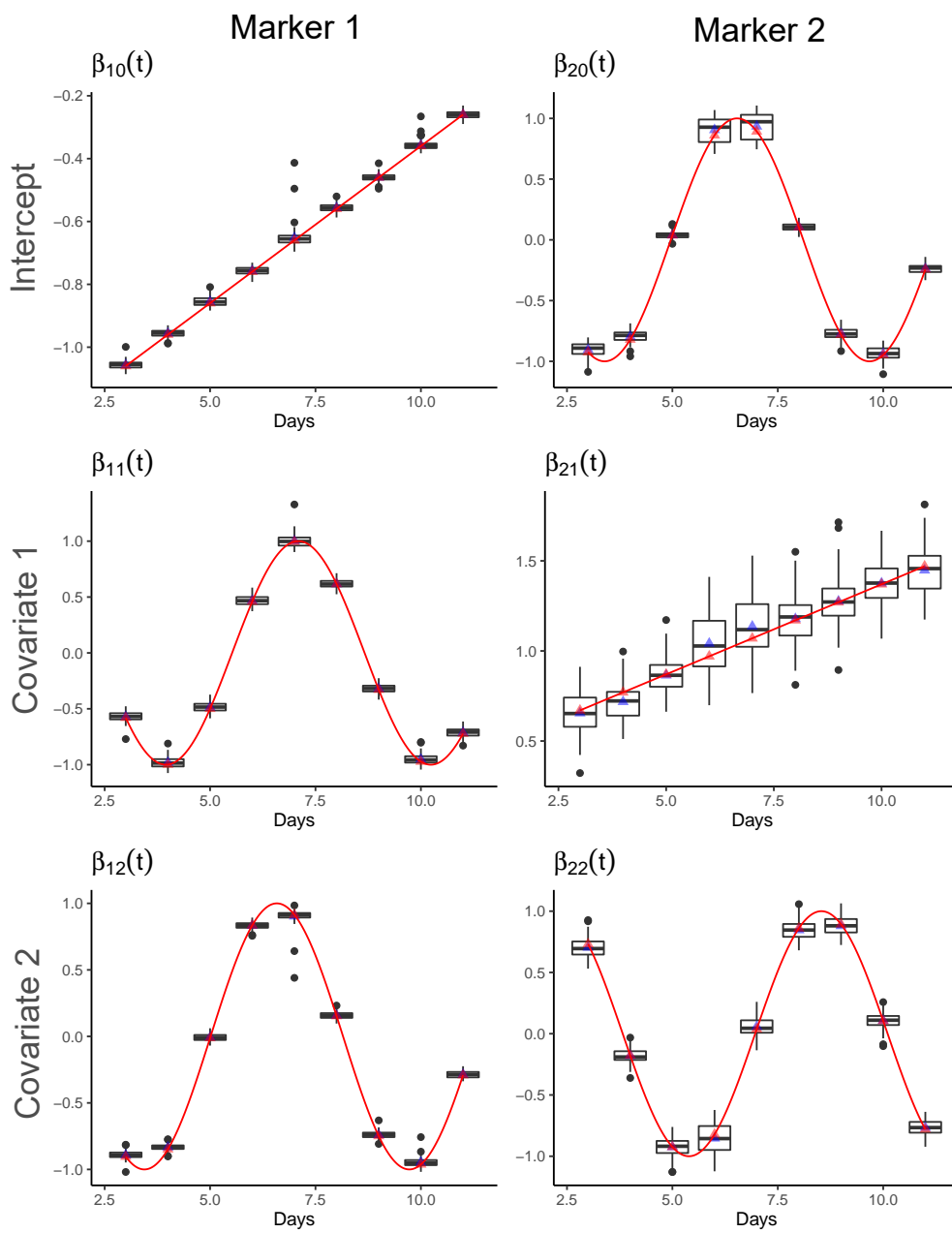
Figure B.2: Estimated regression coefficients $\widehat{\boldsymbol{\beta}}_1(t)$ and $\widehat{\boldsymbol{\beta}}_2(t)$ across 9 time points, based on 100 simulated datasets. Red: true parameters and functions. Blue: estimators.
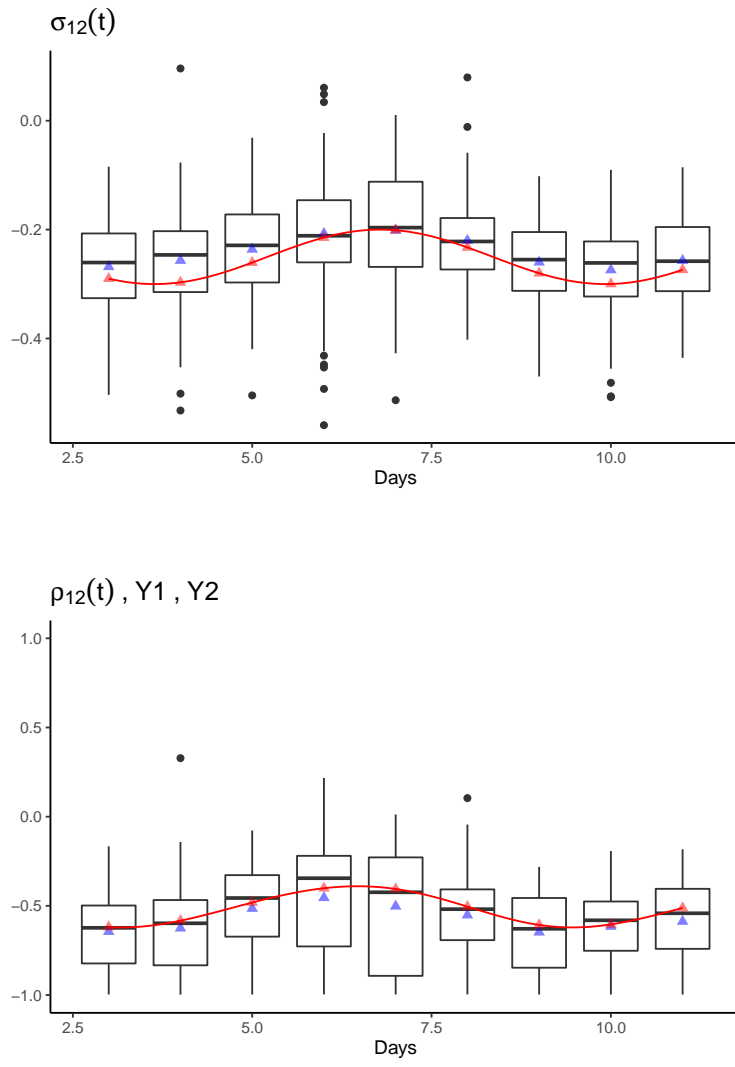
Figure B.3: Top panel: estimated covariance parameters $\widehat{\sigma}_{12}(t)$ across time points, based on 100 simulated datasets. Bottom panel: estimated correlation parameters $\widehat{\rho}_{12}(t)$ across time points, based on 100 simulated datasets. Red: true parameters and functions. Blue: estimators.

and $\rho_{12}(t)$. In both panels of Figure B.3, the mean and median of estimators are close to the true parameters. Also, $\sigma_{12}(t)$ and $\rho_{12}(t)$ are in the interquartile ranges of the estimators. The results in Figure B.3 suggest we can use $\widehat{\sigma}_{12}(t)$ and $\widehat{\rho}_{12}(t)$ to approximate the covariance and correlation structure of the simulated latent processes. Therefore, the proposed estimation method performs consistently and the biases of estimators are fairly acceptable.

## B.2 Results from Latent Models in EHRs Analysis

To estimate the parameters in the joint models, we first implemented the adaptive method of bandwidth selection as stated in Cao et al. (2015) and chose optimal bandwidths among 3, 6, and 9 months. We ended up to select $h = 9$ months for regression coefficients, $h = 3$ months for variances, and $h = 9$ months for correlations as the optimal bandwidths. Using the optimal bandwidths, we estimated $\boldsymbol{\beta}_k(t)$, $\sigma_k^2(t)$, and $\sigma_{kl}^2(t)$ at 25 time points.

The results are presented in Figure B.4, Figure B.5, and Figure B.6, respectively. The salmon-colored ribbons in these two figures are 95% confidence intervals for the parameters based on 100 bootstrap datasets.

Figure B.4 presents the relationships between each pair of health markers and covariates. In general, all health markers exhibit changes over time. Mean HbA1c ($\widehat{\beta}_{20}(t)$) has an increasing trend after about 200 days, which may suggest the difficulty to achieve long-term control of glycemic levels in a chronically ill patient population. Mean SBP ($\widehat{\beta}_{10}(t)$) and BMI ($\widehat{\beta}_{40}(t)$) show decreasing trends over time, suggesting relatively good control of blood pressure and body mass. Mean HDL ($\widehat{\beta}_{30}(t)$) increases during the first one year and has an decreasing trend afterwards, which may suggest the improvement in control of cholesterol in this patient population. The estimated regression coefficients for covariates, i.e., the estimated effects of covariates on health markers, do not show any pattern of drastic changes over time. Instead, the estimated values across time fluctuate around mean values. $\widehat{\beta}_{11}(t)$ and $\widehat{\beta}_{31}(t)$ are positive across time, while $\widehat{\beta}_{21}(t)$, $\widehat{\beta}_{41}(t)$, $\widehat{\beta}_{51}(t)$, and $\widehat{\beta}_{61}(t)$ are negative. Hence, estimators $\widehat{\beta}_{\cdot 1}(t)$ suggest that elder subjects on average have higher SBP and HDL but they have lower HbA1c, HDL, DD, and logMed. Similarly, estimators of sex effect, $\widehat{\beta}_{\cdot 2}(t)$, suggest that compared to men, women tend to have higher expected means of HDL and BMI, but they have lower accesses to DD and logMed. There is no apparent difference in the average values of SBP and HbA1c between elder subjects and younger subjects. For race, the estimators of $\widehat{\beta}_{\cdot 3}(t)$ indicate that white people have lower expected means than other races of people in SBP, HbA1c,
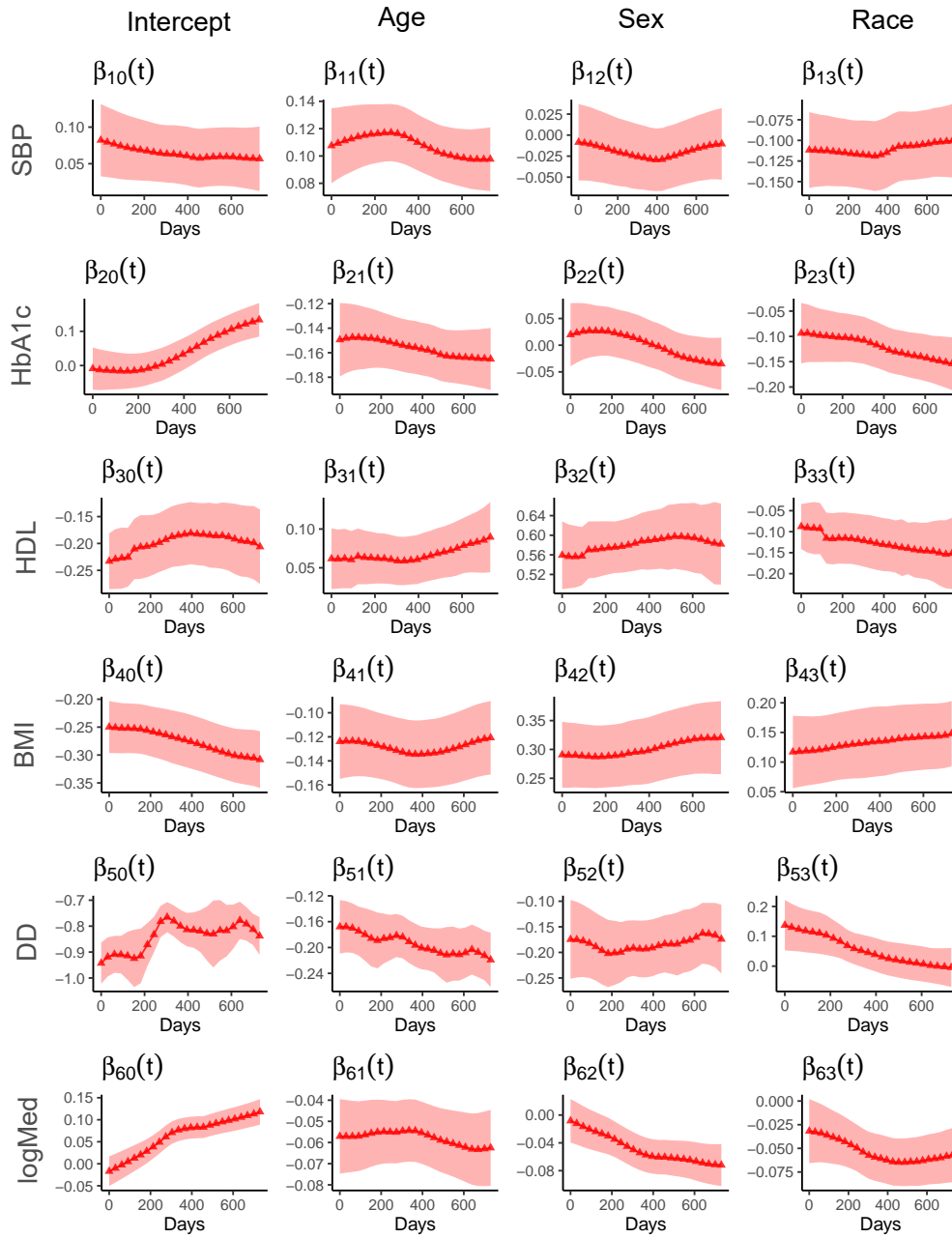
Figure B.4: Estimated regression coefficients $\widehat{\boldsymbol{\beta}}_k(t)$ across 25 time points. Salmon-colored ribbons: 95% confidence intervals for the estimators based on 100 bootstrapped datasets.

HDL, and logMed. However, on average, white people have higher values of BMI and DD, compared with non-white people.
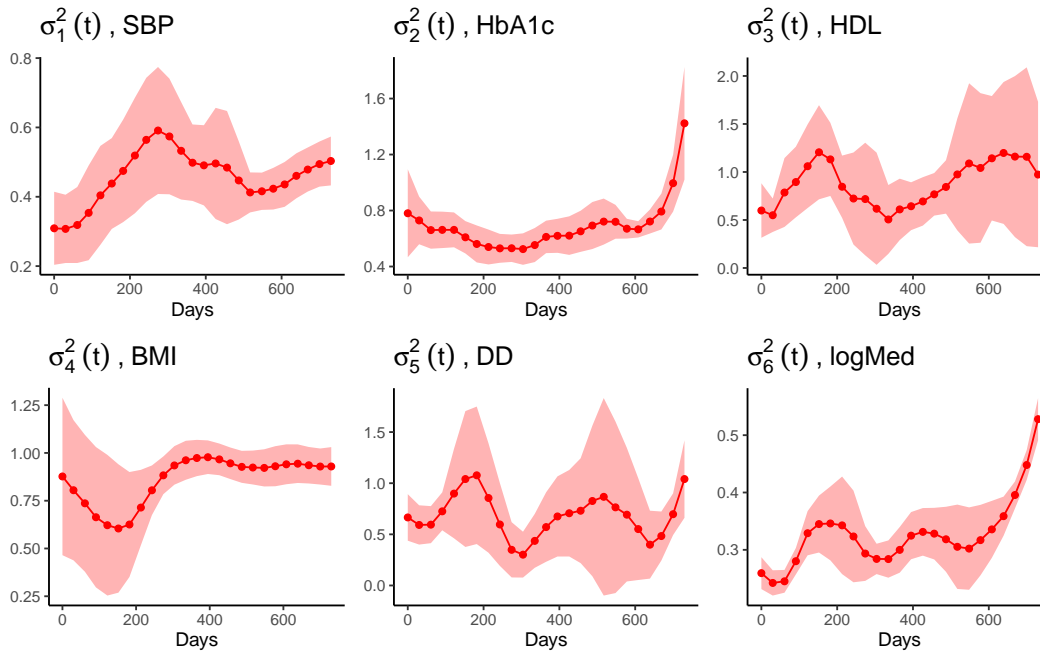


Figure B.5: Estimated variance parameters $\widehat{\sigma}_k^2(t)$ across 25 time points. Salmon-colored ribbons: 95% confidence intervals for the estimators based on 100 bootstrapped datasets.

Figure B.5 presents the estimated variances of the latent process of each health marker. The results suggest, in general, the variances fluctuate across time, and they do not have apparent patterns. However, there is one fact worth mention that the variances of HbA1c and logMed have observable increases after 1.5 years. This phenomenon may reflect the outcome of long-term control of HbA1c varies a lot among the population, so that the number of prescribed medications are different among patients.

Figure B.6 presents the correlations between each pair of latent processes of health markers at three time points. The results suggest the concurrent correlations between SBP and HbA1c, SBP and BMI, HbA1c and DD, DD and logMed are positive and moderate. Moreover, there exist negative and observable concurrent correlations between HDL and BMI, BMI and logMed. One interesting observation from Figure B.6 is that the estimated correlation between DD and logMed is as high as 0.8 and increases from $t = 6$ months to $t = 12$ months. This result is understandable
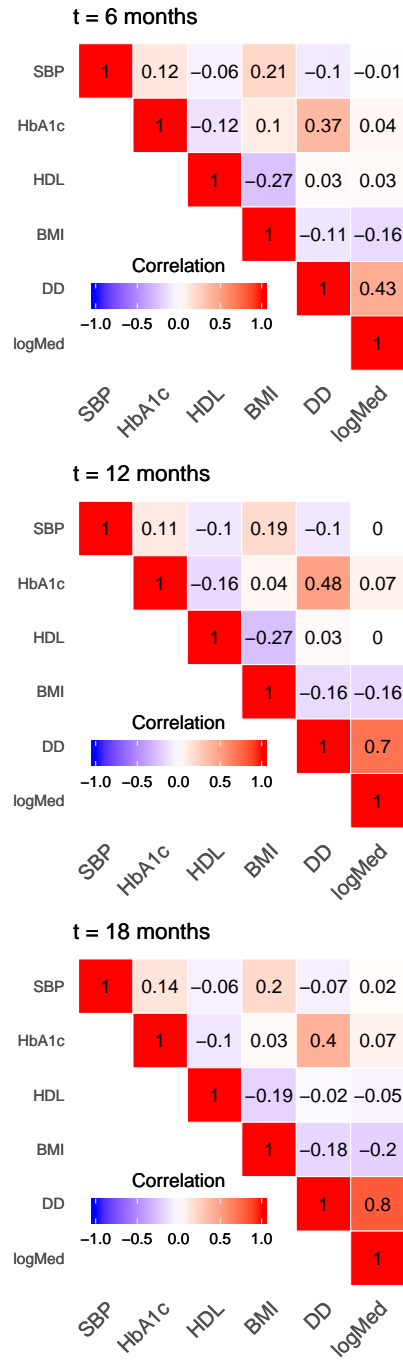
Figure B.6: Estimated correlation parameters $\widehat{\sigma}_{kl}(t)$ at $t = 6, 12, 18$ months prior to baseline treatment dates.

since as a patient receives more medications then the prescription is more likely to contain diabetic drugs. In addition, the estimated correlation between HbA1c and DD fluctuates around 0.4. This may suggest that the patients in this cohort had a greater possibility to receive diabetic drugs if their HbA1c levels are high.
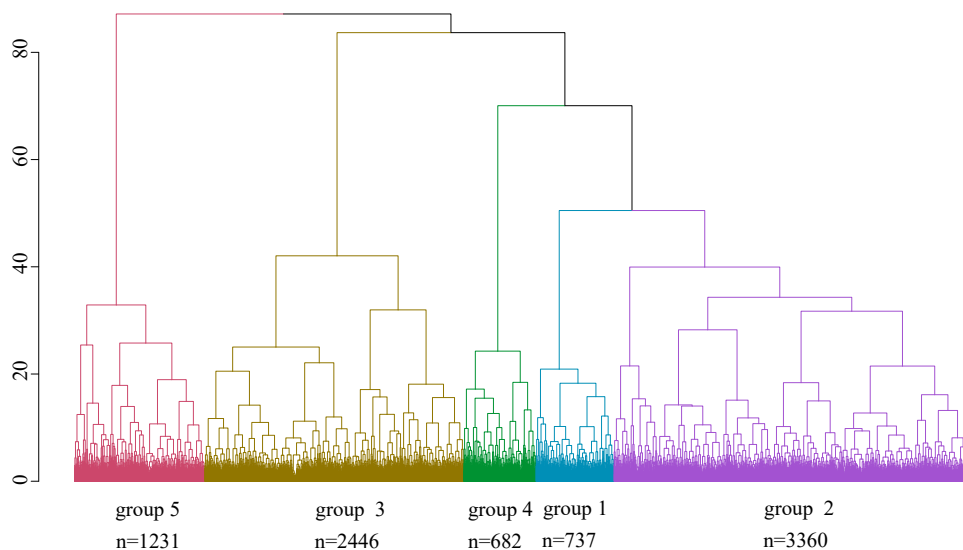


Figure B.7: Dendrogram of Mahalanobis distances for 8,456 patients.

Finally, we computed the similarity between each pair of patients using the Mahalanobis distance defined in (2.7). To compute $\widehat{\epsilon}_i(t)$, we substituted $Y_i(t)$ with the nearest neighbor observation of time $t$ for patient $i$. Using the between-patient similarity matrix, we performed a hierarchical clustering on the 8,456 patients, and the results are given in Figure B.7. We observed 5 clusters within which patients had similar health profiles.

## B.3 Coding Examples for Learning ITRs

In the github repository `https://github.com/jitonglou/EHR_ITR`, we also provide coding examples for learning ITRs as Figure B.8.

Since we were not allowed to release the EHR dataset, we illustrated the proposed method using the same simulated dataset in section B.1. After the estimation of parameters $\beta_k(t)$ and $\Omega(t)$, we calculated the similarity between each pair of subjects and performed a hierarchical clustering
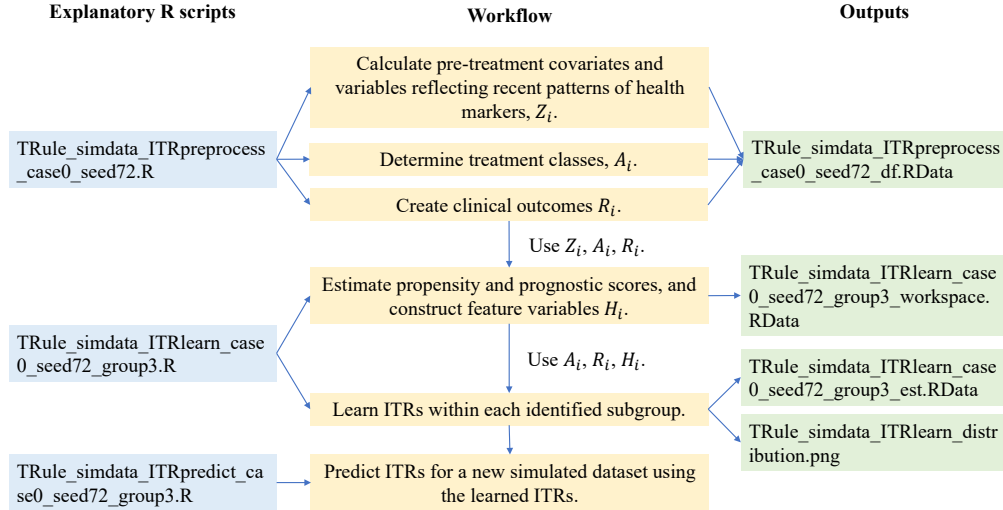
Figure B.8: Flow-chart of learning optimal ITRs in a simulation study.

analysis on the similarity matrix. Some subjects had no measurement for both health markers or had outliers. Thus, we excluded these subjects and the sample size for the dataset was 8,337 subjects. Subsequently, we clustered these subjects into three subgroups. Next, we prepared a dataset for learning ITRs. In particular, for the $i$th subject, we created two variables reflecting recent patterns of health markers by calculating the average value of $Y_{i1}(t)$ and $Y_{i2}(t)$ between $t = 3$ and $t = 11$. We denoted these two variables to $V_{i1}$ and $V_{i2}$. Then, we simulated the outcome reward

$$R_i = 1 + X_{i1} - X_{i2} + 2V_{i1} - 2V_{i2} + e_i,$$

where $e_i \sim \mathcal{N}(0, 1/3)$. Afterwards, the $i$th subject was randomly assigned a treatment $A_i$ in $\{A, B, C\}$ with equal probabilities, and the resulting dataset had 8,330 subjects of three treatments in three subgroups. We let $\boldsymbol{Z}_i = (X_{i1}, X_{i2}, V_{i1}, V_{i2})^T$ and denoted $u$ to a certain treatment selected from $\{A, B, C\}$. Within each subgroup, we estimated the propensity scores $\pi(\boldsymbol{Z}_i) = P(A_i = u|\boldsymbol{Z}_i)$ by a 10-fold cross-validation random forests with 3 repeats. Similarly, we estimated the prognostic scores $\psi(\boldsymbol{Z}_i) = E(R_i|\boldsymbol{Z}_i)$ by a gradient boosting model with 5,000 trees of which the maximum depth was 4. Lastly, we used $\boldsymbol{H}_i = (\boldsymbol{Z}_i, \widehat{\pi}(\boldsymbol{Z}_i), \widehat{\psi}(\boldsymbol{Z}_i))$, $A_i$, and $R_i$ to estimate ITRs. After the estimation, a certain treatment rule $D(\cdot)$ can be evaluated by its empirical value function, which is

defined as

$$\frac{\sum_u \sum_i I\left(A_i = D(\boldsymbol{Z}_i)\right) = u\right) R_i/\hat{\pi}(\boldsymbol{Z}_i)}{\sum_u \sum_i I\left(A_i = D(\boldsymbol{Z}_i)\right) = u\right)/\hat{\pi}(\boldsymbol{Z}_i)}$$

## B.4  Empirical HbA1c Values of Estimated ITRs

Table B.1: Comparison of ITRs using M-learning and Q-learning in terms of the empirical value function for the expected HbA1c level using 2-fold cross-validation with 100 repeats (a lower value means more beneficial).

| Group | Model | Mean (sd) | Median (Q1,Q3) |
|---|---|---|---|
| **Group 1** | **M-learning**, **SVM RBF kernel** | **6.362 (0.070)** | **6.359 (6.322, 6.400)** |
| | Q-learning, random forest | 7.042(0.098) | 7.034 (6.977, 7.106) |
| | Q-learning, SVM RBF kernel | 6.696 (0.073) | 6.572 (6.524, 6.626) |
| Universal rules: MET: 6.911, INS: 7.335, Other: 7.162, Multiple: 7.375 | | | |
| **Group 2** | **M-learning**, **SVM RBF kernel** | **6.397 (0.018)** | **6.398 (6.384, 6.407)** |
| | Q-learning, random forest | 7.284(0.051) | 7.286 (7.251, 7.320) |
| | Q-learning, SVM RBF kernel | 6.776 (0.031) | 6.775 (6.761, 6.794) |
| Universal rules: MET: 6.922, INS: 7.472, Other: 7.329, Multiple: 7.473 | | | |
| **Group 3** | **M-learning**, **SVM RBF kernel** | **6.597 (0.045)** | **6.592 (6.564, 6.629)** |
| | Q-learning, random forest | 7.550(0.068) | 7.554 (7.514, 7.597) |
| | Q-learning, SVM RBF kernel | 6.712 (0.057) | 6.705 (6.673, 6.749) |
| Universal rules: MET: 7.073, INS: 7.711, Other: 7.560, Multiple: 7.792 | | | |
| **Group 4** | **M-learning**, **SVM RBF kernel** | **7.752 (0.117)** | **7.765 (7.670, 7.810)** |
| | Q-learning, random forest | 9.052(0.163) | 9.071 (8.947, 9.166) |
| | Q-learning, SVM RBF kernel | 7.794 (0.131) | 7.793 (7.708, 7.879) |
| Universal rules: MET: 8.479, INS: 9.199, Other: 8.866, Multiple: 9.121 | | | |
| **Group 5** | **M-learning**, **SVM RBF kernel** | **6.464 (0.057)** | **6.462 (6.429, 6.502)** |
| | Q-learning, random forest | 7.393(0.082) | 7.389 (7.338, 7.441) |
| | Q-learning, SVM RBF kernel | 6.697 (0.041) | 6.694 (6.672, 6.724) |
| Universal rules: MET: 6.943, INS: 7.717, Other: 7.378, Multiple: 7.443 | | | |

# REFERENCES

Abdullah, S., Matthews, M., Frank, E., Doherty, G., Gay, G., and Choudhury, T. (2016). Automatic detection of social rhythms in bipolar disorder. *Journal of the American Medical Informatics Association* **23,** 538–543.

Abramowitz, M. and Stegun, I. (1965). *Handbook of Mathematical Functions: With Formulas, Graphs, and Mathematical Tables.* Applied mathematics series. Dover Publications, New York, NY, USA.

American Diabetes Association (2018). Pharmacologic approaches to glycemic treatment: standards of medical care in diabetes–2018. *Diabetes Care* **41,** 73–85.

Andersen, P. K. and Gill, R. D. (1982). Cox's regression model for counting processes: a large sample study. *The Annals of Statistics* **10,** 1100–1120.

Antonelli, J., Cefalu, M., Palmer, N., and Agniel, D. (2018). Doubly robust matching estimators for high dimensional confounding adjustment. *Biometrics* **74,** 1171–1179.

Barnett, I., Torous, J., Staples, P., Sandoval, L., Keshavan, M., and Onnela, J.-P. (2018). Relapse prediction in schizophrenia through digital phenotyping: a pilot study. *Neuropsychopharmacology* **43,** 1660–1666.

Beiwinkel, T., Kindermann, S., Maier, A., Kerl, C., Moock, J., Barbian, G., and Rössler, W. (2016). Using smartphones to monitor bipolar disorder symptoms: a pilot study. *JMIR Mental Health* **3,** e2.

Ben-Zeev, D., Scherer, E. A., Wang, R., Xie, H., and Campbell, A. T. (2015). Next-generation psychiatric assessment: Using smartphone sensors to monitor behavior and mental health. *Psychiatric Rehabilitation Journal* **38,** 218–226.

Bentler, P. M. (1990). Comparative fit indexes in structural models. *Psychological Bulletin* **107,** 238–246.

Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics).* Springer-Verlag, Berlin, Heidelberg.

Borchers, H. W. (2019). *pracma: Practical Numerical Math Functions.* R package version 2.2.9.

Bourla, A., Mouchabac, S., Hage, W. E., and Ferreri, F. (2018). e-PTSD: an overview on how new technologies can improve prediction and assessment of Posttraumatic Stress Disorder (PTSD). *European Journal of Psychotraumatology* **9,** 1424448.

Canivell, S., Mata-Cases, M., Real, J., Franch-Nadal, J., Vlacho, B., Khunti, K., and et al. (2019). Glycaemic control after treatment intensification in patients with type 2 diabetes uncontrolled on two or more non-insulin antidiabetic drugs in a real-world setting. *Diabetes, Obesity and Metabolism* **21,** 1373–1380.

Canzian, L. and Musolesi, M. (2015). Trajectories of depression: Unobtrusive monitoring of depressive states by means of smartphone mobility traces analysis. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, UbiComp'15, pages 1293–1304, New York, NY, USA. Association for Computing Machinery.

Cao, H., Zeng, D., and Fine, J. P. (2015). Regression analysis of sparse asynchronous longitudinal data. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **77,** 755–776.

Cebul, R. D., Love, T. E., Jain, A. K., and Hebert, C. J. (2011). Electronic health records and quality of diabetes care. *The New England journal of medicine* **365,** 825–833.

Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* **20,** 37–46.

Cole, R., Kripke, D., Gruen, W., Mullaney, D., and Gillin, J. (1992). Automatic sleep/wake identification from wrist activity. *Sleep* **15,** 461–469.

Cornelissen, G. (2014). Cosinor-based rhythmometry. *Theoretical biology and medical modelling* **11,** 16.

Cornet, V. P. and Holden, R. J. (2018). Systematic review of smartphone-based passive sensing for health and wellbeing. *Journal of biomedical informatics* **77,** 120–132.

Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine Learning* **20,** 273–297.

DaSilva, A. W., Huckins, J. F., Wang, R., Wang, W., Wagner, D. D., and Campbell, A. T. (2019). Correlates of stress in the college environment uncovered by the application of penalized generalized estimating equations to mobile sensing data. *JMIR mHealth and uHealth* **7,** e12084.

Davidian, M. and Giltinan, D. M. (2003). Nonlinear models for repeated measurement data: An overview and update. *Journal of Agricultural, Biological, and Environmental Statistics* **8,** 387–419.

De Maesschalck, R., Jouan-Rimbaud, D., and Massart, D. L. (2000). The mahalanobis distance. *Chemometrics and Intelligent Laboratory Systems* **50,** 1–18.

Depp, C. A., Bashem, J., Moore, R. C., Holden, J. L., Mikhael, T., Swendsen, J., Harvey, P. D., and Granholm, E. L. (2019). GPS mobility as a digital biomarker of negative symptoms in schizophrenia: a case control study. *NPJ Digital Medicine* **2,** 108.

Drucker, H., Burges, C. J. C., Kaufman, L., Smola, A., and Vapnik, V. (1997). Support vector regression machines. In Mozer, M. C., Jordan, M., and Petsche, T., editors, *Advances in Neural Information Processing Systems*, volume 9. MIT Press.

Drugs.com (2019). Drugs used to treat diabetes, type 2. `https://www.drugs.com/condition/diabetes-mellitus-type-ii.html`[Accessed: June 12, 2019].

Fan, J. and Zhang, W. (2008). Statistical methods with varying coefficient models. *Statistics and its interface* **1,** 179–195.

Faurholt-Jepsen, M., Busk, J., Þórarinsdóttir, H., Frost, M., Bardram, J. E., Vinberg, M., and Kessing, L. V. (2019). Objective smartphone data as a potential diagnostic marker of bipolar disorder. *Australian and New Zealand Journal of Psychiatry* **53,** 119–128.

Faurholt-Jepsen, M., Vinberg, M., Frost, M., Christensen, E. M., Bardram, J. E., and Kessing, L. V. (2015). Smartphone data as an electronic biomarker of illness activity in bipolar disorder. *Bipolar Disorders* **17,** 715–728.

Ferdous, R., Osmani, V., and Mayora, O. (2015). Smartphone app usage as a predictor of perceived stress levels at workplace. In *2015 9th International Conference on Pervasive Computing Technologies for Healthcare (PervasiveHealth)*, pages 225–228.

Forgy, E. (1965). Cluster analysis of multivariate data: efficiency versus interpretability of classification. *Biometrics* **21,** 768–769.

Foster, J. C., Taylor, J. M. G., and Ruberg, S. J. (2011). Subgroup identification from randomized clinical trial data. *Statistics in Medicine* **30,** 2867–2880.

Friedman, J. H. (2001). Greedy function approximation: a gradient boostingmachine. *The Annals of Statistics* **29,** 1189 – 1232.

Fu, H., Zhou, J., and Faries, D. E. (2016). Estimating optimal treatment regimes via subgroup identification in randomized control trials and observational studies. *Statistics in Medicine* **35,** 3285–3302.

Gainesa, B. R., Kim, J., and Zhou, H. (2018). Algorithms for fitting the constrained lasso. *Journal of Computational and Graphical Statistics* **27,** 861–871.

Garcia-Ceja, E., Osmani, V., and Mayora, O. (2016). Automatic stress detection in working environments from smartphones' accelerometer data: a first step. *IEEE Journal of Biomedical and Health Informatics* **20,** 1053–1060.

Ginsburg, G. S. and Phillips, K. A. (2018). Precision medicine: from science to value. *Health Affairs* **37,** 694–701.

Goodday, S. M. and Friend, S. (2019). Unlocking stress and forecasting its consequences with digital technology. *NPJ Digital Medicine* **2,** 75.

Greenwell, B. (2020). *fastshap: Fast Approximate Shapley Values.* R package version 0.0.5.

Gueorguieva, R. V. and Sanacora, G. (2006). Joint analysis of repeatedly observed continuous and ordinal measures of disease severity. *Statistics in medicine* **25,** 1307–1322.

Gunter, T. D. and Terry, N. P. (2005). The emergence of national electronic health record architectures in the United States and Australia: models, costs, and questions. *Journal of medical Internet research* **7,** e3.

Haines-Delmont, A., Chahal, G., Bruen, A. J., Wall, A., Khan, C. T., Sadashiv, R., and Fearnley, D. (2020). Testing suicide risk prediction algorithms using phone measurements with patients in acute mental health settings: feasibility study. *JMIR mHealth and uHealth* **8,** e15901.

Hallensleben, N., Spangenberg, L., Forkmann, T., Rath, D., Hegerl, U., Kersting, A., Kallert, T., and Glaesmer, H. (2017). Investigating the dynamics of suicidal ideation: preliminary findings from a study using ecological momentary assessments in psychiatric inpatients. *The Journal of Crisis Intervention and Suicide Prevention* **39,** 65–69.

Haneuse, S. (2016). Distinguishing selection bias and confounding bias in comparative effectiveness research. *Medical care* **54,** e23.

Haneuse, S. and Daniels, M. (2016). A general framework for considering selection bias in EHR-based studies: what data are observed and why? *EGEMS (Washington DC)* **4,** 1203.

Henao, R., Lu, J. T., Lucas, J. E., Ferranti, J., and Carin, L. (2016). Electronic health record analysis via deep poisson factor models. *The Journal of Machine Learning Research* **17,** 6422–6453.

Henderson, R., Ansell, P., and Alshibani, D. (2010). Regret-regression for optimal dynamic treatment regimes. *Biometrics* **66,** 1192–1201.

Herrin, J., Graca, B., Nicewander, D., Fullerton, C., Aponte, P., Stanek, G., Cowling, T., Collinsworth, A., Fleming, N. S., and Ballard, D. J. (2012). The effectiveness of implementing an electronic health record on diabetes care and outcomes. *Health services research* **47,** 1522–1540.

Ho, J. C., Ghosh, J., Steinhubl, S. R., Stewart, W. F., Denny, J. C., Malin, B. A., and Sun, J. (2014). Limestone: high-throughput candidate phenotype generation via tensor factorization. *Journal of biomedical informatics* **52,** 199–211.

Ho, T. K. (1995). Random decision forests. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, volume 1, pages 278–282.

Hooper, D., Coughlan, J., and Mullen, M. (2008). Structural equation modelling: guidelines for determining model fit. *Electronic Journal of Business Research Methods* **6,** 53–60.

Hripcsak, G., Ryan, P. B., Duke, J. D., Shah, N. H., Park, R. W., Huser, V., and et al. (2016). Characterizing treatment pathways at scale using the ohdsi network. *Proceedings of the National Academy of Sciences of the United States of America* **113,** 7329–7336.

Huang, J. Z., Wu, C. O., and Zhou, L. (2002). Varying-coefficient models and basis function approximations for the analysis of repeated measurements. *Biometrika* **89,** 111–128.

Huang, X., Goldberg, Y., and Xu, J. (2019). Multicategory individualized treatment regime using outcome weighted learning. *Biometrics* **75,** 1216–1227.

Husky, M., Olié, E., Guillaume, S., Genty, C., Swendsen, J., and Courtet, P. (2014). Feasibility and validity of ecological momentary assessment in the investigation of suicide risk. *Psychiatry Research* **220,** 564–570.

Insel, T. R. (2018). Digital phenotyping: a global tool for psychiatry. *World Psychiatry* **17,** 276–277.

Jacobson, N. C., Weingarden, H., and Wilhelm, S. (2019). Using digital phenotyping to accurately detect depression severity. *The Journal of Nervous and Mental Disease* **207,** 893–896.

Jain, S. H., Powers, B. W., Hawkins, J. B., and Brownstein, J. S. (2015). The digital phenotype. *Nature Biotechnology* **33,** 462–463.

Jänig, W. (1989). Autonomic nervous system. In Schmidt, R. F. and Thews, G., editors, *Human Physiology*, pages 333–370. Springer, Berlin, Heidelberg, Germany.

Karatzoglou, A., Smola, A., Hornik, K., and Zeileis, A. (2004). kernlab – an S4 package for kernel methods in R. *Journal of Statistical Software* **11,** 1–20.

Karstoft, K.-I., Galatzer-Levy, I. R., Statnikov, A., Li, Z., Shalev, A. Y., and members of Jerusalem Trauma Outreach and Prevention Study (J-TOPS) group (2015). Bridging a translational gap: using machine learning to improve the prediction of PTSD. *BMC Psychiatry* **15,** 30.

Kassambara, A. and Mundt, F. (2020). *factoextra: extract and visualize the results of multivariate data analyses*. R package version 1.0.7.

Kelly, D., Curran, K., and Caulfield, B. (2017). Automatic prediction of health status using smartphone-derived behavior profiles. *IEEE Journal of Biomedical and Health Informatics* **21,** 1750–1760.

Kessler, R. C. (2000). Posttraumatic stress disorder: the burden to the individual and to society. *The Journal of clinical psychiatry* **61,** 4–14.

Kleiman, E. M., Turner, B. J., Fedor, S., Beale, E. E., Huffman, J. C., and Nock, M. K. (2017). Examination of real-time fluctuations in suicidal ideation and its risk factors: results from two ecological momentary assessment studies. *Journal of Abnormal Psychology* **126,** 726–738.

Kleiman, E. M., Turner, B. J., Fedor, S., Beale, E. E., Picard, R. W., Huffman, J. C., and Nock, M. K. (2018). Digital phenotyping of suicidal thoughts. *Depression and Anxiety* **35,** 601–608.

Kuhn, M. (2020). *caret: Classification and Regression Training.* R package version 6.0-86.

Lambert, P. and Vandenhende, F. (2002). A copula-based model for multivariate non-normal longitudinal data: Analysis of a dose titration safety study on a new antidepressant. *Statistics in Medicine* **21,** 3197–3217.

Lance, G. N. and Williams, W. T. (1966). Computer programs for hierarchical polythetic classification ("similarity analyses"). *The Computer Journal* **9,** 60–64.

Linnstaedt, S. D., Zannas, A. S., McLean, S. A., Koenen, K. C., and Ressler, K. J. (2020). Literature review and methodological considerations for understanding circulating risk biomarkers following trauma exposure. *Molecular psychiatry* **25,** 1986–1999.

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

Liu, Y., Wang, Y., Kosorok, M. R., Zhao, Y., and Zeng, D. (2018). Augmented outcome-weighted learning for estimating optimal dynamic treatment regiments. *Statistics in Medicine* **37,** 3776–3788.

Lo, A., Chernoff, H., Zheng, T., and Lo, S.-H. (2015). Why significant variables aren't automatically good predictors. *Proceedings of the National Academy of Sciences* **112,** 13892–13897.

Lou, Z., Shao, J., and Yu, M. (2018). Optimal treatment assignment to maximize expected outcome with multiple treatments. *Biometrics* **74,** 506–516.

Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems 2017* **2017,** 4766–4775.

McGuire, H., Longson, D., Adler, A., Farmer, A., and Lewin, I. (2016). Management of type 2 diabetes in adults: summary of updated nice guidance. *British Medical Journal* **353,** i1575.

McHugh, M. L. (2012). Interrater reliability: the kappa statistic. *Biochemia Medica* **22,** 276–282.

McLean, S. A., Ressler, K., Koenen, K. C., Neylan, T., Germine, L., Jovanovic, T., Clifford, G. D., Zeng, D., An, X., Linnstaedt, S., et al. (2020). The AURORA study: a longitudinal, multimodal library of brain biology and function after traumatic stress exposure. *Molecular psychiatry* **25,** 283–296.

Mental Health America (2020). Adults with AMI who did not receive treatment 2020. `https://mhanational.org/issues/2020/mental-health-america-access-care-data#adults_ami_no_treatmentt` [Accessed: April 6, 2021].

Mesko, B. (2017). The role of artificial intelligence in precision medicine. *Expert Review of Precision Medicine and Drug Development* **2,** 239–241.

Minassian, A., Maihofer, A. X., Baker, D. G., Nievergelt, C. M., Geyer, M. A., Risbrough, V. B., and Marine Resiliency Study Team (2015). Association of predeployment heart rate variability with risk of postdeployment posttraumatic stress disorder in active-duty marines. *JAMA Psychiatry* **72,** 979–986.

Miscouridou, X., Perotte, A., Elhadad, N., and Ranganath, R. (2018). Deep survival analysis: non-parametrics and missingness. *Proceedings of the 3rd Machine Learning for Healthcare Conference, in PMLR.* **85,** 244–256.

Molenberghs, G. and Verbeke, G. (2005). *Models for Discrete Longitudinal Data.* Springer, New York, NY, USA.

Montvida, O., Shaw, J., Atherton, J., Stringer, F., and Paul, S. (2018). Long-term trends in antidiabetes drug usage in the U.S.: real-world evidence in patients newly diagnosed with type 2 diabetes. *Journal of the American Medical Association* **41,** 69–78.

Moshe, I., Terhorst, Y., Asare, K. O., Sander, L. B., Ferreira, D., Baumeister, H., Mohr, D. C., and Pulkki-Råback, L. (2021). Predicting symptoms of depression and anxiety using smartphone and wearable data. *Frontiers in Psychiatry* **12,** 625247.

Muaremi, A., Arnrich, B., and Tröster, G. (2013). Towards measuring stress with smartphones and wearable devices during workday and sleep. *Bionanoscience* **3,** 172–183.

Murphy, S. A. (2003). Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **65,** 331–366.

Murphy, S. A. (2005). A generalization error for Q-learning. *Journal of Machine Learning Research (JMLR)* **6,** 1073–1097.

National Institute of Mental Health (2021). Research Domain Criteria (RDoC). `https://www.nimh.nih.gov/research/research-funded-by-nimh/rdoc/about-rdoc.shtml` [Accessed: April 6, 2021].

Palmer, S. C., Mavridis, D., and Nicolucci, A. (2016). Comparison of clinical outcomes and adverse events associated with glucose-lowering drugs in patients with type 2 diabetes: a meta-analysis. *Journal of the American Medical Association* **316,** 313–324.

Palmius, N., Tsanas, A., Saunders, K. E. A., Bilderbeck, A. C., Geddes, J. R., Goodwin, G. M., and De Vos, M. (2017). Detecting bipolar depression from geographic location data. *IEEE Transactions on Biomedical Engineering* **64,** 1761–1771.

Peterson, R. A. and Cavanaugh, J. E. (2019). Ordered quantile normalization: a semiparametric transformation built for the cross-validation era. *Journal of Applied Statistics* pages 1–16.

Pincus, S. M., Gladstone, I. M., and Ehrenkranz, R. A. (1991). A regularity statistic for medical data analysis. *Journal of Clinical Monitoring* **7,** 335–345.

Place, S., Blanch-Hartigan, D., Rubin, C., Gorrostieta, C., Mead, C., Kane, J., Marx, B. P., Feast, J., Deckersbach, T., Pentland, A. S., Nierenberg, A., and Azarbayejani, A. (2017). Behavioral indicators on a mobile sensing platform predict clinically validated psychiatric symptoms of mood and anxiety disorders. *Journal of Medical Internet Research* **19,** e75.

Pyne, J. M., Constans, J. I., Wiederhold, M. D., Gibson, D. P., Kimbrell, T., Kramer, T. L., Pitcock, J. A., Han, X., Williams, D. K., Chartrand, D., Gevirtz, R. N., Spira, J., Wiederhold, B. K., McCraty, R., and McCune, T. R. (2016). Heart rate variability: pre-deployment predictor of post-deployment ptsd symptoms. *Biological Psychology* **121,** 91–98.

Qi, Z., Liu, D., Fu, H., and Liu, Y. (2020). Multi-armed angle-based direct learning for estimating optimal individualized treatment rules with various outcomes. *Journal of the American Statistical Association* **115,** 678–691.

Qian, M. and Murphy, S. A. (2011). Performance guarantees for individualized treatment rules. *The Annals of Statistics* **39,** 1180–1210.

Reinertsen, E. and Clifford, G. D. (2018). A review of physiological and behavioral monitoring with digital sensors for neuropsychiatric illnesses. *Physiological Measurement* **39,** 05TR01.

Robins, J. M. (2004). Optimal structural nested models for optimal sequential decisions. In *Proceedings of the Second Seattle Symposium in Biostatistics*, volume 179 of *Lect. Notes Stat.*, pages 189–326, New York, NY, USA. Springer.

Roglic, G. (2016). WHO global report on diabetes: a summary. *International Journal of Noncommunicable Diseases* **1,** 3–8.

Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics* **20,** 53–65.

Saeb, S., Zhang, M., Karr, C. J., Schueller, S. M., Corden, M. E., Kording, K. P., and Mohr, D. C. (2015). Mobile phone sensor correlates of depressive symptom severity in daily-life behavior: an exploratory study. *Journal of Medical Internet Research* **17,** e175.

Sano, A. and Picard, R. W. (2013). Stress recognition using wearable sensors and mobile phones. *2013 Humaine Association Conference on Affective Computing and Intelligent Interaction* pages 671–676.

Sano, A., Taylor, S., McHill, A. W., Andrew Jk Phillips, Barger, L. K., Klerman, E., and Picard, R. (2018). Identifying objective physiological markers and modifiable behaviors for self-reported stress and mental health status using wearable sensors and mobile phones: Observational study. *Journal of Medical Internet Research* **20,** e210.

Shaffer, F. and Ginsberg, J. (2017). An overview of heart rate variability metrics and norms. *Frontiers in Public Health* **5,** 258.

Staples, P., Torous, J., Barnett, I., Carlson, K., Sandoval, L., Keshavan, M., and Onnela, J.-P. (2017). A comparison of passive and active estimates of sleep in a cohort with schizophrenia. *NPJ Schizophrenia* **3,** 37.

Stone, N. J., Robinson, J. G., Lichtenstein, A. H., Bairey Merz, C. N., Blum, C. B., Eckel, R. H., and et al. (2014). 2013 ACC/AHA guideline on the treatment of blood cholesterol to reduce atherosclerotic cardiovascular risk in adults: a report of the American College of Cardiology/American Heart Association task force on practice guidelines. *Journal of the American College of Cardiology* **63,** 2889–2934.

Stütz, T., Kowar, T., Kager, M., Tiefengrabner, M., Stuppner, M., Blechert, J., Wilhelm, F. H., and Ginzinger, S. (2015). Smartphone based stress prediction. In Ricci, F., Bontcheva, K., Conlan, O., and Lawless, S., editors, *User Modeling, Adaptation and Personalization*, UMAP 2015, pages 240–251. Springer, Cham, Switzerland.

Tao, Y. and Wang, L. (2017). Adaptive contrast weighted learning for multi-stage multi-treatment decision-making. *Biometrics* **73,** 145–155.

Task Force Report (1996). Heart rate variability: standards of measurement, physiological interpretation and clinical use. *Circulation* **93,** 1043–1065.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **58,** 267–288.

Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., and Knight, K. (2004). Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **67,** 91–108.

Tibshirani, R., Walther, G., and Hastie, T. (2002). Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **63,** 411–423.

Tucker, L. R. and Lewis, C. (1973). A reliability coefficient for maximum likelihood factor analysis. *Psychometrika* **38,** 1–10.

Tulppo, M. P., Makikallio, T. H., Takala, T. E., Seppanen, T., and Huikuri, H. V. (1996). Quantitative beat-to-beat analysis of heart rate dynamics during exercise. *American Journal of Physiology-Heart and Circulatory Physiology* **271,** H244–H252.

Verbeke, G., Fieuws, S., Molenberghs, G., and Davidian, M. (2014). The analysis of multivariate longitudinal data: a review. *Statistical Methods in Medical Research* **23,** 42–59.

Verbeke, G. and Molenberghs, G. (2000). *Linear Mixed Models for Longitudinal Data.* Springer, New York, NY, USA.

Wahle, F., Kowatsch, T., Fleisch, E., Rufer, M., and Weidt, S. (2016). Mobile sensing and support for people with depression: a pilot trial in the wild. *JMIR mHealth and uHealth* **4,** e111.

Wang, R., Aung, M. S. H., Abdullah, S., Brian, R., Campbell, A. T., Choudhury, T., Hauser, M., Kane, J., Merrill, M., Scherer, E. A., Tseng, V. W. S., and Ben-Zeev, D. (2016). Crosscheck: Toward passive sensing and detection of mental health changes in people with schizophrenia. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, UbiComp'16, pages 886–897, New York, NY, USA. Association for Computing Machinery.

Watkins, C. J. C. H. and Dayan, P. (1992). Q-learning. *Machine Learning* **8,** 279–292.

Whelton, P. K., Carey, R. M., Aronow, W. S., Casey Jr., D. E., Collins, K. J., Dennison Himmelfarb, C., and et al. (2018). 2017 ACC/AHA/AAPA/ABC/ACPM/AGS/APhA/ASH/ASPC/NMA /PCNA guideline for the prevention, detection, evaluation, and management of high blood pressure in adults: a report of the American College of Cardiology/American Heart Association task force on clinical practice guidelines. *Journal of the American College of Cardiology* **71,** e127–e248.

Wu, P., Zeng, D., and Wang, Y. (2020). Matched learning for optimizing individualized treatment strategies using electronic health records. *Journal of the American Statistical Association* **115,** 380–392.

You, K. and Zhu, X. (2018). *ADMM: algorithms using alternating direction method of multipliers.* R package version 0.3.1.

Zeng, D. and Lin, D. (2006). Efficient estimation of semiparametric transformation models for counting processes. *Biometrika* **93,** 627–640.

Zhang, C., Chen, J., Fu, H., He, X., Zhao, Y., and Liu, Y. (2020). Multicategory outcome weighted margin-based learning for estimating individualized treatment rules. *Statistica Sinica* **30,** 1857–1879.

Zhao, Y., Zeng, D., Rush, A. J., and Kosorok, M. R. (2012). Estimating individualized treatment rules using outcome weighted learning. *Journal of the American Statistical Association* **107,** 1106–1118.