

Sparse Multi-Output Gaussian Processes for Medical Time Series Prediction

Li-Fang Cheng

LIFANGC@PRINCETON.EDU

*Department of Electrical Engineering
Princeton University
Princeton, NJ 08544, USA*

Gregory Darnell

GDARNELL@PRINCETON.EDU

*Lewis-Sigler Institute
Princeton University
Princeton, NJ 08544, USA*

Bianca Dumitrascu

BIANCAD@PRINCETON.EDU

*Lewis-Sigler Institute
Princeton University
Princeton, NJ 08544, USA*

Corey Chivers

COREY.CHIVERS@UPHS.UPENN.EDU

*University of Pennsylvania Health System
Philadelphia, PA 19104, USA*

Michael E Draugelis

MICHAEL.DRAUGELIS@UPHS.UPENN.EDU

*University of Pennsylvania Health System
Philadelphia, PA 19104, USA*

Kai Li

LI@CS.PRINCETON.EDU

*Department of Computer Science
Princeton University
Princeton, NJ 08540, USA*

Barbara E Engelhardt

BEE@PRINCETON.EDU

*Department of Computer Science
Center for Statistics and Machine Learning
Princeton University
Princeton, NJ 08540, USA*

Abstract

In the scenario of real-time monitoring of hospital patients, high-quality inference of patients' health status using all information available from clinical covariates and lab tests is essential to enable successful medical interventions and improve patient outcomes. Developing a computational framework that can learn from observational large-scale electronic health records (EHRs) and make accurate real-time predictions is a critical step. In this work, we develop and explore a Bayesian nonparametric model based on Gaussian process (GP) regression for hospital patient monitoring. We propose MedGP, a statistical framework that incorporates 24 clinical and lab covariates and supports a rich reference data set from which relationships between observed covariates may be inferred and exploited for

high-quality inference of patient state over time. To do this, we develop a highly structured sparse GP kernel to enable tractable computation over tens of thousands of time points while estimating correlations among clinical covariates, patients, and periodicity in patient observations. MedGP has a number of benefits over current methods, including (i) not requiring an alignment of the time series data, (ii) quantifying confidence regions in the predictions, (iii) exploiting a vast and rich database of patients, and (iv) inferring interpretable relationships among clinical covariates. We evaluate and compare results from MedGP on the task of online prediction for three patient subgroups from two medical data sets across 8,043 patients. We found MedGP improves online prediction over baseline methods for nearly all covariates across different disease subgroups and studies. The publicly available code is at <https://github.com/bee-hive/MedGP>.

Keywords: Gaussian processes, electronic health records, sparse time series analysis, spectral mixture kernel, kernel density estimation.

1. Introduction

Large-scale collections of electronic health records (EHRs) are becoming useful for understanding disease progress, early diagnosis, and personalized treatments for many clinical diseases (Murdoch and Detsky, 2013; Hripcsak and Albers, 2013; Ghassemi et al., 2015a). EHRs contain rich patient information—disease history, demographics, vital signs, and lab results—that clinicians use to diagnose and treat patients. In this work, we are interested in developing a statistical framework that leverages medical data from a set of reference patients to enable personalized, real-time monitoring of new hospital patients. In particular, we consider data from the Hospitals at the University of Pennsylvania (HUP) containing hospital information for over 260,000 patients, and the public Multiparameter Intelligent Monitoring in Intensive Care (MIMIC-III) data set with more than 53,000 admissions from 38,000 patients in intensive care units (ICUs) (Johnson et al., 2016).

One motivation for monitoring new patients is to characterize patient state to allow the early diagnosis of sepsis or septic shock. Sepsis is one of the leading causes of death in critically ill patients in the United States (Hotchkiss and Karl, 2003). Each year an estimated 750,000 cases of sepsis or septic shock occur in the US. The mortality rate of septic patients ranges from 20% to 30%, and accounts for roughly 9.3% of all US deaths (Angus et al., 2001; Kumar et al., 2011). Sepsis is usually developed during a patient’s stay in the hospital. However, accurate diagnosis of sepsis is difficult due to heterogeneous symptoms across patients (Pierrakos and Vincent, 2010).

One way to reduce the mortality rate of sepsis is to increase the accuracy of early diagnosis of sepsis. To do this, we might develop a model of patient state and fit this model to EHR data from previous hospital patients with sepsis. However, existing EHR data pose several challenges because they have been collected with traditional monitoring methods. Many of the covariates, lab results in particular, are sparsely sampled across patients. That is, there are only a small number of observations of any lab result per patient. For example, vital signs are generally taken once every three to four hours for inpatient data, and once every hour for patients in the intensive care unit (ICU). In contrast, blood tests requiring a blood draw are generally performed at most once a day. We see this sparsity in an example of 24 clinical covariates (Table 1) measured across time for a single patient, including four densely sampled vital signs (respiration rate, heart rate, systolic

Type	Covariate	Sepsis	Neoplasms	Heart Failure	MIMIC-III
Vital	Respiration rate (RR)	87,076	493,964	147,445	291,466
Vital	Heart rate (HR)	96,317	527,989	227,951	294,746
Vital	Systolic blood pressure (SBP)	84,909	447,666	104,129	124,587
Vital	Body temperature (Temp)	80,597	364,286	94,468	56,533
Lab	Blood urea nitrogen (BUN)	12,528	71,825	21,751	25,102
Lab	Carbon dioxide (CO ₂)	12,672	72,784	21,844	20,979
Lab	Calcium level	10,388	66,051	18,867	20,568
Lab	Chloride	10,100	68,534	21,421	26,248
Lab	Creatinine	12,689	72,928	21,889	25,237
Lab	Glucose point-of-care (Glucose POC)	20,444	170,872	54,239	24,196
Lab	Hematocrit (Hct)	12,752	74,060	22,035	24,810
Lab	Hemoglobin (Hgb)	13,005	75,646	27,891	21,226
Lab	Mean cell hemoglobin (MCH)	12,587	69,736	18,379	20,877
Lab	Mean cell hemoglobin concentration (MCHC)	12,577	69,682	18,359	20,885
Lab	Mean cell volume (MCV)	12,587	69,751	18,380	20,875
Lab	International normalization ratio (INR)	5,733	38,810	17,005	15,735
Lab	Prothrombin time (PT)	5,722	38,844	17,007	15,734
Lab	Partial thromboplastin time (PTT)	5,872	41,894	19,596	17,185
Lab	Platelet	12,586	69,945	18,367	21,395
Lab	Potassium level	12,830	77,395	28,470	27,200
Lab	Red blood cell (RBC)	12,600	69,776	18,387	20,876
Lab	Red cell distribution width (RDW)	12,580	69,757	18,381	20,877
Lab	Sodium level	12,848	78,617	28,597	26,383
Lab	White blood cell (WBC)	12,581	69,950	18,384	20,960

Table 1: **The 24 clinical covariates modeled in MedGP.** This table includes the total number of observations for each covariate across patients in three disease groups—sepsis, neoplasms, and heart failure—in the HUP data, and the heart failure patients in the MIMIC-III data.

blood pressure, and body temperature) and 20 sparsely sampled lab covariates (Figure 1). Data missingness is systematic and not at random (Newgard and Lewis, 2015): a doctor will only order a test that will be informative in characterizing patient state relevant to diagnosis. Moreover, these time series data are not aligned across patients to a reference time point or disease onset; instead, patient intake is at time 0 and release is hours or days later. The sparsity over patients and uncalibrated time series make the physiological progression of disease within patients or joint analysis of time series across patients challenging due to substantial uncertainty of patient state and rate of disease progression at any time.

In this work, we build a statistical framework that uses sparse, heterogeneous EHR time series data to monitor and predict vital signs and lab results for each patient in an online way. To do this, we first designed a nonparametric model based on Gaussian process (GP) multivariate regression to explore the correlations both within each clinical covariate across time and across clinical covariates given rich EHR reference data. Our model includes a highly structured GP kernel regularized using sparsity-inducing priors to avoid overfitting, allow interpretability, and ensure computational tractability. Second, we propose a framework based on nonparametric density estimation to tailor the empirical model to a patient-specific model for each new patient. For real-time monitoring, we update the empir-

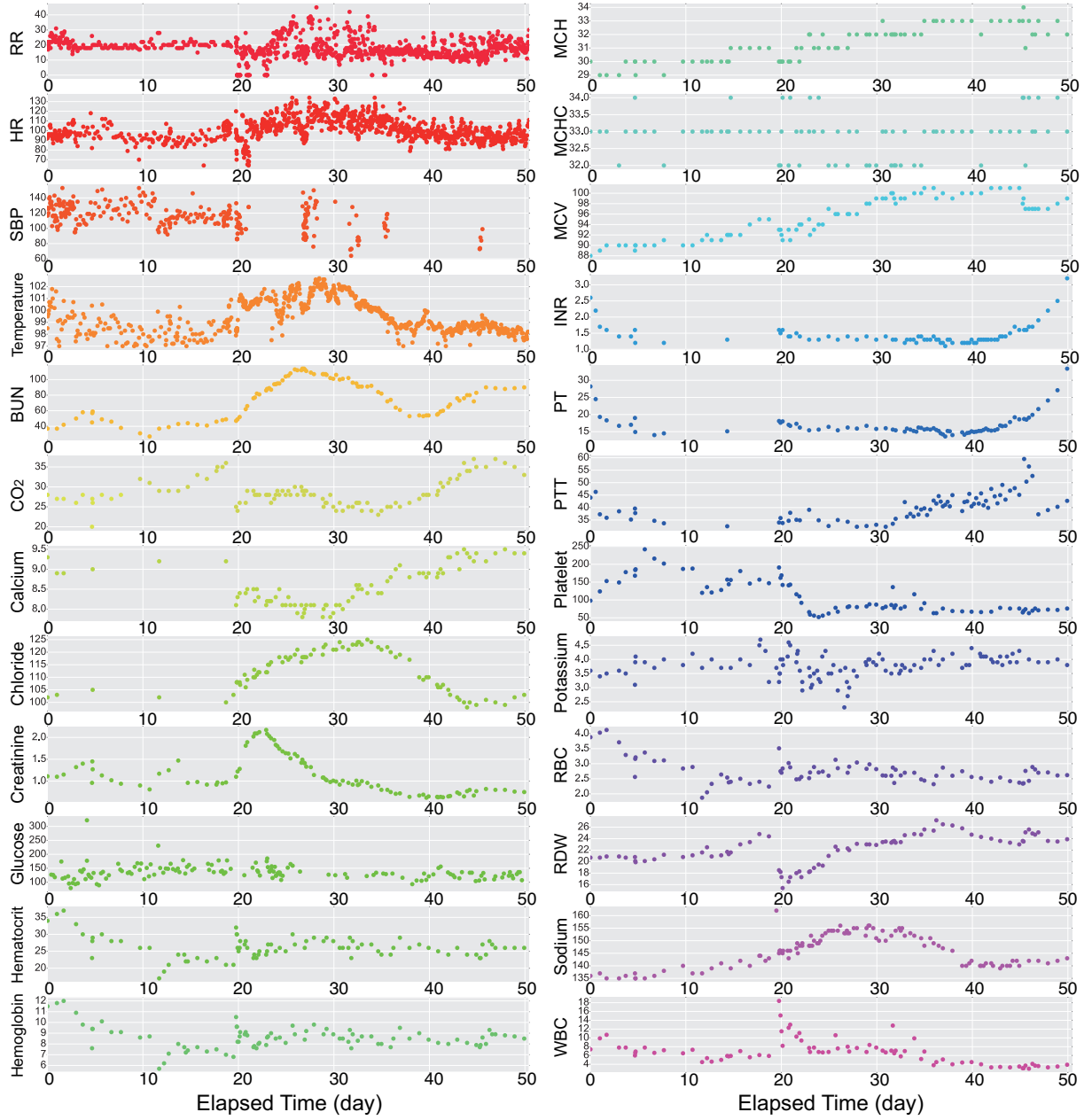


Figure 1: **An example of time series data of 24 clinical covariates for a septic patient in the HUP data.** The 24 covariates include four vital signs—respiration rate (RR), heart rate (HR), systolic blood pressure (SBP), body temperature—and 20 lab results. The time series are aligned by the patient’s admission time. The density of sampling varies widely over the 24 covariates. A full description of these covariates can be found in Table 1.

ical distribution from reference patients with patient-specific observations as measurements are observed. We evaluate our method, MedGP, on over 6,000 patients from three disease groups with more than four million measurements from the HUP data and one disease group from the MIMIC-III data set. We compare results to state-of-the-art approaches for patient online monitoring and investigate differences in correlations among covariates across disease groups.

2. Related Work

Related work falls into three areas of medical time series analysis: (i) incorporating noisy, heterogeneous, irregular, and sparsely sampled time series data; (ii) combining information across multiple time series; and (iii) exploiting reference data in addition to observations about the current patient to enable patient-specific predictions for a new hospital patient.

Most prior work has focused on modeling each clinical covariate separately. Due to the irregularity and temporal sparsity of medical data, conventional time series models, such as hidden Markov models (HMMs), autoregressive (AR) models, state-space models, and linear dynamical systems (LDS), are challenging to apply because of the assumption of regular measurement sampling in time. Recent work has focused on developing methods to compensate for the missing data in order to work with models that assume complete data. In Kim et al. (2010), missing data were imputed by averaging over a time window using a kernel support vector machine (SVM). Methods such as matrix factorization and k -nearest neighbor (KNN) clustering were used for missing data imputation, and improvements in septic shock prediction were reported (Ho et al., 2014). In other work, a hierarchical switching LDS model was used to monitor the physiological signals during neonatal sepsis; the model allows the latent state of a patient to change during periods with fewer observations (Stanulescu et al., 2014). In an alternative approach, noisy and sparse time series data were smoothed temporally by putting Gaussian priors on the mean parameters of the Gaussian mixture model, which is related to a Gaussian process prior, although the distribution is over a finite-dimensional vector (Marlin et al., 2012).

Gaussian processes (GPs) are useful approaches for time series analysis because they can naturally capture irregular time series observations and estimate prediction uncertainties in a probabilistic framework (Roberts et al., 2012). For these reasons, GPs have been applied to the analysis of medical time series data. Previous work used a single-output GP regression model to smooth and impute each covariate independently (Stegle et al., 2008; Lasko et al., 2013). The Probabilistic Subtyping Model (PSM) added patient-specific information for smoothing temporal trajectories of clinical covariates and clustering disease subtypes (Schulam et al., 2015). PSM learns a mixture model based on a B-spline and GPs to impute the clinical measurements for patients with scleroderma. Demographic covariates, including gender, ethnicity, and clinical history, were also incorporated in the model. In an extension of PSM, the authors adapted patient-specific information to forecast specific clinical covariates (Schulam and Saria, 2015); the time series for each covariate was still modeled independently.

The idea of capturing the joint dynamics between vital signs and lab tests has also been explored. Using high-frequency regularly sampled time series, the dynamics between heart rate (HR) and blood pressure (BP) were modeled using a mixture of an LDS model (Nemati

et al., 2012) and a switching vector autoregressive model (SVAR) (Lehman et al., 2015). The joint dynamics estimated across covariates were reported to be associated with hospital mortality. In other work (Rizopoulos and Ghosh, 2011), a multivariate spline-based approach with linear mixed effects was used to predict multiple longitudinal outcomes and time-to-death of patients. Time series graphical models (TGMs) (Dahlhaus, 2000; Tank et al., 2015) have also been studied and applied for analyzing multivariate medical time series of ICU patients (Gather et al., 2002). TGMs model the partial correlations between each dimension of the multivariate time series as an undirected graph. However, both TGMs and SVAR models follow the assumptions of vector autoregressive (VAR) models, and thus assume the sampling interval of the time series is fixed across dimensions. In practice, this means missing data imputation needs to be done in advance (Tank et al., 2015).

Several multi-output GP frameworks have been proposed for other application areas. In the geostatistics literature, the linear model of coregionalization (LMC) characterizes correlations between outputs through a set of kernels and coregionalization matrices that estimate weights for pairwise outputs (Journel and Huijbregts, 1978; Goovaerts, 1997). In the machine learning literature, related models include multi-task GPs (Bonilla et al., 2008), semiparametric latent factor models (Teh et al., 2005), and multi-task kernel learning (Titsias and Lázaro-Gredilla, 2011). These can be viewed as variations of the LMC with different parameterizations and constraints. Convolution processes (CPs) have also been adapted to model multiple correlated outputs through the convolution of smooth kernels and latent processes (Álvarez and Lawrence, 2011). This approach usually has fewer hyperparameters and more efficient computation as compared to LMC, but only squared exponential (SE) kernels have been shown to be computationally tractable. Applying a multi-task GP (MTGP) framework (Bonilla et al., 2008) to clinical time series analysis has also been considered in two studies (Ghassemi et al., 2015b; Dürichen et al., 2015); both studies considered one patient as one task and used the remaining patients as reference training data. Other work adapted the LMC framework with one SE kernel to model three sparsely sampled clinical covariates (intracranial pressure, mean arterial blood pressure, and Pressure-Reactivity Index) jointly (Ghassemi et al., 2015b). The MTGP was shown to outperform a single-task GP (STGP) in prediction error. Both MTGP and CP have also been used with an SE kernel to model three densely sampled vital signs (respiration rate, systolic blood pressure, and heart rate); both methods showed improvements as compared to a single-task GP (Dürichen et al., 2015).

Our work is distinct from previous research in several ways. First, we use the GP regression framework to model multiple irregularly sampled medical time series using a sparse structured multi-output kernel. In contrast to related work (Ghassemi et al., 2015b; Dürichen et al., 2015), our kernel uses a mixture of flexible spectral kernels (Wilson and Adams, 2013), allowing periodic behavior and both short-term and long-term dependencies within and across the clinical covariates over time. Second, we use the LMC framework to enable an interpretable quantification of cross-correlation and sparsity between covariates. Third, we model many more clinical covariates (24) compared with previous studies (at most three); in the online medical setting, efficient and scalable computation in this multi-view model is essential. To the best of our knowledge, this is the first work that uses a sparse and low-rank formulation of the shared covariance matrix across clinical covariates

to estimate and regularize the relationships between covariates in order to learn about covariate relationships specific to patient subgroups and to prevent overfitting.

In our methodology, MedGP, we trained a GP model on each reference patient separately, and used these models to estimate the empirical population-level model using nonparametric density estimation. This approach avoids training procedures that iterate through all reference patients, which is computationally intractable for an online system (Ghassemi et al., 2015b; Dürichen et al., 2015). To speed up training, we optimized the implementation in C++ using multithreading. Finally, in order to personalize the model for a new patient, we update the empirical population-level model on-the-fly to estimate patient specific parameters as measurements from the new patient are observed.

3. Methods

In this section, we describe our method, MedGP, for estimating the underlying dynamic processes jointly across a large number of sparsely sampled clinical covariates. We first describe the design of the Gaussian process kernel for capturing the temporal correlations within and between covariates. Next, we introduce the sparsity-inducing prior to regularize the LMC weight matrix. We then describe estimation of the parameters in the prior and the kernel. Next, we describe how to learn a patient-specific kernel by first building a population-level model from reference patients and then performing online updating of the parameters when observations about a new patient accumulate. Finally, we describe methods to perform computationally tractable online inference in these models, concluding with a discussion of computational complexity.

3.1 Gaussian Processes (GPs)

Gaussian processes (GPs) are distributions over arbitrary functions. By definition, a Gaussian process is a collection of random variables, any finite collection of which have a joint Gaussian distribution. Alternatively, a GP can be described as a distribution on an arbitrary function, defined as

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x}')), \quad (1)$$

where $m(\mathbf{x})$ is the *mean function*:

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})], \quad (2)$$

and $\kappa(\mathbf{x}, \mathbf{x}')$ is the *covariance function* or *kernel*:

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]. \quad (3)$$

Any finite number of function values jointly have a multivariate Gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix \mathbf{K} between any pair of observations, defined by the kernel function,

$$\begin{aligned} [f(x_1), f(x_2), \dots, f(x_T)]^\top &\sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}), \\ \boldsymbol{\mu} &= [m(x_1), m(x_2), \dots, m(x_T)]^\top, \end{aligned} \quad (4)$$

$$\mathbf{K}_{i,j} = \kappa(x_i, x_j).$$

Properties of the function $f(\mathbf{x})$ such as smoothness or periodicity are determined by the kernel function $\kappa(\mathbf{x}, \mathbf{x}')$. One of the most commonly used kernels is the squared exponential (SE) kernel

$$\kappa(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right), \quad (5)$$

which is parameterized by a length scale ℓ and a scale factor σ . The functions generated by a GP with an SE kernel are smooth because the kernel function is infinitely differentiable (Rasmussen and Williams, 2006). The value of the length scale ℓ determines the distribution of changes over the function value with respect to changes in the input \mathbf{x} , encouraging a specific smoothness. Due to its simplicity, SE is used in many applications; however, the properties of the functions that it captures are fairly limited. Periodic functions, for example, are not well modeled by an SE kernel, but instead captured by a periodic kernel

$$\kappa(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left[-\frac{4 \sin^2\left(\frac{\pi \|\mathbf{x} - \mathbf{x}'\|}{p}\right)}{\ell^2}\right], \quad (6)$$

where p is the period of the function. When modeling medical time series, the SE kernel or the periodic kernel are often used in combination to capture the unknown source-specific smoothness and periodicity of the trajectories of clinical covariates (Stegle et al., 2008; Dürichen et al., 2015).

3.2 Gaussian Process Regression with a Structured Multi-Output Kernel

Our first goal is to jointly model multiple clinical covariates—vital signs and lab tests—over time for each patient using GP regression. For the i th patient, we denote the time series of the d th covariate as a vector $\mathbf{x}_{i,d}$, representing the time points that the d th covariate was observed, and the corresponding observation vector $\mathbf{y}_{i,d}$:

$$\mathbf{x}_{i,d}^\top = [x_{i,d,1}, x_{i,d,2}, \dots, x_{i,d,t}, \dots, x_{i,d,T_{i,d}}], \quad (7)$$

$$\mathbf{y}_{i,d}^\top = [y_{i,d,1}, y_{i,d,2}, \dots, y_{i,d,t}, \dots, y_{i,d,T_{i,d}}], \quad (8)$$

where t indexes time, and $T_{i,d}$ is the total number of observations for the d th covariate of the i th patient.

To represent the time series data over all D covariates, we define

$$\mathbf{x}_i^\top = [\mathbf{x}_{i,1}^\top, \mathbf{x}_{i,2}^\top, \dots, \mathbf{x}_{i,D}^\top], \quad (9)$$

$$\mathbf{y}_i^\top = [\mathbf{y}_{i,1}^\top, \mathbf{y}_{i,2}^\top, \dots, \mathbf{y}_{i,D}^\top], \quad (10)$$

where $\mathbf{x}_i, \mathbf{y}_i \in \mathbb{R}^{T_i \times 1}$, $T_i = \left(\sum_{d=1}^D T_{i,d}\right)$. Let \mathcal{F}_i be a multi-output function over time for the i th patient. We capture the relationship between time and clinical observations as a GP regression model:

$$\mathbf{y}_i = \mathcal{F}_i(\mathbf{x}_i) + \boldsymbol{\epsilon}_i, \quad (11)$$

where $\boldsymbol{\epsilon}_i$ is the residual noise vector. Marginally at the t th observation of the d th covariate, the residual noise is modeled as

$$\epsilon_{i,d,t} \sim \mathcal{N}(0, \sigma_{i,d}^2), \quad (12)$$

where $\sigma_{i,d}^2$ is the covariate-specific residual variance for each individual.

We assume that the function \mathcal{F}_i is drawn from a patient-specific Gaussian process \mathcal{GP}_i with mean function $\mu_i(\mathbf{x})$ and kernel $\kappa_i(\mathbf{x}, \mathbf{x}')$:

$$\mathcal{F}_i \sim \mathcal{GP}_i(\mu_i(\mathbf{x}), \kappa_i(\mathbf{x}, \mathbf{x}')). \quad (13)$$

We set $\mu_i(\mathbf{x}) = \mathbf{0}$ (Rasmussen and Williams, 2006).

We designed the kernel $\kappa_i(\mathbf{x}, \mathbf{x}')$ to capture predictive and generalizable covariance structure across medical time series data. Assuming the covariates are correlated across time, we adapted the linear model of coregionalization (LMC) framework (Journel and Huijbregts, 1978; Goovaerts, 1997). We used a set of Q *basis kernels* $\{\kappa_q(\mathbf{x}, \mathbf{x}')\}_{q=1}^Q$ to model D covariates jointly. The kernel for the cross-covariance of any pair of covariate types is modeled by a weighted, structured linear mixture of the Q basis kernels. The full joint kernel is written as a block structured function

$$\kappa_i(\mathbf{x}_i, \mathbf{x}'_i) = \sum_{q=1}^Q \left(\begin{bmatrix} b_{q,(1,1)} \kappa_q(\mathbf{x}_{i,1}, \mathbf{x}'_{i,1}) & \cdots & b_{q,(1,D)} \kappa_q(\mathbf{x}_{i,1}, \mathbf{x}'_{i,D}) \\ b_{q,(2,1)} \kappa_q(\mathbf{x}_{i,2}, \mathbf{x}'_{i,1}) & \cdots & \vdots \\ \vdots & \ddots & \vdots \\ b_{q,(D,1)} \kappa_q(\mathbf{x}_{i,D}, \mathbf{x}'_{i,1}) & \cdots & b_{q,(D,D)} \kappa_q(\mathbf{x}_{i,D}, \mathbf{x}'_{i,D}) \end{bmatrix} \right), \quad (14)$$

where $b_{q,(d,d')}$ scales the covariance (defined by the q th basis kernel) between covariates d and d' , and $\kappa_i(\mathbf{x}_i, \mathbf{x}_i) \in \mathbb{R}^{T_i \times T_i}$. We collapsed $b_{q,(d,d')}$ into a set of weight matrices $\{\mathbf{B}_q\}_{q=1}^Q$, where each \mathbf{B}_q is a symmetric positive definite matrix

$$\mathbf{B}_q = \begin{bmatrix} b_{q,(1,1)} & b_{q,(1,2)} & \cdots & b_{q,(1,D)} \\ b_{q,(1,1)} & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ b_{q,(D,1)} & b_{q,(D,2)} & \cdots & b_{q,(D,D)} \end{bmatrix} \in \mathbb{R}^{D \times D}. \quad (15)$$

If the inputs are the same for all covariates, we can further simplify Eq. (14) with Kronecker product \otimes . That is, if $\mathbf{x}_{i,1} = \mathbf{x}_{i,2} = \cdots = \mathbf{x}_{i,D} \triangleq \mathbf{x}_{i,*}$ and $\mathbf{x}'_{i,1} = \mathbf{x}'_{i,2} = \cdots = \mathbf{x}'_{i,D} \triangleq \mathbf{x}'_{i,*}$:

$$\kappa_i(\mathbf{x}_i, \mathbf{x}'_i) = \sum_{q=1}^Q \mathbf{B}_q \otimes \kappa_q(\mathbf{x}_{i,*}, \mathbf{x}'_{i,*}), \quad (16)$$

although in practice we do not often see this situation in medical time series data.

The properties of the time series observations, such as periodicity and short term dependencies, are captured in the Q basis kernels. For medical covariates, the properties of each patient's time series observations may vary. As a trivial example, when a patient is under age 18, their pulse will be well correlated with their age, height, and weight; above age 18, the correlation among pulse, age, height, and weight is more variable within age than across ages. Furthermore, only a few vital signs, such as heart rate, blood pressure, and body temperature, are known to be periodic with a 24-hour period (i.e., a circadian rhythm), but whether there is a similar period for specific lab results, such as white blood cell count or pressure of carbon dioxide in the blood, is unclear (Widmaier et al., 2004).

To handle the heterogeneity of patterns within covariates and across patients, we selected the spectral mixture (SM) kernel as the basis kernel (Wilson and Adams, 2013). The SM kernel is a general form of a variety of stationary kernels, including the squared exponential (SE) kernel and the periodic kernel, and has also shown good performance in modeling processes generated from more complex kernels through a mixture of kernels approach (Wilson and Adams, 2013). The basis kernel $\kappa_q(x_t, x_{t'})$ is written as

$$\begin{aligned}\kappa_q(x_t, x_{t'}) &= \exp(-2\pi^2\tau^2v_q) \cos(2\pi\tau\mu_q), \\ \tau &= |x_t - x_{t'}| \quad (\text{absolute distance in time}).\end{aligned}\tag{17}$$

In our work, the mixture weights for each basis kernel are encoded in \mathbf{B}_q .

To be used for GP regression, $\kappa_i(\mathbf{x}, \mathbf{x}')$ must be a valid Mercer kernel, i.e., the Gram matrix must be positive definite for all \mathbf{x} and \mathbf{x}' . Since the matrix produced by each basis kernel $\kappa_q(\mathbf{x}, \mathbf{x}')$ is symmetric positive definite, we only need to ensure that every \mathbf{B}_q is positive definite to produce a Mercer kernel. To do this, we parameterized \mathbf{B}_q as

$$\mathbf{B}_q = \mathbf{A}_q \mathbf{A}_q^\top + \begin{bmatrix} \lambda_{q,1} & 0 & \cdots & 0 \\ 0 & \lambda_{q,2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_{q,D} \end{bmatrix} = \mathbf{A}_q \mathbf{A}_q^\top + \text{diag}(\boldsymbol{\lambda}_q),\tag{18}$$

$$\mathbf{A}_q = \begin{bmatrix} a_{q,(1,1)} & \cdots & a_{q,(1,R_q)} \\ \vdots & \ddots & \vdots \\ a_{q,(D,1)} & \cdots & a_{q,(D,R_q)} \end{bmatrix}\tag{19}$$

Here $\mathbf{A}_q \in \mathbb{R}^{D \times R_q}$, $\boldsymbol{\lambda}_q \in \mathbb{R}^{D \times 1}$. We let R_q denote the number of non-zero columns in \mathbf{A}_q , or the rank for \mathbf{B}_q when $\boldsymbol{\lambda}_q = \mathbf{0}$.

For any two observations from the same patient of different covariates at different times, denoted as $x_{i,d,t}$ and $x_{i,d',t'}$, the prior covariance from the GP kernel is

$$\kappa_i(x_{i,d,t}, x_{i,d',t'}) = \sum_{q=1}^Q b_{q,(d,d')} \kappa_q(x_t, x_{t'}).\tag{20}$$

We summarize the parameters and hyperparameters of our SM-LMC kernel in Table 2.

3.3 Sparsity-Inducing Priors on Weight Matrix \mathbf{B}_q

As the number of medical covariates included in the model increases, we need to increase the number of basis kernels Q and corresponding R_q to allow greater representational flexibility. However, too many basis kernels may lead to overfitting and will become computationally intractable. To avoid this, we regularized the elements of each weight matrix \mathbf{B}_q by introducing structured sparsity-inducing priors on each \mathbf{A}_q matrix as follows.

We included two layers of sparsity-inducing priors for flexible, data-adaptive shrinkage behavior, modified from previous work (Polson and Scott, 2010; Gao et al., 2013). First, we put column-wise sparsity-inducing priors to regularize each column in \mathbf{A}_q . This corresponds to regularizing the degree of freedom of the functions, or number of latent processes

notation	size	description
v_q	Q	squared exponential part of q th basis kernel
μ_q	Q	periodicity of q th basis kernel
$a_{q,(d,r)}$	$\sum_{q=1}^Q D \times R_q$	weights of (d, d') for q th basis kernel
$\lambda_{q,(d)}$	$Q \times D$	intra-covariate weights of the d th covariate for q th basis kernel $\mathbf{B}_q = \mathbf{A}_q \mathbf{A}_q^\top + \text{diag}(\boldsymbol{\lambda}_q)$

Table 2: The list of hyperparameters for modeling the $d = 1 : D$ clinical variables and $q = 1 : Q$ mixture kernels.

generated from each basis kernel in the LMC model (Álvarez et al., 2012). Second, we put sparsity-inducing priors on each matrix element $a_{q,(d,r)}$ in \mathbf{A}_q to produce element-wise sparsity. The effect of element-wise sparsity is to perform model selection on the number of basis kernels that each pair of covariates uses for covariance representation. Finally, we put sparsity-inducing priors on the elements of $\boldsymbol{\lambda}_q$ to shrink the covariance for observations from the same covariate.

In practice, we implemented each layer of the prior as a two-layer hierarchical gamma distribution. The generative model is written as

$$\begin{aligned}
\tau_{q,(r)} &\sim \text{Gamma}(d, \eta), \\
\phi_{q,(r)} &\sim \text{Gamma}(\gamma, \tau_{q,(r)}), \\
\delta_{q,(d,r)} &\sim \text{Gamma}(\beta, \phi_{q,(r)}), \\
\psi_{q,(d,r)} &\sim \text{Gamma}(\alpha, \delta_{q,(d,r)}), \\
a_{q,(d,r)} &\sim \mathcal{N}(0, \psi_{q,(d,r)}),
\end{aligned} \tag{21}$$

where each element $a_{q,(d,r)}$ has a Gaussian distribution. Parameters $\phi_{q,(r)}$ and $\tau_{q,(r)}$ control the column-specific shrinkage, while parameters $\psi_{q,(d,r)}$ and $\delta_{q,(d,r)}$ control the local shrinkage of each element in the \mathbf{A}_q matrix. For vector $\boldsymbol{\lambda}_q$, we regularized each element with a local Laplace prior:

$$\lambda_{q,(d)} \sim \text{Laplace}(0, \beta_\lambda). \tag{22}$$

For our results, we set $\alpha = \beta = \gamma = d = 0.5$ to recapitulate two layers of the horseshoe prior, using a statistically equivalent prior represented by a hierarchical gamma with four layers (Carvalho et al., 2010; Armagan et al., 2011; Gao et al., 2013; Zhao et al., 2016). Parameters $\psi_{q,(d,r)}$, $\delta_{q,(d,r)}$, $\phi_{q,(r)}$, and $\tau_{q,(r)}$ were estimated during optimization. We set $\beta_\lambda = 0.01$ to regularize the diagonal terms $\lambda_{q,(d)}$. The hyperparameter η controls the overall shrinkage profile of the hierarchical gamma prior (see Appendix A for more details). We chose η over $\{0.01, 0.1, 1.0\}$ using cross-validation prediction error. Hyperparameter η was chosen using grid search over the range $\{0.01, 0.1, 1.0\}$ using cross-validation prediction error as the objective.

3.4 Parameter Learning

To estimate the parameters for the regularized kernel, we optimized the posterior probability. We denote all parameters that were estimated directly as $\boldsymbol{\theta}$ and hyperparameters in

the sparsity-inducing prior as $\boldsymbol{\theta}_f$:

$$\boldsymbol{\theta} = \left\{ \mu_q, v_q, a_{q,(d,r)}, \lambda_{q,(d)}, \psi_{q,(d,r)}, \delta_{q,(d,r)}, \phi_{q,(r)}, \tau_{q,(r)} \right\}, \quad (23)$$

for $q = 1, \dots, Q \quad d = 1, \dots, D \quad r_{|q} = 1, \dots, R_q$

$$\begin{aligned} \boldsymbol{\theta}_f &= \{\alpha, \beta, \gamma, d, \eta, \beta_\lambda\}, \\ \alpha &= \beta = \gamma = d = 0.5. \end{aligned} \quad (24)$$

The posterior density of our model is then

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}_f) &\propto p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\boldsymbol{\theta}_f) \\ &\propto p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \left[\prod_{q=1}^Q \prod_{d=1}^D \prod_{r=1}^{R_q} p(a_{q,(d,r)}|\psi_{q,(d,r)}) p(\psi_{q,(d,r)}|\alpha, \delta_{q,(d,r)}) p(\delta_{q,(d,r)}|\beta, \phi_{q,(r)}) \right] \\ &\quad \times \left[\prod_{q=1}^Q \prod_{r=1}^{R_q} p(\phi_{q,(r)}|\gamma, \tau_{q,(r)}) p(\tau_{q,(r)}|d, \eta) \right] \left[\prod_{q=1}^Q \prod_{d=1}^D p(\lambda_{q,(d)}|\beta_\lambda) \right] \left[\prod_{q=1}^Q p(v_q) p(\mu_q) \right]. \end{aligned} \quad (25)$$

The term $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$ is found by calculating the GP marginal likelihood given the values of $\boldsymbol{\theta}$ (Rasmussen and Williams, 2006), which is

$$\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{y}^\top (K_{|\boldsymbol{\theta}} + \epsilon I)^{-1} \mathbf{y} - \frac{1}{2} \log |K_{|\boldsymbol{\theta}} + \epsilon I| - \left(\frac{\sum_{d=1}^D T_{i,d}}{2} \right) \log(2\pi). \quad (26)$$

We thus estimated $\boldsymbol{\theta}$ by solving the optimization problem:

$$\arg \max_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}_f) = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}). \quad (27)$$

See Eq. (31) in Appendix B for the derivation of $\mathcal{Q}(\boldsymbol{\theta})$.

Due to the conjugacy of the hierarchical gamma priors, we optimized parameters $\psi_{q,(d,r)}$, $\delta_{q,(d,r)}$, $\phi_{q,(r)}$, $\tau_{q,(r)}$ directly using maximum a posteriori (MAP) estimates of their posterior distribution (or mean when the mode does not exist). Our optimization procedure then consists of two parts. In the first part, we used the update equations to estimate $\psi_{q,(d,r)}$, $\delta_{q,(d,r)}$, $\phi_{q,(r)}$, and $\tau_{q,(r)}$ directly. In the second part, we estimated parameters μ_q , v_q , $a_{q,(d,r)}$, and $\lambda_{q,(d)}$ using a scaled conjugate gradient method to find the local maximum, conditioned on $\hat{\psi}_{q,(d,r)}$, $\hat{\delta}_{q,(d,r)}$, $\hat{\phi}_{q,(r)}$, and $\hat{\tau}_{q,(r)}$. (Details can be found in Appendix B Eq. (32)–(35) and Eq. (36)–(40).) We iterated over the two steps until the change in $\mathcal{Q}(\boldsymbol{\theta})$ reached the convergence criterion (< 0.005) or until the maximum number of iterations (≥ 30).

3.5 Estimating the Population-Level Model and Online Updating

The GP with the structured kernel described above lets us model the patient-specific joint dynamics between covariates within the same patient. We now describe how we built a population-level empirical prior from a set of mixture kernels estimated from all training patients, and how we apply this empirical prior to a new patient.

To estimate the empirical priors across reference patients, we trained one GP kernel for each patient separately, and then we clustered and extracted the distribution of the basis

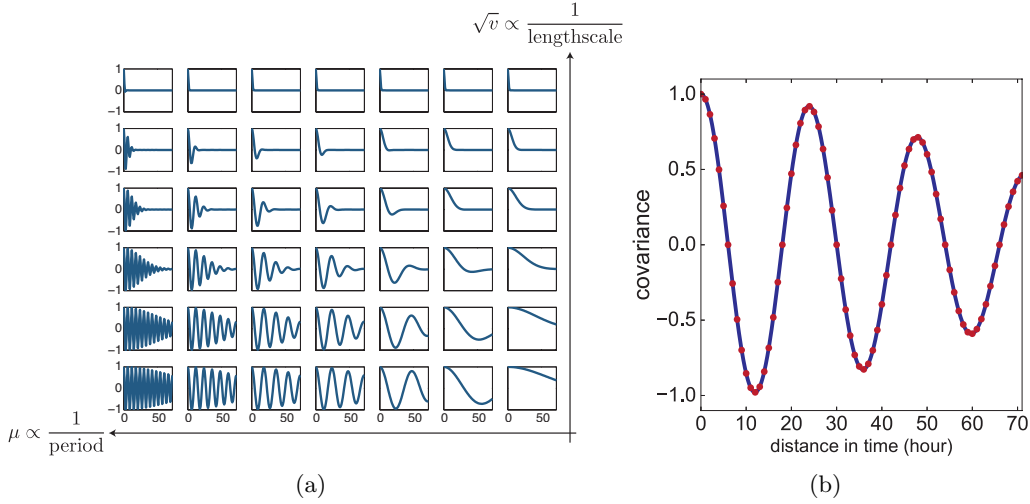


Figure 2: **Illustrations of the basis kernels and the features for kernel clustering.** (a) An example of a discrete set of basis kernels with different μ and v within a 72-hour window. (b) An example of the 72-dim temporal features (shown with red dots) taken from a kernel for GMM clustering.

kernels (defined by hyperparameters μ_q and v_q). The idea here is that, when we observe a set of estimated mixture kernels, we would like to understand the high-level properties of these mixture kernels shared across covariates and patients in the same patient group, and then estimate the distributions of these hyperparameters through observations of basis kernels belonging to this cluster. For instance, a circadian rhythm (24-hour periodicity) may be observed in some covariates for some patients, but the period across patients could vary within a range. Across the space of μ and v the spectral kernels vary substantially (Figure 2a). For each basis kernel that was estimated, the characteristic period is $1/\mu_q$ and the length scale is $1/2\pi\sqrt{v_q}$ (Wilson and Adams, 2013). There are different ways to define the features of a kernel. Here we used the temporal features of the learned kernels directly (Figure 2b). The temporal spacing of two adjacent points is one hour, and we use the kernel values within a window of length 72 hours. We then used a Gaussian mixture model (GMM) model to perform clustering on the kernels, and we chose the best number of kernel clusters Q' ($1 \leq Q' \leq Q$) based on Bayesian information criterion (BIC). For the MedGP implementation, we adapted the open source scikit-learn package (Pedregosa et al., 2011). We used version 0.18.1, with ten random restarts, a maximum of 2,000 iterations, and allowing each mixture component to have its own covariance matrix.

For each identified kernel cluster, we estimated one set of parameters μ_q and v_q for the basis kernel, and the weight coefficients—elements in \mathbf{B}_q matrices, computed using the \mathbf{A}_q matrices and $\boldsymbol{\lambda}_q$ vectors. We do this by building an empirical distribution using kernel density estimation (KDE) with a Gaussian kernel over the GP kernel hyperparameters assigned to that cluster. The bandwidth of the kernel density estimator was chosen based on Silverman’s “rule-of-thumb” (Silverman, 1986). We estimated each new parameter using density weighted means with the density from the univariate KDE as the weights. Note

that if there were multiple kernels in a patient cluster, the estimated \mathbf{B}_q matrices were added based on the additive assumption of our kernel before aggregating to estimate the population-level kernel for that cluster. To allow online updating, we estimated the elements of the new empirical \mathbf{A}_q matrix and $\boldsymbol{\lambda}_q$ vector corresponding to each new \mathbf{B}_q matrix using singular value decomposition (SVD). For the univariate GP regression, we did not use density weighted means because we found them to be unstable; instead we used a grid-based search to identify the hyperparameters with the highest kernel density estimates.

As the number of vital signs and lab measurements of a new patient accumulate, we update the hyperparameters to estimate a patient-specific kernel. Indeed, we update the kernel sequentially every time a new observation arrives. To do this in a computationally tractable way, we used the momentum method (Rumelhart et al., 1988) with a 72-hour window of previous observations to update the kernel hyperparameters when predicting the value of next observation. For all experiments, we chose the momentum as 0.9 and the learning rate as 10^{-5} . For elements in the \mathbf{A}_q matrices, we do not update the values if the elements were regularized to be zero so as to maintain the empirical sparsity structure.

3.6 Efficient Inference in MedGP

The main bottleneck of our method is in learning patient-specific kernel hyperparameters. Let $T_i = \sum_{d=1}^D T_{i,d}$ denote the total number of samples of the i th patient; the computational cost to compute the Gram matrix is $\mathcal{O}(QT_i^2)$, which increases linearly with the chosen number of basis kernels. To find the MAP estimates of the parameters, we need to invert and compute the determinant of the Gram matrix ($K_{|\boldsymbol{\theta}} + \epsilon I$) in Eq. (26). The computational complexity for the full matrix inversion is $\mathcal{O}(T_i^3)$ using Cholesky decomposition. When calculating the gradients for optimizing the hyperparameters, the cost is dominated by $\mathcal{O}(QDRT_i^2)$ after the inverse Gram matrix is pre-computed, which is linear with the total number of the kernel hyperparameters. In practice, the complexity of each iteration is either $\mathcal{O}(T_i^3)$ or $\mathcal{O}(QDRT_i^2)$. That is, the patient with the most measurements is the main bottleneck for training. In our implementation, we mitigate the bottleneck using optimized linear algebra functions in Intel MKL library with multithreading and computing the gradients of the hyperparameters in parallel.

3.7 Medical Data Preprocessing

The HUP medical time series data consist of electronic health records (EHRs) from more than 260,000 patients admitted to a University of Pennsylvania Hospital. For each patient, the data include many heterogeneous clinical covariates, including ICD-9 codes, patient demography, length-of-stay, vital signs, and lab results. We jointly modeled the 24 covariates with the greatest number of observations across patients (Table 1). We selected three groups of discharged patients from these data: 1,365 septic patients, 952 patients with heart failure, and 4,723 patients with neoplasms. Each patient has at least one observation for each of the 24 covariates, and in total over four million observations were evaluated.

For each clinical covariate, we first removed obvious artifacts (e.g., values outside of the possible range in living humans). For the patients with neoplasms or heart failure, we used the full patient length-of-stay in training and testing. For septic patients, the disease progression varies substantially across patients, and the distribution of the covariates

changes dramatically depending on the disease phase. To address this issue, we segmented the time series data into four disjoint partitions based on clinical status: *no sepsis*, *pre-sepsis*, *sepsis*, and *recovery*. To label each stage, we incorporated prior clinical domain knowledge. For instance, we identified sepsis stages using ICD-9 codes and positive blood culture results. Since our model assumes stationarity, to better estimate the temporal correlation across covariates, we chose the *recovery* stage before the patients’ discharge to test our method, since this is a relatively stable stage. We used the bed unit information to identify if the patient is in a stable state. That is, when a patient is transferred to step-down bed, we labeled the time series after the transfer as *recovery*. The median length-of-stay after pre-processing is 140 hours for the sepsis group, 285 hours for the heart failure group, and 197 hours for the neoplasms group.

We applied similar preprocessing procedure to the MIMIC-III data. We selected patients with a heart failure diagnosis that eventually had a routine discharge. We removed artifacts such as out of bounds values for each covariate, and applied the criteria to each patient that at least five measurements were taken for all 24 selected covariates. We extracted 1,004 heart failure under these criteria and used 1,003 of them, excluding one patient with more than 50K measurements due to memory constraints.

3.8 Experimental Setup

We applied MedGP to the three selected groups of patients separately, and evaluated characteristics and performance of MedGP under two different experimental settings. In the first analysis, we evaluated the model’s ability to learn the covariance between a pair of highly correlated clinical covariates, and we measured the imputation performance in an online setting. In the second analysis, we follow the same online setting, but instead jointly model all 24 clinical covariates, including four vital signs and 20 lab covariates. In both settings, we evaluated our method using 10-fold cross-validation at the patient level. That is, for each fold we ran the kernel clustering step on the kernels from the training patients to estimate a set of population-level basis kernels and \mathbf{B}_q matrices. This set of kernels was then applied to the held-out patients to predict the value of each covariate using observations from all other covariates measured at the same time as, or earlier than, the test observation (i.e., no future information included). After each prediction, we updated the patient-specific kernel parameters using the new observations from the test patient.

We compared our method to several univariate methods that modeled each covariate separately: (i) a naive one-lag prediction procedure, which predicts an observation equal to the last observation available from the same patient; (ii) an independent GP with squared exponential (SE) or spectral mixture (SM) kernels fitting each covariate separately (we tested with $Q = 1$ for SM); (iii) the multi-resolution Probability Subtyping Model (PSM) combining linear regression, B-splines, and independent GPs (Schulam et al., 2015). To estimate the spectral kernel parameters, for each patient we initialized 1,000 random kernels by drawing uniformly from a length scale range (between 6 and 72 hours) and period range (between 24 and 72 hours). We computed the marginal likelihood of all random kernels for each patient, and then initialize optimization using the kernels with the highest marginal likelihood. The elements in the \mathbf{A}_q matrices are initialized randomly between -1.5 and 1.5 .

We compared results from MedGP to these various methods using two metrics: (i) mean absolute error (MAE) of the predicted observations with the true observations, and (ii) 95% coverage, the percentage of true observations that fell within the predictive 95% confidence region. We quantified and reported the improvements with respect to both metrics compared to all three baselines (naive prediction, univariate GP, and PSM). To test if the differences in prediction results from different approaches were statistically significant, we performed paired t-tests for the results of each covariate and compared the p -values with a Bonferroni corrected threshold (dependent on the number of jointly modeled covariates in each experiment).

We note that the original PSM was designed to model scleroderma disease (Schulam et al., 2015). Thus, to make it applicable to our different patient groups, several adjustments were made. First, we omitted the population and environmental factors selected for their relevance to scleroderma. Second, we chose the knots of the B-spline basis by sampling every hour for vital signs and every 24 hours for lab results between zero and the longest length-of-stay for patients in each disease group. Third, to make PSM training feasible on the scale of our data set, we limited the maximum number of subtypes to ten for the sepsis and heart failure groups, and 20 for the neoplasms group.

4. Experimental Results

We analyzed the performance of the method, MedGP—multi-output GP with a sparse SLMC kernel and online updating—by applying it to time series data from the Hospital of the University of Pennsylvania (HUP) and the public MIMIC-III data set (Johnson et al., 2016). We ran two types of experiments—one on two correlated lab covariates and the other with 24 covariates jointly. The results were compared against the baseline methods for prediction accuracy and 95% coverage calibration.

4.1 Results of Two Lab Covariates

As a proof of principle, we jointly modeled two well correlated lab covariates, prothrombin time (PT) and international normalization ratio (INR) on three HUP subgroups. PT measures the time it takes for the plasma in the blood to clot, and is often ordered to check bleeding problems. INR is an international standard for PT to account for possible variations across different labs. For the same patient, the two covariates usually have similar trajectories over time (Figure 1).

We trained the kernels for one patient’s INR and PT time series data both with and without the structured sparse prior (Figure 3). Both \mathbf{A}_q and \mathbf{B}_q matrices estimated using the sparse prior have higher levels of sparsity versus those estimated without using the sparse prior. We observed that for both methods, one of the estimated basis kernels κ_1 captures long-term (around one month) dependencies. However, with the sparse prior, the estimated weights associated with this long term kernel \mathbf{A}_1 are rank one instead of rank two. This means the trajectories of the two covariates are similar enough to be explained by one instead of two functions, and thus fewer hyperparameters. Moreover, two basis kernels were found with zeros weights \mathbf{A}_2 and \mathbf{A}_5 (Figure 3b), suggesting that the prespecified number of basis kernels may be reduced. We also found that the off-diagonal elements in the \mathbf{B}_q matrices in both cases have nonzero values, suggesting a nonzero covariance between PT and INR.

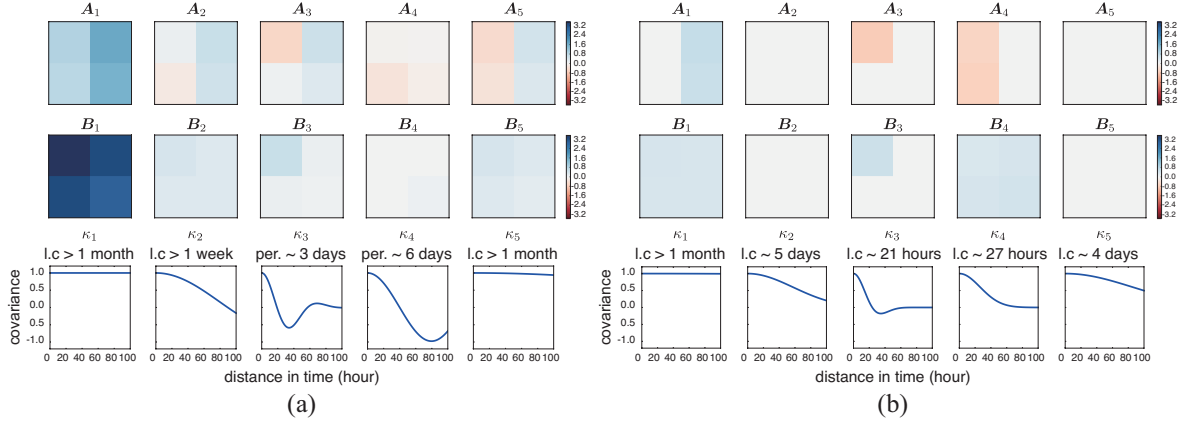


Figure 3: **The trained kernel for one patient jointly modeling PT and INR.** For both the (a) SM-LMC kernel and (b) sparse SM-LMC kernel, the \mathbf{A}_q matrices (upper row), \mathbf{B}_q matrices (middle row), and the basis kernel κ_q (bottom row) are illustrated. The zero elements are colored in light grey. Here *l.c.* denotes length scale for each basis kernel, and *per.* denotes period. The length-of-stay for this patient was over 90 days.

observations. In particular, two basis kernels captured the covariance between PT and INR: one with a greater than one-month trend (Figure 3b, \mathbf{B}_1 and κ_1), and one with a 27-hour trend (Figure 3b, \mathbf{B}_4 and κ_4). Here, the sparse kernel has 18 non-zero hyperparameters, whereas there are 40 for the non-sparse kernel. We can compare the two fitted kernels using both log marginal likelihoods and model selection scores. The log marginal likelihoods of the two kernels are -118.16 (SM-LMC) and -128.50 (sparse SM-LMC), indicating a better fit for the SM-LMC model without sparsity. However, the Bayesian information criterion (BIC) values, which take into account the number of parameters in a model, were 353.63 (SM-LMC) and 309.79 (sparse SM-LMC), where values closer to zero reflect better models. Thus, using a sparse prior has the advantage of a more compact kernel representation.

We then ran our model on all three disease groups separately, and compared our method with the univariate baselines described in Section 3.8 under the scenario of online imputation of the same two well-correlated clinical covariates. For independent GPs, we used gradient descent to optimize the hyperparameters. For PSM, we performed grid search for the parameters of the B-spline and the independent GP kernel. For our method, we set $Q = 5$ and $R_q = 2$ for the \mathbf{A}_q matrices for training. In the sepsis and heart failure groups, three nonzero basis kernel functions ($Q' = 3$) were found for the model using the SM-LMC kernel, while only two nonzero basis kernel functions ($Q' = 2$) were found using the sparse SM-LMC kernel; the number of nonzero hyperparameters were 18 and 12 respectively. In the neoplasms group, the number of nonzero basis kernels were the same as the pre-specified number ($Q' = Q = 5$). With 10-fold cross-validation, we found that results using the SM-LMC kernel showed smaller imputation error than those using the baselines for both PT and INR (Figure 4). The mean absolute errors (MAEs) showed that the non-sparse SM-LMC kernels perform imputation the best among the related approaches. On the other

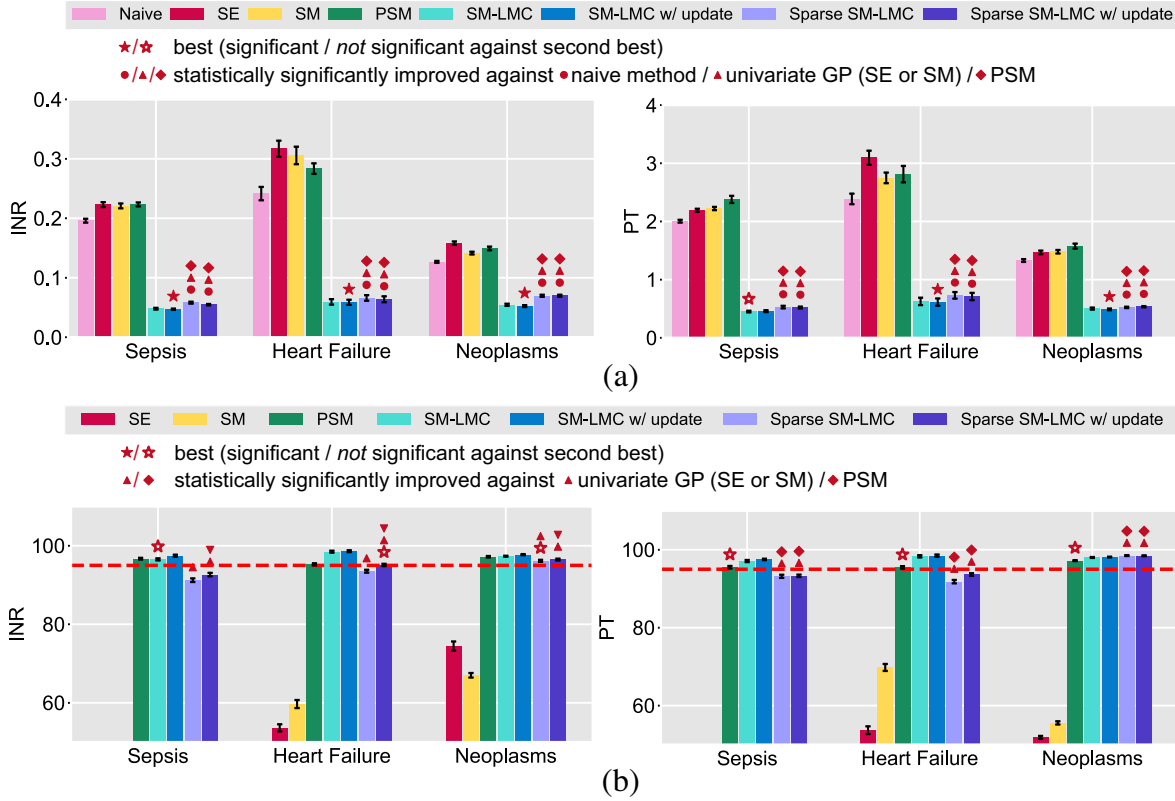


Figure 4: **The results of prediction when jointly modeling INR and PT.** The figure illustrates (a) mean absolute error (MAE), and (b) 95% coverage (the dashed red line indicates 95%). The error bars denote ± 1 standard error.

hand, looking at the 95% coverage, results using non-sparse or sparse SM-LMC kernel were well calibrated with respect to the confidence region compared with independent GPs, although sometimes slightly worse than PSM. Note that in this experiment we used a p-value threshold $p < 0.005$ to detect statistical significance, which reflects the Bonferroni correction. The results indicate that the sparse prior finds models with sparse structure while maintaining prediction performance in this two covariate case.

4.2 Results of a Joint Model Including 24 Vital Signs and Lab Covariates

In the second experimental setting, we jointly modeled 24 vital signs and clinical covariates ($D = 24$) for all three disease groups (Table 1). We set the number of basis kernels $Q = 5$ and the number of nonzero columns in A_q as $R_q = 8$ in this experiment for the three HUP subsets. For the MIMIC-III heart failure subset, we set $Q = 4$. More detailed results of the best setup as well as the results with different Q could be found in Appendix C and Appendix D.

4.2.1 ESTIMATING POPULATION-LEVEL KERNELS

We first visualized the population-level kernels estimated from the three patient groups of the HUP data (Figure 5–7) and the MIMIC-III patient subgroup (Figure 8). We observed shared patterns in the basis kernels κ_q and the weight matrices \mathbf{B}_q across all patient groups. Comparing the estimated population-level kernels, we found at least one long-term smoothing basis kernel with length scale longer than three days, and one 24- to 25-hour periodic basis kernel, which indicates the existence of circadian rhythms in specific covariates as expected. Furthermore, in the neoplasms group, which consists of more patients than the other two groups, we found additional short-term smoothing basis kernels and one 12- to 13-hour periodic basis kernel, which may correspond to known circasemidian rhythm of clinical covariates, such as body temperature. We also observed an 11-hour periodic kernel in the MIMIC-III subset.

In addition to the characteristics of the basis kernels, our model with the sparse prior also showed interpretable cross-covariate patterns (Figure 5b, Figure 6b, and Figure 7b and Figure 8b). Based on the \mathbf{B}_q matrices, we identified groups of well correlated covariates. For instance, lab covariates hematocrit (Hct), hemoglobin (Hgb), and red blood cell (RBC) count showed the highest levels of correlation. Since both Hct and Hgb are known to be proportional to the number of red blood cells, this positive correlation was encouraging (Widmaier et al., 2004). The pair of lab covariates studied in the previous section, INR and PT, also showed substantial positive correlation. We found that the four vital signs—respiration rate (RR), heart rate (HR), systolic blood pressure (SBP), and body temperature (Temp)—had substantial correlations with each other as well as weak correlations with some lab covariates. Another identifiable set of well-correlated covariates includes lab measurements of carbon dioxide (CO_2), calcium, chloride, potassium, and sodium. The three lab covariates related to the concentration of hemoglobin—mean cell hemoglobin (MCH), mean cell volume (MCV), and mean cell hemoglobin concentration (MCHC)—appeared to have substantial correlation (Figure 5). The correlations modeled in these covariance matrices are exploited for accurate prediction and imputation in the MedGP framework.

To learn more about the importance of each kernel type across all subsets, we visualized the percent coverage of each type of kernel clusters found in all the subsets we have worked on (Figure 9). The coverage of each kernel type is computed as the ratio of patients that have non-zero \mathbf{B}_q matrix corresponding to it. We found that the kernel clusters of long-term (length scale > 3 days) and short-term (length scale < 12 hours) smooth dependencies have the highest coverage across four subsets. In the MIMIC-III subset, the coverages of the short-term kernel, and the 12-hour and 24-hour periodic kernels are higher than that of in the HUP subsets. We think this is because the higher sampling frequency in the MIMIC-III subset enables more accurate estimation of the short-term and periodic dependencies.

4.2.2 RESULTS FOR ONLINE IMPUTATION

Next, we used the trained kernels to perform online imputation for each patient subgroup, where the goal is to predict the next observation for each covariate given the observations at previous time points. Across these methods, we used the percentage of improvement in MAE over three types of baselines—naïve prediction, univariate GP (with SE or SM

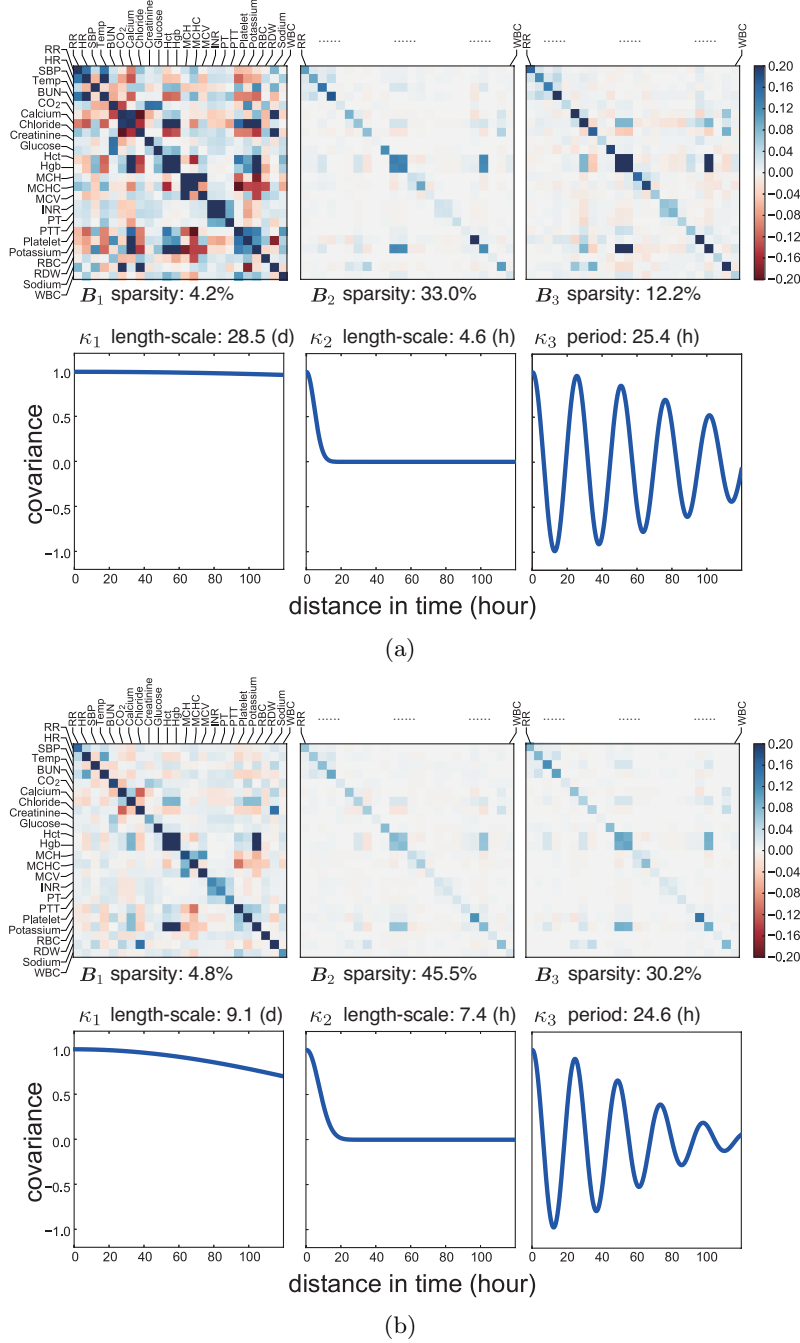


Figure 5: **The estimated population-level basis kernels and corresponding B_q matrices for septic patients.** We show the kernels estimated (a) without a sparse prior ($Q' = 3$) and (b) with a sparse prior ($Q' = 3$). The sparsity of the B_q matrices is calculated as the percentage of nearly zero entries (i.e., values $\leq 10^{-3}$). The units for length scale or period are (d) for days and (h) for hours.

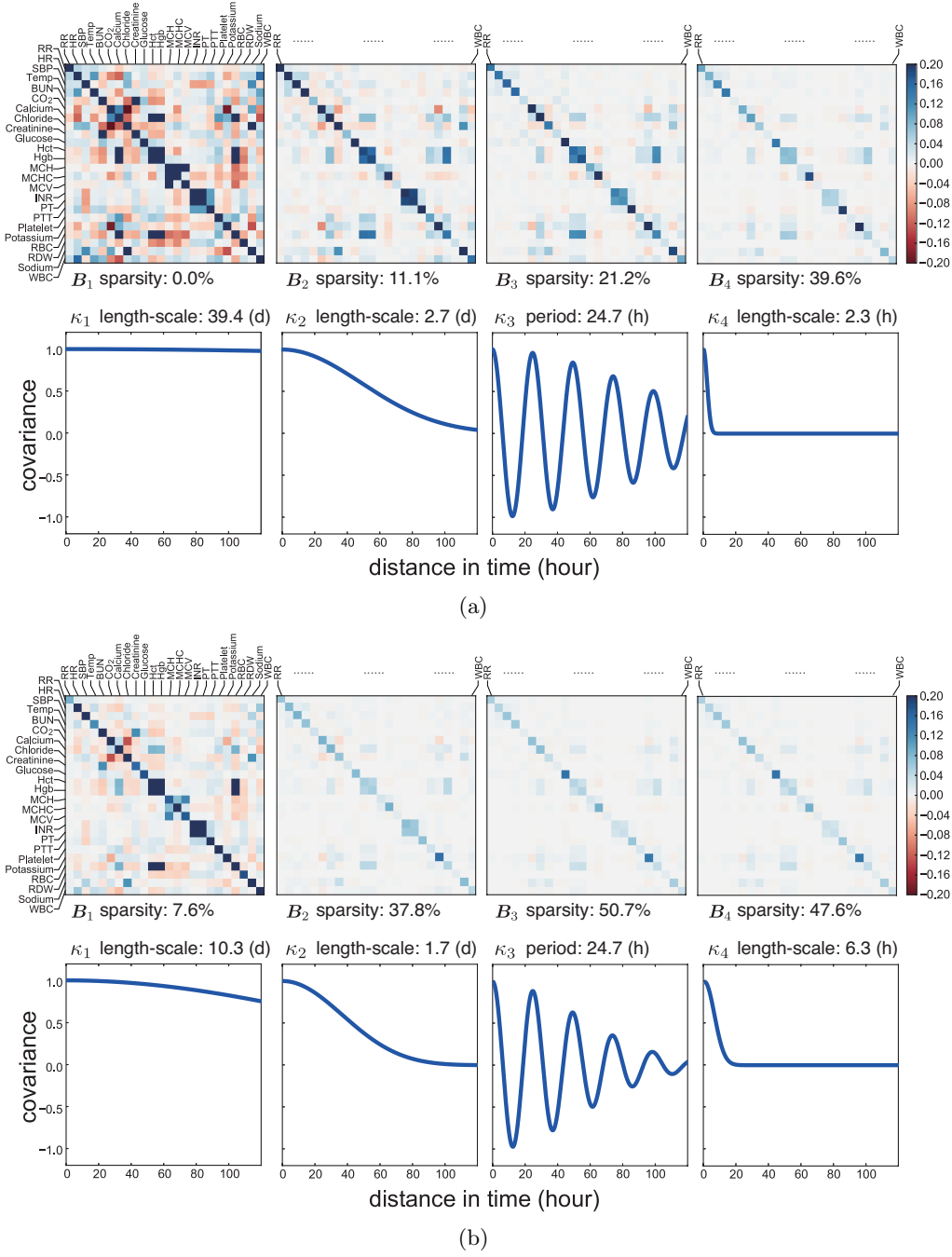
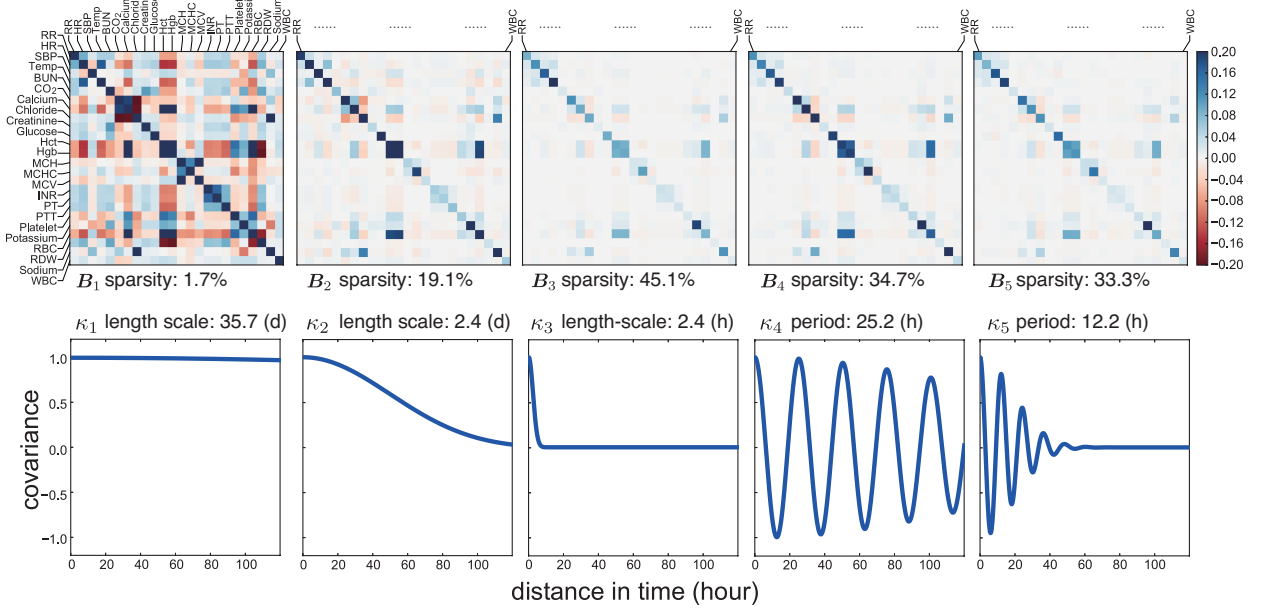
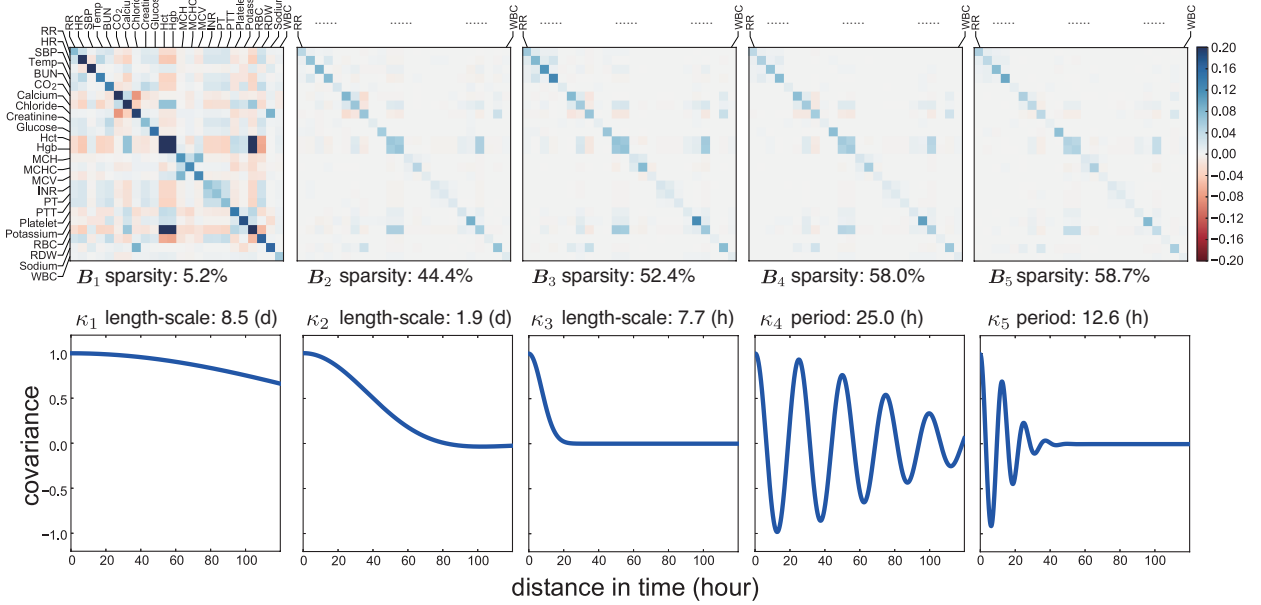


Figure 6: **The estimated population-level basis kernels and corresponding B_q matrices for patients with heart failure.** We show the kernels estimated (a) without a sparse prior ($Q' = 4$) and (b) with a sparse prior ($Q' = 4$). The sparsity of the B_q matrices are calculated as the percentage of nearly zero entries (i.e., values $\leq 10^{-3}$). The units for length scale or period are (d) for days and (h) for hours.



(a)



(b)

Figure 7: **The estimated population-level basis kernels and corresponding B_q matrices for patients with neoplasms.** We show the kernels estimated (a) without a sparse prior ($Q' = 5$) and (b) with a sparse prior ($Q' = 5$). The sparsity of the B_q matrices are calculated as the percentage of nearly zero entries (i.e., values $\leq 10^{-3}$). The units for length scale or period are (d) for days and (h) for hours.

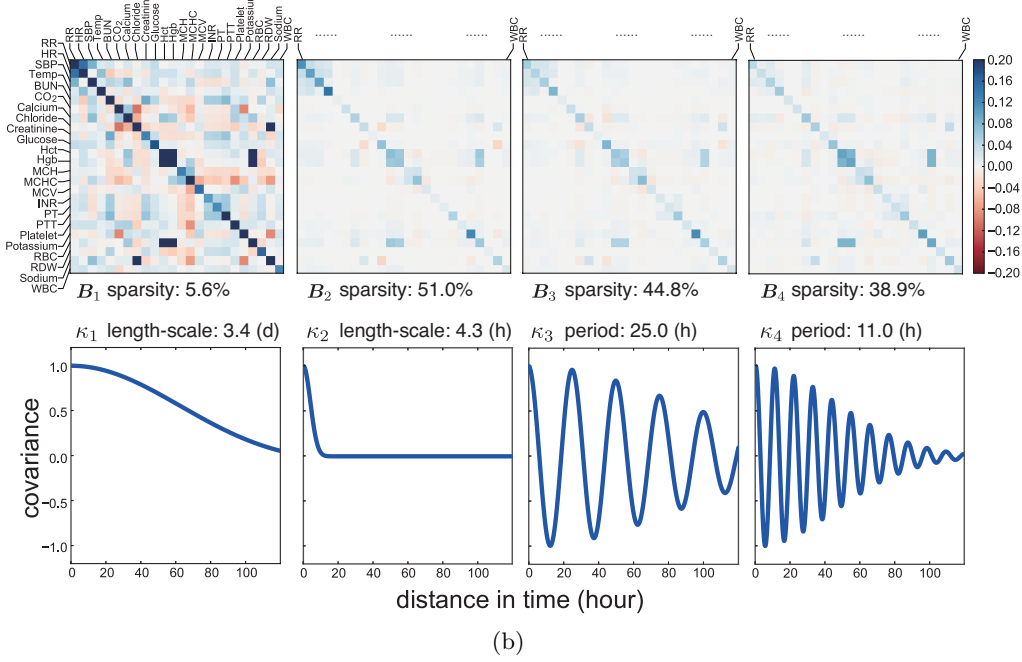
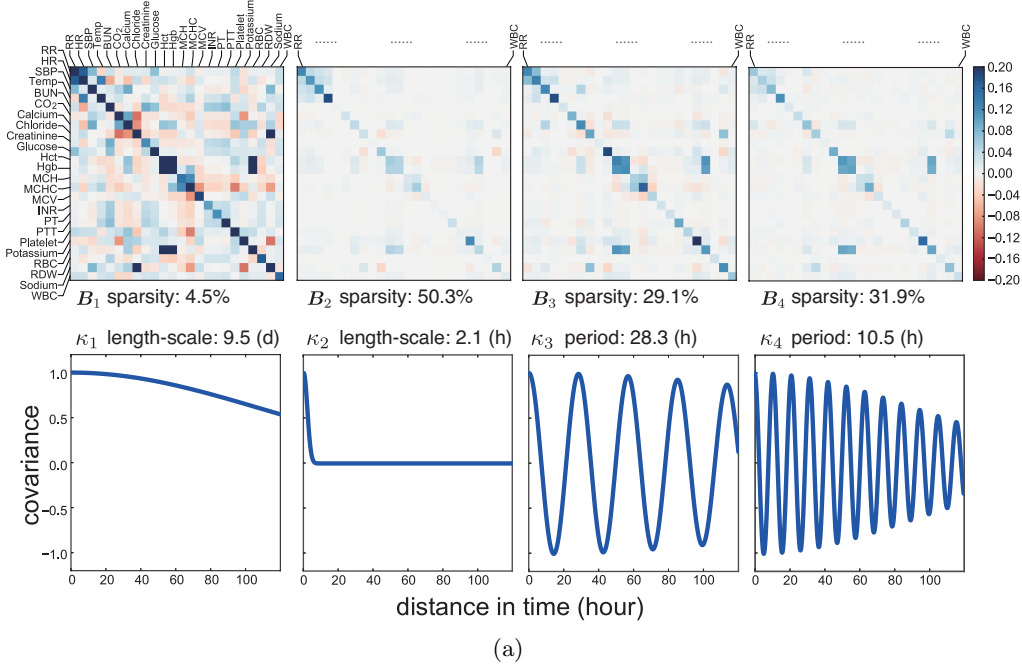


Figure 8: **The estimated population-level basis kernels and corresponding B_q matrices for 1003 patients with heart failure in MIMIC-III data set.** We show the kernels estimated (a) without a sparse prior ($Q' = 4$) and (b) with a sparse prior ($Q' = 4$). The sparsity of the B_q matrices is calculated as the percentage of nearly zero entries (i.e., values $\leq 10^{-3}$). The units for length scale or period are (d) for days and (h) for hours.

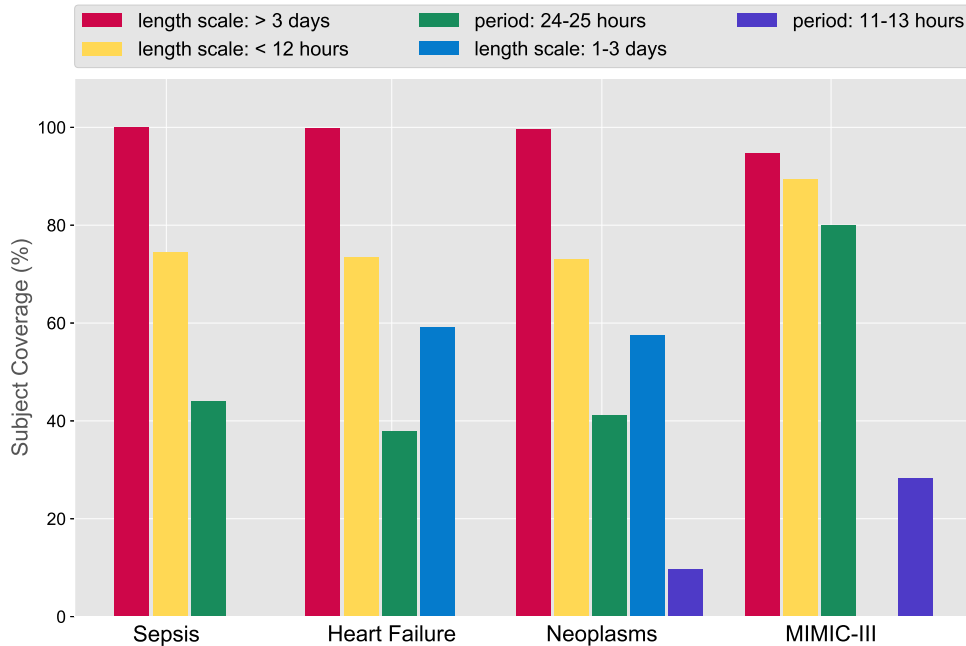


Figure 9: **The coverage over subject for each discovered kernel.** We illustrated the proportion of subjects that have non-zero \mathbf{B} matrix of a kernel that is in the same cluster of the population-level kernel clusters.

kernel), and PSM—to compare results for each of the 24 clinical covariates; we visualized the results separately (Figure 10–12; Figure 14–17 in Appendix C; Figure 18–33 in Appendix D). We also show the results of variations of our method for comparison (with or without the proposed sparse prior; with or without online updating). We performed paired t-tests on predictions from MedGP and each baseline to quantify the improvements, and statistical significance was evaluated using Bonferroni-corrected $p < 4.17 \times 10^{-4}$.

Comparing results with the independent GP model—specifically, selecting the best results from the SE or SM kernel, we found that MedGP, and in particular sparse SM-LMC with online updating, outperformed the independent GP model on the online imputation task for most covariates across the four patient groups (Figure 10). In the HUP data, we found 18, 21, 22 covariates significantly improved by MedGP in the sepsis, heart failure, and neoplasms subgroups respectively. In the MIMIC-III subset, we found 19 covariates were improved. For all four groups, the number of covariates that were improved significantly by MedGP is greater than using SM-LMC kernels without the sparse prior. We found that the covariates that were well correlated in \mathbf{B}_q usually showed significant positive improvements over independent GPs; Hct, Hgb, and RBC are notable examples. Similar observations could be made for INR and PT, the pair of lab covariates studied previously (Figure 4). Across 24 covariates, the MAEs for INR and PT were slightly worse compared with only modeling these two covariates. However, we also observed that using the sparse prior with the SM-LMC kernel led to better performance as compared to not using the sparse prior,

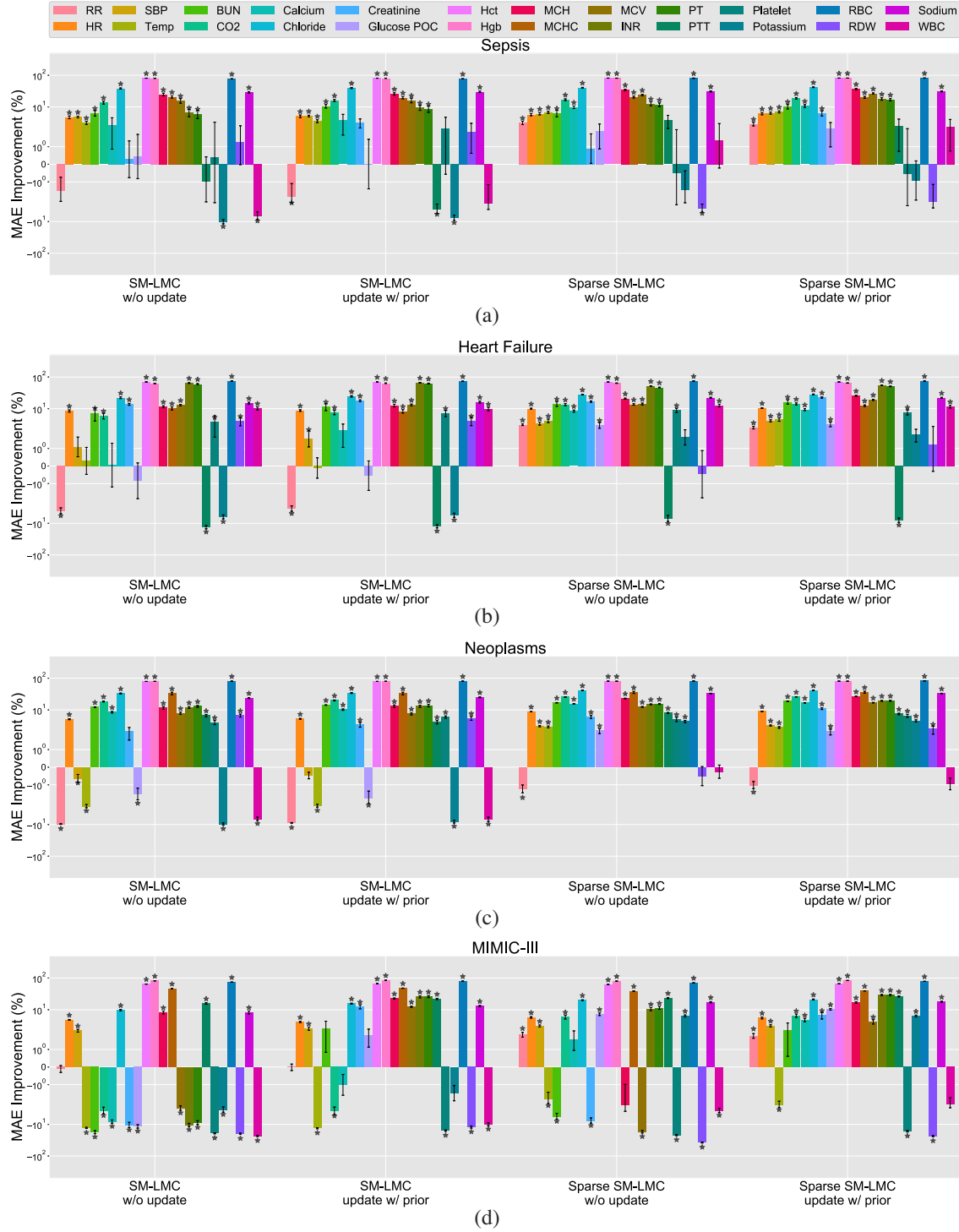


Figure 10: **The percent improvement using MedGP for online imputation compared to independent (univariate) GPs.** The figures depicts the results of 24 covariates for the (a) sepsis, (b) heart failure, and (c) neoplasms and (d) MIMIC-III heart failure subgroups. The y -axis is on log scale. The error bars denote ± 1 standard error. The \star indicates statistical significance.

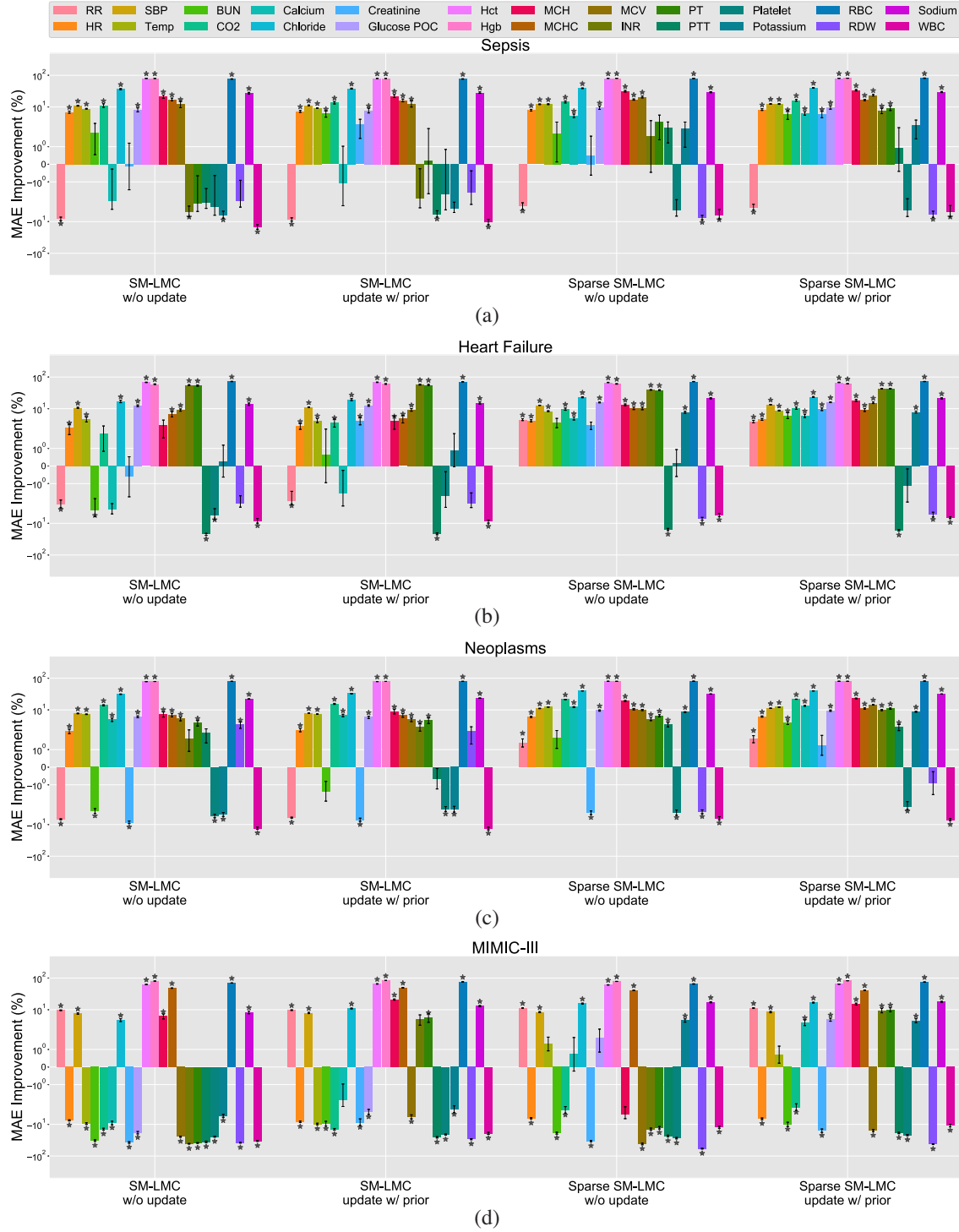


Figure 11: **The percent improvement using MedGP for online imputation compared to the naive method.** The figures depicts the results of 24 covariates for the (a) sepsis, (b) heart failure, and (c) neoplasms and (d) MIMIC-III heart failure subgroups. The y -axis is on log scale. The error bars denote ± 1 standard error. The \star indicates statistical significance.

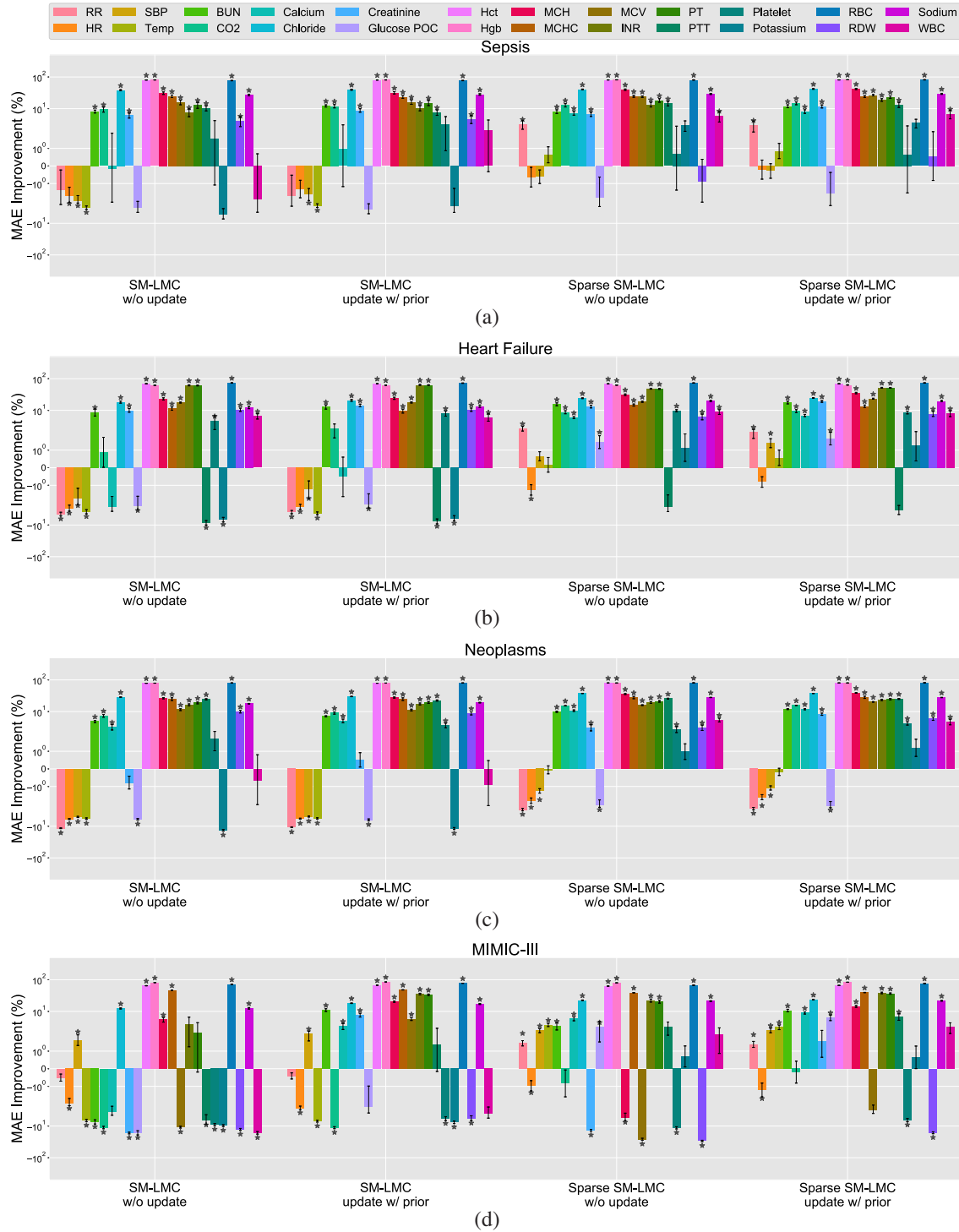


Figure 12: **The percent improvement using MedGP for online imputation compared to PSM.** The figure depicts the results of 24 covariates for the (a) sepsis, (b) heart failure, and (c) neoplasms and (d) MIMIC-III heart failure subgroups. The y -axis is on log scale. The error bars denote ± 1 standard error. The \star indicates statistical significance.

indicating that sparse regularization is helpful when jointly modeling heterogeneous covariates. Finally, there were some covariates for which MedGP did not improve over univariate GPs in two or more disease groups, including red cell distribution width (RDW), white blood cell count (WBC) and platelets.

When the baseline method is the naive one-lag method, for all three disease groups, we found fewer covariates with significant improvements compared with improvements over univariate GPs (Figure 11). In particular, the covariates for which the naive method had an advantage were lab covariates that have piece-wise linear behavior, such as mean cell hemoglobin (MCH) and mean cell hemoglobin concentration (MCHC) (Figure 1). In the case of piece-wise linear behavior, our kernel does not improve the performance compared with the naive approach since the time series are neither smooth nor periodic. Moreover, we also found that the naive method performed better in respiration rate, PTT, platelet, RDW, and white blood cell (WBC) count. Overall, however, our method improved online prediction results for 18, 20, 20 of the 24 covariates in sepsis, heart failure, and neoplasms groups, respectively. In the MIMIC-III subset, we found 14 covariates were improved significantly over the naive method.

When the baseline method is PSM (Schulam et al., 2015), we found that our method outperformed PSM for most of the lab covariates, but PSM outperformed MedGP in imputation of vital signs and two lab covariates: glucose point-of-care (Glucose POC) and potassium (Figure 12). For vital signs and glucose level, PSM has an advantage because of a higher sampling rate in those covariates and the highly structured mean function in PSM in the HUP subsets. The sampling rates are usually every 4 hours for vital signs and every 8 hours for glucose, which is more frequent than other lab covariates. Since PSM uses a B-spline basis function to capture the empirical mean, it may tolerate non-stationarity better. However, in the MIMIC-III subset, we observed that our method improved in imputing glucose and three vital signs (RR, SBP, temperature) over PSM significantly. We think this reflects the higher sampling rate of the covariates that allows better estimation of the short-term temporal dependencies. Overall, MedGP significantly improved the imputation of 17, 20, 18 covariates in sepsis, heart failure, neoplasms subsets respectively in the HUP data set, and 16 covariates in the MIMIC-III subset when compared with PSM. We contrast the PSM approach of structuring the mean function with our approach of structuring the kernel function, which leads to different types of gains in this problem.

Next, we looked at the calibration of the 95% coverage (Figure 16 and Figure 17 in Appendix C; Figure 26–33 in Appendix D). We found that MedGP outperformed independent GPs in terms of calibration of the 95% confidence region for all covariates. For this evaluation, the values closer to 95% are better. We observed that the coverage using the non-sparse SM-LMC kernel was usually higher than the coverage using the sparse SM-LMC kernel in the three HUP subgroups, indicating that MedGP may slightly underestimate covariate-specific noise. In contrast, in the MIMIC-III subset, we observed that MedGP gave consistently more accurate 95% coverage than without regularization in most covariates. We also found that, in all patient subsets, online updating significantly improves the accuracy of the 95% coverage. Among all tested methods, PSM tended to overestimate the 95% confidence region. We think this is because PSM assumes that the input time series are aligned by patient status, and this alignment is not the case in our data. With unaligned data, PSM learned large marginal variance parameters due to high empirical variance of

Implementation	Sequential	Multithreading
Computing Gram matrix	11	2
Inverting Gram matrix	13	3
Computing gradients	2497	97
Total per iteration	2521	102

Table 3: **Training time (in seconds) for a single iteration under different implementations of MedGP.** The total number of observations across time for this patient is 6,679. The sequential test used a single CPU, while the multithreading test used 35 CPUs—one thread per CPU.

the observations across patients at the same elapsed time. In contrast, the estimation of marginal covariance parameters in MedGP is not affected by alignment because estimates are patient-specific. We also observed that for either MedGP or PSM, the coverage was lower for some covariates in the MIMIC-III subset than in HUP subsets, such as temperature, CO₂, and PTT. This potentially reflects greater non-stationarity in the MIMIC-III subset, whose records were from intensive care units (ICUs) instead of regular hospital beds.

Finally, we compared the prediction performance of MedGP compared with the version without patient-specific online updating. We observed that online updating significantly improves the imputation errors of at least 12 out of 24 covariates in sepsis, heart failure, neoplasms, and the MIMIC-III subset (Figure 14 and Figure 15 in Appendix C; Figure 18–25 in Appendix D). Similarly, evaluating the 95% coverage, all 24 covariates were improved by the online updating across the three diseases groups in HUP, and 18 covariates were improved in the MIMIC-III subset (Figure 16 and Figure 17 in Appendix C; Figure 26–33 in Appendix D). This improvement highlights the importance of updating the empirical priors with patient-specific observations for this problem.

4.3 Computational Efficiency and Scalability

In this section, we compare computational speed between different implementations of our method. For patients with only a few observations, an existing implementation using conventional GP inference is sufficient for computationally tractable online inference. However, since our data include a large number of patients with potentially thousands of observations each, we implemented an exact inference algorithm in C++ and optimized it through Intel MKL libraries and customized multithreading blocks. In the experimental setting of $Q = 5$, $D = 24$, and $R_q = 8$, there are 1114 hyperparameters to estimate. We summarized the runtime under different implementations for one patient with 2,028 unique time points and 6,679 observations (Table 3); the tests were performed using a machine with Intel(R) Xeon(R) CPUs running at 2.40GHz. Using our optimized implementation, for patients with large number of observations ($T_i \geq 5000$), we accelerated training by a factor of 10 to 25 on average as compared with the sequential approach. We also compared our implementation with the standard GPy (GPy, since 2012) implementation under different sample sizes and Q , and reached empirically at least three times speed up. We provide these results in Appendix E.

The proposed framework can be parallelized at the patient level and is suitable for analysis when patient data are observed in a streaming form. For each reference patient, we distributed the optimized training process on a computing cluster to estimate the patient-specific hyperparameters in parallel. In addition, the population-level kernels could be updated sequentially; the computationally expensive GP training procedure does not need to be applied to patient data in bulk. That is, when we receive more data from new patients, we only need to update the kernel density estimators. Our framework provides better computational efficiency compared to models designed for smaller collections of observations (e.g., approximately two hundred observations for each patient) as in most previous work. Those approaches are computationally intractable when working on a set of rich patient observations of the magnitude of the HUP data due to large matrix inversions and summing marginal likelihoods across patients at each iteration.

5. Discussion

In this paper, we propose a flexible and efficient framework for estimating the temporal dependencies across multiple sparse and irregularly sampled medical time series data. We developed a model with multi-output Gaussian process regression with a highly structured kernel. We fit this model using an optimized implementation of exact GP inference to three different disease groups in the HUP medical data set and the MIMIC-III ICU data set. We showed that our method, MedGP, improves performance for online prediction of 24 clinical covariates as compared with independent univariate GPs, a naive method of propagating the previous observation, and an earlier state-of-the-art approach, PSM (Schulam et al., 2015). We found that, for well-correlated covariates, our method improves online imputation performance substantially over the related methods in most tested covariates. The improvements over the naive one-lag prediction and univariate GPs were significant in both vital signs and lab covariates. We found that PSM was, in general, better at predicting vital signs with more densely sampled observations. However, our approach does not require patient time series alignment and shows better calibration of the 95% confidence region as compared to PSM.

There are several directions that will be explored using the MedGP framework motivated by the present results. The first direction is to allow time-varying covariances by specifically modeling non-stationarity. Some possible approaches to explore include incorporating state-space models or change point detection (Adams and MacKay, 2007; Saatçi et al., 2010), and extending those methods to work on multivariate scenarios. Another direction of interest is to consider latent subpopulation-level structured kernels through multivariate medical time series. We expect that our results could be further improved through incorporating hierarchical methods with proper features or metrics to represent the differences between patients within the same disease group and across disease groups more carefully. For instance, the original PSM used three levels of hierarchy based on the subpopulations of patients with scleroderma, including population level, subpopulation level, and individual level. Our model may benefit from such an approach, but more efficient inference procedures are needed to train on our large data set (Feinberg et al., 2017). We should point out that this is possible through, for instance, deriving corresponding stochastic variation inference (SVI) algorithm. For example, previous work develops the SVI algorithm for semiparamet-

ric latent factor model (SLFM) with $R_q = 1$ (Nguyen and Bonilla, 2014), which could be generalized to apply to MedGP.

For future applications, we will use the framework to monitor the health status of patients in a hospital setting and identify those patients at high risk for acute diseases in order to assist with decision making in treatment plans. Specifically, MedGP can impute latent state in patients at any time point, including confidence region around those estimates; this latent state can be used for a number of downstream analyses which require complete knowledge of patient state at specific time points. For instance, the changes of dynamics and temporal correlations between two vital signs have been found to be useful for disease detection given high-frequency regularly sampled time series (Nemati et al., 2012; Lehman et al., 2015). We demonstrated that MedGP accurately estimates the temporal correlations in the presence of sparse, unaligned time series data for up to 24 covariates, and we would expect to further associate the cross-covariate dynamics to more complicated diseases, such as septic shock (Henry et al., 2015), where the interactions of multiple covariates are jointly taken into consideration for diagnosis.

Appendix A. Details of the Hierarchical Gamma Prior

In this appendix, we provide more background and visualization for the hierarchical gamma prior we used for regularization. For the convenience, we use $\Gamma(\cdot)$ to denote the gamma function, and $\mathcal{G}(a, b)$ to represent a gamma distribution with shape parameter a and rate parameter b .

Following Proposition 1 in Armagan et al. (2011), for a random variable x drawn from a normal distribution with two-layered gamma priors on variance

$$x \sim \mathcal{N}(0, \psi_1), \quad \psi_1 \sim \mathcal{G}(\alpha, \delta), \quad \delta \sim \mathcal{G}(\beta, \nu), \quad (28)$$

is equivalent to the hierarchy

$$x \sim \mathcal{N}(0, 1/\rho - 1), \quad \rho \sim \mathcal{TPB}(\alpha, \beta, \nu), \quad (29)$$

where $\mathcal{TPB}(\alpha, \beta, \nu)$ denotes the three-parameter beta distribution. The probability density function of ρ is given as

$$f(\rho; \alpha, \beta, \nu) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \nu^\beta \rho^{\beta-1} (1 - \rho)^{\alpha-1} [1 + (\nu - 1)\rho]^{-(\alpha+\beta)}. \quad (30)$$

In Figure 13, we visualized the density of ρ in Equation (29) for $\alpha = \beta = 0.5$, and under different values of ν (Armagan et al., 2011). In this case, the prior distribution of x is equivalent to a horseshoe prior, and ρ can be interpreted as the shrinkage coefficient (Carvalho et al., 2010).

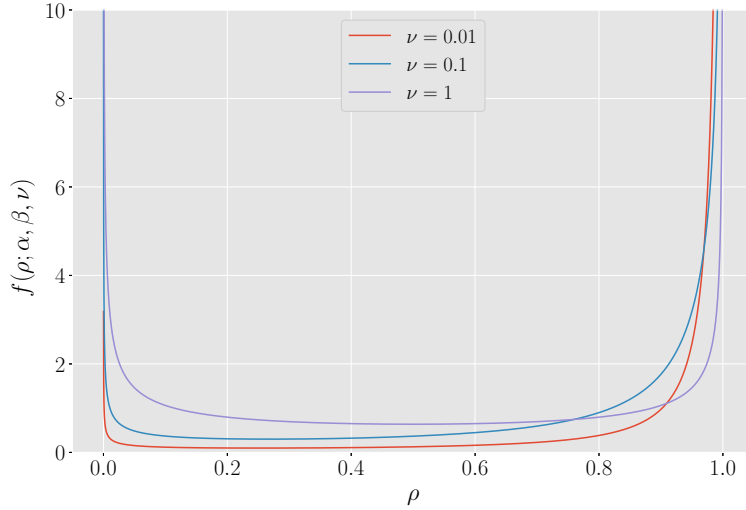


Figure 13: **The density of ρ drawn from a three parameter beta prior with different values of ν . For all values of ν , we set $\alpha = \beta = 0.5$.**

Specifically, for the case with four layers of gamma prior used in our work,

$$x \sim \mathcal{N}(0, \psi_2), \quad \psi_2 \sim \mathcal{G}(\alpha, \delta), \quad \delta \sim \mathcal{G}(\beta, \phi), \quad \phi \sim \mathcal{G}(\gamma, \tau), \quad \tau \sim \mathcal{G}(d, \eta),$$

is equivalent to

$$x \sim \mathcal{N}(0, 1/\rho - 1), \quad \rho \sim \mathcal{TPB}(\alpha, \beta, 1/\zeta - 1), \quad \zeta \sim \mathcal{TPB}(\gamma, d, \eta).$$

In our case, we set $\alpha = \beta = \gamma = d = 0.5$ so both ρ and ζ recapitulate horseshoe priors (Armagan et al., 2011; Gao et al., 2013; Zhao et al., 2016).

Appendix B. Details of Gradient Computation and Update Equations

In this appendix, the equations for the objective function during optimization, update equations for the parameters in the sparsity inducing prior and the gradients for the hyper-parameters of the GP kernel are listed as reference.

The objective function to optimize for training one patient, $\mathcal{Q}(\boldsymbol{\theta})$, is

$$\begin{aligned}
\mathcal{Q}(\boldsymbol{\theta}) \propto & \left[-\frac{1}{2} \mathbf{y}^\top (K_{|\boldsymbol{\theta}} + \epsilon I)^{-1} \mathbf{y} - \frac{1}{2} \log |K_{|\boldsymbol{\theta}} + \epsilon I| - \left(\frac{\sum_{d=1}^D T_{i,d}}{2} \right) \log 2\pi \right] \\
& + \sum_{q=1}^Q \sum_{d=1}^D \sum_{r=1}^{R_q} \left(-\frac{1}{2} \log \psi_{q,(d,r)} - \frac{a_{q,(d,r)}^2}{2\psi_{q,(d,r)}} \right) \\
& + \sum_{q=1}^Q \sum_{d=1}^D \sum_{r=1}^{R_q} [\alpha \log \delta_{q,(d,r)} + (\alpha - 1) \log \psi_{q,(d,r)} - \delta_{q,(d,r)} \psi_{q,(d,r)}] \\
& + \sum_{q=1}^Q \sum_{d=1}^D \sum_{r=1}^{R_q} [\beta \log \phi_{q,(r)} + (\beta - 1) \log \delta_{q,(d,r)} - \phi_{q,(r)} \delta_{q,(d,r)}] \\
& + \sum_{q=1}^Q \sum_{r=1}^{R_q} [\gamma \log \tau_{q,(r)} + (\gamma - 1) \log \phi_{q,(r)} - \tau_{q,(r)} \phi_{q,(r)}] \\
& + \sum_{q=1}^Q \sum_{r=1}^{R_q} (d \log \eta + (d - 1) \log \tau_{q,(r)} - \eta \tau_{q,(r)}) \\
& + \sum_{q=1}^Q \sum_{d=1}^D \left(-\log 2\beta_\lambda - \frac{|\lambda_{q,(d)}|}{\beta_\lambda} \right). \tag{31}
\end{aligned}$$

For update equations, we quoted from Zhao et al. (2016):

$$\hat{\psi}_{q,(d,r)} = \frac{(2\alpha - 3) + \sqrt{(2\alpha - 3)^2 + 8a_{q,(d,r)}^2 \delta_{q,(d,r)}}}{4\delta_{q,(d,r)}} \tag{32}$$

$$\hat{\delta}_{q,(d,r)} = \frac{\alpha + \beta}{\psi_{q,(d,r)} + \phi_{q,(r)}} \tag{33}$$

$$\hat{\phi}_{q,(r)} = \frac{D\beta + \gamma - 1}{\sum_{d=1}^D \delta_{q,(d,r)} + \tau_{q,(r)}} \tag{34}$$

$$\hat{\tau}_{q,(r)} = \frac{\gamma + d}{\phi_{q,(r)} + \eta} \tag{35}$$

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2} \text{tr} \left(\left(\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - K_{|\boldsymbol{\theta}}^{-1} \right) \frac{\partial K_{|\boldsymbol{\theta}}}{\partial \theta_j} \right) \quad \text{where } \boldsymbol{\alpha} = K_{|\boldsymbol{\theta}}^{-1} \mathbf{y}, \theta_j \in \boldsymbol{\theta} \tag{36}$$

$$\begin{aligned}
\frac{\partial \mathcal{Q}(\boldsymbol{\theta})}{\partial a_{q,(d,r)}} &= \frac{1}{2} \text{tr} \left(\left(\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - K_{|\boldsymbol{\theta}}^{-1} \right) \frac{\partial K_{|\boldsymbol{\theta}}}{\partial a_{q,(d,r)}} \right) - \frac{a_{q,(d,r)}}{\psi_{q,(d,r)}}, \\
\text{where } \frac{\partial K_{|\boldsymbol{\theta}}}{\partial a_{q,(d,r)}} &= B'_q \otimes k_q(\mathbf{x}, \mathbf{x}'), \\
B'_{q,(i,j)} &= \begin{cases} 2a_{q,(d,r)} & , \text{ for } i = j = d, \\ a_{q,(j,r)} & , \text{ for } i = d, j \neq d, \\ a_{q,(i,r)} & , \text{ for } i \neq d, j = d, \\ 0 & , \text{ otherwise.} \end{cases}
\end{aligned} \tag{37}$$

For partial gradients used for optimization:

$$\begin{aligned}
\frac{\partial \mathcal{Q}(\boldsymbol{\theta})}{\partial \lambda_{q,(d)}} &= \frac{1}{2} \text{tr} \left(\left(\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - K_{|\boldsymbol{\theta}}^{-1} \right) \frac{\partial K_{|\boldsymbol{\theta}}}{\partial \lambda_{q,(d)}} \right) - \frac{\text{sign}(\lambda_{q,(d)})}{\beta_\lambda}, \\
\text{where } \frac{\partial K_{|\boldsymbol{\theta}}}{\partial \lambda_{q,(d)}} &= \text{diag}(\boldsymbol{\lambda}'_q) \otimes k_q(\mathbf{x}, \mathbf{x}'), \\
\lambda'_{q,(i)} &= \begin{cases} 1 & , \text{ for } i = d, \\ 0 & , \text{ otherwise.} \end{cases}
\end{aligned} \tag{38}$$

$$\begin{aligned}
\frac{\partial \mathcal{Q}(\boldsymbol{\theta})}{\partial v_q} &= \frac{1}{2} \text{tr} \left(\left(\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - K_{|\boldsymbol{\theta}}^{-1} \right) \frac{\partial K_{|\boldsymbol{\theta}}}{\partial v_q} \right), \\
\text{where } \frac{\partial K_{|\boldsymbol{\theta}}}{\partial v_q} &= B_q \otimes k_{qv}(\mathbf{x}, \mathbf{x}'),
\end{aligned} \tag{39}$$

$$\begin{aligned}
k_{qv}(\mathbf{x}, \mathbf{x}') &= -2\pi^2 \tau^2 \exp(-2\pi^2 \tau^2 v_q) \cos(2\pi \tau \mu_q), \\
\frac{\partial \mathcal{Q}(\boldsymbol{\theta})}{\partial \mu_q} &= \frac{1}{2} \text{tr} \left(\left(\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - K_{|\boldsymbol{\theta}}^{-1} \right) \frac{\partial K_{|\boldsymbol{\theta}}}{\partial \mu_q} \right), \\
\text{where } \frac{\partial K_{|\boldsymbol{\theta}}}{\partial \mu_q} &= B_q \otimes k_{q\mu}(\mathbf{x}, \mathbf{x}'),
\end{aligned} \tag{40}$$

$$k_{q\mu}(\mathbf{x}, \mathbf{x}') = -2\pi \tau \exp(-2\pi^2 \tau^2 v_q) \sin(2\pi \tau \mu_q).$$

Appendix C. Detailed Results of Imputation Error and 95% Coverage

We organized the detailed results of online imputation on all 24 covariates under the best number of basis kernel ($Q = 5$ for HUP subsets and $Q = 4$ for the MIMIC-III subset) in Figure 14 to Figure 17. For Figure 14 and Figure 15, the mean absolute errors (MAEs) for each covariate is shown (in the original unit of measure). In Figure 16 and Figure 17, we showed the percentage for the prediction lied within the 95% confidence region (i.e. 95% coverage). We put markers in the figures to indicate the best among all methods, and the comparison of MedGP (sparse SM-LMC with online updating) against other methods. The statistical significance were tested using paired t-tests on patient-level results.

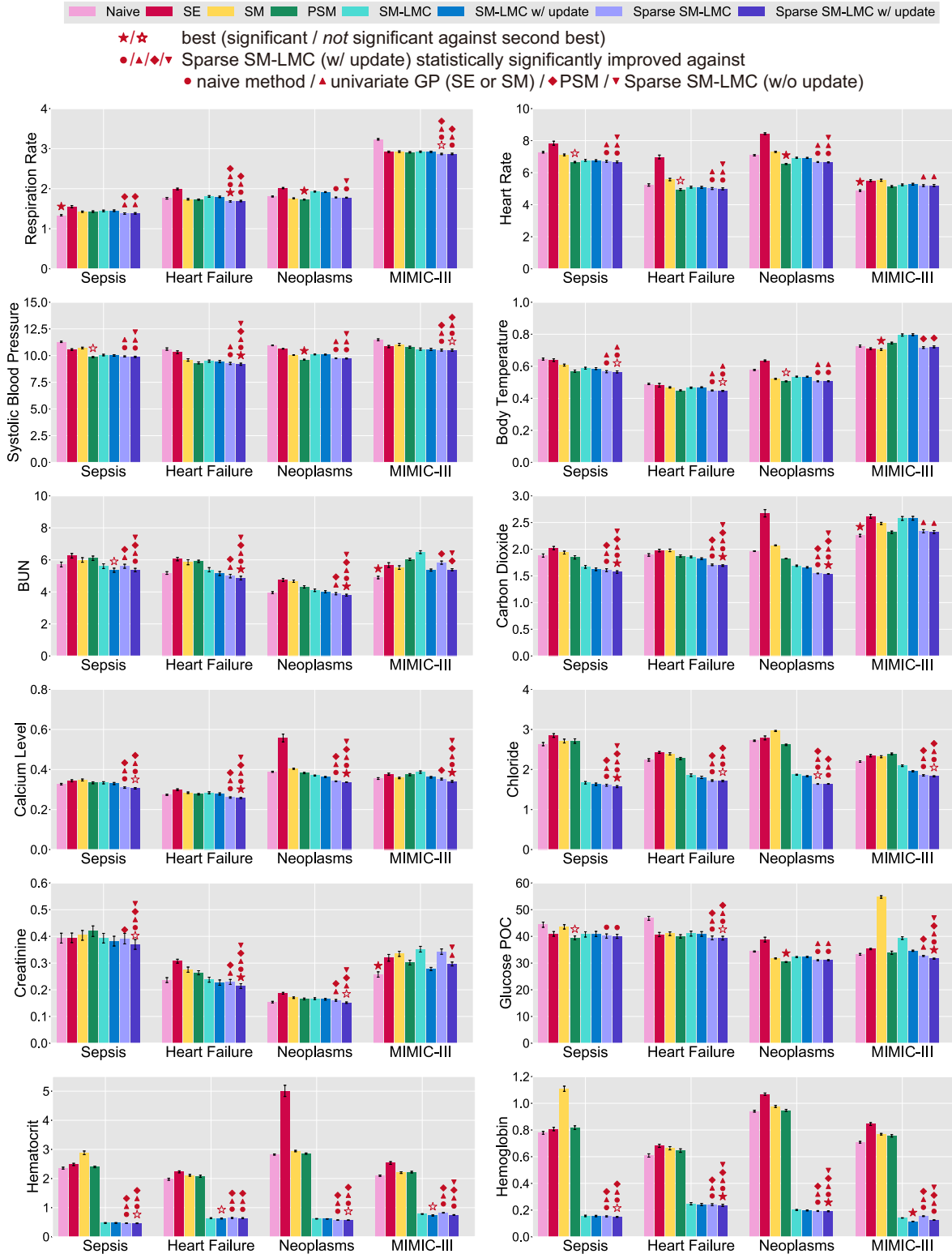


Figure 14: Mean absolute error (MAE) for 12 out of 24 covariates tested. The error bars denote ± 1 standard error.

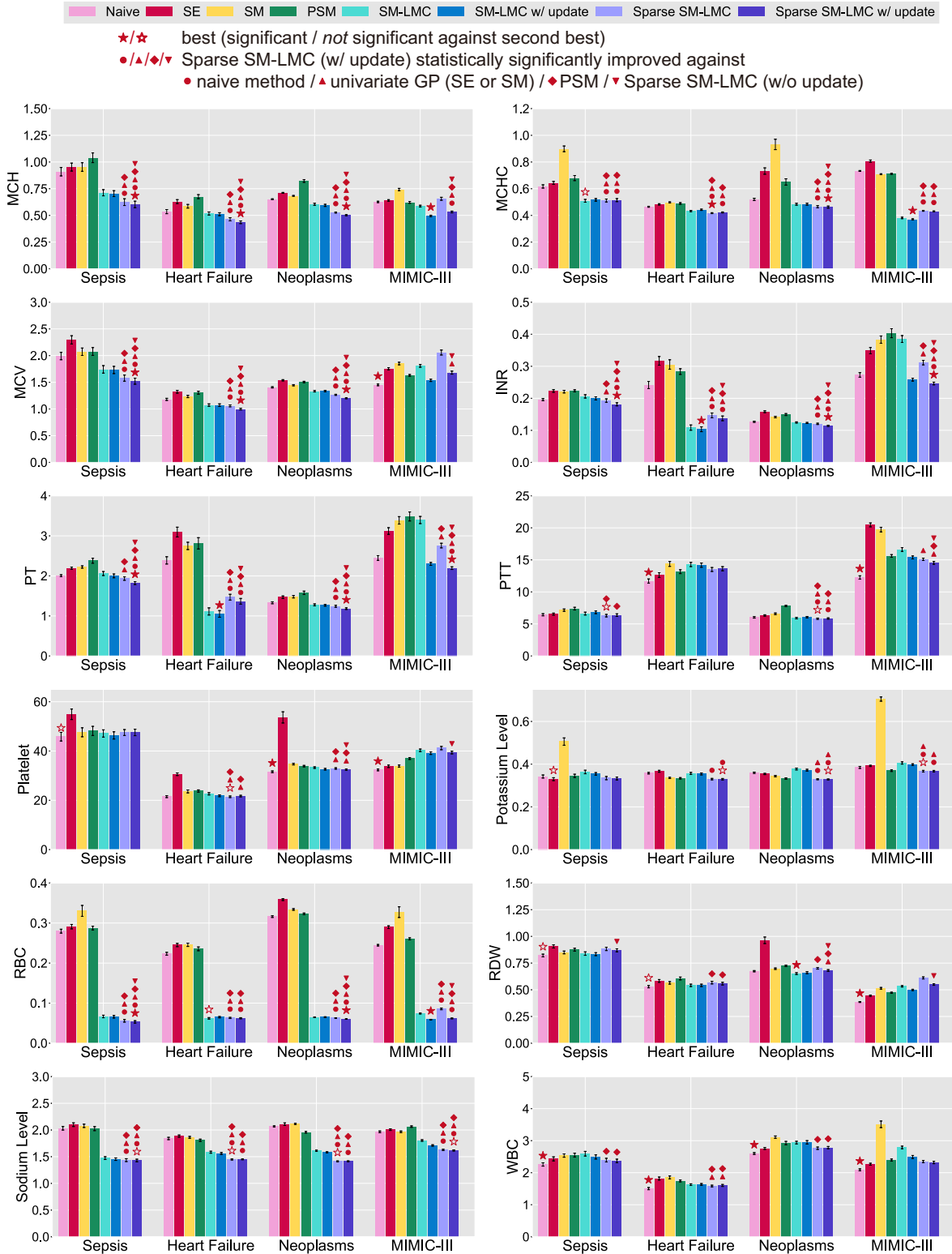


Figure 15: Mean absolute error (MAE) for 12 out of 24 covariates tested. The error bars denote ± 1 standard error.

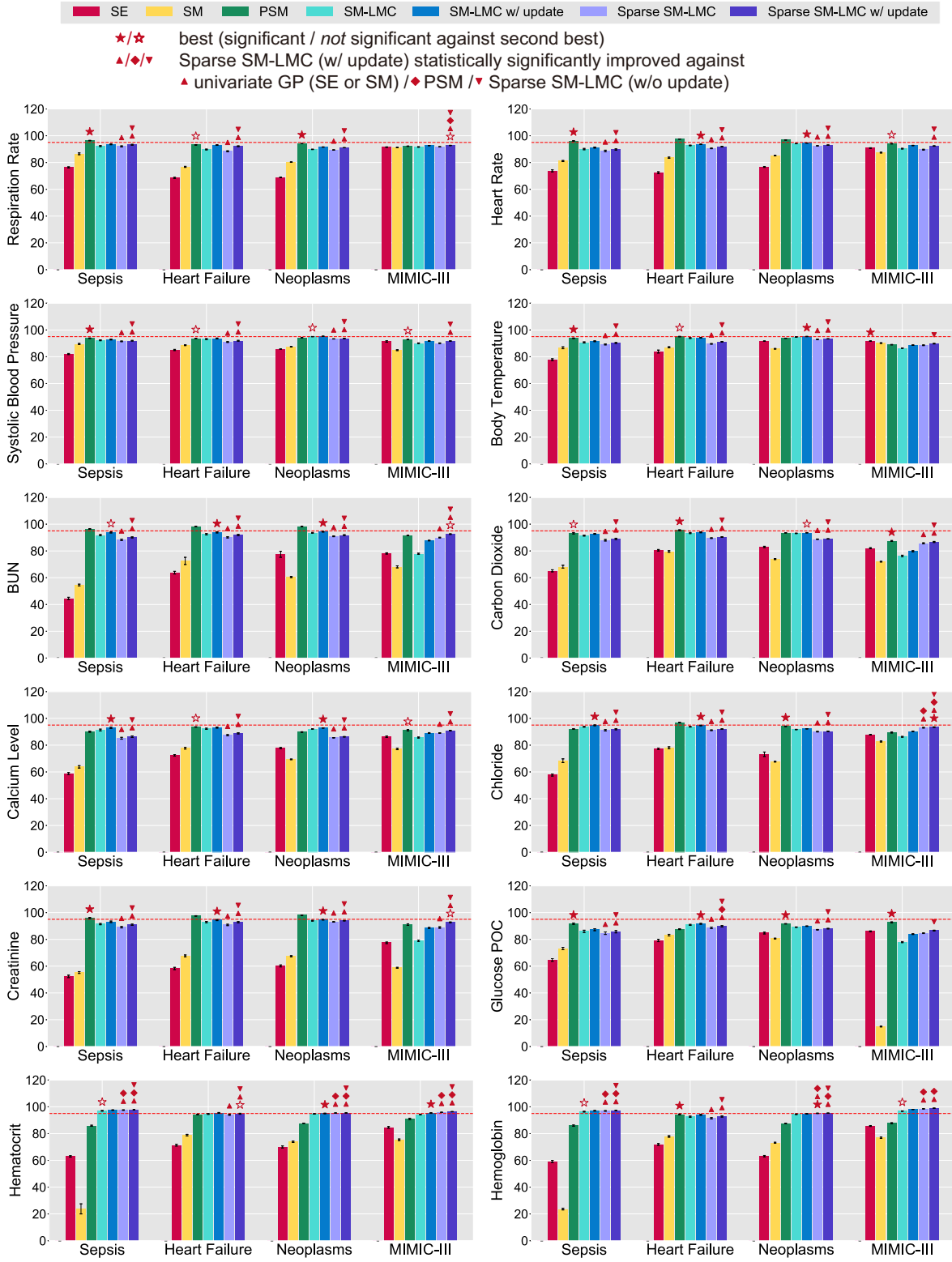


Figure 16: **The 95% coverage for 12 out of 24 covariates tested.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

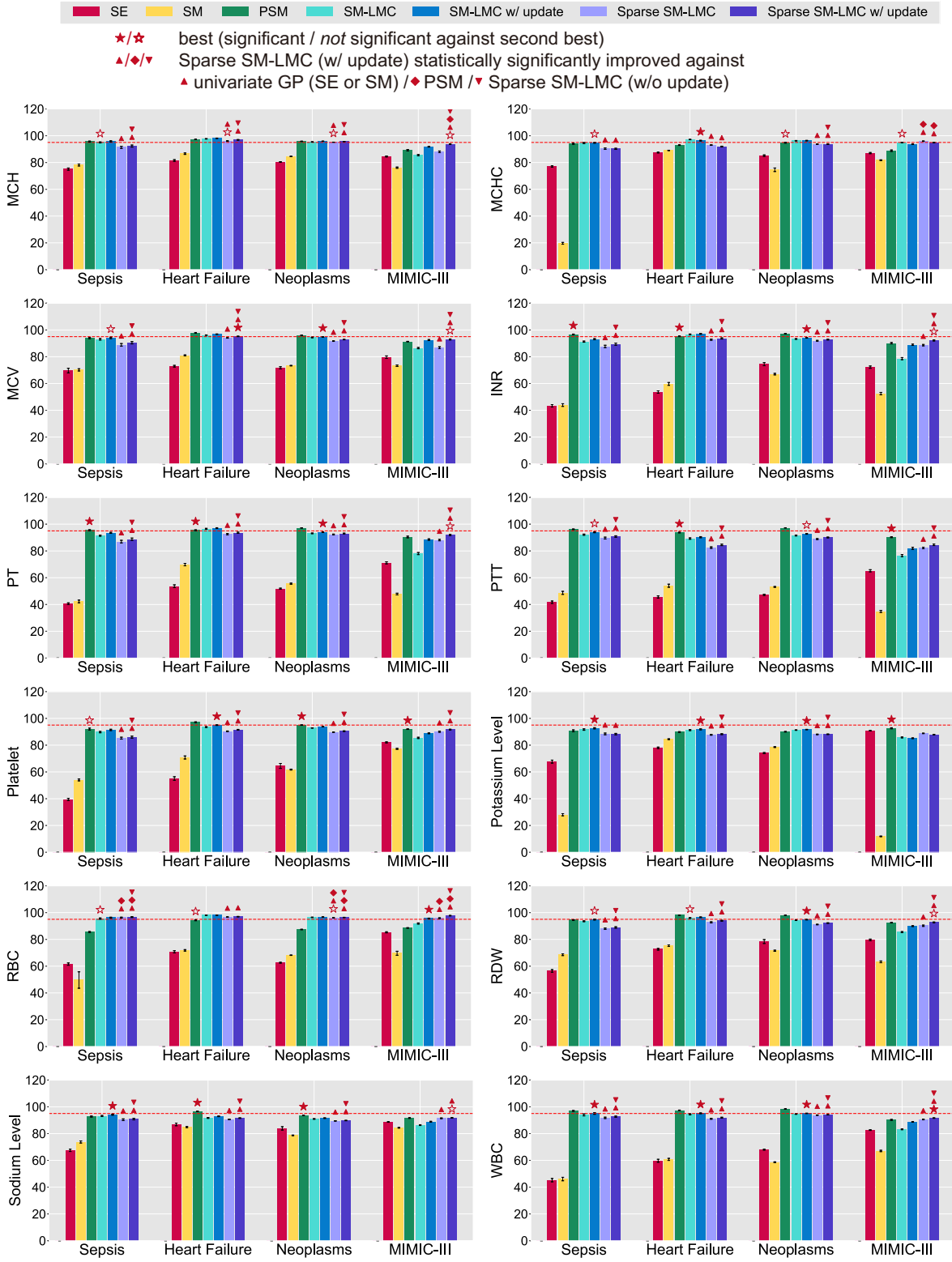


Figure 17: **The 95% coverage for 12 out of 24 covariates tested.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

Appendix D. Results under Different Number of Basis Kernels

In this appendix, we showed more detailed results of the experiments using different number of basis kernels. We ran experiments with for $Q = 1, \dots, 5$ on all four subsets. The results include all three subgroups in the HUP data set and the MIMIC-III heart failure subset. We visualized the results in Figure 18–33. We noticed that for most of the covariates, the imputation performance (both MAE and 95% coverage) improves as the number of Q increases. We also observed that the best number of Q varies across covariates under different metrics. For instance, for lab covariates INR and PT, we observed that setting $Q = 1$ or $Q = 2$ reduces MAE compared with $Q = 5$, but the coverage still improves after $Q = 2$. Allowing more numbers of basis kernels increases the flexibility for customization, but also increases complexity and thus the risk of overfitting for some covariates or patients. Overall $Q = 5$ for HUP subsets and $Q = 4$ for the MIMIC-III subset reached the largest number of covariates improved over the best of baselines using imputation error as the performance metric. How to improve the performance for a specific clinical covariate at patient-level would be one future direction of interest.

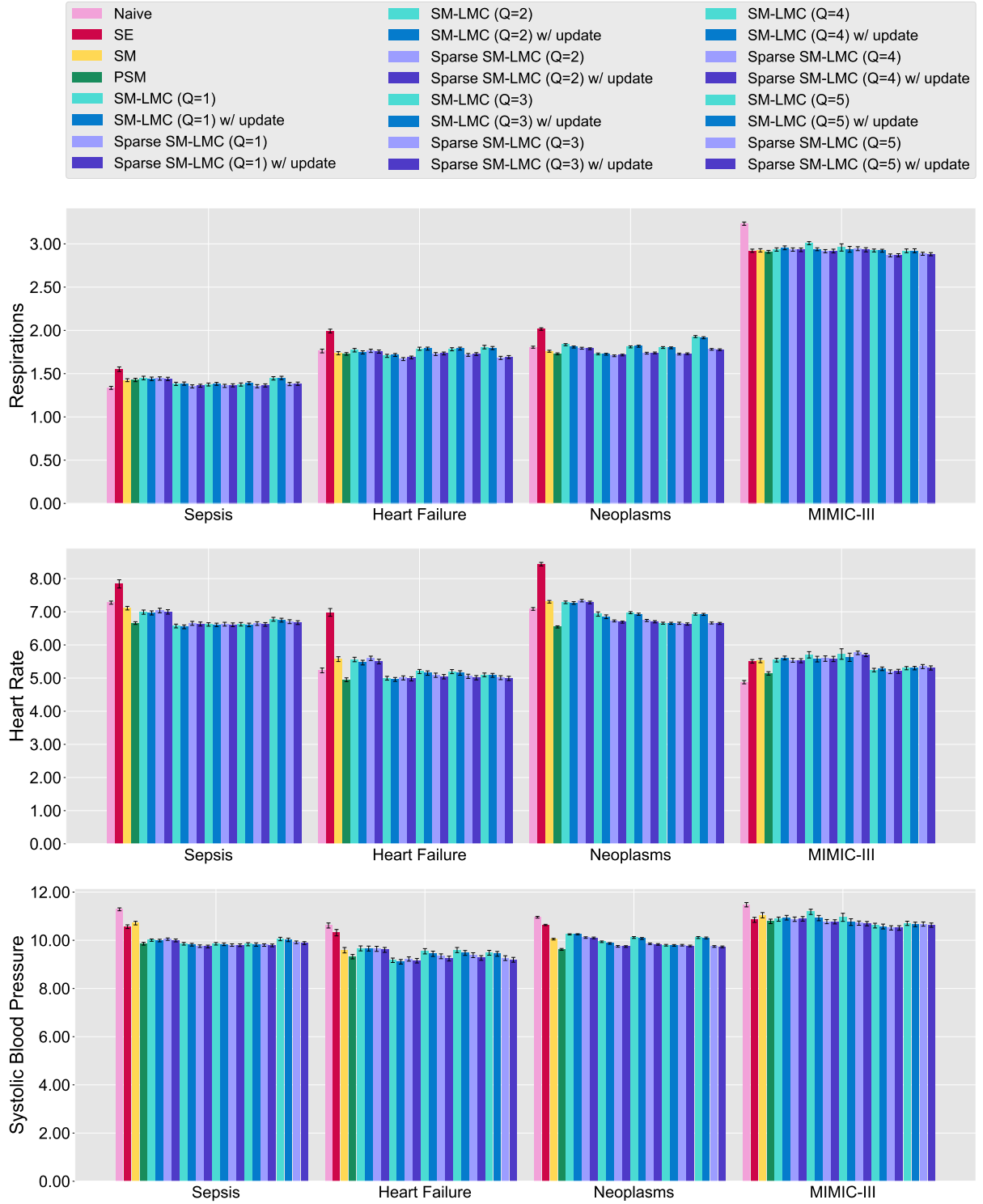


Figure 18: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

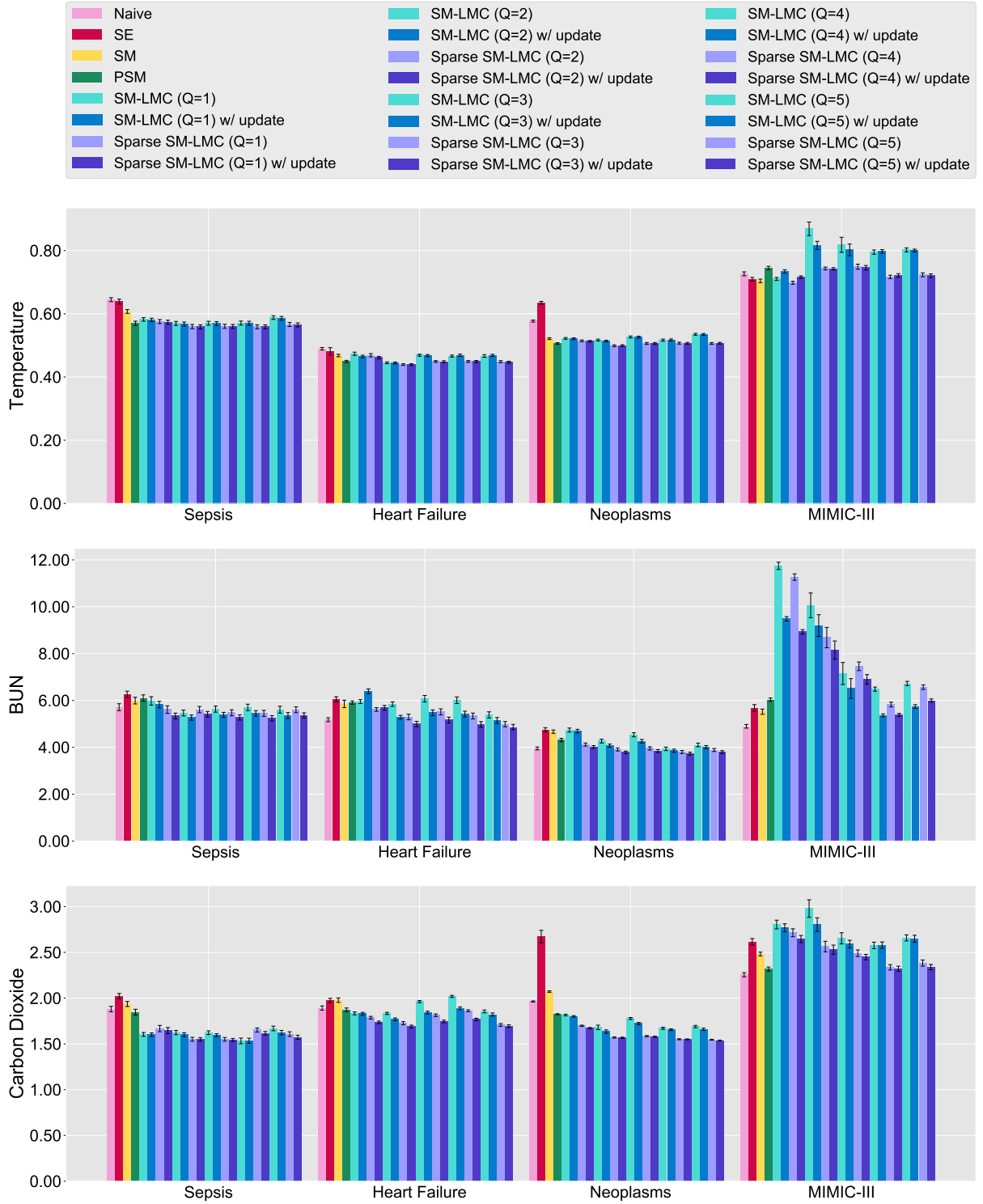


Figure 19: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

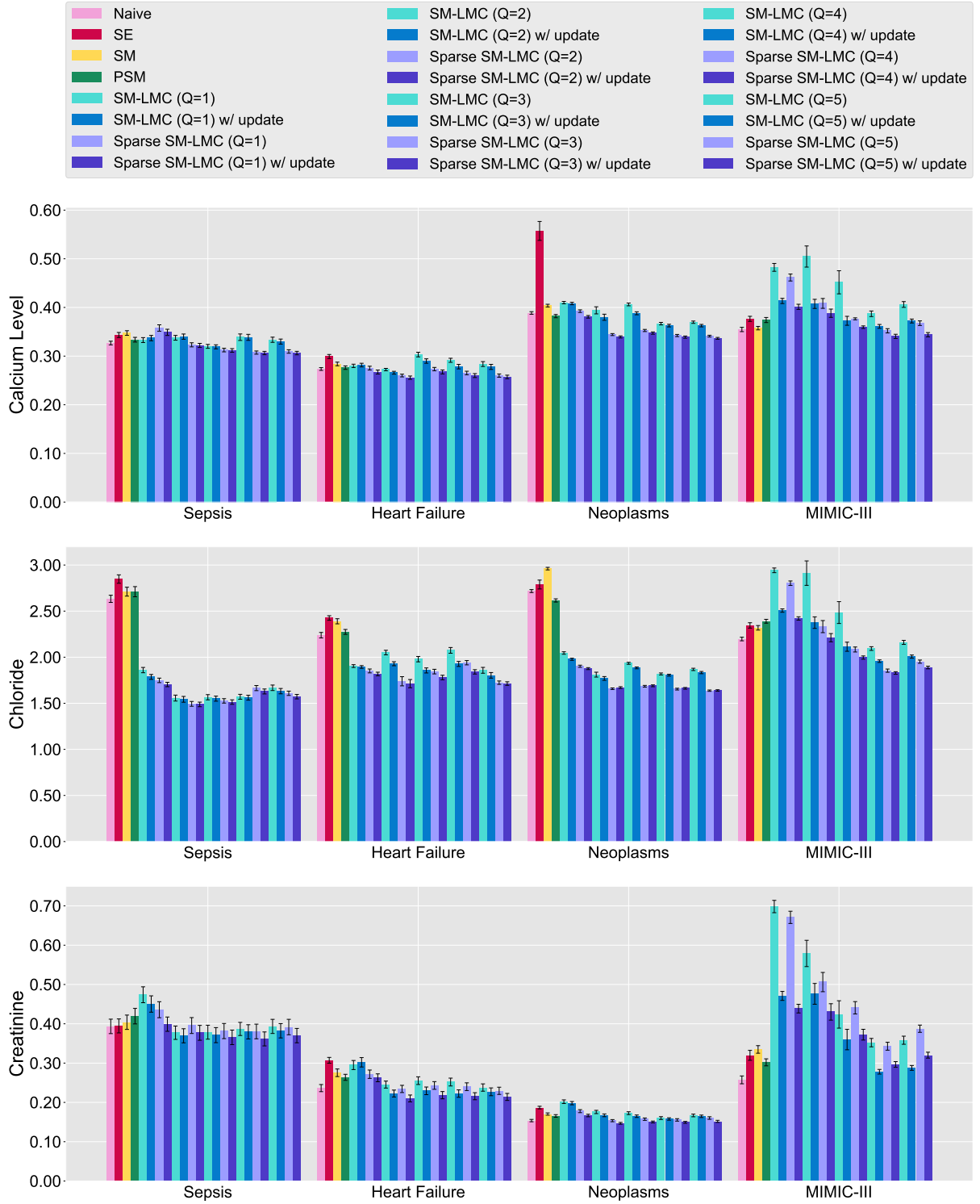


Figure 20: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

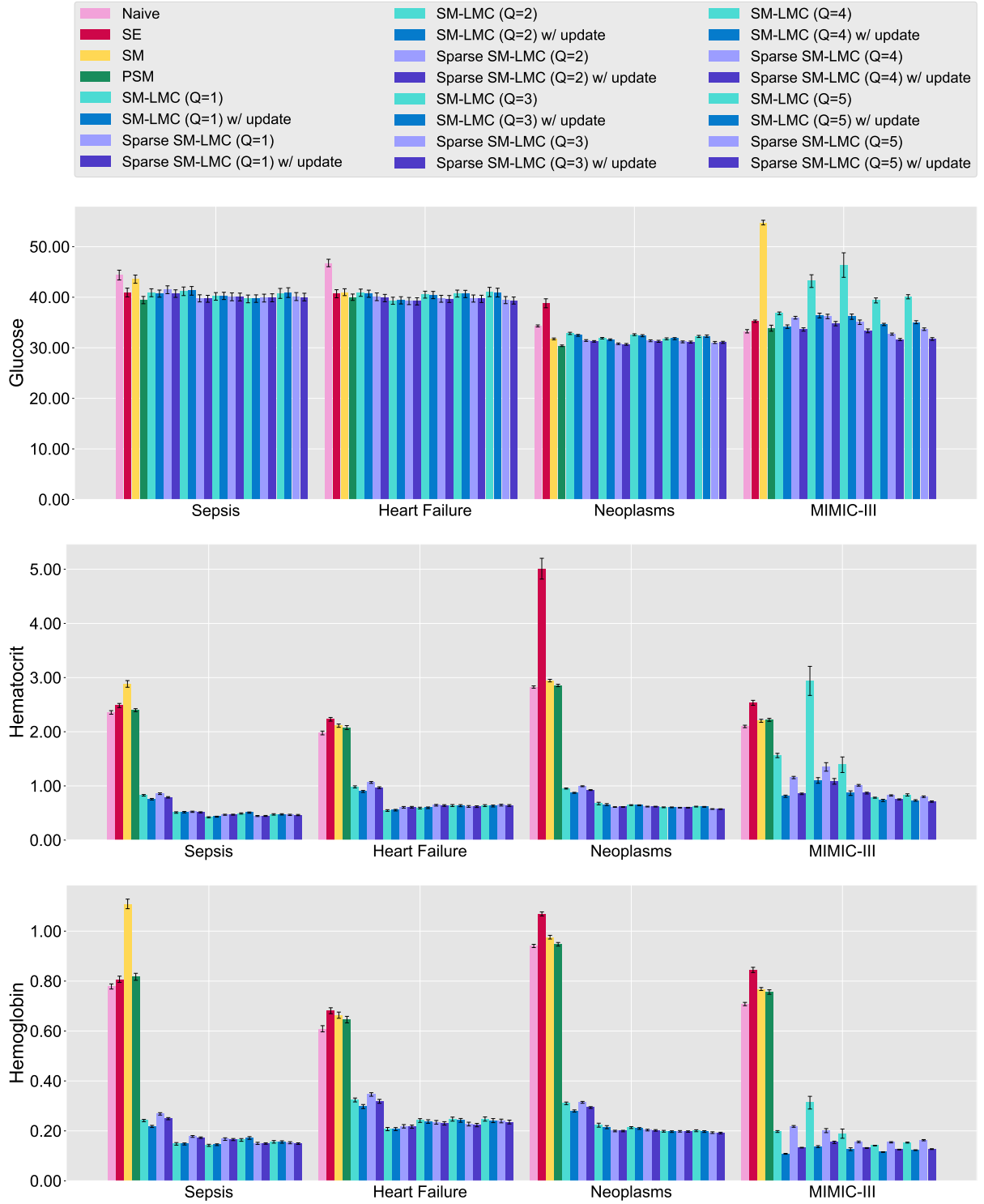


Figure 21: **The mean absolute error (MAE) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error.

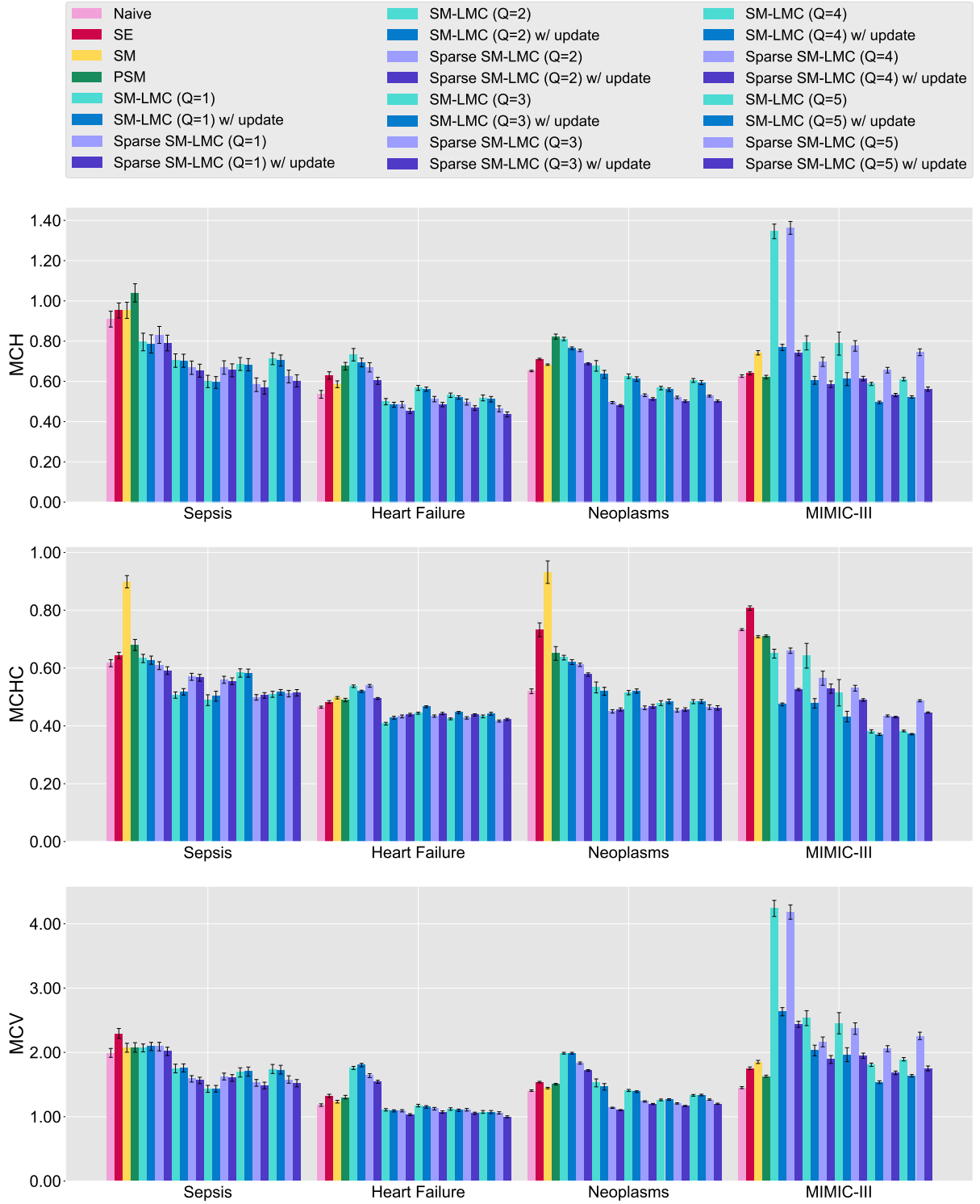


Figure 22: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

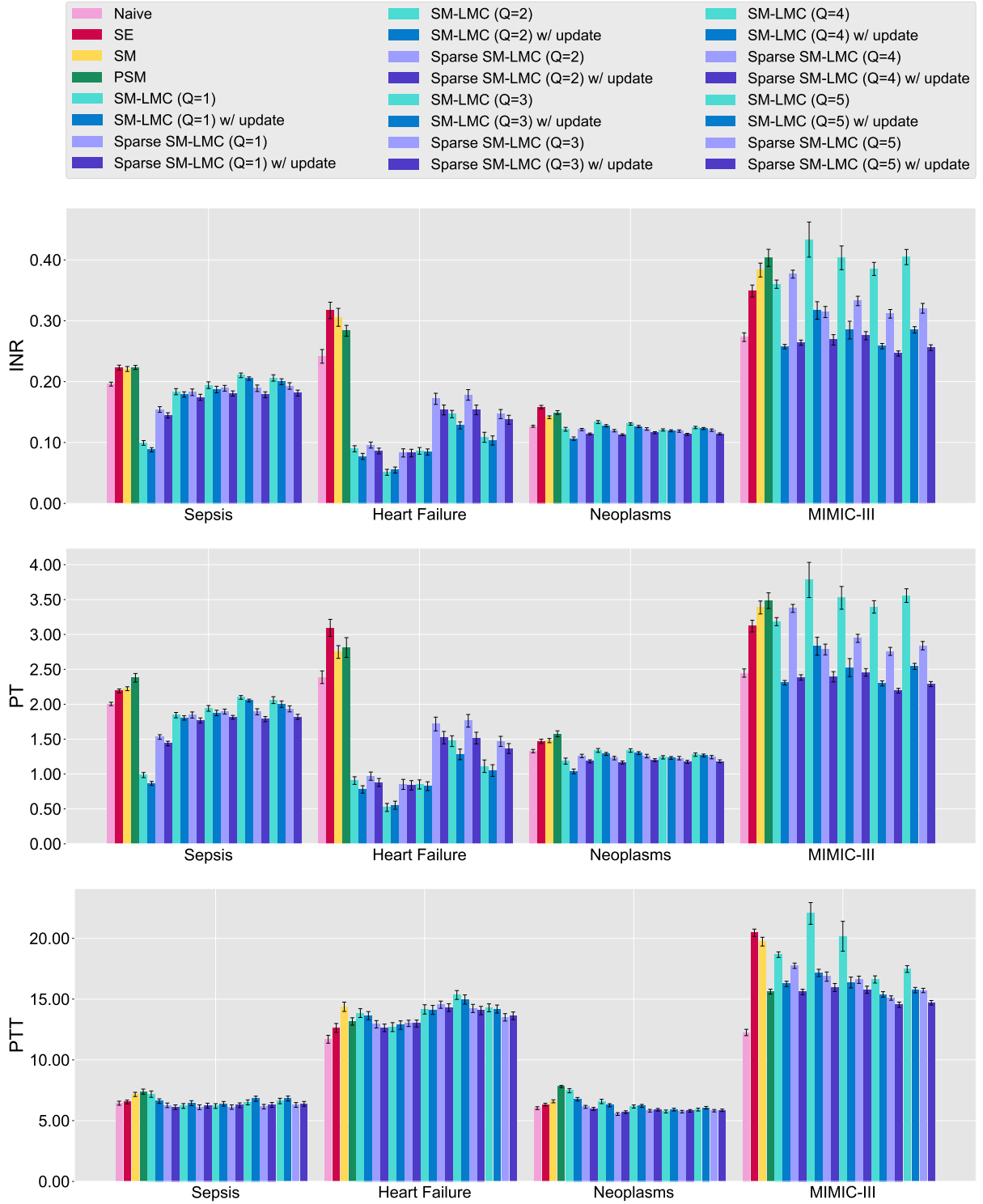


Figure 23: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

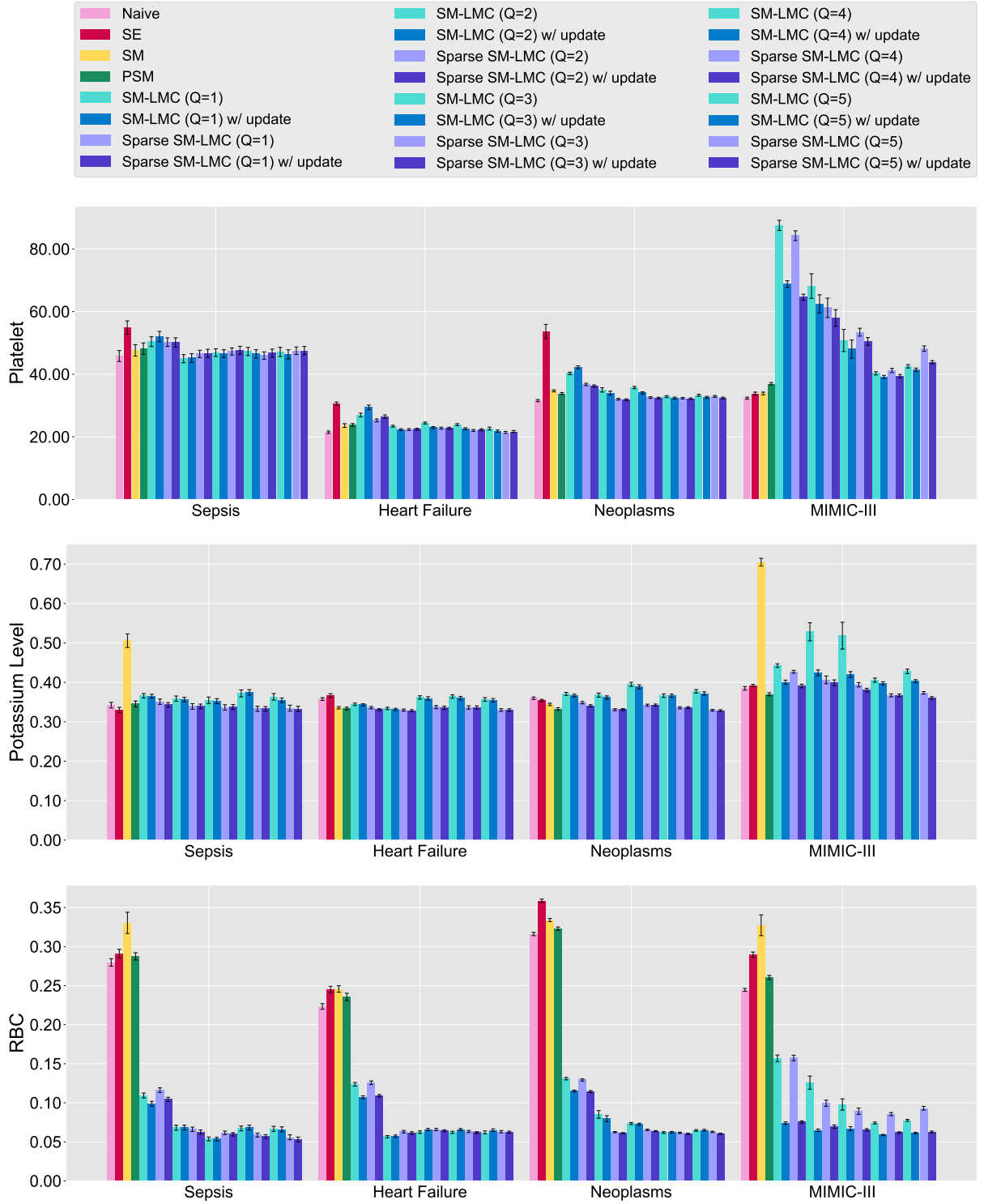


Figure 24: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

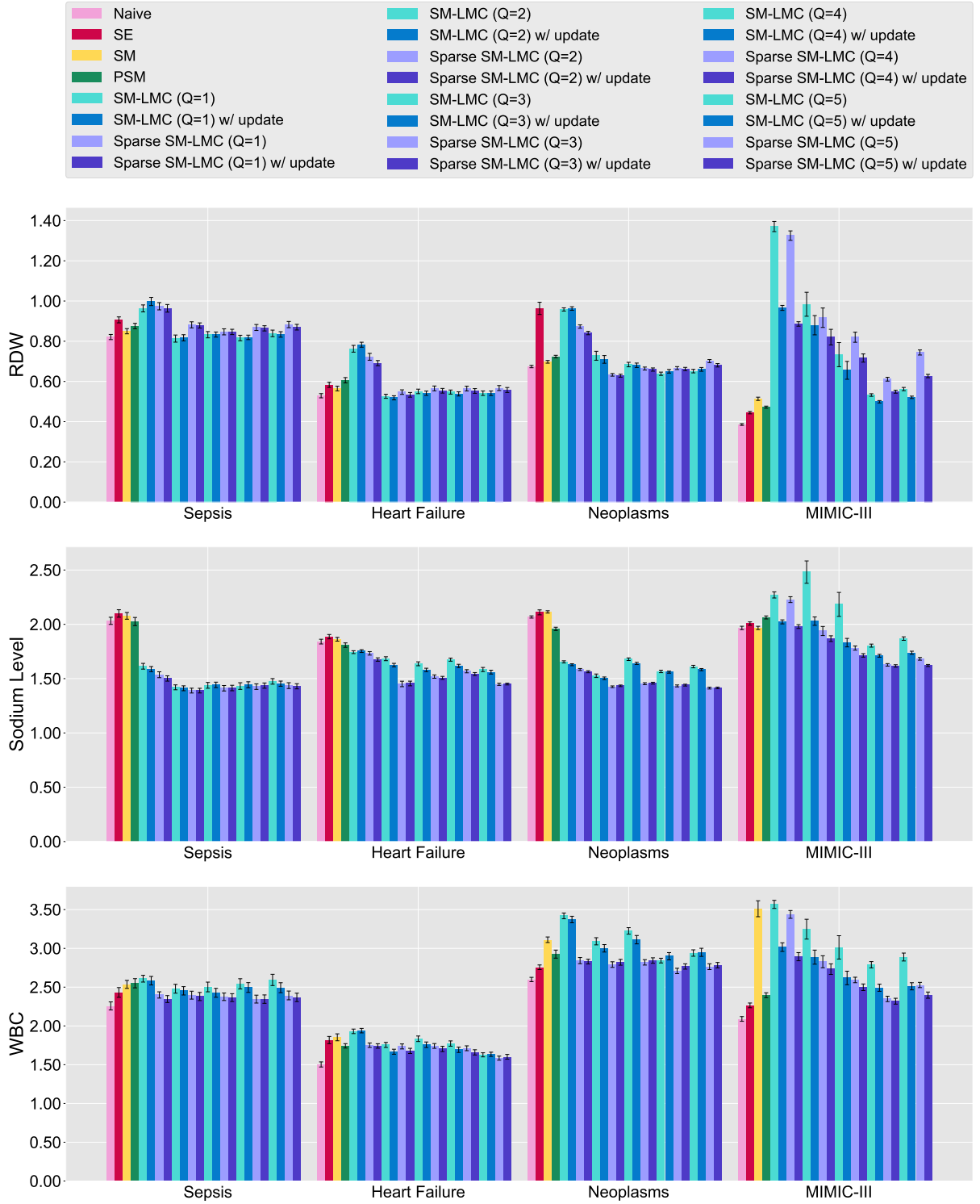


Figure 25: The mean absolute error (MAE) of online imputation under different Q for all cohorts. The error bars denote ± 1 standard error.

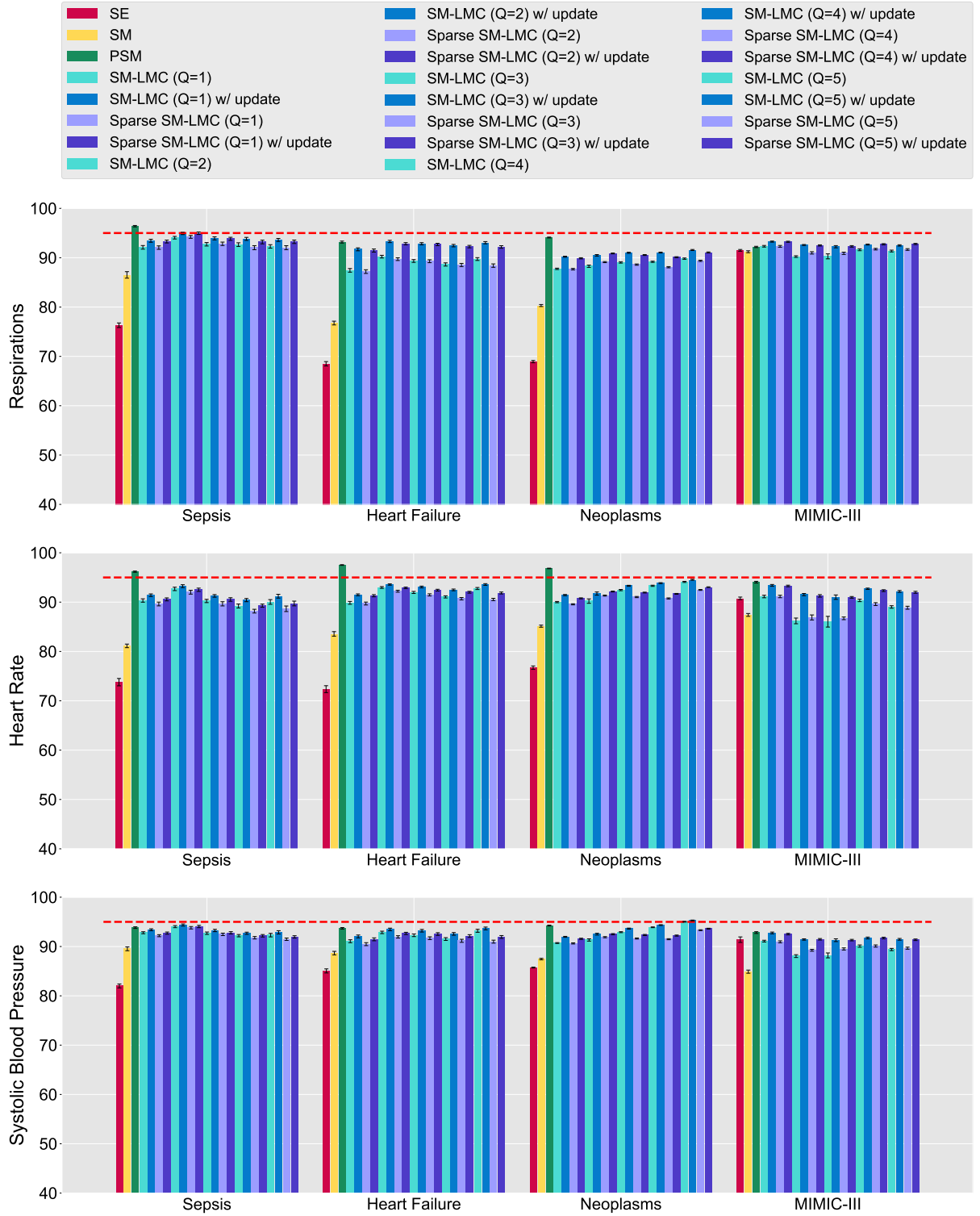


Figure 26: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

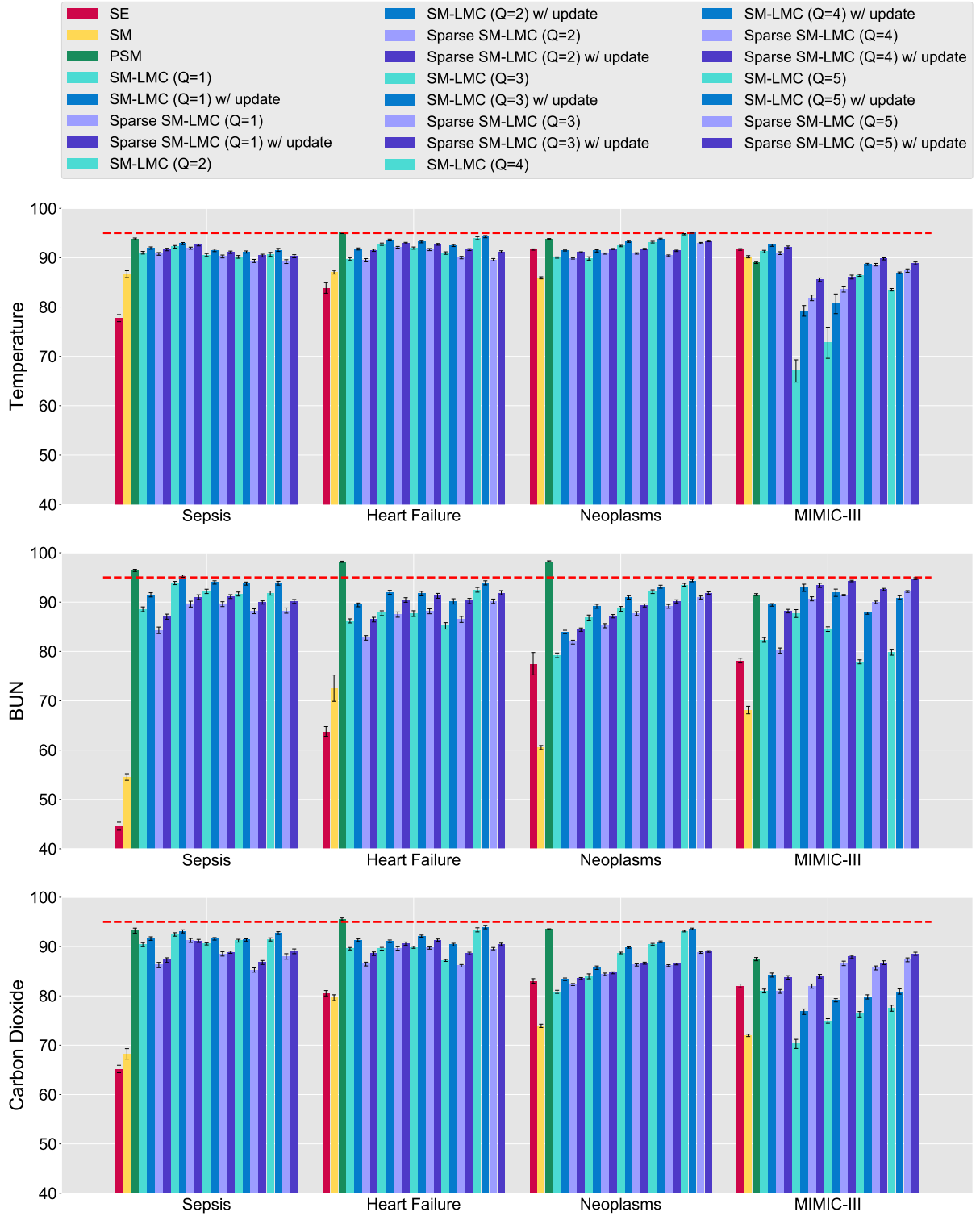


Figure 27: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

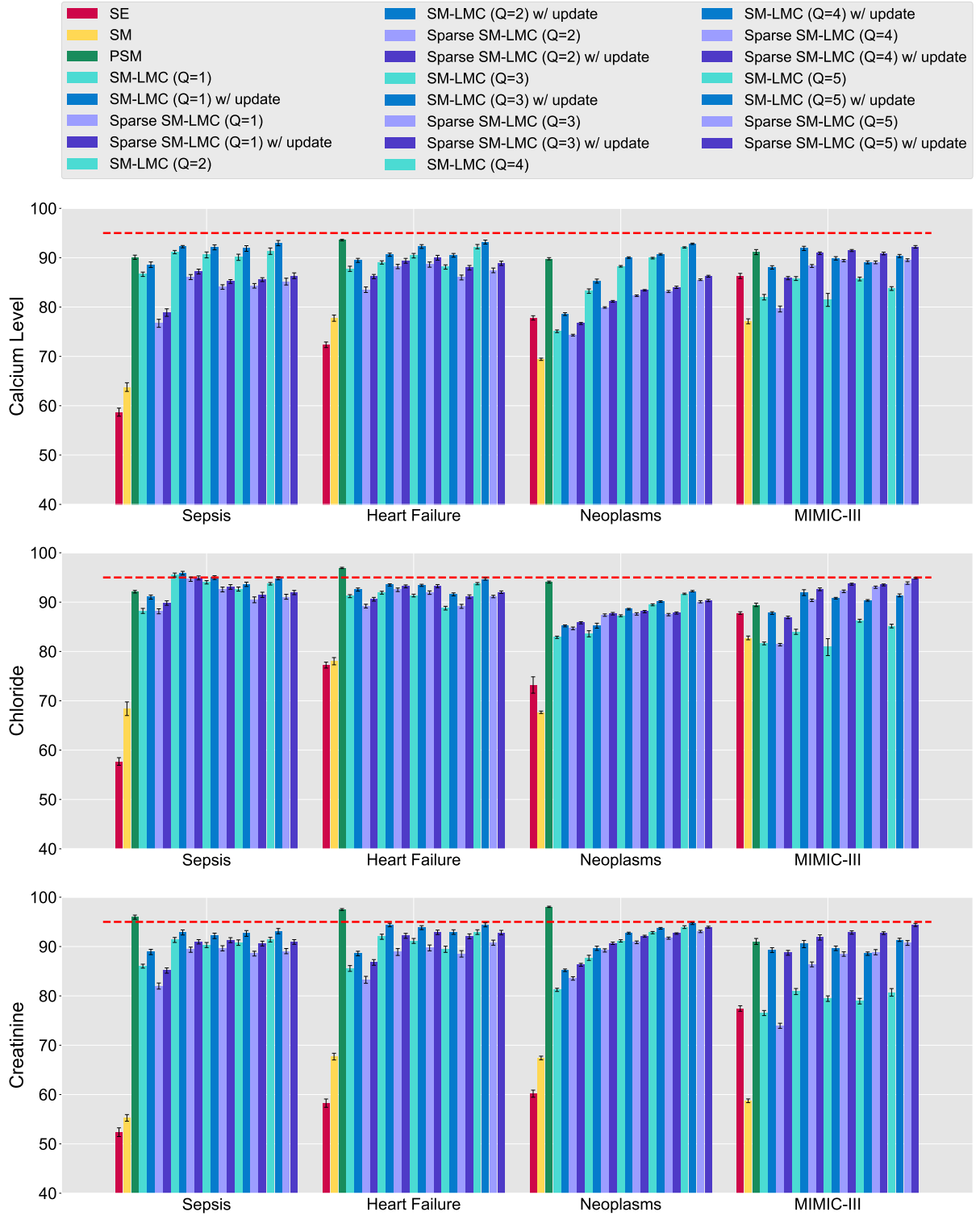


Figure 28: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

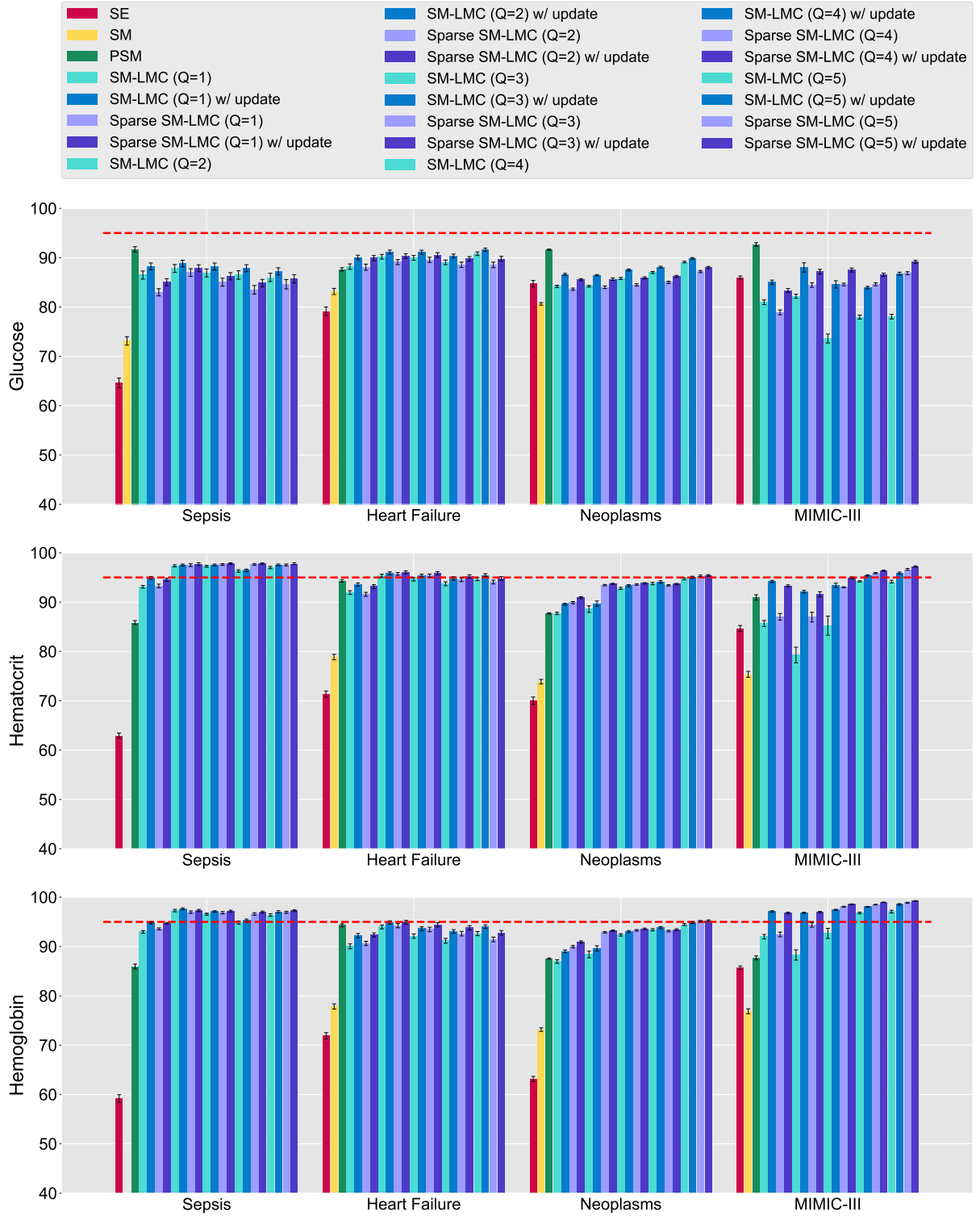


Figure 29: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

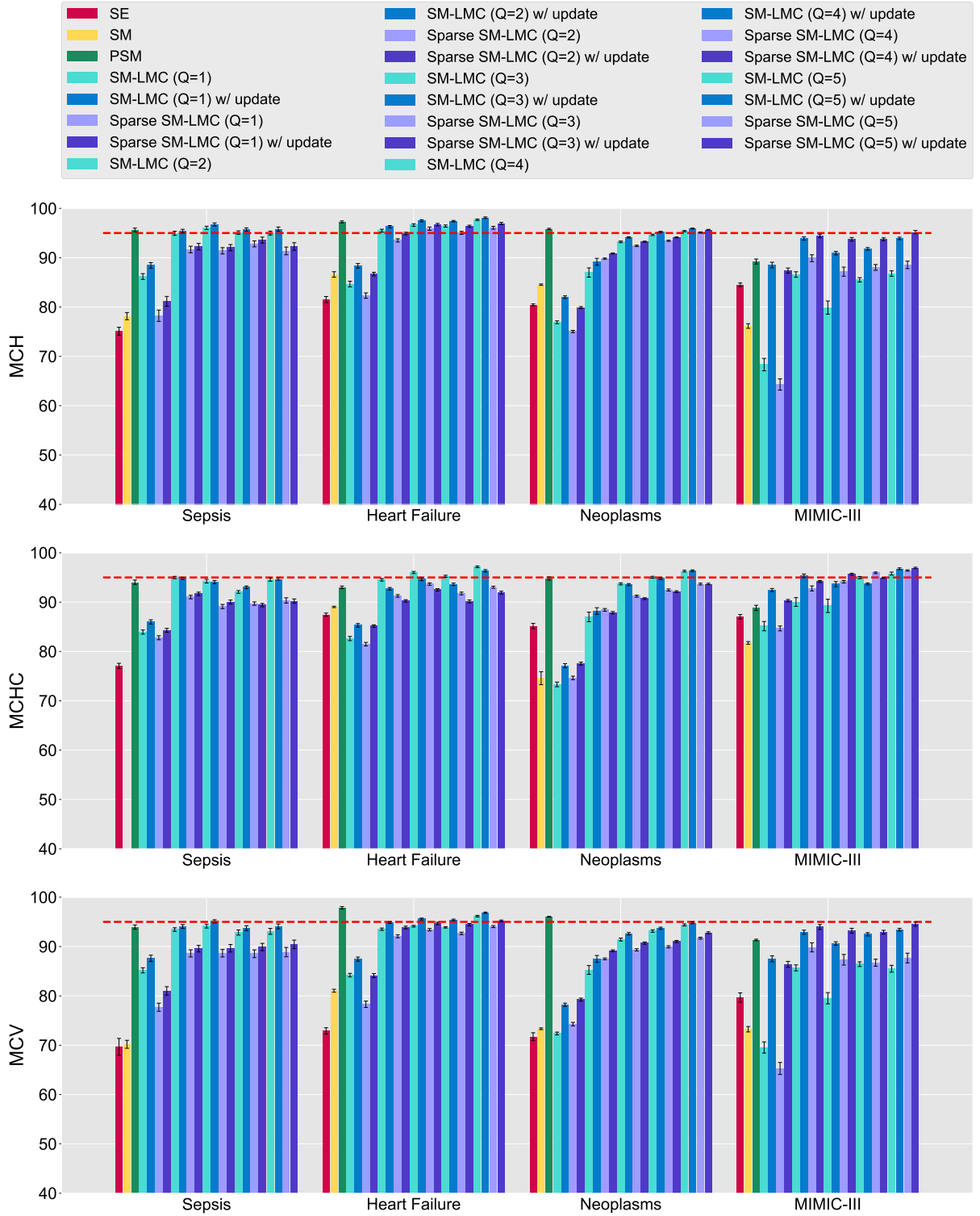


Figure 30: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

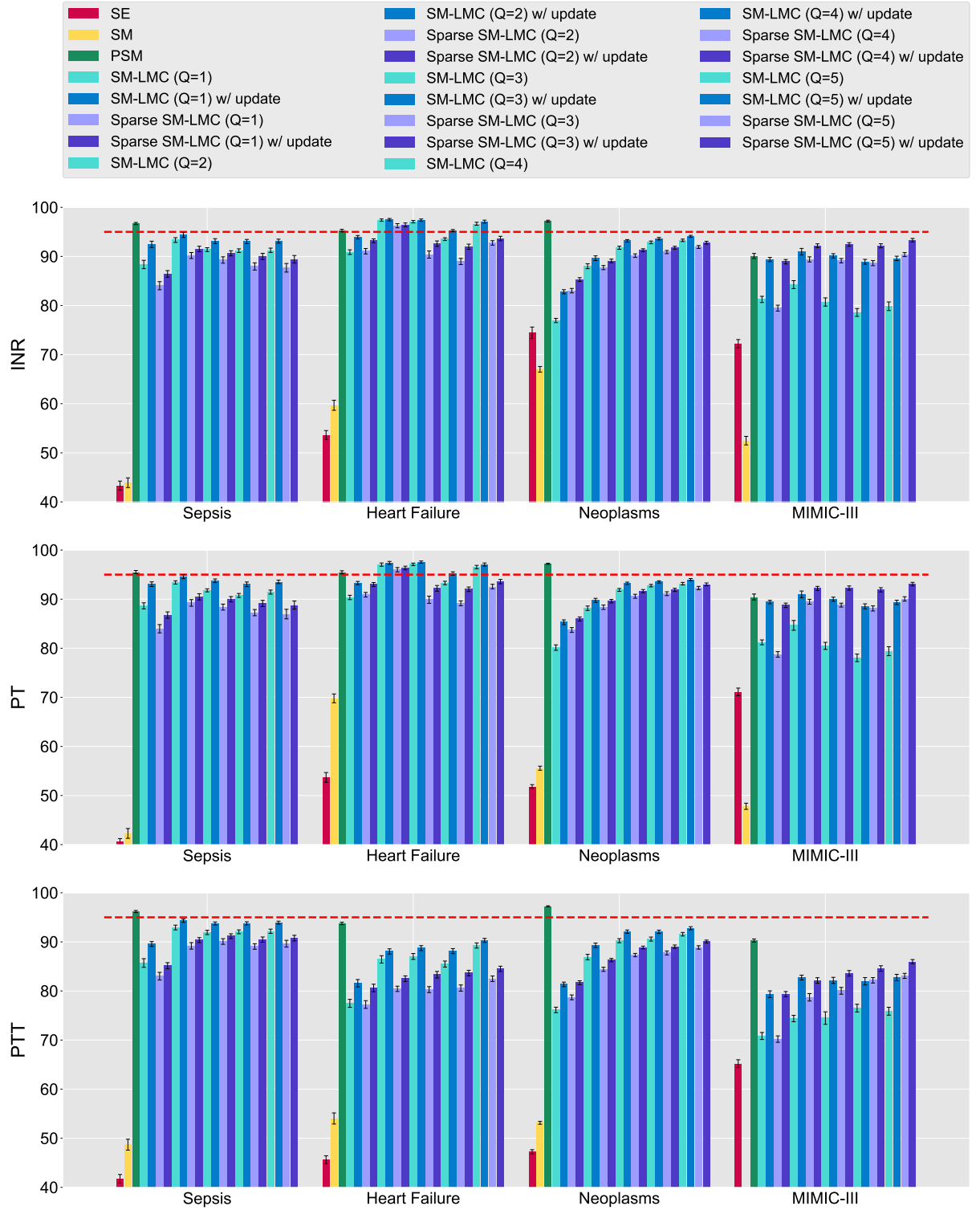


Figure 31: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

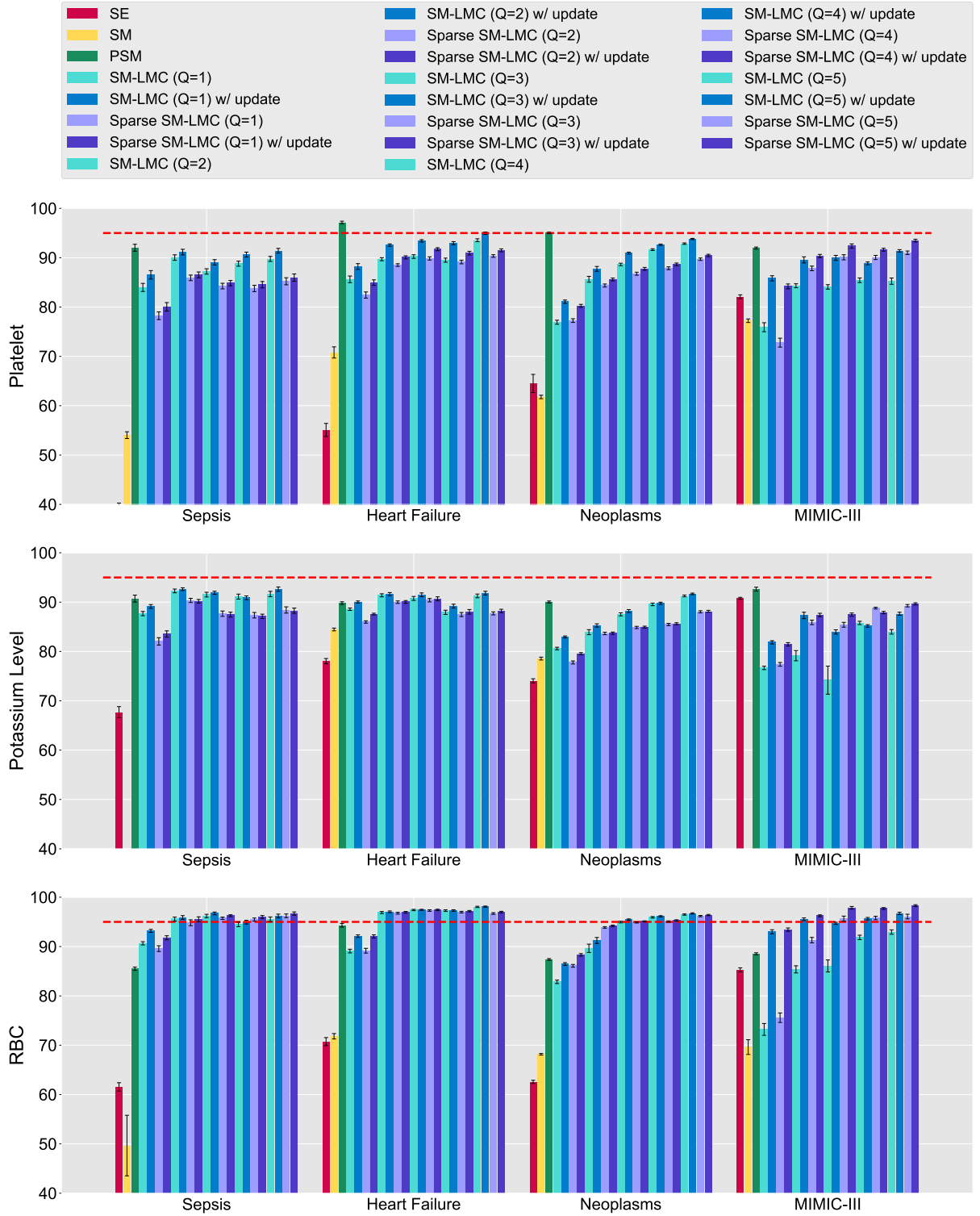


Figure 32: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

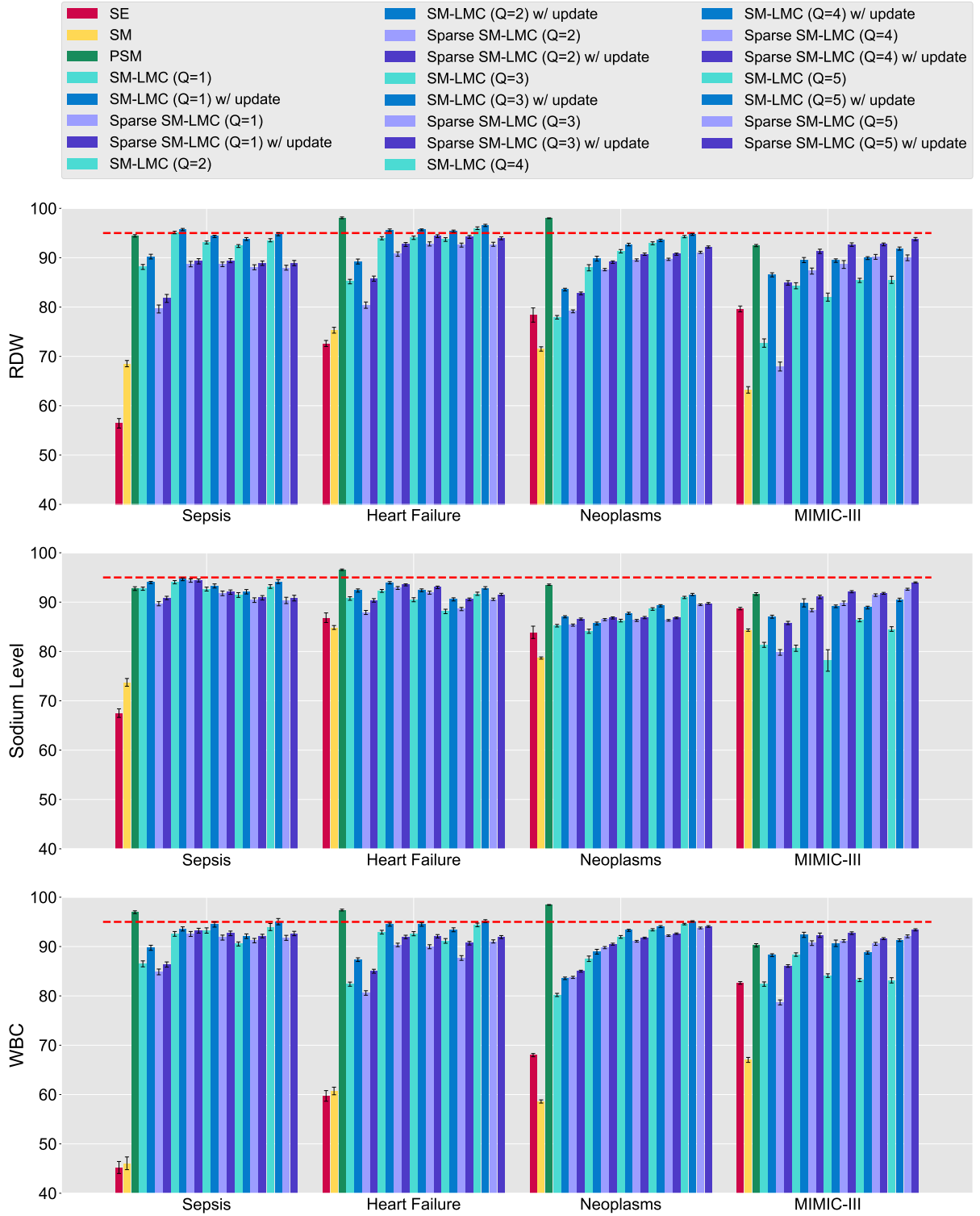


Figure 33: **The 95% coverage (in percentage) of online imputation under different Q for all cohorts.** The error bars denote ± 1 standard error. The red dashed line indicates 95%.

Appendix E. Improvements in Empirical Runtime

In this appendix, we provide the comparisons in runtime with GPy (GPy, since 2012), a state-of-the-art optimized Python library for GPs. We selected few benchmark cases from the MIMIC-III subset, and profiled the runtime for performing one iteration when using gradient-based optimizers. That is, the runtime for computing the gram matrix, log marginal likelihood, and gradients of all parameters. The experiments were performed on the machine with 20 Intel(R) Xeon(R) CPUs running at 2.50GHz (no GPUs were used). For GPy implementation, we also allowed multithreading and the access to MKL optimization for matrix operations, provided by Anaconda with academic license. In Figure 34, we show the average runtime for a single iteration under different number of basis kernels: $Q = 1$ and $Q = 5$, corresponding to 242 and 1114 parameters ($D = 24, R = 8$). We found that for training cases smaller than 10^4 observations, GPy with multithreading is comparable to our implementation. However, for the cases larger than 10^4 observations, our implementation speeds up by up to 2.5 times. The largest case we tested here includes 29,525 observations.

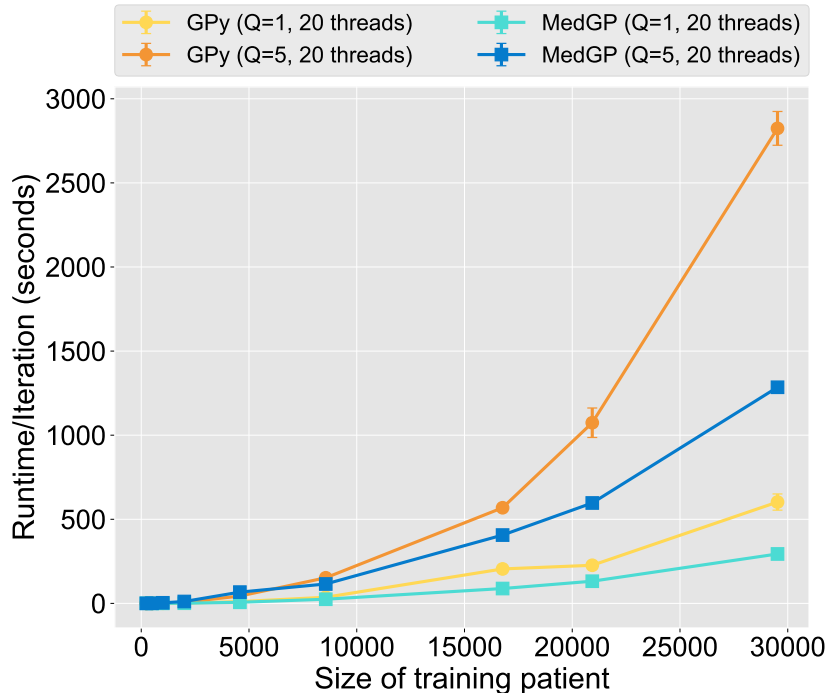


Figure 34: **The empirical runtime of our implementation.** A comparison of the average runtimes for one iteration (including computation of gradients) for MedGP and optimized baseline GPy.

References

- Ryan P. Adams and David J. C. MacKay. Bayesian online changepoint detection. *arXiv preprint arXiv:0710.3742*, 2007.
- Mauricio A. Álvarez and Neil D. Lawrence. Computationally efficient convolved multiple output Gaussian processes. *Journal of Machine Learning Research*, 12:1459–1500, 2011.
- Mauricio A. Álvarez, Lorenzo Rosasco, and Neil D. Lawrence. Kernels for vector-valued functions: a review. *Foundations and Trends in Machine Learning*, 4(3):195–266, 2012.
- Derek C. Angus, Walter T. Linde-Zwirble, Jeffrey Lidicker, Gilles Clermont, Joseph Carcillo, and Michael R. Pinsky. Epidemiology of severe sepsis in the United States: analysis of incidence, outcome, and associated costs of care. *Critical care medicine*, 29(7):1303–1310, 2001.
- Artin Armagan, Merlise Clyde, and David B. Dunson. Generalized beta mixtures of Gaussians. In *Advances in Neural Information Processing Systems 24*, pages 523–531. 2011.
- Edwin V. Bonilla, Kian Ming Chai, and Christopher K. I. Williams. Multi-task Gaussian process prediction. In *Advances in Neural Information Processing Systems 20*, pages 153–160. 2008.
- Carlos M. Carvalho, Nicholas G. Polson, and James G. Scott. The horseshoe estimator for sparse signals. *Biometrika*, 97(2):465–480, 2010.
- Rainer Dahlhaus. Graphical interaction models for multivariate time series. *Metrika*, 51(2):157–172, 2000.
- Robert Dürichen, Marco A. F. Pimentel, Lei Clifton, Achim Schweikard, and David A. Clifton. Multitask Gaussian processes for multivariate physiological time-series analysis. *IEEE Transactions on Biomedical Engineering*, 62(1):314–322, 2015.
- Vladimir Feinberg, Li-Fang Cheng, Kai Li, and Barbara E Engelhardt. Large linear multi-output gaussian process learning for time series. *arXiv preprint arXiv:1705.10813*, 2017.
- Chuan Gao, Christopher D. Brown, and Barbara E. Engelhardt. A latent factor model with a mixture of sparse and dense factors to model gene expression data with confounding effects. *arXiv preprint arXiv:1310.4792*, 2013.
- Ursula Gather, Michael Imhoff, and Roland Fried. Graphical models for multivariate time series from intensive care monitoring. *Statistics in medicine*, 21(18):2685–2701, 2002.
- Marzyeh Ghassemi, Leo Anthony Celi, and David J. Stone. State of the art review: the data revolution in critical care. *Critical Care*, 19(1):118, 2015a.
- Marzyeh Ghassemi, Marco A. F. Pimentel, Tristan Naumann, Thomas Brennan, David A. Clifton, Peter Szolovits, and Mengling Feng. A multivariate timeseries modeling approach to severity of illness assessment and forecasting in ICU with sparse, heterogeneous clinical data. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 446–453, 2015b.

- Pierre Goovaerts. *Geostatistics for natural resources evaluation*. Oxford university press, 1997.
- GPy. GPy: A gaussian process framework in python. <http://github.com/SheffieldML/GPy>, since 2012.
- Katharine E. Henry, David N. Hager, Peter J. Pronovost, and Suchi Saria. A targeted real-time early warning score (TREWScore) for septic shock. *Science Translational Medicine*, 7(299):299ra122–299ra122, 2015.
- Joyce C. Ho, Cheng H. Lee, and Joydeep Ghosh. Septic shock prediction for patients with missing data. *ACM Transactions on Management Information Systems*, 5(1):1:1–1:15, 2014.
- Richard S. Hotchkiss and Irene E. Karl. The pathophysiology and treatment of sepsis. *New England Journal of Medicine*, 348(2):138–150, 2003.
- George Hripcsak and David J. Albers. Next-generation phenotyping of electronic health records. *Journal of the American Medical Informatics Association*, 20(1):117–121, 2013.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3, 2016.
- Andre G. Journel and Charles J. Huijbregts. *Mining geostatistics*. Academic Press, 1978.
- JooSeuk Kim, James M. Blum, and Clayton D. Scott. Temporal features and kernel methods for predicting sepsis in postoperative patients. 2010.
- Gagan Kumar, Nilay Kumar, Amit Taneja, Thomas Kaleekal, Sergey Tarima, Emily McGinley, Edgar Jimenez, Anand Mohan, Rumi Ahmed Khan, Jeff Whittle, Elizabeth Jacobs, and Rahul Nanchal. Nationwide trends of severe sepsis in the 21st century (2000-2007). *Chest*, 140(5):1223–1231, 2011.
- Thomas A. Lasko, Joshua C. Denny, and Mia A. Levy. Computational phenotype discovery using unsupervised feature learning over noisy, sparse, and irregular clinical data. *PloS One*, 8(6):1–13, 2013.
- Li-Wei H. Lehman, Ryan P. Adams, Louis Mayaud, George B. Moody, Atul Malhotra, Roger G. Mark, and Shamim Nemati. A physiological time series dynamics-based approach to patient monitoring and outcome prediction. *IEEE Journal of Biomedical and Health Informatics*, 19(3):1068–1076, May 2015.
- Benjamin M. Marlin, David C. Kale, Robinder G. Khemani, and Randall C. Wetzel. Unsupervised pattern discovery in electronic health care data using probabilistic clustering models. In *Proceedings of the Second ACM SIGHIT International Health Informatics Symposium*, pages 389–398, 2012.
- Travis B. Murdoch and Allan S. Detsky. The inevitable application of big data to health care. *JAMA*, 309(13):1351–1352, 2013.

- Shamim Nemati, Li-Wei H. Lehman, Ryan P. Adams, and Ahana Malhotra. Discovering shared cardiovascular dynamics within a patient cohort. In *Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 6526–6529, 2012.
- Craig D. Newgard and Roger J. Lewis. Missing data: How to best account for what is not known. *JAMA*, 314(9):940–941, 2015.
- Trung V. Nguyen and Edwin V. Bonilla. Collaborative multi-output Gaussian processes. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, 2014.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Charalampos Pierrakos and Jean-Louis Vincent. Sepsis biomarkers: a review. *Critical Care*, 14(1):R15, 2010.
- Nicholas G. Polson and James G. Scott. Shrink globally, act locally: Sparse Bayesian regularization and prediction. *Bayesian Statistics*, 9:501–538, 2010.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- Dimitris Rizopoulos and Pulak Ghosh. A Bayesian semiparametric multivariate joint model for multiple longitudinal outcomes and a time-to-event. *Statistics in Medicine*, 30(12):1366–1380, 2011.
- S. Roberts, M. Osborne, M. Ebden, S. Reece, N. Gibson, and S. Aigrain. Gaussian processes for time-series modelling. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 371(1984), 2012.
- David E Rumelhart, Geoffrey E Hinton, Ronald J Williams, et al. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1, 1988.
- Yunus Saatçi, Ryan Turner, and Carl Edward Rasmussen. Gaussian process change point models. In *Proceedings of the 27th International Conference on Machine Learning*, pages 927–934, 2010.
- Peter Schulam and Suchi Saria. A framework for individualizing predictions of disease trajectories by exploiting multi-resolution structure. In *Advances in Neural Information Processing Systems 28*, pages 748–756. 2015.
- Peter Schulam, Fredrick Wigley, and Suchi Saria. Clustering longitudinal clinical marker trajectories from electronic health data: Applications to phenotyping and endotype discovery. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- Bernard W. Silverman. *Density estimation for statistics and data analysis*. Chapman and Hall, 1986.

- Ioan Stanculescu, Christopher K. I. Williams, and Yvonne Freer. A hierarchical switching linear dynamical system applied to the detection of sepsis in neonatal condition monitoring. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, 2014.
- Oliver Stegle, Sebastian V. Fallert, David J. C. MacKay, and Søren Brage. Gaussian process robust regression for noisy heart rate data. *IEEE Transactions on Biomedical Engineering*, 55(9):2143–2151, 2008.
- Alex Tank, Nicholas Foti, and Emily Fox. Bayesian structure learning for stationary time series. In *Proceedings of the Thirty-first Conference on Uncertainty in Artificial Intelligence*, 2015.
- Yee Whye Teh, Matthias Seeger, and Michael I. Jordan. Semiparametric latent factor models. In *Proceedings of the International Workshop on Artificial Intelligence and Statistics*, volume 10, 2005.
- Michalis K. Titsias and Miguel Lázaro-Gredilla. Spike and slab variational inference for multi-task and multiple kernel learning. In *Advances in Neural Information Processing Systems 24*, pages 2339–2347. 2011.
- Eric P. Widmaier, Hershel Raff, and Kevin T. Strang. *Vander, Sherman, Luciano’s Human Physiology: the Mechanisms of Body Function. 9th edition*. Boston: McGraw-Hill Higher Education, 2004.
- Andrew G. Wilson and Ryan P. Adams. Gaussian process kernels for pattern discovery and extrapolation. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1067–1075, 2013.
- Shiwen Zhao, Chuan Gao, Sayan Mukherjee, and Barbara E. Engelhardt. Bayesian group factor analysis with structured sparsity. *Journal of Machine Learning Research*, 17(196): 1–47, 2016.