# Longitudinal Variational Autoencoder: Supplementary Materials

# 1 Covariance functions

#### Squared exponential CF

Let  $\mathbf{x}^{(r)} = x \in \mathcal{X}_j$  denote an univariate continuous-valued covariate. The squared exponential (SE) CF is defined as

$$k_{\rm se}(\boldsymbol{x}^{(r)}, \boldsymbol{x}^{(r)\prime} | \boldsymbol{\theta}_{\rm se}) = \sigma_{\rm se}^2 \exp\left(-\frac{(x-x')^2}{2\ell_{\rm se}^2}\right), \qquad \boldsymbol{\theta}_{\rm se} = (\sigma_{\rm se}^2, \ell_{\rm se})$$

where  $\sigma_{se}^2$  is the magnitude parameter (also called scale) and  $\ell_{se} \ge 0$  is the length-scale. The magnitude controls the marginal variance of the GP and length-scale controls its smoothness [Rasmussen and Williams, 2006].

#### **Categorical CF**

Let  $\mathbf{x}^{(r)} = x \in \mathcal{X}_j$  denote a categorical or discrete covariate. The categorical CF is defined as:

$$k_{\rm ca}(\boldsymbol{x}^{(r)}, \boldsymbol{x}^{(r)\prime}) = \begin{cases} 1, & \text{if } x = x' \\ 0, & \text{otherwise} \end{cases} \qquad \theta_{\rm ca} = \emptyset$$

#### **Binary CF**

Let  $x^{(r)} = x \in \mathcal{X}_i$  denote an arbitrary univariate covariate. The binary CF is defined as:

$$k_{\rm bi}(\boldsymbol{x}^{(r)}, \boldsymbol{x}^{(r)\prime}) = \begin{cases} 1, & \text{if } x = x' = 1\\ 0, & \text{otherwise} \end{cases} \qquad \theta_{\rm bi} = \emptyset$$

#### Interaction CF

Let  $\boldsymbol{x}^{(r)} = [\boldsymbol{x}^{(a)^T}, \boldsymbol{x}^{(b)^T}]^T$ , where  $\boldsymbol{x}^{(a)}$  and  $\boldsymbol{x}^{(b)}$  are arbitrary sub-vectors of  $\boldsymbol{x}$ . We define the interaction CF for these subsets as the product of two CFs defined over  $\boldsymbol{x}^{(a)}$  and  $\boldsymbol{x}^{(b)}$  respectively:

$$k_{\rm in}(\boldsymbol{x}^{(r)}, \boldsymbol{x}^{(r)\prime} | \boldsymbol{\theta}^{(r)}) = k^{(a)}(\boldsymbol{x}^{(a)}, \boldsymbol{x}^{(a)\prime} | \boldsymbol{\theta}^{(a)})k^{(b)}(\boldsymbol{x}^{(b)}, \boldsymbol{x}^{(b)\prime} | \boldsymbol{\theta}^{(b)}),$$

where  $\theta^{(r)} = \theta^{(a)} \cup \theta^{(b)}$ . The interaction CF enables us to combine any combination of univariate SE, categorical, and binary CFs in a product CF. However, in practice we restrict such combinations to include no more than a single SE CF. As Cheng et al. [2019] stated, such a condition affords the SE GP flexibility with random intercept/slope constructions similar to the linear mixed effect modelling framework, without sacrificing interpretability. For example, in longitudinal studies an instance-specific auto-correlated temporal deviation from the population-level temporal mean can be captured by an interaction term between a categorical CF over the instance identifiers and a SE CF over the temporal covariate. Furthermore, Cheng et al. [2019] adapted an approach to handle missing covariates based on an interaction CF containing a missing-ness mask.

### 2 Efficient KL divergence computation

As the main text states, optimising the variational objective of L-VAE involves the computation of L KL divergences  $D_{\text{KL}} = D_{\text{KL}}(\mathcal{N}(\bar{\boldsymbol{\mu}}_l, W_l) || \mathcal{N}(\boldsymbol{0}, \Sigma_l))$ , where  $\bar{\boldsymbol{\mu}}_l = [\mu_{\phi,l}(\boldsymbol{y}_1), \dots, \mu_{\phi,l}(\boldsymbol{y}_N)]^T$ ,  $W_l = \text{diag}(\sigma_{\phi,l}^2(\boldsymbol{y}_1), \dots, \sigma_{\phi,l}^2(\boldsymbol{y}_N))$ , and  $\Sigma_l = \sum_{r=1}^R K_{XX}^{(l,r)} + \sigma_{zl}^2 I_N$ . Henceforth, we drop the index l for notational simplicity. Each of the KL divergences is available in closed form using the well-known expression for the KL divergence between two multivariate normal distributions:

$$D_{\mathrm{KL}} = \frac{1}{2} \left( \mathrm{tr}(\Sigma^{-1}W) + \bar{\boldsymbol{\mu}}^T \Sigma^{-1} \bar{\boldsymbol{\mu}} - N + \log |\Sigma| - \log |W| \right),$$

but its exact computation requires  $\mathcal{O}(N^3)$  flops, which makes it impractical when N exceeds a few thousands. In this section, we provide a derivation of a novel strategy to approximately compute this KL divergence at a reduced computational cost.

**KL** divergence and evidence lower bound We start by exploiting the diagonal structure of W and establish the connection between the upper bound for  $D_{\text{KL}}$  and the evidence lower bound for the marginal log likelihood (MLL) of a Gaussian process:

$$D_{\mathrm{KL}}(\mathcal{N}(\bar{\boldsymbol{\mu}}, W)||\mathcal{N}(\mathbf{0}, \Sigma)) \triangleq \int_{\boldsymbol{z}} \mathcal{N}(\boldsymbol{z}|\bar{\boldsymbol{\mu}}, W) \log\left(\frac{\mathcal{N}(\boldsymbol{z}|\bar{\boldsymbol{\mu}}, W)}{\mathcal{N}(\boldsymbol{z}|\mathbf{0}, \Sigma)}\right) d\boldsymbol{z}$$
  
$$= -\int_{\boldsymbol{z}} \log\left(\mathcal{N}(\boldsymbol{z}|\mathbf{0}, \Sigma)\right) \mathcal{N}(\boldsymbol{z}|\bar{\boldsymbol{\mu}}, W) d\boldsymbol{z} - \frac{1}{2} \log\left|(2\pi e)W\right|$$
  
$$\leq -\int_{\boldsymbol{z}} \mathcal{L}(\boldsymbol{z}; \Sigma) \mathcal{N}(\boldsymbol{z}|\bar{\boldsymbol{\mu}}, W) d\boldsymbol{z} - \frac{1}{2} \log\left|(2\pi e)W\right|$$
(1)  
for any function  $\mathcal{L}(\boldsymbol{z}; \Sigma) : \mathcal{L}(\boldsymbol{z}; \Sigma) < \log\left(\mathcal{N}(\boldsymbol{z}|\mathbf{0}, \Sigma)\right) \ \forall \boldsymbol{z}.$ 

Hence, for any lower bound of GP MLL  $\mathcal{L}(\boldsymbol{z}; \Sigma)$ , the corresponding expression would provide an upper bound for KL divergence between the considered multivariate normal distributions. However, the lower bound of GP MLL is a much more common and well-studied problem. Please note that in the expression (1) and further throughout this section, we specifically use the *lower bound* and *ELBO* terms for the GP MLL lower bound  $\mathcal{L}(\boldsymbol{z}; \Sigma)$ , and not for the lower bound of the deep generative model as in the main text.

Variational learning of inducing variables in sparse Gaussian processes One of the most fundamental works on ELBOs for GP MLL was done by Titsias [2009], and is based on the paradigm of low-rank inducing point approximations of GPs. We briefly recap their key results here, and later build upon their derivation by introducing modifications that would suit the specific structure of matrices  $\Sigma$  in our problem setting.

We denote the set of inducing points locations in  $\mathcal{X}$  as  $S = [\mathbf{s}_1^T, \ldots, \mathbf{s}_M^T]^T$ , and the value of the Gaussian process at the inducing locations as  $\mathbf{u} = [u_1, \ldots, u_M]^T$ . We recall that  $\Sigma = K_{XX} + \sigma_z^2 I_N$  and therefore,  $\mathbf{z}$  can be represented as the sum of noise-free GP,  $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, K_{XX})$  and i.i.d. Gaussian noise. The following identities always stand:

$$p(\boldsymbol{z}|\boldsymbol{f}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{f}, \sigma_{\boldsymbol{z}}^{2}I_{N})$$

$$p(\boldsymbol{f}|\boldsymbol{u}) = \mathcal{N}(\boldsymbol{f}|K_{XS}K_{SS}^{-1}\boldsymbol{u}, \tilde{K}),$$

$$\tilde{K} = K_{XX} - K_{XS}K_{SS}^{-1}K_{SX}$$

$$p(\boldsymbol{u}) = \mathcal{N}(\boldsymbol{u}|\boldsymbol{0}, K_{SS})$$

$$p(\boldsymbol{z}) = \int_{\boldsymbol{u}}\int_{\boldsymbol{f}} p(\boldsymbol{z}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{u})p(\boldsymbol{u})d\boldsymbol{f}d\boldsymbol{u}.$$
(2)

Applying the Jensen inequality on the conditional log-probability  $p(\boldsymbol{z}|\boldsymbol{u})$  leads to:

$$\log p(\boldsymbol{z}|\boldsymbol{u}) = \log \int_{\boldsymbol{f}} p(\boldsymbol{z}|\boldsymbol{f}) p(\boldsymbol{f}|\boldsymbol{u}) d\boldsymbol{f}$$
(3)

$$\geq \int_{\boldsymbol{f}} \log\left(p(\boldsymbol{z}|\boldsymbol{f})\right) p(\boldsymbol{f}|\boldsymbol{u}) d\boldsymbol{f} = \sum_{i=1}^{N} \left[ \log \mathcal{N}(z_i|\mu_i, \sigma_{\boldsymbol{z}}^2) - \frac{\tilde{K}_{ii}}{2\sigma_{\boldsymbol{z}}^2} \right],\tag{4}$$

where  $\boldsymbol{\mu} = [\mu_1, \ldots, \mu_N]^T = K_{XS} K_{SS}^{-1} \boldsymbol{u}$  and  $\tilde{K}_{ii}$  denotes the *i*<sup>th</sup> diagonal element of  $\tilde{K}$ . The inequation reduces to identity *iff*  $\boldsymbol{u}$  is a sufficient statistic of  $\boldsymbol{f}$ , so that all elements of  $\tilde{K}$  are zero. However, the inequation remains tight and the approximation is justified as long as the  $\tilde{K}_{ii}$  elements remain small, which is achieved by setting Mto be sufficiently high and optimising the inducing point locations S. After integrating out  $\boldsymbol{u}$ , this approximation leads to the collapsed representation of the variational evidence lower bound,

$$\mathcal{L}_{1}(\boldsymbol{z};\boldsymbol{\Sigma}) \triangleq \log \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, K_{XS}K_{SS}^{-1}K_{SX} + \sigma_{\boldsymbol{z}}^{2}I_{N}) - \frac{1}{2\sigma_{\boldsymbol{z}}^{2}}\mathrm{tr}(\tilde{K}).$$
(5)

**Divergence upper bound for longitudinal Gaussian process** The free-form bound is known to be tight when M is sufficiently high and the covariance function is sufficiently smooth. However, longitudinal studies, by definition, always contain a categorical covariate corresponding to instances, which makes the covariance function non-continuous. By separating the additive component that corresponds to the interaction between instances and time (or age) from the other additive components, the covariance matrix has the following general form  $\Sigma = K_{XX}^{(A)} + \hat{\Sigma}$ , where  $\hat{\Sigma} = \text{diag}(\hat{\Sigma}_1, \dots, \hat{\Sigma}_P)$ ,  $\hat{\Sigma}_P = K_{X_PX_P}^{(R)} + \sigma_z^2 I_{n_P}$ , and  $K_{XX}^{(A)} = \sum_{r=1}^{R-1} K_{XX}^{(r)}$  contains all the other R-1 components. In order to keep the  $\tilde{K}_{ii}$  tight, the bound from eq. (5) would mandate that  $M \geq P$ , thus rendering the bound of Titsias [2009] either computationally inefficient once P is large (due to high M) or insufficiently tight. Since the interaction covariance function is essential for accurate longitudinal modelling, we devised a novel free-form divergence upper bound for such a class of GPs. Similar to eq. (2), we exploit the opportunity to represent z as a sum of noise-free GP  $f_A \sim \mathcal{N}(\mathbf{0}, K_{XX}^{(A)})$  and structured noise  $\hat{f} \sim \mathcal{N}(\mathbf{0}, \hat{\Sigma})$ , and we assign the locations as well as values of inducing points for  $f_A$ . We denote  $\tilde{K}^{(A)} = K_{XX}^{(A)} - K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)}$  and  $\tilde{K}_{XPXP}^{(A)} = K_{XPX}^{(A)} - K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)}$ .

$$\begin{split} \log p(\boldsymbol{z}|\boldsymbol{u}) &= \log \int_{\boldsymbol{f}_A} p(\boldsymbol{z}|\boldsymbol{f}_A) p(\boldsymbol{f}_A|\boldsymbol{u}) d\boldsymbol{f}_A \\ &= \log \int_{\boldsymbol{f}_A} \mathcal{N}(\boldsymbol{z}|\boldsymbol{f}_A, \hat{\Sigma}) \mathcal{N}(\boldsymbol{f}_A|K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\boldsymbol{u}, \tilde{K}^{(A)}) d\boldsymbol{f}_A \\ &\geq \int_{\boldsymbol{f}_A} \log \left( \mathcal{N}(\boldsymbol{z}|\boldsymbol{f}_A, \hat{\Sigma}) \right) \mathcal{N}(\boldsymbol{f}_A|K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\boldsymbol{u}, \tilde{K}^{(A)}) d\boldsymbol{f}_A \\ &= \sum_{p=1}^P \left[ \log \mathcal{N}\left(\boldsymbol{z}_p|\boldsymbol{\mu}_p, \hat{\Sigma}_p\right) - \frac{1}{2} \mathrm{tr}\left(\hat{\Sigma}_p^{-1}\tilde{K}_{X_pX_p}^{(A)}\right) \right], \end{split}$$

where  $\boldsymbol{z}_p$  is the sub-vector of  $\boldsymbol{z}$  that corresponds to the  $p^{\text{th}}$  individual,  $\boldsymbol{\mu}_p$  is the subvector of  $K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\boldsymbol{u}$ , and  $\text{tr}(\cdot)$  denotes the matrix trace operator. After integrating out  $\boldsymbol{u}$ , we obtain a novel variational evidence lower bound,

$$\mathcal{L}_{2}(\boldsymbol{z};\boldsymbol{\Sigma}) \triangleq \log \mathcal{N}\left(\boldsymbol{z}|\boldsymbol{0}, K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\boldsymbol{\Sigma}}\right) - \frac{1}{2} \sum_{p=1}^{P} \operatorname{tr}\left(\hat{\boldsymbol{\Sigma}}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)}\right).$$
(6)

We substitute this lower bound in eq. (1), which links the ELBO and upper bound for KLD:

$$D_{\mathrm{KL}} = D_{\mathrm{KL}}(\mathcal{N}(\bar{\mu}, W || \mathcal{N}(\mathbf{0}, \Sigma)) \leq -\int_{\boldsymbol{z}} \mathcal{L}_{2}(\boldsymbol{z}; \Sigma) \mathcal{N}(\boldsymbol{z} | \bar{\mu}, W) d\boldsymbol{z} - \frac{1}{2} \log |(2\pi e)W|$$

$$= -\int_{\boldsymbol{z}} \log \mathcal{N}\left(\boldsymbol{z} |K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\Sigma}\right) \mathcal{N}(\boldsymbol{z} | \bar{\mu}, W) d\boldsymbol{z} - \frac{1}{2} \log |(2\pi e)W| + \frac{1}{2} \sum_{p=1}^{P} \operatorname{tr}\left(\hat{\Sigma}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)}\right)$$

$$= D_{\mathrm{KL}}\left(\mathcal{N}(\bar{\mu}, W || \mathcal{N}(\mathbf{0}, K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\Sigma}\right) + \frac{1}{2} \sum_{p=1}^{P} \operatorname{tr}\left(\hat{\Sigma}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)}\right)$$

$$= \frac{1}{2} \left(\operatorname{tr}(\bar{\Sigma}^{-1}W) + \bar{\mu}^{T} \bar{\Sigma}^{-1} \bar{\mu} - N + \log |\bar{\Sigma}| - \log |W| + \sum_{p=1}^{P} \operatorname{tr}\left(\hat{\Sigma}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)}\right)\right) \triangleq \mathcal{D}_{2}, \quad (7)$$
where  $\bar{\Sigma} = K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\Sigma}.$ 

**Theorem 1.** For any set of inducing points S and  $\overline{S}$  where  $\overline{S} \subseteq S$ , such that  $\operatorname{rank}(K_{SS}) = \operatorname{rank}(K_{\overline{SS}}^{(A)})$ , non-strict inequality  $\mathcal{L}_1(\boldsymbol{z}; \Sigma, S) \leq \mathcal{L}_2(\boldsymbol{z}; \Sigma, \overline{S})$  holds.

Proof. Based on the univariate additive GP  $f(\boldsymbol{x}) = \sum_{r=1}^{R} f^{(r)}(\boldsymbol{x})$ , we construct another GP  $g(v_1, v_2, \boldsymbol{x}) = \mathbb{1}(v_1) \sum_{r=1}^{R-1} f^{(r)}(\boldsymbol{x}) + \mathbb{1}(v_2) f^{(R)}(\boldsymbol{x})$ , where  $\mathbb{1}(v)$  is the indicator function that equals 1 when v = 1 and 0 otherwise. Accordingly, we define augmented sets  $\hat{X} = [\mathbf{1}_N, \mathbf{1}_N, X]$  and  $\hat{S} = [\mathbf{1}_M, \mathbf{1}_M, S]$ , where  $\mathbf{1}_N$  is a column vector of ones that has length N, which effectively adds the two additional covariates,  $v_1$  and  $v_2$ , into X and S. Then, the marginal covariance of  $g(v_1, v_2, \boldsymbol{x})$  for  $(\hat{X}, \hat{X})$  is  $K_{XX}$ , for  $(\hat{X}, \hat{S})$  is  $K_{XS}$ , and for  $(\hat{S}, \hat{S})$  is  $K_{SS}$ . Thus,  $\mathcal{L}_1(\boldsymbol{z}; \Sigma, S) = \mathcal{L}_1(\boldsymbol{z}; \Sigma, \hat{S})$ .

Since the collapsed lower bound, eq. (5), is obtained as a closed-form solution to the optimisation problem of the variational parameters of the distribution over the inducing points, expanding the set of the inducing points can only expand the variational family and, therefore, never decreases the  $\mathcal{L}_1(\boldsymbol{z}; \Sigma, \hat{S})$  bound [Titsias, 2009]. We consider an expanded set of inducing points  $\tilde{S} = [\hat{S}^T, \hat{S}_A^T, \hat{S}_B^T, \hat{X}_B^T]^T$ , where  $\hat{S}_A = [\mathbf{1}_M, \mathbf{0}_M, S]$  ( $\mathbf{0}_M$ is a column vector of zeros that has length M),  $\hat{S}_B = [\mathbf{0}_M, \mathbf{1}_M, S]$ , and  $\hat{X}_B = [\mathbf{0}_N, \mathbf{1}_N, X]$ . For any  $\boldsymbol{s} \in \hat{S}$ ,  $g(1, 1, \boldsymbol{s}) = g(1, 0, \boldsymbol{s}) + g(0, 1, \boldsymbol{s})$ , so that the values of  $g(v_1, v_2, \boldsymbol{x})$  over  $\hat{S}$  are linearly dependent on the values over  $\hat{S}_A$  and  $\hat{S}_B$ . Therefore, the variational family remains the same for the reduced set of inducing points  $\check{S} = [\hat{S}_A^T, \hat{S}_B^T, \hat{X}_B^T]^T$ , and  $\mathcal{L}_1(\boldsymbol{z}; \Sigma, \hat{S}) \leq \mathcal{L}_1(\boldsymbol{z}; \Sigma, \tilde{S}) = \mathcal{L}_1(\boldsymbol{z}; \Sigma, \check{S})$ .

We will next write the  $\mathcal{L}_1$  bound from equation (5) for  $\check{S}$ :

$$\mathcal{L}_{1}(\boldsymbol{z};\boldsymbol{\Sigma},\check{S}) = \log \mathcal{N}(\boldsymbol{z}|\boldsymbol{0},\check{K}_{X\check{S}}\check{K}_{\check{S}\check{S}}^{-1}\check{K}_{\check{S}X} + \sigma_{\boldsymbol{z}}^{2}I_{N}) - \frac{1}{2\sigma_{\boldsymbol{z}}^{2}}\operatorname{tr}(K_{XX} - \check{K}_{X\check{S}}\check{K}_{\check{S}\check{S}}^{-1}\check{K}_{\check{S}X}),$$

where (from the definition of  $\check{S}$ )

$$\check{K}_{X\check{S}} = \begin{bmatrix} K_{XS}^{(A)}, K_{XS}^{(R)}, K_{XX}^{(R)} \end{bmatrix} \text{ and } \check{K}_{\check{S}\check{S}} = \begin{bmatrix} K_{SS}^{(A)} & 0 & 0\\ 0 & K_{SS}^{(R)} & K_{SX}^{(R)} \\ 0 & K_{XS}^{(R)} & K_{XX}^{(R)} \end{bmatrix}.$$

We recall the 2-by-2 block-matrix inverse formula for the second and third blocks of  $\check{K}_{\check{S}\check{S}}$ :

$$\begin{bmatrix} K_{SS}^{(R)} & K_{SX}^{(R)} \\ K_{XS}^{(R)} & K_{XX}^{(R)} \end{bmatrix}^{-1} = \begin{bmatrix} Q^{-1} & -Q^{-1}K_{SX}^{(R)}K_{XX}^{(R)-1} \\ -K_{XX}^{(R)^{-1}}K_{XS}^{(R)}Q^{-1} & K_{XX}^{(R)^{-1}} + K_{XX}^{(R)^{-1}}K_{XS}^{(R)}Q^{-1}K_{XX}^{(R)^{-1}} \end{bmatrix}$$

where  $Q = K_{SS}^{(R)} - K_{SX}^{(R)} K_{XX}^{(R)^{-1}} K_{XS}^{(R)}$ . Then, we substitute this expression into the matrix-product term  $\check{K}_{X\check{S}}\check{K}_{\check{S}\check{S}}^{-1}\check{K}_{\check{S}X}$ , alongside utilising the zero-blocks in first row/column of  $\check{K}_{\check{S}\check{S}}$ . Such a manipulation yields a very simple expression, as all the terms involving Q are cancelled out:

$$\begin{split} \check{K}_{X\check{S}}\check{K}_{\check{S}\check{S}}^{-1}\check{K}_{\check{S}X} &= K_{XS}^{(A)}K_{SS}^{(A)}{}^{-1}K_{SX}^{(A)} + \begin{bmatrix} K_{XS}^{(R)}, K_{XX}^{(R)} \end{bmatrix} \begin{bmatrix} K_{SS}^{(R)} & K_{SX}^{(R)} \\ K_{XS}^{(R)} & K_{XX}^{(R)} \end{bmatrix}^{-1} \begin{bmatrix} K_{XS}^{(R)}, K_{XX}^{(R)} \end{bmatrix}^{-1} \\ &= K_{XS}^{(A)}K_{SS}^{(A)}{}^{-1}K_{SX}^{(A)} + K_{XX}^{(R)}. \end{split}$$

The  $\mathcal{L}_1$  bound for the inducing point set  $\check{S}$  can then be written as:

$$\mathcal{L}_{1}(\boldsymbol{z};\boldsymbol{\Sigma},\check{S}) = \log \mathcal{N}\left(\boldsymbol{z}|\boldsymbol{0}, K_{XS}^{(A)}K_{SS}^{(A)}^{-1}K_{SX}^{(A)} + K_{XX}^{(R)} + \sigma_{\boldsymbol{z}}^{2}I_{N}\right) - \frac{1}{2\sigma_{\boldsymbol{z}}^{2}} \operatorname{tr}(K_{XX} - K_{XS}^{(A)}K_{SS}^{(A)}^{-1}K_{SX}^{(A)} - K_{XX}^{(R)})$$
$$= \log \mathcal{N}\left(\boldsymbol{z}|\boldsymbol{0}, K_{XS}^{(A)}K_{SS}^{(A)}^{-1}K_{SX}^{(A)} + \hat{\boldsymbol{\Sigma}}\right) - \frac{1}{2\sigma_{\boldsymbol{z}}^{2}} \operatorname{tr}(K_{XX}^{(A)} - K_{XS}^{(A)}K_{SS}^{(A)}^{-1}K_{SX}^{(A)}). \tag{8}$$

Focusing on the last trace term we have,

where  $\hat{L}$  is the Cholesky decomposition of  $(K_{XX}^{(R)}{}^{-1} + \sigma_z^{-2}I_N)^{-1}$ . As matrix  $K_{XX}^{(R)}{}^{-1} + \sigma_z^{-2}I_N$  is always positive definite, its inverse and corresponding Cholesky decomposition always exist. Additionally, we used the matrix trace cyclic rotation property and the fact that the trace of a positive semidefinite matrix is  $\hat{L}^T \tilde{K}_{XX}^{(A)} \hat{L}$ , which is always non-negative. Combining equations (8) and (9) we have

$$\mathcal{L}_{1}(\boldsymbol{z}; \Sigma, \check{S}) = \log \mathcal{N}\left(\boldsymbol{z} | \boldsymbol{0}, K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\Sigma}\right) - \frac{1}{2\sigma_{\boldsymbol{z}}^{2}} \operatorname{tr}(K_{XX}^{(A)} - K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)})$$
  
$$\leq \log \mathcal{N}\left(\boldsymbol{z} | \boldsymbol{0}, K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\Sigma}\right) - \frac{1}{2} \operatorname{tr}\left(\hat{\Sigma}^{-1} \tilde{K}_{XX}^{(A)}\right)$$
  
$$= \mathcal{L}_{2}(\boldsymbol{z}; \Sigma, S).$$

Consequently, the non-strict inequalities chain gives the final result:

$$\mathcal{L}_1(\boldsymbol{z}; \boldsymbol{\Sigma}, S) = \mathcal{L}_1(\boldsymbol{z}; \boldsymbol{\Sigma}, \hat{S}) \leq \mathcal{L}_1(\boldsymbol{z}; \boldsymbol{\Sigma}, \tilde{S}) = \mathcal{L}_1(\boldsymbol{z}; \boldsymbol{\Sigma}, \check{S}) \leq \mathcal{L}_2(\boldsymbol{z}; \boldsymbol{\Sigma}, S)$$

**Computational complexity** The computational complexity of the bound described in eq. (7), is  $\mathcal{O}(\sum_{p=1}^{P} n_p^3 + NM^2)$  flops, which leads to an approximately similar computational complexity as the Titsias [2009] bound, when  $n_p \simeq M \ll N$ , but is significantly tighter. Here, we elucidate in detail how to perform the computations in eq. (7) to achieve such complexity. The first term of eq. (7) can be written as:

$$\operatorname{tr}(\bar{\Sigma}^{-1}W) = \operatorname{tr}(\hat{\Sigma}^{-1}W) - \operatorname{tr}\left(\hat{\Sigma}^{-1}K_{XS}^{(A)}\left[K_{SS}^{(A)} + K_{SX}^{(A)}\hat{\Sigma}^{-1}K_{XS}^{(A)}\right]^{-1}K_{SX}^{(A)}\hat{\Sigma}^{-1}W\right)$$

$$= \sum_{p=1}^{P} \left(\operatorname{diag}(\hat{\Sigma}^{-1}) \cdot \operatorname{diag}(W_p)\right) - \operatorname{tr}\left(V^{-1}\left(K_{SX}^{(A)}\hat{\Sigma}^{-1}W\hat{\Sigma}^{-1}K_{XS}^{(A)}\right)\right),$$

where  $V = K_{SS}^{(A)} + K_{SX}^{(A)} \hat{\Sigma}^{-1} K_{XS}^{(A)}$ . We have used the Woodbury matrix identity to obtain the first equality and the cyclic rotation property of the matrix trace to obtain the second equality [Press et al., 2007]. Since  $\hat{\Sigma}$  is a block-diagonal matrix, obtaining  $\hat{\Sigma}^{-1}$  takes  $\mathcal{O}(\sum_{p=1}^{P} n_p^3)$  flops and obtaining products  $K_{SX}^{(A)} \hat{\Sigma}^{-1} K_{XS}^{(A)}$  or  $K_{SX}^{(A)} (\hat{\Sigma}^{-1} W \hat{\Sigma}^{-1}) K_{XS}^{(A)}$  takes  $\mathcal{O}(\sum_{p=1}^{P} n_p M^2) = \mathcal{O}(NM^2)$  flops. Moreover, inverting V takes  $\mathcal{O}(M^3)$  flops.

We can use the Woodbury matrix identity to write the second term as,

$$\bar{\mu}^T \bar{\Sigma}^{-1} \bar{\mu} = \bar{\mu}^T \hat{\Sigma}^{-1} \bar{\mu} - \bar{\mu}^T \hat{\Sigma}^{-1} K_{XS}^{(A)} V^{-1} K_{SX}^{(A)} \hat{\Sigma}^{-1} \bar{\mu}.$$

Obtaining  $\hat{\Sigma}^{-1}\bar{\mu}$  takes  $\mathcal{O}(\sum_{p=1}^{P} n_p^2)$  flops and obtaining  $K_{SX}^{(A)}\left(\hat{\Sigma}^{-1}\bar{\mu}\right)$  takes  $\mathcal{O}(NM^2)$  flops.

For the forth term we use the generalised determinant lemma so that,

$$|\bar{\Sigma}| = |\hat{\Sigma}| |K_{SS}^{(A)}|^{-1} |V|.$$

These determinant computations take  $\mathcal{O}(\sum_{p=1}^{P} n_p^3)$ ,  $\mathcal{O}(M^3)$ , and  $\mathcal{O}(M^3)$ , respectively.

The fifth term is trivially  $\mathcal{O}(N)$ . We can again use the cyclic rotation property of the matrix trace to write the last term as:

$$\operatorname{tr}\left(\hat{\Sigma}_{p}^{-1}\tilde{K}_{X_{p}X_{p}}^{(A)}\right) = \operatorname{tr}\left(\hat{\Sigma}_{p}^{-1}K_{X_{p}X_{p}}^{(A)}\right) - \operatorname{tr}\left(\hat{\Sigma}_{p}^{-1}K_{X_{p}S}^{(A)}K_{SS}^{(A)}K_{SX_{p}}^{(A)}\right)$$
$$= \operatorname{tr}\left(\hat{\Sigma}_{p}^{-1}K_{X_{p}X_{p}}^{(A)}\right) - \operatorname{tr}\left(\left(K_{SX_{p}}^{(A)}\hat{\Sigma}_{p}^{-1}K_{X_{p}S}^{(A)}\right)K_{SS}^{(A)}\right).$$

The first trace takes  $\mathcal{O}(n_p^2)$  to compute once the  $\hat{\Sigma}^{-1}$  is available. It takes  $\mathcal{O}(n_p M^2)$  to compute the product in parenthesis in the second trace, and  $\mathcal{O}(M^2)$  for the trace itself once  $K_{SS}^{(A)^{-1}}$  is available. Since we need to compute these traces for  $p = 1, \ldots, P$ , the overall complexity of the last term in eq. (7) is  $\mathcal{O}(\sum_{p=1}^{P} n_p^2 + NM^2)$ . Combining the terms we get the final time complexity:

Complexity = 
$$\mathcal{O}(\sum_{p=1}^{P} n_p^3) + \mathcal{O}(NM^2) + \mathcal{O}(M^3) + \mathcal{O}(\sum_{p=1}^{P} n_p^2) + \mathcal{O}(NM^2)$$
  
+ $\mathcal{O}(\sum_{p=1}^{P} n_p^3) + \mathcal{O}(M^3) + \mathcal{O}(M^3) + \mathcal{O}(N) + \mathcal{O}(\sum_{p=1}^{P} n_p^2 + NM^2)$   
=  $\mathcal{O}(\sum_{p=1}^{P} n_p^3 + NM^2).$  (10)

### 3 Stochastic Variational Inference for longitudinal Gaussian process

A fundamental drawback of the methods described in the previous section is that they require using the full training dataset to compute the loss and perform a gradient step. This can be an issue for many problems with large data, such as sequences of images, electronic health records, etc. A common machine learning technique to tackle such problems is based on training the model using mini-batches. The mini-batch approach makes use of unbiased stochastic estimates of the loss and its gradients, which are computed based on a subset of the data. The subsets are chosen such that all training data points are used within an epoch. In this section, we first recall earlier work on how to adjust the bound of eq. (5) to become compatible with stochastic variational inference and mini-batching. Then, we modify this bound to account for the specifics of the GP covariance structure in L-VAE.

In contrast to Titsias [2009], Hensman et al. [2013] proposed to avoid analytical marginalisation of inducing values  $\boldsymbol{u}$  in eq. (3). Instead, Hensman et al. [2013] proposed to explicitly keep track of its distribution, which is assumed to be Gaussian  $\boldsymbol{u} \sim \mathcal{N}(\boldsymbol{m}, H)$ . Then, the authors derived an alternative evidence lower bound:

$$\log p(\boldsymbol{z}) = \int_{\boldsymbol{u}} \int_{\boldsymbol{f}} p(\boldsymbol{z}|\boldsymbol{f}) p(\boldsymbol{f}|\boldsymbol{u}) p(\boldsymbol{u}) d\boldsymbol{f} d\boldsymbol{u} = \log \int_{\boldsymbol{u}} p(\boldsymbol{z}|\boldsymbol{u}) \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} q(\boldsymbol{u}) d\boldsymbol{u} \ge \int_{\boldsymbol{u}} \log p(\boldsymbol{z}|\boldsymbol{u}) q(\boldsymbol{u}) d\boldsymbol{u} - D_{KL}(q(\boldsymbol{u})||p(\boldsymbol{u}))$$

$$\geq \int_{\boldsymbol{u}} \log \mathcal{N}(\boldsymbol{z}|K_{XS}K_{SS}^{-1}\boldsymbol{u}, \sigma_{z}^{2}I_{N}) \mathcal{N}(\boldsymbol{u}|\boldsymbol{m}, H) d\boldsymbol{u} - \frac{1}{2\sigma_{z}^{2}} \operatorname{tr}(\tilde{K}) - D_{KL}(\mathcal{N}(\boldsymbol{m}, H)||\mathcal{N}(\boldsymbol{0}, K_{SS}))$$

$$= \log \mathcal{N}(\boldsymbol{z}|K_{XS}K_{SS}^{-1}\boldsymbol{m}, \sigma_{z}^{2}I_{N}) - \frac{1}{2\sigma_{z}^{2}} \operatorname{tr}(HK_{SS}^{-1}K_{SX}K_{XS}K_{SS}^{-1}) - \frac{1}{2\sigma_{z}^{2}} \operatorname{tr}(\tilde{K}) - D_{KL}(\mathcal{N}(\boldsymbol{m}, H)||\mathcal{N}(\boldsymbol{0}, K_{SS})))$$

$$\triangleq \mathcal{L}_{3}.$$
(11)

Substituting this bound in eq. (1) yields:

$$\begin{split} D_{\mathrm{KL}}(\mathcal{N}(\bar{\mu},W)||\mathcal{N}(\mathbf{0},\Sigma)) &\leq -\int_{\mathbf{z}} \mathcal{L}_{3}\mathcal{N}(\mathbf{z}|\bar{\mu},W)d\mathbf{z} - \frac{1}{2}\log|(2\pi e)W| \\ &= -\int_{\mathbf{z}}\log\mathcal{N}(\mathbf{z}|K_{XS}K_{SS}^{-1}\mathbf{m},\sigma_{z}^{2}I_{N})\mathcal{N}(\mathbf{z}|\bar{\mu},W)d\mathbf{z} + \frac{1}{2\sigma_{z}^{2}}\mathrm{tr}(\tilde{K}) \\ &+ \frac{1}{2\sigma_{z}^{2}}\mathrm{tr}(HK_{SS}^{-1}K_{SX}K_{XS}K_{SS}^{-1}) + D_{KL}(\mathcal{N}(\mathbf{m},H)||\mathcal{N}(\mathbf{0},K_{SS})) - \frac{1}{2}\log|(2\pi e)W| \\ &= \frac{1}{2}(K_{XS}K_{SS}^{-1}\mathbf{m} - \bar{\mu})^{T}(\sigma_{z}^{2}I_{N})^{-1}(K_{XS}K_{SS}^{-1}\mathbf{m} - \bar{\mu}) + \frac{1}{2}\mathrm{tr}((\sigma_{z}^{2}I_{N})^{-1}W) + \frac{N}{2}\log\sigma_{z}^{2} \\ &+ \frac{N}{2}\log2\pi + \frac{1}{2\sigma_{z}^{2}}\mathrm{tr}(\tilde{K}) + \frac{1}{2\sigma_{z}^{2}}\mathrm{tr}(HK_{SS}^{-1}K_{SX}K_{XS}K_{SS}^{-1}) \\ &+ D_{KL}(\mathcal{N}(\mathbf{m},H)||\mathcal{N}(\mathbf{0},K_{SS})) - \frac{1}{2}\log|(2\pi e)W| \\ &= \frac{1}{2}\left(\sigma_{z}^{-2}\sum_{i=1}^{N}(K_{\mathbf{x}_{iS}}K_{SS}^{-1}\mathbf{m} - \bar{\mu}_{i})^{2} + \sigma_{z}^{-2}\sum_{i=1}^{N}\sigma_{\phi}^{2}(\mathbf{y}_{i}) + N\log\sigma_{z}^{2} + \sigma_{z}^{-2}\sum_{i=1}^{N}\tilde{K}_{ii} \\ &+ \sigma_{z}^{-2}\sum_{i=1}^{N}\mathrm{tr}\left((K_{SS}^{-1}HK_{SS}^{-1})(K_{S\mathbf{x}_{i}}K_{\mathbf{x}_{i}S})\right) - \sum_{i=1}^{N}\log\sigma_{\phi}^{2}(\mathbf{y}_{i}) - N\right) \\ &+ D_{KL}(\mathcal{N}(\mathbf{m},H)||\mathcal{N}(\mathbf{0},K_{SS})) \\ &\triangleq \mathcal{D}_{3} \end{split}$$

Each term, except the last one, is additive over i = 1, ..., N. Therefore, replacing the sum over all i = 1, ..., N with a batch-normalised partial sum over a subset of indices,  $\mathcal{I} \subset \{1, ..., N\}$  of size  $|\mathcal{I}| = \hat{N}$ :

$$\hat{\mathcal{D}}_{3} = \frac{1}{2} \frac{N}{\hat{N}} \sum_{i \in \mathcal{I}} \left( \sigma_{z}^{-2} (K_{\boldsymbol{x}_{i}S} K_{SS}^{-1} \boldsymbol{m} - \bar{\boldsymbol{\mu}}_{i})^{2} + \sigma_{z}^{-2} \sigma_{\phi}^{2} (\boldsymbol{y}_{i}) + \sigma_{z}^{-2} \tilde{K}_{ii} + \sigma_{z}^{-2} \operatorname{tr} \left( \left( K_{SS}^{-1} H K_{SS}^{-1} \right) (K_{S\boldsymbol{x}_{i}} K_{\boldsymbol{x}_{i}S}) \right) - \log \sigma_{\phi}^{2} (\boldsymbol{y}_{i}) \right) \\ + \frac{N}{2} \log \sigma_{z}^{2} - \frac{N}{2} + D_{KL} (\mathcal{N}(\boldsymbol{m}, H) || \mathcal{N}(\boldsymbol{0}, K_{SS})),$$
(13)

is an unbiased estimate of the KL divergence upper bound  $E_{\mathcal{I}\sim\mathfrak{S}\{1,\ldots,N\}}(\hat{\mathcal{D}}_3) = \mathcal{D}_3 \geq D_{\mathrm{KL}}(\mathcal{N}(\bar{\boldsymbol{\mu}},W)||\mathcal{N}(\mathbf{0},\Sigma))$ . Here,  $\mathfrak{S}\{1,\ldots,N\}$  is a uniform distribution over elements of arbitrary fixed partitions of set  $\{1,\ldots,N\}$ . This property enables us to use the mini-batching technique for the approximate computation of the KL divergence term of L-VAE and its gradients.

Similar to our criticism of the ELBO  $\mathcal{L}_1$ , the  $\mathcal{L}_3$  is not well suited to the typical properties of GP covariance structures used for longitudinal modelling. Following the same notation as in the previous section, we introduce a modification to eq. (11):

$$\log p(\boldsymbol{z}) = \int_{\boldsymbol{u}} \int_{\boldsymbol{f}} p(\boldsymbol{z}|\boldsymbol{f}) p(\boldsymbol{f}|\boldsymbol{u}) p(\boldsymbol{u}) d\boldsymbol{f} d\boldsymbol{u} = \log \int_{\boldsymbol{u}} p(\boldsymbol{z}|\boldsymbol{u}) \frac{p(\boldsymbol{u})}{q(\boldsymbol{u})} q(\boldsymbol{u}) d\boldsymbol{u} \ge \int_{\boldsymbol{u}} \log p(\boldsymbol{z}|\boldsymbol{u}) q(\boldsymbol{u}) d\boldsymbol{u} - D_{KL}(q(\boldsymbol{u})||p(\boldsymbol{u}))$$

$$\geq \int_{\boldsymbol{u}} \log \mathcal{N}(\boldsymbol{z}|K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\boldsymbol{u}, \hat{\Sigma}) \mathcal{N}(\boldsymbol{u}|\boldsymbol{m}, H) d\boldsymbol{u} - \frac{1}{2} \sum_{p=1}^{P} \operatorname{tr} \left( \hat{\Sigma}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)} \right) - D_{KL}(\mathcal{N}(\boldsymbol{m}, H)||\mathcal{N}(\boldsymbol{0}, K_{SS}^{(A)}))$$

$$= \log \mathcal{N}(\boldsymbol{z}|K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\boldsymbol{m}, \hat{\Sigma}) - \frac{1}{2} \operatorname{tr}(HK_{SS}^{(A)^{-1}}K_{SX}^{(A)} \hat{\Sigma}^{-1}K_{XS}^{(A)}K_{SS}^{(A)^{-1}}) - \frac{1}{2} \sum_{p=1}^{P} \operatorname{tr} \left( \hat{\Sigma}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)} \right)$$

$$- D_{KL}(\mathcal{N}(\boldsymbol{m}, H)||\mathcal{N}(\boldsymbol{0}, K_{SS}^{(A)}))$$

$$\triangleq \mathcal{L}_{4}$$
(14)

Substituting this novel ELBO in eq. (1) yields:

$$\begin{split} D_{\mathrm{KL}}(\mathcal{N}(\bar{\mu},W)||\mathcal{N}(\mathbf{0},\Sigma)) &\leq -\int_{\mathbf{z}}^{\mathcal{L}} \mathcal{L}_{4}\mathcal{N}(\mathbf{z}|\bar{\mu},W)d\mathbf{z} - \frac{1}{2}\log|(2\pi e)W| \\ &= -\int_{\mathbf{z}}^{1}\log\mathcal{N}(\mathbf{z}|K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\mathbf{m},\hat{\Sigma})\mathcal{N}(\mathbf{z}|\bar{\mu},W)d\mathbf{z} + \frac{1}{2}\sum_{p=1}^{P}\mathrm{tr}\left(\hat{\Sigma}_{p}^{-1}\tilde{K}_{X_{p}X_{p}}^{(A)}\right) \\ &+ \frac{1}{2}\mathrm{tr}(HK_{SS}^{(A)^{-1}}K_{SX}^{(A)}\hat{\Sigma}^{-1}K_{XS}^{(A)}K_{SS}^{(A)^{-1}}) + D_{KL}(\mathcal{N}(\mathbf{m},H)||\mathcal{N}(\mathbf{0},K_{SS}^{(A)})) \\ &- \frac{1}{2}\log|(2\pi e)W| \\ &= \frac{1}{2}(K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\mathbf{m} - \bar{\mu})^{T}\hat{\Sigma}^{-1}(K_{XS}^{(A)}K_{SS}^{(A)^{-1}}\mathbf{m} - \bar{\mu}) + \frac{1}{2}\mathrm{tr}(\hat{\Sigma}^{-1}W) + \frac{1}{2}\log|\hat{\Sigma}| \\ &+ \frac{N}{2}\log2\pi + \frac{1}{2}\sum_{p=1}^{P}\mathrm{tr}\left(\hat{\Sigma}_{p}^{-1}\hat{K}_{X_{p}X_{p}}^{(A)}\right) + \frac{1}{2}\mathrm{tr}(HK_{SS}^{(A)^{-1}}K_{SX}^{(A)}\hat{\Sigma}^{-1}K_{XS}^{(A)}K_{SS}^{(A)^{-1}}) \\ &+ D_{KL}(\mathcal{N}(\mathbf{m},H)||\mathcal{N}(\mathbf{0},K_{SS}^{(A)})) - \frac{1}{2}\log|(2\pi e)W| \\ &= \frac{1}{2}\left(\sum_{p=1}^{P}(K_{X_{p}S}^{(A)}K_{SS}^{(A)^{-1}}\mathbf{m} - \hat{\mu}_{p})^{T}\hat{\Sigma}_{p}^{-1}(K_{X_{p}S}^{(A)}K_{SS}^{(A)^{-1}}\mathbf{m} - \hat{\mu}_{p}) \\ &+ \sum_{p=1}^{N}\sum_{i=1}^{n_{p}}(\hat{\Sigma}_{p}^{-1})_{ii}\sigma_{\phi}^{2}(\mathbf{y}_{I_{p}i}) + \sum_{p=1}^{P}\log|\hat{\Sigma}_{p}| + \frac{1}{2}\sum_{p=1}^{P}\mathrm{tr}\left(\hat{\Sigma}_{p}^{-1}\tilde{K}_{X_{p}X_{p}}\right) \\ &+ \sum_{p=1}^{P}\mathrm{tr}\left(\left(K_{SS}^{(A)^{-1}}\mathbf{H}K_{SS}^{(A)^{-1}}\right)\left(K_{SX_{p}}^{(A)}\hat{\Sigma}_{p}^{-1}K_{X_{p}S}^{(A)}\right)\right) - \sum_{i=1}^{N}\log\sigma_{\phi}^{2}(\mathbf{y}_{i}) - N\right) \\ &+ D_{KL}(\mathcal{N}(\mathbf{m},H)||\mathcal{N}(\mathbf{0},K_{SS}^{(A)})) \end{aligned}$$

where  $\mathcal{I}_{pi}$  is the index of the *i*<sup>th</sup> sample for the *p*<sup>th</sup> patient and  $\hat{\mu}_p = [\bar{\mu}_{\mathcal{I}_{p1}}, \dots, \bar{\mu}_{\mathcal{I}_{pnp}}]^T$  is a sub-vector of  $\bar{\mu}$  that corresponds to the *p*<sup>th</sup> patient. Each term, except for the last one, is additive over  $p = 1, \dots, P$ . Therefore, replacing the sum over all  $p = 1, \dots, P$  with a batch-normalised partial sum over a subset of indices  $\mathcal{P} \subset \{1, \dots, P\}$  of size  $|\mathcal{P}| = \hat{P}$ :

$$\hat{\mathcal{D}}_{4} = \frac{1}{2} \frac{P}{\hat{P}} \sum_{p \in \mathcal{P}} \left( (K_{X_{pS}}^{(A)} K_{SS}^{(A)^{-1}} \boldsymbol{m} - \hat{\boldsymbol{\mu}}_{p})^{T} \hat{\Sigma}_{p}^{-1} (K_{X_{pS}}^{(A)} K_{SS}^{(A)^{-1}} \boldsymbol{m} - \hat{\boldsymbol{\mu}}_{p}) + \sum_{i=1}^{n_{p}} (\hat{\Sigma}_{p}^{-1})_{ii} \sigma_{\phi}^{2} (\boldsymbol{y}_{\mathcal{I}_{pi}}) + \log |\hat{\Sigma}_{p}| \right. \\ \left. + \operatorname{tr} \left( \hat{\Sigma}_{p}^{-1} \tilde{K}_{X_{p}X_{p}}^{(A)} \right) + \operatorname{tr} \left( \left( K_{SS}^{(A)^{-1}} H K_{SS}^{(A)^{-1}} \right) \left( K_{SX_{p}}^{(A)} \hat{\Sigma}_{p}^{-1} K_{X_{pS}}^{(A)} \right) \right) - \sum_{i=1}^{n_{p}} \log \sigma_{\phi}^{2} (\boldsymbol{y}_{\mathcal{I}_{pi}}) \right) \\ \left. - \frac{N}{2} + D_{KL} (\mathcal{N}(\boldsymbol{m}, H) || \mathcal{N}(\mathbf{0}, K_{SS}^{(A)})), \tag{16}$$

is an unbiased estimate of the KL divergence upper bound  $E_{\mathcal{P}\sim\mathfrak{S}\{1,\ldots,P\}}(\hat{\mathcal{D}}_4) = \mathcal{D}_4 \geq D_{\mathrm{KL}}(\mathcal{N}(\bar{\boldsymbol{\mu}}, W)||\mathcal{N}(\mathbf{0}, \Sigma))$ . This property enables us to use the mini-batching technique for a more precise approximate computation of the KL divergence term of L-VAE and its gradients, by approximately splitting equal number of patients to each batch.

The learning of variational parameters m and H can be done either by explicitly parameterising and updating them within the overall optimisation scheme, or by using the natural gradients approach similar to Hensman et al. [2013]. Similar to the claims of the original paper, using the natural gradients can mitigate the challenge of finding a robust unconstrained parameterisation for the variational parameters. The gradients of KL divergence unbiased estimate  $\hat{\mathcal{D}}_4$  w.r.t. the variational parameters have the following form:

$$\frac{\partial \hat{\mathcal{D}}_4}{\partial \boldsymbol{m}} = -\sum_{p \in \mathcal{P}} K_{SS}^{(A)^{-1}} K_{SX_p}^{(A)} \hat{\Sigma}_p^{-1} \hat{\boldsymbol{\mu}}_p + \left( \sum_{p \in \mathcal{P}} K_{SS}^{(A)^{-1}} K_{SX_p}^{(A)} \hat{\Sigma}_p^{-1} K_{X_pS}^{(A)} K_{SS}^{(A)^{-1}} + K_{SS}^{(A)^{-1}} \right) \boldsymbol{m}$$
(17)

$$\frac{\partial \hat{D}_4}{\partial H} = -\frac{1}{2}H^{-1} + \frac{1}{2}\sum_{p\in\mathcal{P}} K_{SS}^{(A)^{-1}} K_{SX_p}^{(A)} \hat{\Sigma}_p^{-1} K_{X_pS}^{(A)} K_{SS}^{(A)^{-1}} + \frac{1}{2} K_{SS}^{(A)^{-1}}$$
(18)

Making use of the chain rule, we can write the following:

$$\frac{\partial \mathcal{D}_4}{\partial \boldsymbol{\eta}} = \frac{\partial \mathcal{D}_4}{\partial [m,H]} \frac{\partial [m,H]}{\partial \boldsymbol{\eta}} = \begin{bmatrix} \frac{\partial \mathcal{D}_4}{\partial m} & \frac{\partial \mathcal{D}_4}{\partial H} \end{bmatrix} \begin{bmatrix} \frac{\partial m}{\partial \eta_1} & \frac{\partial m}{\partial \eta_2} \\ \frac{\partial H}{\partial \eta_1} & \frac{\partial H}{\partial \eta_2} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{D}_4}{\partial m} & \frac{\partial \mathcal{D}_4}{\partial H} \end{bmatrix} \begin{bmatrix} I & 0 \\ -2m & I \end{bmatrix}$$

Then, using the update rule and following notation similar to Hensman et al. [2013], we get:

$$\theta_{2(t+1)} = \theta_{2(t)} - l \frac{\partial \hat{\mathcal{D}}_4}{\partial \eta_2}$$

$$\theta_{2(t+1)} = \theta_{2(t)} - l \left( \frac{\partial \hat{\mathcal{D}}_4}{\partial m} \frac{\partial m}{\partial \eta_2} + \frac{\partial \hat{\mathcal{D}}_4}{\partial H} \frac{\partial H}{\partial \eta_2} \right)$$

$$\theta_{2(t+1)} = \theta_{2(t)} - l \frac{\partial \hat{\mathcal{D}}_4}{\partial H}$$

$$-\frac{1}{2} H_{(t+1)}^{-1} = -\frac{1}{2} H_{(t)}^{-1} - l \frac{\partial \hat{\mathcal{D}}_4}{\partial H}$$

$$H_{(t+1)} = \left( H_{(t)}^{-1} + 2l \frac{\partial \hat{\mathcal{D}}_4}{\partial H} \right)^{-1}$$
(19)

and

$$\boldsymbol{\theta}_{1(t+1)} = \boldsymbol{\theta}_{1(t)} - l \frac{\partial \hat{\mathcal{D}}_4}{\partial \eta_1}$$
$$\boldsymbol{\theta}_{1(t+1)} = \boldsymbol{\theta}_{1(t)} - l \left( \frac{\partial \hat{\mathcal{D}}_4}{\partial \boldsymbol{m}} \frac{\partial \boldsymbol{m}}{\partial \eta_1} + \frac{\partial \hat{\mathcal{D}}_4}{\partial H} \frac{\partial H}{\partial \eta_1} \right)$$
$$\boldsymbol{\theta}_{1(t+1)} = \boldsymbol{\theta}_{1(t)} - l \left( \frac{\partial \hat{\mathcal{D}}_4}{\partial \boldsymbol{m}} - 2 \frac{\partial \hat{\mathcal{D}}_4}{\partial H} \boldsymbol{m} \right)$$
$$\boldsymbol{H}_{(t+1)}^{-1} \boldsymbol{m}_{(t+1)} = H_{(t)}^{-1} \boldsymbol{m}_{(t)} - l \left( \frac{\partial \hat{\mathcal{D}}_4}{\partial \boldsymbol{m}} - 2 \frac{\partial \hat{\mathcal{D}}_4}{\partial H} \boldsymbol{m}_{(t)} \right)$$
$$\boldsymbol{m}_{(t+1)} = H_{(t+1)} \left( H_{(t)}^{-1} \boldsymbol{m}_{(t)} - l \left( \frac{\partial \hat{\mathcal{D}}_4}{\partial \boldsymbol{m}} - 2 \frac{\partial \hat{\mathcal{D}}_4}{\partial H} \boldsymbol{m}_{(t)} \right) \right)$$

### 4 Predictive distribution

The problem of obtaining the predictive distribution for the high-dimensional, out-of-sample data  $\boldsymbol{y}_*$  given covariates  $\boldsymbol{x}_*$  with L-VAE can be split into two parts: obtaining the predictive distribution of the latent representation  $\boldsymbol{z}_*$  and propagating the obtained distribution through the probabilistic decoder  $p(\boldsymbol{y}_*|\boldsymbol{z}_*)$ . Given the training samples Y, covariate information X, the learnt parameters of the generative model  $\boldsymbol{\omega} = \{\psi, \theta\}$ , and the inference model  $\phi$ , the predictive distribution follows:

$$p_{\omega}(\boldsymbol{y}_{*}|\boldsymbol{x}_{*},Y,X) = \int_{\boldsymbol{z}_{*}} p_{\omega}(\boldsymbol{y}_{*}|\boldsymbol{z}_{*},\boldsymbol{x}_{*},Y,X) p_{\omega}(\boldsymbol{z}_{*}|\boldsymbol{x}_{*},Y,X) d\boldsymbol{z}_{*}$$

$$= \int_{\boldsymbol{z}_{*},Z} \underbrace{p_{\psi}(\boldsymbol{y}_{*}|\boldsymbol{z}_{*})}_{\text{decode GP prediction GP predictive posterior posterior of } Z$$

$$\approx \int_{\boldsymbol{z}_{*},Z} \underbrace{p_{\psi}(\boldsymbol{y}_{*}|\boldsymbol{z}_{*})}_{\text{decode GP prediction GP predictive posterior encode training samples}} \underbrace{p_{\phi}(Z|Y,X)}_{\boldsymbol{q}_{*}dZ} d\boldsymbol{z}_{*} d$$

The true posterior  $p_{\omega}(Z|Y, X)$  is intractable and, similar to the model inference or learning problem, it is replaced with the variational approximation defined by the inference model  $q_{\phi}(Z|Y, X)$ . Given such a substitution, the (approximate) predictive GP distribution for the latent representation is available in closed form,

$$\hat{p}_{\omega}(\boldsymbol{z}_{*}|\boldsymbol{x}_{*},Y,X) = \int_{Z} p_{\theta}(\boldsymbol{z}_{*}|\boldsymbol{x}_{*},Z,X) q_{\phi}(Z|Y,X) dZ = N(\boldsymbol{z}_{*}|\boldsymbol{\mu}_{*},\boldsymbol{\Sigma}_{*})$$

Since, in the current work, we only consider multi-output GPs with diagonal cross-covariance functions and the output of probabilistic encoder  $p_{\psi}(\boldsymbol{y}_n | \boldsymbol{z}_n)$  is restricted to be a multivariate normal distribution with diagonal covariance matrix, the predictive distribution in the latent space also factorises across the latent dimensions  $N(\boldsymbol{z}_* | \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) = \prod_{l=1}^L N(\boldsymbol{z}_{*l} | \boldsymbol{\mu}_{*l}, \sigma_{*l}^2)$  with:

$$\mu_{*l} = K_{\boldsymbol{x}_{*}X}^{(l)} \Sigma_{l}^{-1} \bar{\boldsymbol{\mu}}_{l}$$
  
$$\sigma_{*l}^{2} = k_{\boldsymbol{x}_{*}\boldsymbol{x}_{*}}^{(l)} - K_{\boldsymbol{x}_{*}X}^{(l)} \Sigma_{l}^{-1} K_{X\boldsymbol{x}_{*}}^{(l)} + K_{\boldsymbol{x}_{*}X}^{(l)} \Sigma_{l}^{-1} W_{l} \Sigma_{l}^{-1} K_{X\boldsymbol{x}_{*}}^{(l)} + \sigma_{zl}^{2}$$

where  $k_{\boldsymbol{x}*\boldsymbol{x}*}^{(l)} = K_{\boldsymbol{x}*\boldsymbol{x}*}^{(l)} = \sum_{r=1}^{R} K_{\boldsymbol{x}*\boldsymbol{x}*}^{(l,r)}(\theta_{l}^{(r)})$  and  $K_{\boldsymbol{x}*\boldsymbol{X}}^{(l)} = K_{\boldsymbol{X}\boldsymbol{x}*}^{(l)T} = \sum_{r=1}^{R} K_{\boldsymbol{x}*\boldsymbol{X}}^{(l,r)}(\theta_{l}^{(r)})$  (see eq. (4) in the main text), and  $\Sigma_{l} = \sum_{r=1}^{R} K_{\boldsymbol{X}\boldsymbol{X}}^{(l,r)}(\theta_{l}^{(r)}) + \sigma_{zl}^{2}I$ ,  $W_{l} = \text{diag}(\sigma_{\phi,l}^{2}(\boldsymbol{y}_{1}), \dots, \sigma_{\phi,l}^{2}(\boldsymbol{y}_{N}))$  and  $\bar{\boldsymbol{\mu}}_{l} = [\mu_{\phi,l}(\boldsymbol{y}_{1}), \dots, \mu_{\phi,l}(\boldsymbol{y}_{N})]^{T}$  (as in eq. (8) in the main text). Incorporating this property with eq. (21) leads to,

$$p_{\omega}(\boldsymbol{y}_{*}|\boldsymbol{x}_{*},Y,X) \approx \int_{\boldsymbol{z}_{*}} \prod_{d=1}^{D} \mathcal{N}\left(y_{*d}|g_{\psi,d}(\boldsymbol{z}_{*}),\sigma_{yd}^{2}\right) \prod_{l=1}^{L} \mathcal{N}(z_{*l}|\mu_{*l},\sigma_{*l}^{2}) d\boldsymbol{z}_{*}.$$

#### 5 Scalable predictive distribution

Computing the predictive distribution, as described above, requires performing cubic operations over the  $N \times N$  matrices, which makes it practically infeasible for problems with large training data. In this section, we exploit the same inducing points paradigm that was used for  $\mathcal{D}_2$  in eq. (7) to alleviate the cubic complexity. For simplicity of notation, we omit the latent dimension index l in the rest of the section. We use the same notation as in eq. (6) and set  $U = K_{SS}^{(A)} + K_{SX}^{(A)} \hat{\Sigma}^{-1} K_{XS}^{(A)}$ . Let  $X_* = [\mathbf{x}_{1*}, \ldots, \mathbf{x}_{N'*}]$  denote a collection of N' test data points from P' many subjects. Then,

$$\bar{\boldsymbol{\mu}}_{*} = K_{X_{*}X} \Sigma^{-1} \bar{\boldsymbol{\mu}} = \left( K_{X_{*}S}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + K_{X_{*}X}^{(R)} \right) \left( K_{XS}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + \hat{\Sigma} \right)^{-1} \bar{\boldsymbol{\mu}} = \left( K_{X_{*}S}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} + K_{X_{*}X}^{(R)} \right) \left( \hat{\Sigma}^{-1} \bar{\boldsymbol{\mu}} - \hat{\Sigma}^{-1} K_{XS}^{(A)} U^{-1} K_{SX}^{(A)} \hat{\Sigma}^{-1} \bar{\boldsymbol{\mu}} \right)$$
(22)

In practice, the efficient computations are executed in the following order:

- 1. Compute  $\hat{\Sigma}^{-1}\bar{\mu}$ ; exploiting the block-diagonal property of  $\hat{\Sigma}^{-1}$ .
- 2. Compute  $K_{XS}^{(A)}U^{-1}K_{SX}^{(A)}(\hat{\Sigma}^{-1}\bar{\mu})$  using the low-rank properties of all matrices involved.
- 3. Compute  $\tilde{\mu} = \hat{\Sigma}^{-1} \bar{\mu} \hat{\Sigma}^{-1} K_{XS}^{(A)} U^{-1} K_{SX}^{(A)} \hat{\Sigma}^{-1} \bar{\mu}$ .
- 4. Compute  $K_{X*S}^{(A)}K_{SS}^{(A)^{-1}}K_{SX}^{(A)}\tilde{\mu}$  using the low-rank properties of all matrices involved.

- 5. Compute  $K_{X_*X}^{(R)}\tilde{\mu}$ . This relies on the fact that the  $K_{X_*X}^{(R)}$  matrix will be very sparse, which can be exploited either via sparse matrix operations or simply by cycling over the sets of rows that correspond to each subject and performing the multiplication only for the non-zero part of  $K_{X_*X}^{(R)}$  (if none, this product will be zero for the given subject).
- 6. Finally, sum up  $\bar{\mu}_* = K_{X_*S}^{(A)} K_{SS}^{(A)^{-1}} K_{SX}^{(A)} \tilde{\mu} + K_{X_*X}^{(R)} \tilde{\mu}$

In practice, the predictive distribution will be subject to the inducing point locations S and is guaranteed to be exact if and only if u is a sufficient statistic of  $f_A$ .

The predictive mean can also be expressed in terms of the variational parameters:

$$K_{X_*S}^{(A)} K_{SS}^{(A)^{-1}} \boldsymbol{m} + K_{X_*X}^{(R)} \hat{\Sigma}^{-1} (\bar{\boldsymbol{\mu}} - K_{XS}^{(A)} K_{SS}^{(A)^{-1}} \boldsymbol{m})$$

Also, the predictive covariance can be expressed as follows:

$$\left(K_{X_*S}^{(A)} - K_{X_*X}^{(R)}\hat{\Sigma}^{-1}K_{XS}^{(A)}\right)K_{SS}^{(A)^{-1}}HK_{SS}^{(A)^{-1}}\left(K_{X_*S}^{(A)} - K_{X_*X}^{(R)}\hat{\Sigma}^{-1}K_{XS}^{(A)}\right)^T + \sigma_z^2 I_{N'} + K_{X_*X_*}^{(R)} - K_{X_*X}^{(R)}\hat{\Sigma}^{-1}K_{XX_*}^{(R)}$$

#### 6 Downstream classification task in the healthcare data experiment

Given the test set, our objective is to predict the patient *mortality*. Since our proposed model is generative in nature, to allow such an objective, we shall include the *mortality* covariate in the additive GP prior. The downstream classification task is done by computing the probability that *mortality* = 0 and *mortality* = 1. Given the learnt L-VAE model with fixed parameters  $\phi, \psi, \theta$ , we can obtain the probability for the binary mortality event for each patient  $(X_*, Y_*)$  in the test set as,

$$P(\text{mortality} = 0) = \frac{\exp(\mathcal{L}(\phi, \psi, \theta; Y_*, X_*, \text{mortality} = 0))}{\sum_{i=0}^{1} \exp(\mathcal{L}(\phi, \psi, \theta; Y_*, X_*, \text{mortality} = i))},$$
  
$$P(\text{mortality} = 1) = 1 - P(\text{mortality} = 0).$$

Furthermore, we also introduce a time to mortality covariate (or *mortalityTime*) that is based on the survival time. This covariate is relevant for individuals whose *mortality* = 1 and is treated as missing for the individuals that survive. However, in the testing phase, the exact mortality time is not known as our objective is to perform classification based on *mortality*. To overcome this problem, we estimate the distribution of the *mortalityTime* covariate based on the values of the covariate in the training set and, in the test, we compute the expectation of the ELBO described in eq. (5) of the main manuscript w.r.t. to that distribution.



Figure 1: Histogram of survival times in training and test set

We approximate the expectation with a finite weighted average. Concretely, we first allocate the *mortalityTime* covariate values in the training set into B bins based on a logarithmic scale. Let  $\alpha_i$  be the bin count and  $t_i$  be the average value of the *mortalityTime* in the *i*<sup>th</sup> bin. The proportion of *mortalityTime* values in bin *i* is  $w_i = \alpha_i / (\sum_{i=1}^{B} \alpha_i)$ . The weighted ELBO is then computed as:

$$\mathcal{L}(\phi, \psi, \theta; Y, X, \text{mortality} = 1) = \sum_{i=1}^{B} w_i \cdot \mathcal{L}(\phi, \psi, \theta; Y, X, \text{mortality} = 1, t_i),$$

where the ELBO terms are now explicitly conditioned by the *mortalityTime*,  $t_i$ . We use B = 6 in our analysis. The histograms of the survival times in the training and test data are shown in Fig. 1.

We follow the data preprocessing steps described in Luo et al. [2018], and standardise the measurements of the 35 different attributes.

### 7 Optimisation and practical considerations

We make use of a suitable stochastic optimisation technique to minimise the ELBO in eq. (4) of the main manuscript. The parameters that we need to optimise include the neural network weights ( $\phi$  and  $\psi$ ) and kernel parameters ( $\theta$ ) of the multi-output additive GPs. In particular, the optimisation is done using the Adam optimiser [Kingma and Ba, 2015], which is an adaptive learning rate method that maintains an exponentially decaying average of past gradients as well as squared gradients. In case of mini-batch training, the Adam steps are conducted interchangeably with natural gradient-based updates of the variational parameters. For the inference implementation, we make use of PyTorch [Paszke et al., 2019] which allows the computation of derivatives using automatic differentiation.

In all the experiments, we first pre-train the neural networks with a standard normal distribution as the prior on the latent space (standard VAE [Kingma and Welling, 2014]) for 1000 epochs. This is followed by training the L-VAE model using the pre-trained encoder and decoder networks as initial values for  $\phi$  and  $\psi$ , respectively. While training the L-VAE model, we monitor the loss on the validation (independent) dataset as a performance metric. Similar to the strategy of early stopping, we save the weights of the model that has the best performance on our defined metric. These model weights are chosen to perform predictions and other downstream tasks. However, to handle the possibility of a local minimum, we do not specify a stopping criterion and continue the training procedure till a predefined number of epochs has been performed. In the Rotated MNIST and Health MNIST experiments, L-VAE is trained for a maximum of 2000 epochs. Moreover, in the Healthcare data experiment, L-VAE is trained for a maximum of 1000 epochs.

#### 8 Supplementary tables

Table 1 describes the neural network architecture used for the Rotated MNIST experiment. The hyperparameter choices are similar to Casale et al. [2018]. The architecture used for the Health MNIST experiment is described in table 2. We have tried to replicate the hyperparameter choices from Fortuin et al. [2020] for this experiment. For the Physionet Challenge 2012 dataset, we did not make use of a convolutional neural network (CNN) as was done for the Rotated MNIST as well as Health MNIST experiments because CNNs are more appropriate for image based (visual) data where the regional correlation (receptive field) of the measured values is important [Goodfellow et al., 2016]. Table 3 describes the architecture for the multi layered perceptron (MLP) that we used for the Physionet Challenge 2012 dataset. It is similar to the architecture used in Fortuin et al. [2020].

	Hyperparameter	Value
Inference network	Dimensionality of input	$28 \times 28$
	Number of convolution layers	3
	Number of filters per convolution layer	72
	Kernel size	$3 \times 3$
	Stride	2
	Number of feedforward layers	1
	Width of feedforward layers	128
	Dimensionality of latent space	L
	Activation function of layers	ELU
	Dimensionality of input	L
	Number of transposed convolution layers	3
	Number of filters per transposed convolution layer	72
Generative	Kernel size	3  imes 3
network	Stride	1
	Number of feedforward layers	1
	Width of feedforward layers	128
	Activation function of layers	ELU

Table 1: Neural network architectures used in the Rotated MNIST dataset.

	Hyperparameter	Value
Inference network	Dimensionality of input	$36 \times 36$
	Number of convolution layers	2
	Number of filters per convolution layer	144
	Kernel size	3 imes 3
	Stride	2
	Pooling	Max pooling
	Pooling kernel size	$2 \times 2$
	Pooling stride	2
	Number of feedforward layers	2
	Width of feedforward layers	300, 30
	Dimensionality of latent space	L
	Activation function of layers	RELU
Generative network	Dimensionality of input	L
	Number of transposed convolution layers	2
	Number of filters per transposed convolution layer	256
	Kernel size	$4 \times 4$
	Stride	2
	Number of feedforward layers	2
	Width of feedforward layers	30,  300
	Activation function of layers	RELU

Table 2: Neural network architectures used in the Health MNIST dataset.

	Hyperparameter	Value
Inference network	Dimensionality of input	35
	Number of feedforward layers	2
	Number of elements in each feedforward layer	128,64
	Dimensionality of latent space	L
	Activation function of layers	RELU
Generative network	Dimensionality of input	L
	Number of feedforward layers	2
	Number of elements in each feedforward layer	64, 128
	Activation function of layers	RELU

Table 3: Neural network architectures used in the Physionet Challenge 2012 dataset.

## 9 Supplementary images



Figure 2: Plate diagram of the model. The shaded circle refers to an observed variable, the partially shaded circle refers to a partially observed variable (due to missing values), and the un-shaded circle refers to an un-observed variable. (a) Represents the inference (or encoder) model and (b) Represents the generative (or decoder) model.



Figure 3: Comparison of the resulting latent space using VAE, PCA, and L-VAE on the Health MNIST dataset. The L-VAE model is fit using  $\boldsymbol{f}_{ca}(id) + \boldsymbol{f}_{se}(age) + \boldsymbol{f}_{ca\times se}(id \times age) + \boldsymbol{f}_{ca\times se}(sex \times age) + \boldsymbol{f}_{ca\times se}(diseasePresence \times diseaseAge)$  as the multi-output additive GP prior. The number of latent dimensions is set to 2. The points are coloured according to the diseaseAge as shown in the colour bar.



Figure 4: GP model fittings of L-VAE in the latent space with dimension 2 on the Health MNIST dataset.



Figure 5: AUROC scores for the patient mortality prediction task on the test set of the Physionet Challenge 2012 dataset. This is an extension to Fig. 4 in the main manuscript. In this figure, we can observe L-VAE's performance with different latent dimensions as well as all the patient-specific general auxiliary covariates. Higher AUROC score is better.

### References

- F. P. Casale, A. V. Dalca, L. Saglietti, J. Listgarten, and N. Fusi. Gaussian process prior variational autoencoders. In Advances in Neural Information Processing Systems, NeurIPS, 2018.
- L. Cheng, S. Ramchandran, T. Vatanen, N. Lietzén, R. Lahesmaa, A. Vehtari, and H. Lähdesmäki. An additive gaussian process regression model for interpretable non-parametric analysis of longitudinal data. *Nature Communications*, 2019.
- V. Fortuin, D. Baranchuk, G. Rätsch, and S. Mandt. GP-VAE: deep probabilistic time series imputation. In The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS. PMLR, 2020.

- I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016.
- J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence, UAI*. AUAI Press, 2013.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR, 2015.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. In 2nd International Conference on Learning Representations, ICLR, 2014.
- Y. Luo, X. Cai, Y. Zhang, J. Xu, and X. Yuan. Multivariate time series imputation with generative adversarial networks. In Advances in Neural Information Processing Systems, NeurIPS, 2018.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems, NeurIPS, 2019.
- W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery. Numerical Recipes 3rd Edition: The Art of Scientific Computing. Cambridge University Press, USA, 3 edition, 2007.
- C. E. Rasmussen and C. K. I. Williams. Gaussian processes for machine learning. MIT Press, 2006.
- M. K. Titsias. Variational learning of inducing variables in sparse gaussian processes. In Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics, AISTATS, 2009.