

---

# Bayesian Sequential Batch Design in Functional Data

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many longitudinal studies are hindered by noisy observations sampled at irregular  
2 and sparse time points. In handling such data and optimizing the design of a study,  
3 most of the existing functional data analysis focuses on the frequentist approach  
4 that bears the uncertainty of model parameter estimation. While the Bayesian  
5 approach as an alternative takes into account the uncertainty, little attention has  
6 been given to sequential batch designs that enable information update and cost  
7 efficiency. To fill the gap, we propose a Bayesian hierarchical model with Gaussian  
8 processes which allows us to propose a new form of the utility function based on  
9 the Shannon information between posterior predictive distributions. The proposed  
10 procedure sequentially identifies optimal designs for new subject batches, opening  
11 a new way for incorporating the Bayesian approach in finding the optimal design  
12 and enhancing model estimation and the quality of analysis with sparse data.

## 13 1 Introduction

14 Many of the longitudinal studies suffer from noisy observations. It is often the case that only a small  
15 number of irregularly spaced observations can be taken for each subject, making it a sparse dataset  
16 for the subsequent analysis (Zeger and Diggle, 1994; Brumback and Rice, 1998; Guo, 2004; Yao  
17 et al., 2005). In light of this issue, functional data analysis (FDA) has been developed as one of the  
18 most popular methods to handle such data and enhance the quality of estimation. In particular, as the  
19 sparse observations can only provide limited information for recovering the underlying trajectory,  
20 FDA offers an effective way to optimize the design of a study by judiciously selecting optimal time  
21 points for taking observations.

22 Existing FDA literature has mostly focused on rather a frequentist approach that considers the “best  
23 guess” of parameters to find an optimal design (Ji and Müller, 2017; Park et al., 2018; Rha et al.,  
24 2020). However, this approach oftentimes bears uncertainty of the model parameter estimation and  
25 can possibly hinder the quality of analysis. A Bayesian approach, on the other hand, takes into  
26 account this uncertainty and conducts the analysis based on a prior distribution of the parameters  
27 (Chaloner and Verdinelli, 1995). Specifically, a Bayesian hierarchical model assumes a common mean  
28 function for the underlying subject trajectories, enabling us to borrow the strength of all observations  
29 across subjects to recover the trajectories. (Yang et al., 2016)

30 Ryan et al. (2015) proposed the fully Bayesian static design for mixed effect model to determine  
31 sampling time points for precise estimation of the model parameters. Nevertheless, the static design  
32 uses the same design throughout the experimental process without accounting for any incoming  
33 information that may be collected during the experiment (Ryan et al., 2016). In this regard, a  
34 sequential design may offer more efficient and flexible design schemes as it updates the optimal  
35 design at each stage with new information provided from the previous stages. (Chaloner, 1986;  
36 Müller et al., 2007). Yet scant work has been done on constructing Bayesian hierarchical models

37 with a sequential design that considers the uncertainty of model parameter estimation and updates the  
 38 optimal design with newly acquired information at each stage.

39 To fill this gap, in this study, we propose a Bayesian hierarchical model that sequentially identifies  
 40 optimal designs for new batches of subjects by (1) providing information for updating the posterior  
 41 mean function of the underlying trajectories of existing subjects and (2) offering sufficient information  
 42 for accurate estimation of new subject trajectories. Particularly, we first obtain the posterior distribu-  
 43 tions of underlying trajectories from our Bayesian hierarchical model, and update the distributions  
 44 with new observations. Then based on the posterior distributions, we find the optimal design by the  
 45 simulated annealing (SA) algorithm proposed by Van Laarhoven et al. (1987), which is widely known  
 46 for its strengths in search in large space and computational efficiency.

47 In sum, our study is expected to open a new way for incorporating the Bayesian approach in handling  
 48 noisy observations with sequential batch designs and further enhance model estimations with new  
 49 information update. The rest of paper is organized as follows: Section 2 introduces our Bayesian  
 50 hierarchical model that is used to obtain the posterior distributions of underlying trajectories. Section  
 51 3 formulates the utility function as the design criterion for finding the optimal design. Section 4  
 52 details the implementation of the simulated annealing algorithm on the search of the optimal design.  
 53 A discussion can be found in Section 5.

## 54 2 Bayesian Hierarchical Model

55 In longitudinal studies, it is not uncommon to have observations that are sampled at sparse and  
 56 irregular time points. The collected samples are viewed as functional observations and are often  
 57 contaminated with unknown noises. Assuming each subject following their independent stochastic  
 58 process, we consider the Bayesian hierarchical model proposed by Yang et al. (2016) as follows:

$$\begin{aligned} \mathbf{Y}_i(\mathbf{t}_i) &= \mathbf{X}_i(\mathbf{t}_i) + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \stackrel{i.i.d.}{\sim} N(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}), \\ \mathbf{X}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma} &\stackrel{i.i.d.}{\sim} GP(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad i = 1, \dots, n, \\ \boldsymbol{\mu} &\sim GP\left(\boldsymbol{\mu}_0, \frac{1}{c} \boldsymbol{\Sigma}\right), \end{aligned}$$

59 where  $\mathbf{Y}_i(\mathbf{t}_i) = \{Y_i(t_{i,1}), \dots, Y_i(t_{i,n_i})\}$  are the noisy observations of the underlying trajectory  $\mathbf{X}_i$   
 60 at time  $\mathbf{t}_i = (t_{i,1}, \dots, t_{i,n_i})'$ . We consider the additive error vector  $\boldsymbol{\epsilon}_i$  that follows i.i.d. normal  
 61 with mean vector  $\mathbf{0}$  and variance  $\sigma_\epsilon^2 \mathbf{I}$  and is independent of  $\mathbf{X}_i$ . We assume each  $\mathbf{X}_i$  follows i.i.d.  
 62 Gaussian process with a prespecified mean function  $\boldsymbol{\mu}$  and covariance kernel  $\boldsymbol{\Sigma}$ . The universal mean  
 63 function  $\boldsymbol{\mu}$  is assumed unknown and is assigned with a Gaussian process as  $\boldsymbol{\mu} \sim GP(\boldsymbol{\mu}_0, (1/c)\boldsymbol{\Sigma})$   
 64 with the mean function  $\boldsymbol{\mu}_0$  and the covariance kernel  $\boldsymbol{\Sigma}$  scaled by some  $c > 0$ . For simplicity, we  
 65 denote  $\mathbf{Y}_i(\mathbf{t}_i)$  by  $\mathbf{Y}_{i,t_i}$ ,  $\mathbf{X}_i(\mathbf{t}_i)$  by  $\mathbf{X}_{i,t_i}$ ,  $\boldsymbol{\mu}(\mathbf{t}_i)$  by  $\boldsymbol{\mu}_{t_i}$ , and  $\boldsymbol{\Sigma}(\mathbf{t}_i, \mathbf{t}_i)$  by  $\boldsymbol{\Sigma}_{t_i,t_i}$ . Given time grid  $\{\mathbf{t}_i\}$ ,  
 66 we have the following hierarchical structure in multivariate forms for subject  $i$ :

$$\begin{aligned} \mathbf{Y}_{i,t_i} | \mathbf{X}_{i,t_i} &\sim MVN(\mathbf{X}_{i,t_i}, \sigma_\epsilon^2 \mathbf{I}), \\ \mathbf{X}_{i,t_i} | \boldsymbol{\mu}_{t_i}, \boldsymbol{\Sigma}_{t_i,t_i} &\sim MVN(\boldsymbol{\mu}_{t_i}, \boldsymbol{\Sigma}_{t_i,t_i}), \\ \boldsymbol{\mu}_{t_i} | \boldsymbol{\mu}_0, \boldsymbol{\Sigma} &\sim MVN(\boldsymbol{\mu}_{0t_i}, \frac{1}{c} \boldsymbol{\Sigma}_{t_i,t_i}). \end{aligned} \quad (1)$$

67 For simplicity, we assume that the error variance is fixed and the covariance kernel to follow a  
 68 pre-specified structure as squared exponential kernel. The scaling constant  $c$  for the covariance kernel  
 69 of the mean function is set to 1 and thus does not require posterior update in the estimation step. For  
 70 the hyperparameter  $\boldsymbol{\mu}_0$ , We set it to be the smoothed sample mean of  $\{\mathbf{Y}_{i,t_i}\}$ .

71 Different from previous approaches in functional data analysis that mainly focus on smoothing each  
 72 curve individually, the hierarchical GP model borrows the strength of all observations and smooth  
 73 the entire functional observations at once by assuming a common mean function  $\boldsymbol{\mu}$  (Yang et al.,  
 74 2016). In addition, two layers of GPs with the same covariance kernel function provide important  
 75 insights and computational efficiency to our design problem. Assigning a GP on  $\boldsymbol{\mu}$  allows the model  
 76 to share information across the subjects and to predict the trajectories at unobserved time grids for  
 77 all of the subjects based on the collected observations of only a portion of subjects. Besides, the  
 78 hierarchical structure of GPs still gives us a closed form of the predictive distribution which reduces  
 79 the computational cost in evaluating the optimal design criterion significantly. In the next section, we  
 80 will detail the design problem and propose a utility function for the corresponding optimal design.

### 81 3 Utility Function and Optimal Sequential Batch Design

82 Conventional sequential design approach adopts one-step-look-ahead method that only considers the  
 83 next subject, which is often not optimal. Static design approach determines the optimal design in  
 84 a holistic view but uses the same fixed protocol throughout the experiments. To combine the best  
 85 of two worlds, we adopt a sequential batch scheme. We consider the problem of multistage design  
 86 that sequentially finds optimal sampling times for a new batch of subjects based on the information  
 87 obtained from observations of existing subjects from previous stages. For demonstration purposes,  
 88 we only display the utility function for one future stage. However, by including new observations  
 89 with the obtained optimal design at the current stage, one is able to update the optimal design criterion  
 90 and acquire new optimal designs for all future stages in a sequential manner.

91 For the experiments, we assume that observations can be taken on an equally-spaced common grid that  
 92 has  $T_0$  time points. Yet, for each subject, only  $k (< T_0)$  observations can be taken. Before stage 1, we  
 93 assume that an experiment is already conducted and observations  $\mathbf{Y}_0 = \{\mathbf{Y}_1(t_i), \dots, \mathbf{Y}_N(t_i)\}$   
 94 for  $N$  subjects are taken based on a fixed design  $\mathbf{D}_0 = \{t_1, \dots, t_N\}$ . Suppose we are now  
 95 at stage 1 and we are to recruit a new batch of  $M (> 1)$  subjects and take observations  $\mathbf{Y}_1 =$   
 96  $\{\mathbf{Y}_{N+1}(t_i), \dots, \mathbf{Y}_{N+M}(t_i)\}$  from these subjects according to a design  $\mathbf{D}_1 = \{t_{N+1}, \dots, t_{N+M}\}$ .  
 97 Here, we consider the batch size  $M$  and the number of observations per subject  $k$  to be fixed. Our  
 98 attempt is to find the optimal design  $\mathbf{D}_1$  that achieves two goals: (1) the newly-added observations  
 99 based on  $\mathbf{D}_1$  should provide more information to update the posterior mean function so as to improve  
 100 the recovery of underlying trajectories  $\mathbf{X}_0$  for the existing subjects  $1, \dots, N$ ; (2) the observations  
 101 based on  $\mathbf{D}_1$  should also provide sufficient information for the estimation of new batch of subject  
 102 trajectories.

103 Specifically, when recovering the trajectories of existing and new subjects, we focus on the trajectory  
 104 values at unobserved time points, denoted by  $\mathbf{X}^c$ . We would like to compare the posterior predictive  
 105 distributions  $p(\mathbf{X}_0^c, \mathbf{X}_1^c | \mathbf{Y}_0)$  of  $\mathbf{X}_0^c$  and  $\mathbf{X}_1^c$  given the information from existing subjects to the  
 106 posterior predictive distributions  $p(\mathbf{X}_0^c, \mathbf{X}_1^c | \mathbf{Y}_0, \mathbf{Y}_{1, \mathbf{D}_1})$  of  $\mathbf{X}_0^c$  and  $\mathbf{X}_1^c$  given the information from  
 107 existing subjects and the new batch of subjects. That is, we would like to maximize the improvement  
 108 in prediction of  $\mathbf{X}^c$  before and after including the new batch of subjects.

109 We consider an information-based approach and measure the improvement by Kullback-Leibler (KL)  
 110 divergence, which is a classic metric in information theory that measures the difference between two  
 111 distributions. Therefore, we propose the following utility function as the optimal design criterion:

$$U(\mathbf{D}_1, \mathbf{Y}_0) = D_{KL}(p_1 || p_0) = \int \log \left( \frac{p_1}{p_0} \right) dp_1, \quad (2)$$

112 where we denote by  $p_0 = p(\mathbf{X}_0^c, \mathbf{X}_1^c | \mathbf{Y}_0)$  and  $p_1 = p(\mathbf{X}_0^c, \mathbf{X}_1^c | \mathbf{Y}_0, \mathbf{Y}_{1, \mathbf{D}_1})$ , which are both multivari-  
 113 ate normal distributions under our model framework.

114 To evaluate the above utility function, we consider a combination of implementing the predictive  
 115 formula of Gaussian process and using empirical Bayes procedure for the rest of model parameters  
 116 to obtain a closed-form solution for the utility function. Concerning the page limit, we refer the  
 117 readers to Appendix A for the detailed derivation. This closed-form solution facilitates computational  
 118 efficiency by avoiding the evaluation of intractable marginal likelihood in the utility function as  
 119 commonly seen in many optimal Bayesian design problems.

### 120 4 Computation

121 Because of the closed-form solution of the utility function in Section 3, it is easy to evaluate the  
 122 utility function with a given design. Yet, the design space remains large as we are exploring optimal  
 123 designs for a batch of subjects simultaneously. Therefore, we implement a simulated annealing (SA)  
 124 algorithm (Van Laarhoven et al., 1987) that enables efficient exploration of large and complex design  
 125 spaces and easy implementation. Specifically, the SA algorithm is used at every stage such that it  
 126 incorporates existing and new information from all previous and current stages and finds optimal  
 127 design for the next stage in a sequential manner.

128 The SA algorithm starts with an initial “temperature”  $T_{initial}$  and a randomly generated design  
 129  $\mathbf{D}_{initial}$ . The “energy”  $e$  of this design is then computed based on the utility function defined in  
 130 Equation (2). Then the algorithm generates another candidate design  $\mathbf{D}_{test}$  from the “neighborhood”

131 of  $D_{initial}$  and calculates its energy  $e_{test}$ . If the difference between two energies  $\Delta e = e - e_{test} \leq 0$ ,  
132 the candidate design  $D_{test}$  is accepted and the algorithm will continue to compare it to other  
133 neighborhood designs. At the current temperature  $T$ , if  $\Delta e > 0$ , the candidate design is accepted  
134 with a probability of  $\exp(\Delta e/T)$ . This process is repeated until no further improvements can be  
135 made within a maximum number of iterations. Then the temperature will be lowered according to a  
136 “cooling schedule” and the whole procedure will be repeated again. Finally, we follow the approach  
137 proposed by Aragon et al. (1991) to terminate the algorithm if the acceptance probability is smaller  
138 than some threshold  $P_{threshold}$ .

139 In the algorithm, a number of parameters, initial temperature, cooling schedule, neighborhood of a  
140 design, maximum number of iterations, and acceptance threshold, require initial values. Nevertheless,  
141 as the SA algorithm is a heuristic algorithm, the parameter values heavily depend on the problem  
142 settings and experiment setup. Therefore, we also set the parameter values in a heuristic way so as  
143 to be able to adapt to different scenarios. Based on suggestions in Van Laarhoven et al. (1987), we  
144 set the initial temperature  $T_{initial}$  to be  $\Delta e / \log(0.7)$  so that the initial acceptance probability for  
145 designs with  $\Delta e > 0$  is 0.6. This is to limit the time spent at high temperatures. The cooling schedule  
146 is an exponential decaying function of the temperature  $T_{new} = 0.95 \times T_{old}$ .

147 For the neighborhood of a design, there are many choices, such as changing only one time point for  
148 one subject in the batch or changing one set of time points for one subject in the batch. However, the  
149 candidate set for the former can easily increase exponentially with different time grid and observation  
150 sizes and it is also suspected that a single time point can make much difference on the trajectory  
151 recovery of all subjects. Thus, considering computation efficiency, we define the neighborhood of  
152 a design by changing one set time points from one subject in the batch. Here we propose to set  
153 the maximum number of iterations to be 10 and the acceptance threshold to be 0.2, as suggested  
154 in Aragon et al. (1991). As noted before, since the SA algorithm is a heuristic approach that is  
155 contingent upon a specific problem, empirical tuning on the initial parameters is necessary when  
156 conducting different experiments. A pseudo code that illustrates the structure of the algorithm can be  
157 found in Appendix B.

## 158 5 Discussion

159 To handle the noisy observations in many fields such as longitudinal studies, extant FDA literature  
160 mostly adopts rather a frequentist approach and bears the uncertainty of parameter estimation. As  
161 an alternative to improve the quality of model estimation, a Bayesian approach naturally takes into  
162 account the uncertainty in estimation and produces posterior predictive distribution. In this study, we  
163 adopt a Bayesian hierarchical model of Gaussian processes for the underlying trajectories, which  
164 enables us to obtain the trajectory predictive distributions with closed-form expressions at reduced  
165 computational cost. We propose an optimal Bayesian sequential batch design scheme that sequentially  
166 finds optimal design for a batch of subjects based on the information obtained from all previous  
167 and current stages. Specifically, its sequential feature helps update the optimal design criterion with  
168 new information at each stage, whereas its batch feature controls for a small number of stages and  
169 maintains the overall cost effectiveness. Combining these two features, this scheme is designed to  
170 improve the trajectory recovery of current subjects and achieve accurate estimation of future subject  
171 trajectories. Finally, in the optimization step, we implement a simulated annealing algorithm that  
172 takes in empirically-tuned parameters and outputs a final design with computational efficiency.

173 Further refinement of this study can be done by altering the assumptions made in our analysis.  
174 Particularly in the design setup, we assume that the batch size  $M$  of the optimal design is small. This  
175 is established as  $M$  should not be too large to only have too few updates on the design optimality  
176 criterion. Nonetheless, in practice,  $M$  is often contingent upon the size of the initial data set and  
177 the number of design stages. The interactions between these factors may change the optimal size  
178 of the batch. To account for this, there are two potential approaches to find the optimal  $M$ . One  
179 is to iteratively test different values of  $M$  from 1 to the existing subject size  $N$ . Yet additional  
180 consideration will need to be put in to reduce its computational expensiveness. Another is to include  
181  $M$  as a random variable and incorporate it inside the utility function. That is, the optimal design and  
182 the optimal batch size are obtained in each stage.

183 **References**

- 184 Aragon, C. R., Johnson, D., McGeoch, L., and Schevon, C. (1991). Optimization by simulated  
185 annealing: An experimental evaluation; part ii, graph coloring and number partitioning. *Operations*  
186 *research*, 39(3):378–406.
- 187 Brumback, B. A. and Rice, J. A. (1998). Smoothing spline models for the analysis of nested and  
188 crossed samples of curves. *Journal of the American Statistical Association*, 93(443):961–976.
- 189 Chaloner, K. (1986). Optimal bayesian design for non-linear estimation.
- 190 Chaloner, K. and Verdinelli, I. (1995). Bayesian experimental design: A review. *Statistical science*,  
191 pages 273–304.
- 192 Guo, W. (2004). Functional data analysis in longitudinal settings using smoothing splines. *Statistical*  
193 *methods in medical research*, 13(1):49–62.
- 194 Ji, H. and Müller, H.-G. (2017). Optimal designs for longitudinal and functional data. *Journal of the*  
195 *Royal Statistical Society Series B: Statistical Methodology*, 79(3):859–876.
- 196 Müller, P., Berry, D. A., Grieve, A. P., Smith, M., and Krams, M. (2007). Simulation-based sequential  
197 bayesian design. *Journal of statistical planning and inference*, 137(10):3140–3150.
- 198 Park, S. Y., Xiao, L., Willbur, J. D., Staicu, A.-M., and Jumbe, N. (2018). A joint design for functional  
199 data with application to scheduling ultrasound scans. *Computational Statistics & Data Analysis*,  
200 122:101–114.
- 201 Rha, H., Kao, M.-H., and Pan, R. (2020). Design optimal sampling plans for functional regression  
202 models. *Computational Statistics & Data Analysis*, 146:106925.
- 203 Ryan, E. G., Drovandi, C. C., McGree, J. M., and Pettitt, A. N. (2016). A review of modern  
204 computational algorithms for bayesian optimal design. *International Statistical Review*, 84(1):128–  
205 154.
- 206 Ryan, E. G., Drovandi, C. C., and Pettitt, A. N. (2015). Simulation-based fully bayesian experimental  
207 design for mixed effects models. *Computational Statistics & Data Analysis*, 92:26–39.
- 208 Van Laarhoven, P. J., Aarts, E. H., van Laarhoven, P. J., and Aarts, E. H. (1987). *Simulated annealing*.  
209 Springer.
- 210 Yang, J., Zhu, H., Choi, T., and Cox, D. D. (2016). Smoothing and mean–covariance estimation of  
211 functional data with a bayesian hierarchical model. *Bayesian Analysis*, 11(3):649.
- 212 Yao, F., Müller, H.-G., and Wang, J.-L. (2005). Functional data analysis for sparse longitudinal data.  
213 *Journal of the American statistical association*, 100(470):577–590.
- 214 Zeger, S. L. and Diggle, P. J. (1994). Semiparametric models for longitudinal data with application to  
215 cd4 cell numbers in hiv seroconverters. *Biometrics*, pages 689–699.

216 **A Derivation of Utility Function**

217 Let  $A$  be  $(\mathbf{X}_0^c, \mathbf{X}_1^c) | \mathbf{Y}_0$  with distribution  $p_0$  and let  $B$  be  $(\mathbf{X}_0^c, \mathbf{X}_1^c) | (\mathbf{Y}_0, \mathbf{Y}_{1, D_1})$  with distribution  $p_1$ .  
 218 We first derive the distribution  $p_1$  of  $B$ , then the distribution  $p_0$  of  $A$  follows by omitting  $\mathbf{Y}_{1, D_1}$ . For  
 219 notation simplicity, let  $\mathbf{Y}_B$  be  $(N + M) \times k$  dimensional vector containing the observations from  
 220 existing and new batch of subjects, let  $\mathbf{X}^c$  be  $(N + M) \times (T_0 - k)$  dimensional vector containing  
 221 the underlying trajectory values evaluated at unobserved time points. And let  $\mathbf{t}$  be the time points that  
 222 have observations for subjects  $1, \dots, N + M$ , and let  $\mathbf{t}^c$  be the time points that have missing values  
 223 for subjects  $1, \dots, N + M$ .

224 Recall in the Bayesian hierarchical model (1) in Section 2, we assume multivariate distributions for  
 225 the finite observations and underlying trajectory values. We may obtain the joint distribution of  $\mathbf{Y}_B$   
 226 and  $\mathbf{X}_c$  given the hyperparameter  $\boldsymbol{\mu}_0$  as follows:

$$\begin{aligned} \begin{pmatrix} \mathbf{Y}_B \\ \mathbf{X}_c \end{pmatrix} \Big| \boldsymbol{\mu}_0 &\sim MVN \left( \begin{pmatrix} \boldsymbol{\mu}_0(\mathbf{t}) \\ \boldsymbol{\mu}_0(\mathbf{t}^c) \end{pmatrix}, \begin{pmatrix} (1+c)\boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}) + \sigma_\epsilon^2 \mathbf{I} & \boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}^c) \\ \boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}) & (1+c)\boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}^c) \end{pmatrix} \right), \\ \text{where } \mathbf{X}^c &= \begin{pmatrix} \mathbf{X}_0^c \\ \mathbf{X}_1^c \end{pmatrix}, \mathbf{Y}_B = \begin{pmatrix} \mathbf{Y}_{0, D_0} \\ \mathbf{Y}_{1, D_1} \end{pmatrix}. \end{aligned}$$

227 Then with the joint distribution, we may derive the conditional distribution of  $\mathbf{X}_c | \mathbf{Y}_B$  by the condi-  
 228 tional expectation property of multivariate normal distribution. Therefore, we get the distribution of  
 229  $B$  as

$$\begin{aligned} B = \mathbf{X}_c | \mathbf{Y}_B &\sim MVN(\mathbf{m}_B, \boldsymbol{\nu}_B), \\ \text{where } \mathbf{m}_B &= \boldsymbol{\mu}_0(\mathbf{t}) + \boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}^c) ((1+c)\boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}^c))^{-1} (\mathbf{y} - \boldsymbol{\mu}_0(\mathbf{t}^c)), \\ \boldsymbol{\nu}_B &= ((1+c)\boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}) + \sigma_\epsilon^2 \mathbf{I}) - \boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}^c) ((1+c)\boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}^c))^{-1} \boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}). \end{aligned}$$

230 One thing worth noting is that the error variance  $\sigma_\epsilon^2$  is unknown. To keep computation simplicity, we  
 231 adopt the empirical Bayes method that uses the maximum likelihood estimator  $\hat{\sigma}_\epsilon^2$  as the estimated  
 232 value of  $\sigma_\epsilon^2$ .

233 After getting the distribution of  $B$ , we may obtain the distribution of  $A$  by letting  $\mathbf{Y}_A = (\mathbf{Y}_{0, D_0})$  to  
 234 be  $N \times k$  dimensional vector. Then we substitute  $\mathbf{Y}_B$  with  $\mathbf{Y}_A$  and obtain the joint distribution of  $\mathbf{Y}_A$   
 235 and  $\mathbf{X}_c$  given the hyperparameter  $\boldsymbol{\mu}_0$  as:

$$\begin{aligned} \begin{pmatrix} \mathbf{Y}_A \\ \mathbf{X}_c \end{pmatrix} \Big| \boldsymbol{\mu}_0 &\sim MVN \left( \begin{pmatrix} \boldsymbol{\mu}_0(\mathbf{t}) \\ \boldsymbol{\mu}_0(\mathbf{t}^c) \end{pmatrix}, \begin{pmatrix} (1+c)\boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}) + \sigma_\epsilon^2 \mathbf{I} & \boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}^c) \\ \boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}) & (1+c)\boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}^c) \end{pmatrix} \right), \\ \text{where } \mathbf{X}^c &= \begin{pmatrix} \mathbf{X}_0^c \\ \mathbf{X}_1^c \end{pmatrix}, \mathbf{Y}_A = (\mathbf{Y}_{0, D_0}). \end{aligned}$$

236 Similarly, by the conditional expectation property of multivariate normal distribution, we get the  
 237 distribution of  $A$  as

$$\begin{aligned} A = \mathbf{X}_c | \mathbf{Y}_A &\sim MVN(\mathbf{m}_A, \boldsymbol{\nu}_A), \\ \text{where } \mathbf{m}_A &= \boldsymbol{\mu}_0(\mathbf{t}) + \boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}^c) ((1+c)\boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}^c))^{-1} (\mathbf{y} - \boldsymbol{\mu}_0(\mathbf{t}^c)), \\ \boldsymbol{\nu}_A &= ((1+c)\boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}) + \sigma_\epsilon^2 \mathbf{I}) - \boldsymbol{\Sigma}(\mathbf{t}, \mathbf{t}^c) ((1+c)\boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}^c))^{-1} \boldsymbol{\Sigma}(\mathbf{t}^c, \mathbf{t}). \end{aligned}$$

238 Lastly, since both  $A$  and  $B$  follow multivariate normal distributions, the closed-form of the KL  
 239 divergence between two multivariate normal distributions is

$$D_{KL}(p_1 || p_0) = \frac{1}{2} \left[ \log \left( \frac{|\boldsymbol{\nu}_A|}{|\boldsymbol{\nu}_B|} \right) - k + \text{tr} \left\{ \boldsymbol{\nu}_A^{-1} \boldsymbol{\nu}_B \right\} + (\mathbf{m}_A - \mathbf{m}_B)^T \boldsymbol{\nu}_A^{-1} (\mathbf{m}_A - \mathbf{m}_B) \right].$$

---

**Algorithm 1** Simulated-Annealing Algorithm

---

```
 $D \leftarrow D_{initial}$   
 $e \leftarrow Energy(D_{initial})$   
 $T \leftarrow T_{initial}$   
while  $\exp(\Delta e/T) > 0.2$  do  
   $D_{test} \leftarrow neighborhood(D_{initial})$   
   $e_{test} = Energy(D_{test})$   
   $\Delta e = e - e_{test}$   
  if  $\Delta e \leq 0$  then  
     $D \leftarrow D_{test}$   
     $e \leftarrow e_{test}$   
  else  
     $q \leftarrow Random(0, 1)$   
    if  $q < \exp(\Delta e/T)$  then  
       $D \leftarrow D_{test}$   
       $e \leftarrow e_{test}$   
    end if  
  end if  
   $T = 0.95 \times T$   
end while
```

---