

# Moment Matching Based Conjugacy Approximation for Bayesian Ranking and Selection

Qiong Zhang\*

Yongjia Song†

June 17, 2021

## Abstract

We study the conjugacy approximation method in the context of Bayesian ranking and selection with unknown correlations. Under the assumption of normal-inverse-Wishart prior distribution, the posterior distribution remains a normal-inverse-Wishart distribution thanks to the conjugacy property when all alternatives are sampled at each step. However, this conjugacy property no longer holds if only one alternative is sampled at a time, an appropriate setting when there is a limited budget on the number of samples. We propose two new conjugacy approximation methods based on the idea of moment matching. Both of them yield closed-form Bayesian prior updating formulas. This updating formula can then be combined with the knowledge gradient algorithm under the “value of information” framework. We conduct computational experiments to show the superiority of the proposed conjugacy approximation methods, including applications in wind farm placement and computer model calibration.

## Keywords:

Bayesian learning; ranking and selection; moment matching; approximate conjugacy

## 1 Introduction

In this work, we are concerned about selecting the best among a finite set of alternatives. We consider the scenario where we are given a budget to perform a limited number of measurements to evaluate the performances of these alternatives, before the final selection is made. In many real-world applications, the performances of the alternatives may have an underlying but unknown correlation structure, which could be exploited to improve learning for the whole set of alternatives while only a small number of measurements are performed. This situation arises in a variety of applications. One such example is computer model parameter calibration where one aims at selecting parameters that best matches the original physical system. Another example is the optimal wind farm placement [14], where one selects a candidate location that has the highest expected wind power output. In these applications, it is usually too costly to first measure all the alternatives multiple times and then select the best according to the estimated expected performances. We need to wisely allocate the measurement budget among these alternatives.

In the literature, this type of problem has been studied under the methodology known as ranking and selection. The basic idea of ranking and selection is to replicate more on “promising” candidates. More specifically, ranking and selection first builds a statistical model that quantifies the decision maker’s estimation of the expected performances of the alternatives, and then solves

---

\*Virginia Commonwealth University, USA, qzhang4@vcu.edu

†Virginia Commonwealth University, USA, ysong3@vcu.edu

an optimization problem to allocate measurement budget among all alternatives. In the literature, ranking and selection is studied under two different streams. In the frequentists’ perspective, ranking and selection is based on the indifferent-zone approach [9, 10, 11, 8]. From the Bayesian perspective, ranking and selection is studied under the “value of information” framework, see, e.g., [1], and [13] for overviews of the framework. See [2], [17], and [14] for recent development of this approach.

We focus on problems where the performances of different alternatives are likely to be correlated, but such a correlation structure is unknown apriori. If a good approximation of this correlation structure is available, it will help to prevent wasting costly measurements on alternatives that are highly correlated, since one may take advantage of the correlation information to learn about other alternatives using measurement results from a single alternative. Classical ranking and selection methods are well-developed for cases where the performances of different alternatives are assumed to be independent (e.g., [13]). Recently, approaches that exploit the underlying correlation structure have been developed. [5] study Bayesian ranking and selection for correlated normal beliefs. [15] build a Gaussian process model to incorporate the correlation information, and their numerical studies show that the correlation matrix can be accurately approximated by a parametric model (based on kernel function or other known structures). However, in many situations it is a luxury to obtain such an accurate approximation of correlation matrix, and we may only be able to gradually learn the correlation structure while making more measurements. Along this line, [14] recently propose a Bayesian sequential learning procedure based on the normal-inverse-Wishart distribution (e.g., [7]) to address the issue of unknown correlation matrix. The normal-inverse-Wishart distribution provides a very convenient way of updating the prior distribution (i.e., beliefs about the alternatives) using a simple closed-form updating formula, a property known as “conjugacy” in Bayesian statistics. The full conjugacy condition requires that all alternatives should be sampled simultaneously. However, in the context of fully sequential ranking and selection, if only a single alternative is measured in one step, this conjugacy property no longer holds, i.e., the posterior distribution is no longer normal-inverse-Wishart. [14] approximate the posterior distribution as a normal-inverse-Wishart distribution by minimizing their Kullback-Leibler divergence. This conjugacy approximation approach still gives rise to a closed-form prior updating formula. They provide extensive computational results to show the superiority of their method compared to many existing methods. We follow this idea of conjugacy approximation proposed by [14] and propose alternative approximation methods.

Specifically, we propose a different approximation scheme to match the posterior distribution with a normal-inverse-Wishart distribution, using the idea of matching their first moments. This different approximation scheme is motivated by the fact that matching two distributions by minimizing their Kullback-Leibler divergence, a distance measure of two distributions over all the moments, may be unnecessarily strong and induce some over-fitting issue. In contrast, the parameters required in the updating formula only involve the first and second order moments. Therefore, a complete matching of two distributions over all moments may not be necessary. Along this line, we develop two moment matching based conjugacy approximation for sequential ranking and selection under a normal-inverse-Wishart Bayesian model.

The contribution of this paper is two-folds. From the methodology perspective, we provide two new alternative conjugacy approximation methods for Bayesian ranking and selection under a normal-inverse-Wishart Bayesian model, both of which also yield closed-form prior updating formulas. We also show that they are superior to the Kullback-Leibler based approximation in [14] in certain cases according to our numerical study. From the application perspective, this paper is the first one that applies the methodology of Bayesian ranking and selection to calibration of computer models.

The rest of the paper is organized as follows. In Section 2, we review the normal-inverse-Wishart Bayesian model for sequential learning, and the idea of approximating conjugacy for the Bayesian framework proposed by [14]. In Section 3, we propose two new methods for updating the prior information in the Bayesian framework. In Section 4, we briefly review the knowledge gradient method used to select the alternative to sample at each step based on the value of information. We show our computational experiment results in Section 5. Proofs of theoretical results are deferred in the appendix.

An extended abstract of this paper appeared in a conference proceeding [18]. This full version of the paper presents an additional conjugacy approximation method that combines the ideas of moment matching and Kullback-Leibler divergence. We also present additional computational experiments motivated by applications in wind farm placement and computer model calibration.

## 2 Problem Setup

We aim to select the best alternative from a candidate set  $\{1, \dots, K\}$  according to their performances. For example, in the wind farm placement application, we choose the location with the highest wind power; in the computer experiments calibration, we choose the parameter setting which best matches the physical system. To be specific, let  $\mu = (\mu_1, \dots, \mu_K)^\top$  be their true performances, our goal is to find

$$k^* \in \operatorname{argmax}_{i=1}^K \mu_k.$$

However,  $\mu_k$ 's are unknown, so we can only choose the best alternative based on our belief about  $\mu$ . Following the standard assumptions in Bayesian ranking and selection [5, 14], we assume that our belief about  $\mu$  follows a multivariate normal distribution (note that  $\mu$  is used to denote both true performance and our belief.)

$$\mu|\Sigma \sim N_K(\theta^0, (q^0)^{-1}\Sigma), \quad \Sigma \sim IW_K(\mathbf{B}^0, b^0), \quad (1)$$

where given  $\Sigma$ , the conditional distribution of  $\mu$  is a multivariate normal distribution with mean vector  $\theta^0$  and covariance matrix  $(q^0)^{-1}\Sigma$ , and  $\Sigma$  follows an inverse-Wishart distribution with parameter  $\mathbf{B}^0$  and degree of freedom  $b^0$ . In the literature of Bayesian statistics, the joint distribution of  $\mu$  and  $\Sigma$  is also referred to as the normal-inverse-Wishart distribution. The expectation of  $\Sigma$  is  $\mathbf{B}^0/(b^0 - K - 1)$ , which quantifies the correlation between the performances of different alternatives.

We update our belief based on random measurements of the performances. Let  $\hat{\mathbf{Y}} = (y_1, \dots, y_K)^\top$  be a sample of the random measurement, which follows a multivariate distribution

$$\hat{\mathbf{Y}} \sim N_K(\mu, \Sigma). \quad (2)$$

Our belief about  $\mu$  can be updated sequentially as samples  $\hat{\mathbf{Y}}_1, \hat{\mathbf{Y}}_2, \dots$  are collected in a sequence. The Bayesian sequential selection is very efficient using the Bayesian model in (1) and (2). This is due to the conjugacy property of the normal-inverse-Wishart distribution in (1), which allows us to update the prior information after each new sample in a computationally tractable way [4]. Specifically, suppose that the parameters in (1) and (2) have been updated to  $\theta^n$ ,  $\mathbf{B}^n$ ,  $q^n$  and  $b^n$  at the  $n$ -th step, i.e.,

$$\mu|\Sigma \sim N_K(\theta^n, (q^n)^{-1}\Sigma), \quad \Sigma|\hat{\mathbf{Y}}^n \sim IW_K(\mathbf{B}^n, b^n). \quad (3)$$

Given a new sample  $\hat{\mathbf{Y}}^{n+1}$ , the posterior density function of  $\mu$  and  $\Sigma$  can be computed by combining the density functions of  $\hat{\mathbf{Y}}^{n+1}|\mu, \Sigma$ ,  $\mu|\Sigma$ , and  $\Sigma$ :

$$p^{n+1}(\mu, \Sigma|\hat{\mathbf{Y}}^{n+1}) \propto p^n(\hat{\mathbf{Y}}^{n+1}|\mu, \Sigma)p^n(\mu|\Sigma)p^n(\Sigma), \quad (4)$$

where

$$p^n(\hat{\mathbf{Y}}^{n+1}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{q^n}{2}(\hat{\mathbf{Y}}^{n+1} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\hat{\mathbf{Y}}^{n+1} - \boldsymbol{\mu}) \right\}, \quad (5)$$

$$p^n(\boldsymbol{\mu}|\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{q^n}{2}(\boldsymbol{\mu} - \boldsymbol{\theta}^n)^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \boldsymbol{\theta}^n) \right\}, \quad (6)$$

and

$$p^n(\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-\frac{b^n+K+1}{2}} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{B}^n \boldsymbol{\Sigma}^{-1}) \right\}. \quad (7)$$

By combining the above terms, it can be shown that  $p^n(\boldsymbol{\mu}, \boldsymbol{\Sigma}|\hat{\mathbf{Y}}^{n+1})$  follows a normal-inverse-Wishart distribution with parameters

$$\begin{aligned} q^{n+1} &= q^n + 1 \\ b^{n+1} &= b^n + 1 \\ \boldsymbol{\theta}^{n+1} &= \frac{q^n \boldsymbol{\theta}^n + \hat{\mathbf{Y}}^{n+1}}{q^n + 1} \\ \mathbf{B}^{n+1} &= \mathbf{B}^n + \frac{q^n}{q^n + 1}(\boldsymbol{\theta}^n - \hat{\mathbf{Y}}^{n+1})(\boldsymbol{\theta}^n - \hat{\mathbf{Y}}^{n+1})^\top. \end{aligned} \quad (8)$$

The normal-inverse-Wishart distribution provides a very convenient way to update the prior. However, this update requires a sample of all alternatives  $\hat{\mathbf{Y}}$  at each step, which could be too expensive when the number of alternatives is large or sampling is costly. [13] and [5] show that it is computationally advantageous to choose the most promising alternative to sample at each step. However, the flexibility of choosing only one alternative at a time will cause a significant challenge: the updating formula (8) cannot be applied in this case. The reason is that  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma}|\hat{y}_k^{n+1})$  ( $\neq p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma}|\hat{\mathbf{Y}}^{n+1})$ ) no longer follows a normal-inverse-Wishart distribution. To address this challenge, two important questions need to be answered: first, how to update the prior information in (1) when only alternative is sampled at each step; and second, how to choose the most “promising” alternative at each step. For the first question, we review an existing method in the rest of this section, and propose two new methods in Section 3. For the second question, we show in Section 4 how the proposed new methods can be used in the knowledge gradient algorithm [5], where the alternative to sample at each step is chosen by maximizing the value of information.

We now review an existing prior updating method proposed by [14] for the Bayesian model in (1)-(2) using the idea of approximate conjugacy. For the convenience of presentation, we introduce notations that will be used throughout the rest of the paper.

**Notation** For any vector  $\mathbf{x} \in \mathbb{R}^K$ , we denote the  $k$ -th element of  $\mathbf{x}$  as  $\mathbf{x}_k$ , and we denote the vector consisting of all elements of  $\mathbf{x}$  except  $\mathbf{x}_k$  as  $\mathbf{x}_{-k} \in \mathbb{R}^{K-1}$ . For any  $K \times K$  symmetric matrix  $\mathbf{X}$ , we let  $\mathbf{X}_{kk}$  be the  $k$ -th diagonal element of  $\mathbf{X}$ ,  $\mathbf{X}_{\cdot,k}$  be the  $k$ -th column of  $\mathbf{X}$ ,  $\mathbf{X}_{-k,k} \in \mathbb{R}^{K-1}$  be the subvector of  $\mathbf{X}_{\cdot,k}$  whose  $k$ -th element is excluded, and  $\mathbf{X}_{-k,-k}$  be the submatrix of  $\mathbf{X}$  constructed by removing the  $k$ -th row and the  $k$ -th column of  $\mathbf{X}$ . We also define:

$$\mathbf{X}_{-k|k} := \mathbf{X}_{-k,-k} - \frac{\mathbf{X}_{-k,k}\mathbf{X}_{k,-k}}{\mathbf{X}_{kk}}.$$

We first consider how to update the prior information from the  $n$ -th step, given that  $k$  is the alternative chosen to be sampled in the  $(n+1)$ -th step. Given  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , the new update  $\hat{y}_k^{n+1}$

follows a normal distribution,  $\hat{y}_k^{n+1} \sim N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_{kk})$ . Using the Bayes' rule, the posterior distribution of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  given  $\hat{y}_k^{n+1}$  is:

$$\begin{aligned} p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{y}_k^{n+1}) &\propto |\boldsymbol{\Sigma}|^{-\frac{b^n + K + 1}{2}} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{B}^n \boldsymbol{\Sigma}^{-1}) \right\} \\ &\cdot |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{q^n}{2} (\boldsymbol{\mu} - \boldsymbol{\theta}^n)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\theta}^n) \right\} \\ &\cdot \boldsymbol{\Sigma}_{kk}^{-1/2} \exp \left\{ -\frac{(\hat{y}_k^{n+1} - \boldsymbol{\mu}_k)^2}{2\boldsymbol{\Sigma}_{kk}} \right\}. \end{aligned} \quad (9)$$

We see that the posterior distribution is no longer a normal-inverse-Wishart distribution. Therefore, the conjugacy property of normal-inverse-Wishart distribution cannot be applied.

To address this issue, [14] proposed to use the ‘‘optimal approximation of conjugacy’’ based on minimizing the Kullback-Leibler divergence between the posterior distribution (9) and a normal-inverse-Wishart distribution, which also leads to a closed-form updating formula as follows:

$$\begin{aligned} q^{n+1} &= q^n + \frac{1}{K} \\ b^{n+1} &= b^n + \Delta b^n \\ \boldsymbol{\theta}^{n+1} &= \boldsymbol{\theta}^n + \frac{\hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n}{\frac{b^{n+1}(q^n+1)-K+1}{b^{n+1}-K+1} \mathbf{B}_{kk}^n} \mathbf{B}_{\cdot,k}^n \\ \mathbf{B}^{n+1} &= \frac{b^{n+1}}{b^n} \mathbf{B}^n + \frac{b^{n+1}}{b^n + 1} \left( \frac{q^n(b^{n+1} - K + 1)(\hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n)^2}{b^{n+1}(q^n + 1) - K + 1} - \frac{\mathbf{B}_{kk}^n}{b^n} \right) \\ &\quad \cdot \frac{\mathbf{B}_{\cdot,k}^n \mathbf{B}_{k,\cdot}^n}{\mathbf{B}_{kk}^2} \end{aligned} \quad (10)$$

where  $\Delta b^n$  is a number that can be numerically computed by a bisection algorithm, or approximated by  $K^{-1}$  according to [14].

Although this framework works well in the numerical experiments shown by [14], matching two distributions using Kullback-Leibler divergence is a very strong requirement. The Kullback-Leibler divergence of two distributions is equivalent to a distance measure of two distributions over the moments of all orders. When the true distribution is far away from normal-inverse-Wishart distribution, it may generate over-fitting issues. Therefore, a complete matching of two distributions may not necessarily lead to more accurate approximation. To address this issue, we propose two alternative methods to match the posterior distribution with a normal-inverse-Wishart distribution.

### 3 Moment Matching based Approximate Conjugacy

In this section, we consider two alternative methods to approximate the posterior distribution (9) to a normal-inverse-Wishart distribution using the idea of moment matching. The first method employs the first-order moment matching, and the second method combines the idea of moment matching and Kullback-Leibler divergence minimization. Same as the method in [14], both our new proposed methods yield closed-form updating formulas, which make the Bayesian sequential ranking and selection procedure computationally tractable. A preliminary version of the first approximation method has appeared in a conference proceeding [18].

### 3.1 Conjugacy approximation based on first-order moment matching

We consider how to update the prior information in (1) in each step given a new observation  $\hat{y}_k^{n+1}$ . Following [14], we set the increase of number of samples as  $K^{-1}$  at each step, since only one among  $K$  alternatives is sampled. Therefore, we update  $q^{n+1}$  and  $b^{n+1}$  by  $q^{n+1} = q^n + K^{-1}$  and  $b^{n+1} = b^n + K^{-1}$ .

We now consider how to update  $\boldsymbol{\theta}^{n+1}$  and  $\mathbf{B}^{n+1}$ . Let us first recall how this is done when we obtain a sample of all alternatives  $\hat{\mathbf{Y}}^{n+1}$  at the  $n$ th step. Notice that,  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{\mathbf{Y}}^{n+1})$  in (4) matches the density function of a normal-inverse-Wishart distribution with parameters  $q^{n+1}$ ,  $b^{n+1}$ ,  $\boldsymbol{\theta}^{n+1}$  and  $\mathbf{B}^{n+1}$  in (8). Meanwhile, this implies that

$$\boldsymbol{\theta}^{n+1} = \mathbb{E}(\boldsymbol{\mu} | \hat{\mathbf{Y}}^{n+1}) \quad (11)$$

and

$$\mathbf{B}^{n+1} = (b^{n+1} - K - 1) \mathbb{E} \left\{ \boldsymbol{\Sigma} | \hat{\mathbf{Y}}^{n+1} \right\} = (b^{n+1} - K - 1) \mathbb{E} \left\{ q^{n+1} \text{Var}(\boldsymbol{\mu} | \boldsymbol{\Sigma}, \hat{\mathbf{Y}}^{n+1}) | \hat{\mathbf{Y}}^{n+1} \right\}. \quad (12)$$

That is, the updated parameters  $\boldsymbol{\theta}^{n+1}$  and  $\mathbf{B}^{n+1}$  in (8) match the first-order posterior moments of  $\boldsymbol{\mu}$  and  $q^{n+1} \text{Var}(\boldsymbol{\mu} | \boldsymbol{\Sigma}, \hat{\mathbf{Y}}^{n+1})$ . When only one alternative  $k$  is sampled, the posterior distribution  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{y}_k^{n+1})$  does not follow a normal-inverse-Wishart distribution. The updating formula (10) in [14] is developed by minimizing the Kullback-Leibler divergence between  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{y}_k^{n+1})$  and the density function of a normal-inverse-Wishart distribution. Instead of matching the density functions, we develop updating formulas by matching the first-order posterior moments as in (11) and (12). To do this, we need to compute the posterior moments with regard to  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{y}_k^{n+1})$ . In Proposition 1 (a) and (b), we decompose  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{y}_k^{n+1})$  into a few parts whose first-order moments can be obtained easily. They will then be used to calculate the first-order posterior moments in (11) and (12).

**Proposition 1.** (a) Given  $\boldsymbol{\Sigma}$  and  $\hat{y}_k^{n+1}$ ,  $\boldsymbol{\mu}$  follows a multivariate normal distribution

$$\boldsymbol{\mu} | \boldsymbol{\Sigma}, \hat{y}_k^{n+1} \sim N_K(\tilde{\boldsymbol{\theta}}, (q^{n+1})^{-1} \tilde{\boldsymbol{\Sigma}}), \quad (13)$$

where

$$\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^n + \frac{(\hat{y}_k^{n+1} - \theta_k^n) \boldsymbol{\Sigma}_{\cdot, k}}{(q^n + 1) \boldsymbol{\Sigma}_{kk}},$$

and

$$\tilde{\boldsymbol{\Sigma}} = \frac{q^{n+1}}{q^n + 1} \begin{pmatrix} \frac{q^{n+1}}{q^n} \boldsymbol{\Sigma}_{-k|k} + \frac{\boldsymbol{\Sigma}_{-k, k} \boldsymbol{\Sigma}_{k, -k}}{\boldsymbol{\Sigma}_{k, k}} & \boldsymbol{\Sigma}_{-k, k} \\ \boldsymbol{\Sigma}_{k, -k} & \boldsymbol{\Sigma}_{kk} \end{pmatrix}$$

(b) Let  $A = \tilde{\boldsymbol{\Sigma}}_{-k|k}$ ,  $a = \tilde{\boldsymbol{\Sigma}}_{k, k}^{-1} \tilde{\boldsymbol{\Sigma}}_{-k, k}$ ,  $\tilde{a} = \tilde{\boldsymbol{\Sigma}}_{-k, k}$ , and  $c = \tilde{\boldsymbol{\Sigma}}_{k, k}$ , we have

$$a | A, \hat{y}_k^{n+1} \sim N_{K-1} \left( \frac{\mathbf{B}_{-k, k}^n}{\mathbf{B}_{k, k}^n}, \frac{q^n A}{q^{n+1} \mathbf{B}_{k, k}^n} \right) \quad (14)$$

$$\tilde{a} | A, c, \hat{y}_k^{n+1} \sim N_{K-1} \left( \frac{c \mathbf{B}_{-k, k}^n}{\mathbf{B}_{k, k}^n}, \frac{q^n c^2 A}{q^{n+1} \mathbf{B}_{k, k}^n} \right) \quad (15)$$

$$A | \hat{y}_k^{n+1} \sim IW_{K-1} \left( b^n, \frac{q^{n+1}}{q^n} \mathbf{B}_{-k|k}^n \right) \quad (16)$$

$$c | \hat{y}_k^{n+1} \sim IW_1 \left( b^n - K + 2, \frac{q^{n+1}}{(q^n + 1)} \left[ \mathbf{B}_{kk}^n + \frac{q^n}{q^n + 1} (\hat{y}_k^{n+1} - \theta_k^n)^2 \right] \right), \quad (17)$$

and  $A$  and  $c$  are independent.

As mentioned earlier, we update  $\boldsymbol{\theta}^{n+1}$  and  $\mathbf{B}^{n+1}$  to be

$$\boldsymbol{\theta}^{n+1} = \mathbb{E}(\boldsymbol{\mu}|\hat{y}_k^{n+1}),$$

and

$$\mathbf{B}^{n+1} = (b^{n+1} - K - 1)\mathbb{E}\left\{q^{n+1}\text{Var}(\boldsymbol{\mu}|\boldsymbol{\Sigma}, \hat{\mathbf{Y}}^{n+1})|\hat{y}_k^{n+1}\right\} = (b^{n+1} - K - 1)\mathbb{E}\left\{\tilde{\boldsymbol{\Sigma}}|\hat{y}_k^{n+1}\right\}.$$

The expectations can be calculated using the distributions given in Proposition 1. Proposition 2 summarizes the results.

**Proposition 2.** *Given  $q^{n+1}$  and  $b^{n+1}$ , the updating formulas of  $\boldsymbol{\theta}^{n+1}$  and  $\mathbf{B}^{n+1}$  based on moment matching are given by:*

$$\boldsymbol{\theta}^{n+1} = \boldsymbol{\theta}^n + \frac{\mathbf{B}_{:,k}^n \hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n}{\mathbf{B}_{kk}^n q^n + 1}, \quad (18)$$

$$\mathbf{B}_{-k,-k}^{n+1} = \frac{q^{n+1}(b^{n+1} - K - 1)}{b^n - K} \left\{ \frac{\mathbf{B}_{-k|k}^n}{q^n} + \frac{\tilde{q}}{q^n + 1} \left[ \frac{\mathbf{B}_{-k|k}^n}{b^n - K} + \frac{\mathbf{B}_{-k,k}^n \mathbf{B}_{k,-k}^n}{\mathbf{B}_{kk}^n} \right] \right\}, \quad (19)$$

$$\mathbf{B}_{-k,k}^{n+1} = \frac{q^{n+1}(b^{n+1} - K - 1)\tilde{q}}{(q^n + 1)(b^n - K)} \mathbf{B}_{-k,k}^n, \quad (20)$$

and

$$\mathbf{B}_{kk}^{n+1} = \frac{q^{n+1}(b^{n+1} - K - 1)\tilde{q}}{(q^n + 1)(b^n - K)} \mathbf{B}_{kk}^n, \quad (21)$$

where

$$\tilde{q} = \left[ 1 + \frac{q^n(\hat{y}_{k_{n+1}} - \boldsymbol{\theta}_k^n)^2}{(q^n + 1)\mathbf{B}_{kk}^n} \right].$$

**Remark.** *As shown in Proposition 1–2, the moment matching method contains two folds of moment matching. In the first fold, we match*

$$\tilde{\boldsymbol{\mu}} = \mathbb{E}(\boldsymbol{\mu}|\boldsymbol{\Sigma}, \hat{y}_k^{n+1}) \quad (22)$$

and

$$\tilde{\boldsymbol{\Sigma}} = q^{n+1}\text{Var}(\boldsymbol{\mu}|\boldsymbol{\Sigma}, \hat{y}_k^{n+1}). \quad (23)$$

*In the second fold, we set  $\boldsymbol{\theta}^{n+1} = \mathbb{E}(\tilde{\boldsymbol{\mu}}|\hat{y}_k^{n+1})$  and  $\mathbf{B}^{n+1} = (b^{n+1} - K - 1)\mathbb{E}(\tilde{\boldsymbol{\Sigma}}|\hat{y}_k^{n+1})$ . According to (13), the moment matching in the first fold also guarantees that the distribution of  $\boldsymbol{\mu}|\boldsymbol{\Sigma}, \hat{y}_k^{n+1}$  exactly matches a multivariate normal distribution with parameters  $\tilde{\boldsymbol{\mu}}$  and  $\tilde{\boldsymbol{\Sigma}}$ . However, the moment matching in the second fold does not exactly match two distributions.*

Indicated in Remark 3.1, the moment matching in the second fold does not exactly match two distributions. We next consider an alternative conjugacy approximation by combining the ideas of moment matching and Kullback-Leibler divergence minimization. Specifically, we use moment matching in the first fold of approximation (which is exact), but use Kullback-Leibler divergence minimization in the second fold.

### 3.2 Conjugacy approximation by combining moment matching and Kullback-Leibler divergence minimization

We now present an alternative conjugacy approximation method that combines the idea of moment matching and minimization of the Kullback-Leibler divergence. According to the proof of Proposition 1 (available in the appendix), the posterior distribution  $p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{\mathbf{y}}_k^{n+1})$  can be decomposed to

$$p^{n+1}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \hat{\mathbf{y}}_k^{n+1}) \propto p^{n+1}(\boldsymbol{\mu}, |\tilde{\boldsymbol{\Sigma}}, \hat{\mathbf{y}}_k^{n+1}) p^{n+1}(\tilde{\boldsymbol{\Sigma}} | \hat{\mathbf{y}}_k^{n+1}).$$

Since  $\boldsymbol{\mu}, |\tilde{\boldsymbol{\Sigma}}, \hat{\mathbf{y}}_k^{n+1}$  follows a multivariate normal distribution, the moment matching and distribution matching give same results as indicated in Remark 3.1. However,  $p^{n+1}(\tilde{\boldsymbol{\Sigma}} | \hat{\mathbf{y}}_k^{n+1})$  is not the density function of an inverse-Wishart distribution. Unlike the method in Section 3.1, we consider minimizing the Kullback-Leibler divergence between  $p^{n+1}(\tilde{\boldsymbol{\Sigma}} | \hat{\mathbf{y}}_k^{n+1})$  and an inverse-Wishart distribution to find  $\mathbf{B}^{n+1}$ . The Kullback-Leibler divergence between  $p^{n+1}(\tilde{\boldsymbol{\Sigma}} | \hat{\mathbf{y}}_k^{n+1})$  and the density function of an inverse-Wishart distribution  $\xi(\tilde{\boldsymbol{\Sigma}})$  with parameter  $\mathbf{B}$  and degree of freedom  $b^{n+1}$  is given by:

$$D_{KL}^n(\mathbf{B}) = \mathbb{E}_\xi \left\{ \log \frac{\xi(\tilde{\boldsymbol{\Sigma}})}{p^{n+1}(\tilde{\boldsymbol{\Sigma}} | \hat{\mathbf{y}}_k^{n+1})} \right\}. \quad (24)$$

$\mathbf{B}^{n+1}$  is then obtained by solving  $\min_{\mathbf{B} > 0} D_{KL}^n(\mathbf{B})$ , which has a closed-form that we show in Proposition 3.

**Proposition 3.**  $\mathbf{B}^{n+1} = \operatorname{argmin}_{\mathbf{B} > 0} D_{KL}^n(\mathbf{B})$ , is given by:

$$\mathbf{B}_{k,k} = \frac{q^{n+1}(b^{n+1} - K + 1) \left[ \mathbf{B}_{k,k}^n + \frac{q^n}{q^{n+1}} (\hat{\mathbf{y}}_k^{n+1} - \boldsymbol{\theta}_k^n)^2 \right]}{(b^n + 1)(q^n + 1)}, \quad (25)$$

$$\mathbf{B}_{-k,k}^{n+1} = \frac{\mathbf{B}_{k,k} \mathbf{B}_{-k,k}^n}{\mathbf{B}_{k,k}^n}, \quad (26)$$

and

$$\mathbf{B}_{-k,-k}^{n+1} = \frac{b^{n+1} q^{n+1}}{b^n q^n} \mathbf{B}_{-k|k}^n + \frac{\mathbf{B}_{-k,k} \mathbf{B}_{k,-k}}{\mathbf{B}_{k,k}}. \quad (27)$$

By combining (18) with the formula of  $\mathbf{B}^{n+1}$  in Proposition 3, we obtain a new sequential prior updating procedure.

**Remark.** When all alternatives  $\hat{\mathbf{Y}}$  are sampled in each step, both moment matching and Kullback-Leibler divergence minimization will lead to the same updating formula as in (8).

In Section 5, we compare the performances of the proposed two new conjugacy approximation methods with the one proposed in [14] through experiments motivated by two applications.

## 4 Computation of the Value of Information

In this section, we follow the knowledge gradient framework developed in [5], [13], and [14] on how to sequentially select the alternative to sample in each step. According to this framework, suppose we have obtained  $n$  samples, the alternative to be sampled in the next step,  $k^{n+1}$ , is the one that maximizes the value of information [13]:

$$V_n(k) = \mathbb{E} \left[ \max_{k'=1,2,\dots,K} \theta_{k'}^{n+1} | k^{n+1} = k \right] - \max_{k'=1,2,\dots,K} \theta_{k'}^n, \quad (28)$$



where the expectation is taken with regard to the predictive distribution of  $\theta_k^{n+1}$  given all collected data points.

When the updating formula (10) derived from minimizing the Kullback-Leibler divergence is used, we have:

$$\boldsymbol{\theta}^{n+1} = \boldsymbol{\theta}^n + \mathbf{s}^n(k)T^n, \quad (29)$$

where

$$\mathbf{s}^n(k) = \frac{\sqrt{\frac{q^{n+1}}{q^n(b^n-K+1)}}}{\left(\frac{q^n b^{n+1}}{b^{n+1}-K+1} + 1\right)\sqrt{\mathbf{B}_{kk}^n}} \mathbf{B}_{:,k}^n, \quad (30)$$

and

$$T^n = \frac{\hat{y}_{k_{n+1}} - \theta_k^n}{\sqrt{\frac{q^{n+1}}{q^n(b^n-K+1)} \mathbf{B}_{kk}^n}}. \quad (31)$$

According to [14], the predictive distribution of  $T^n$  is a  $t$ -distribution with degree of freedom  $b^n - K + 1$ . Thus, the expectation in (28) can be computed using the properties of the  $t$ -distribution.

Similarly, when the updating formula (18) is derived based on moment matching, or the combination between moment matching and Kullback-Leibler divergence minimization as described in Section 3,  $\mathbf{s}^n(k)$  in (29) is defined as:

$$\mathbf{s}^n(k) = \frac{\mathbf{B}_{:,k}^n}{\sqrt{q^n(q^n + 1)(b^n - K + 1)\mathbf{B}_{kk}^n}},$$

and  $T^n$  is the same as in (31). Therefore, we can also use the predictive distribution of  $T^n$ , i.e., a  $t$ -distribution with degree of freedom  $b^n - K + 1$ , to compute the expectation in (28).

For all three conjugacy approximation methods, according to the above analysis, the optimization problem that maximizes (28) can be written as:

$$\max_{k=1,2,\dots,K} \{V_n(k)\}, \quad (32)$$

where

$$V_n(k) := \mathbb{E} \left[ \max_{k'=1,2,\dots,K} (\theta_{k'}^n + s_{k'}^n(k^{n+1})T^n) \mid k^{n+1} = k \right] - \max_{k'=1,2,\dots,K} \theta_{k'}^n,$$

and  $s_{k'}^n(k^{n+1})$  is the  $k'$ -th element of vector  $\mathbf{s}^n(k^{n+1})$ . A closed-form solution of (32) has been provided by [13] and [14].

## 5 Numerical Experiments

We present numerical results to compare three approximate conjugacy methods. In particular, we compare the performances of the proposed methods and the one proposed in [14] based on minimizing the Kullback-Leibler divergence using various test cases. The three methods are labeled as:

1. KL: Minimizing the Kullback-Leibler divergence as in [14] (as described in Section 2).
2. Moment: Matching the first-order moments as described in Section 3.1.
3. Moment-KL: Combination of moment matching and Kullback-Leibler divergence minimization as described in Section 3.2.

The performances of the three conjugacy approximation methods are compared using their corresponding opportunity costs at each step. As in [14], the opportunity cost of each method in step  $n$  is defined by

$$C_n = \max_k \mu_k - \mu_{\arg\max_k \theta_k^n}, \quad (33)$$

where  $\mu_k$  is the true performance of the  $k$ th alternative,  $\theta_k^n$  is the posterior mean given by a certain method at step  $n$ , and  $\mu_{\arg\max_k \theta_k^n}$  is the true performance of the best alternative selected by a certain method at the  $n$ th step. A smaller opportunity cost indicates that the method is more accurate in selecting the best alternative. We would also expect that  $C_n$  decreases with  $n$ . For all cases shown below, we replicate the overall procedure 500 times, and report the average results.

### 5.1 Data generated from a multivariate normal distribution

We first consider an example where the samples  $\hat{\mathbf{Y}}$  are generated from a multivariate normal distribution. We consider nine alternatives, and let their corresponding true mean values be  $\frac{1}{9}, \frac{2}{9}, \dots, 1$ , respectively. The true covariance matrix  $A$  is given as:  $A_{ij} = (-\rho)^{|i-j|}$ , and we consider three different values for  $\rho$ , 0.1, 0.5 and 0.9, which indicate three different levels of correlation strength, low, median and high, respectively.

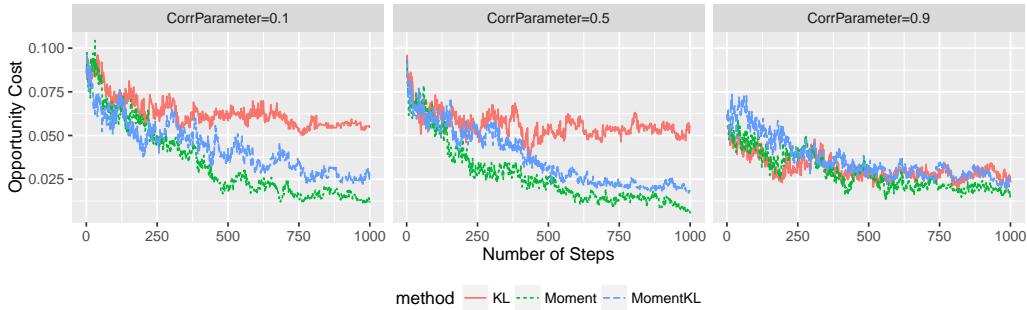


Figure 1: Average opportunity cost at each step for 1000 steps for each method in the multivariate normal case (Section 5.1) over 500 replications.

Figure 1 shows the performances of the three methods on this example in terms of their opportunity costs (33) at each step for 1000 steps. The prior parameters are estimated by the sample mean and sample covariance of 25 samples from all alternatives. We can see from Figure 1 that, as the number of steps increases, the opportunity cost decreases for all three methods. In general, method “Moment” gives the smallest opportunity cost for the low and medium correlation cases. However, when the number of steps is small, the performance of method “Moment-KL” is comparable with and sometimes better than method “Moment”. In the high correlation case, the performances of all three methods are close.

We next consider how the performances of three methods vary using different numbers of samples for prior estimation. Table 1 shows the means and standard deviations of the final results (at the 1000-th step) on the multivariate normal example with 5, 15, and 25 samples for the prior estimation. In Table 1, columns labeled as “Opp. cost” and “Error” show the mean and standard deviation of the opportunity cost, respectively (the same labels are also used in Table 2 and Table 4). We can see from Table 1 that, in terms of the final opportunity cost, “Moment” and “Moment-KL” perform better than “KL” in almost all cases considered. We also see that, when the prior information is more accurate (when a larger number of samples are used), the opportunity cost is significantly lower for all three methods in most cases.

Table 1: The mean and standard deviation of the final (at the 1000-th step) opportunity cost for methods “KL”, “Moment”, and “Moment-KL” on the multivariate normal example (Section 5.1) with various number of samples for prior estimation, and correlation.

| Corr | # prior | KL        |        | Moment    |        | Moment-KL |        |
|------|---------|-----------|--------|-----------|--------|-----------|--------|
|      |         | Opp. cost | Error  | Opp. cost | Error  | Opp. cost | Error  |
| 0.1  | 5       | 0.1767    | 0.0023 | 0.1407    | 0.0023 | 0.1553    | 0.0021 |
|      | 15      | 0.0689    | 0.0011 | 0.0225    | 0.0005 | 0.0201    | 0.0004 |
|      | 25      | 0.0554    | 0.0010 | 0.0123    | 0.0003 | 0.0286    | 0.0007 |
| 0.5  | 5       | 0.1286    | 0.0018 | 0.1650    | 0.0025 | 0.1418    | 0.0020 |
|      | 15      | 0.0844    | 0.0011 | 0.0149    | 0.0004 | 0.0195    | 0.0005 |
|      | 25      | 0.0557    | 0.0009 | 0.0053    | 0.0002 | 0.0172    | 0.0005 |
| 0.9  | 5       | 0.1085    | 0.0017 | 0.0476    | 0.0008 | 0.0472    | 0.0010 |
|      | 15      | 0.0347    | 0.0007 | 0.0084    | 0.0003 | 0.0199    | 0.0005 |
|      | 25      | 0.0233    | 0.0004 | 0.0149    | 0.0004 | 0.0238    | 0.0006 |

As observed from both Figure 1 and Table 1, all three methods have similar results when the alternatives are highly correlated. This can be explained by Remark 3.2. When the correlation is high, and the number of alternatives is small (say,  $K = 9$  in this case), a sample from a single alternative can indicate the performances of other alternatives with high probability. In this sense, sampling a single alternative has a similar effect as sampling all alternatives, in which case the three methods are equivalent as discussed in Remark 3.2.

## 5.2 Wind farm placement using wind speed historical data

We next study the three methods of conjugacy approximation, “KL”, “Moment”, and “Moment-KL” on the application of wind farm placement problem using real-world data. This application is borrowed from [14], where method “KL” is compared with several other alternative approaches in the literature. In this problem, the goal is to select the best site among a set of candidate sites for installing new wind farms, in terms of average wind power output. We use the publicly available historical wind speed data in the United States from [3]. To be consistent with the results shown in [14], we use exactly the same setting described in that paper. However, we may have used a different time period from the wind database [3]: we collected hourly wind speed data from June 30th, 2008 to December 31st, 2011, whereas [14] did not report the range of dates where the data was collected.

As in [14], we choose from 64 candidate sites distributed on an  $8 \times 8$  grid from the state of Washington. We use three different levels of latitude and longitude resolutions, that is, 0.125 degrees (High), 0.25 degrees (Medium) and 0.375 degrees (Low). A higher resolution means more spatial correlations between different locations, and less differences between their true means. Figure 2 shows the average performances of the three conjugacy approximation methods in each step for 200 steps over 500 replications. Table 2 shows the average means and standard deviations of the final opportunity cost (at the 200-th step) of the three methods over 500 replications. Similar to what we have observed in Section 5.1, the proposed methods “Moment” and “Moment-KL” perform better than “KL” in most scenarios. Furthermore, we observe in Figure 2 that the performances under three resolutions are significantly different from each other. This can be explained by the different resolutions considered in the three cases. For the low resolution case (the distance between two alternatives is large), the true performances of different alternatives are significantly

different from each other, therefore, it is easy to distinguish among these alternatives, and make the correct selection, which ends up with a small opportunity cost. For the high resolution case (the distance between two alternatives is small), the true performances of different alternatives are similar, therefore, even if a wrong selection is made, it does not lead to a large opportunity cost. The medium resolution case does not enjoy the advantages in either low or high resolution case, and it gives the worst results in terms of opportunity cost among the three cases.

Table 2: The mean and standard deviation of the final opportunity cost (at the 200-th step) for methods “KL”, “Moment”, and “Moment-KL” on the wind farm example (Section 5.2) with different resolutions.

| Resolution | KL        |        | Moment    |        | Moment-KL |        |
|------------|-----------|--------|-----------|--------|-----------|--------|
|            | Opp. cost | Error  | Opp. cost | Error  | Opp. cost | Error  |
| Low        | 0.0892    | 0.0093 | 0.0443    | 0.0069 | 0.0613    | 0.0076 |
| Medium     | 0.1518    | 0.0089 | 0.0962    | 0.0068 | 0.1190    | 0.0074 |
| High       | 0.0236    | 0.0061 | 0.0032    | 0.0023 | 0.0412    | 0.0078 |

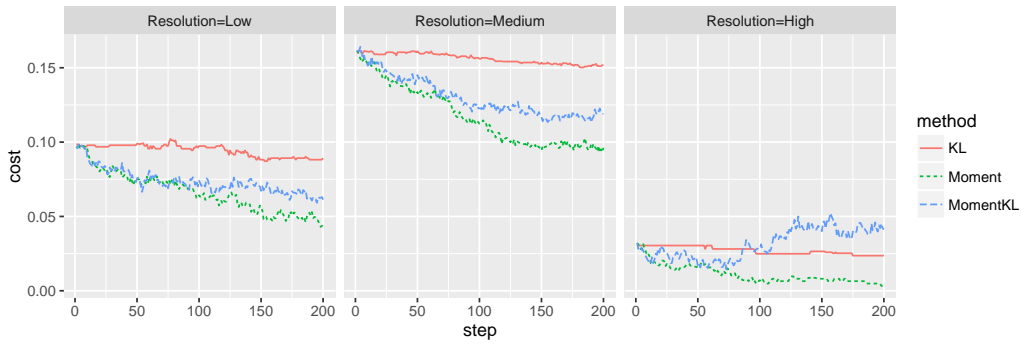


Figure 2: Average opportunity cost at each step for 200 steps for three conjugacy approximation method in the Bayesian sequential ranking and selection in the wind farm placement test case (Section 5.2) over 500 replications.

### 5.3 Computer model calibration

In this section, we formulate the computer model calibration problem as a Bayesian ranking and selection problem. Consider a physical system with  $\mathbf{x} \in \mathcal{X} \subset \mathbf{R}^d$  being the control variables. The response of the system can be seen as a real-valued stochastic function, denoted by  $\eta(\mathbf{x})$ . When running physical experiments is costly, a statistical predictor  $\hat{\eta}(\mathbf{x})$  (such as the interpolator in [16]) can be used to model the unknown true response  $\eta(\mathbf{x})$  based on a set of observations.

Computer experiments are usually used to mimic costly physical experiments. Input parameters of a computer model include control variables  $\mathbf{x}$  in the physical system, as well as a calibration variable  $\lambda$ , which describes some inherent features of the physical system. Let the response of this computer model be  $f(\mathbf{x}, \lambda)$ , the goal of calibrating this computer experiment is to reduce the gap between  $f(\mathbf{x}, \lambda)$  and  $\eta(\mathbf{x})$  by choosing an appropriate  $\lambda$ . We consider the case where the calibration variable  $\lambda$  is a qualitative parameter with  $K$  different qualitative levels. In cases where multiple qualitative variables exist, we let  $\lambda$  be an aggregate qualitative parameter whose qualitative levels correspond to all level combinations of these variables. Similar to [16], the calibration variable  $\lambda$  is

chosen by minimizing the mean squared error (MSE) between the physical model and the computer model:

$$\text{MSE}(\lambda) = \text{E} \{ \eta(\mathbf{x}) - f(\mathbf{x}, \lambda) \}^2, \quad (34)$$

where the expectation is taken with regard to the randomness of  $\mathbf{x}$ , and the randomness of the physical system and/or the computer model (depending on whether or not the computer model is stochastic). The mean squared error in (34) measures the model discrepancy of  $f(\mathbf{x}, \lambda)$ . Since  $f(\mathbf{x}, \lambda)$  and  $\eta(\mathbf{x})$  are unknown and only available at a few design points, the MSE in (34) is not readily available for each  $\lambda$ . By surrogating  $\eta(\mathbf{x})$  with  $\hat{\eta}(\mathbf{x})$ , the MSE of each qualitative level could be estimated empirically. By setting  $\lambda$  at its  $k$ -th qualitative level, we run the computer model on a design of control variables  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , and denote the outputs as  $f_k(\mathbf{x}_i)$  for  $i = 1, \dots, m$ . The empirical estimation of (34) is given by:

$$\hat{y}_k = m^{-1} \sum_{i=1}^m \{ \hat{\eta}(\mathbf{x}_i) - f_k(\mathbf{x}_i) \}^2. \quad (35)$$

We arrange  $\hat{y}_k$ 's from all qualitative levels in a single vector

$$\hat{\mathbf{Y}} = (\hat{y}_1, \dots, \hat{y}_K)^\top, \quad (36)$$

which estimates the model discrepancy of  $f(\mathbf{x}, \lambda)$  at all  $K$  qualitative levels of  $\lambda$ . Hence, we have formulated this problem into a Bayesian ranking and selection problem: the samples are the estimates of MSEs, and the alternatives are the qualitative levels of the calibration parameter  $\lambda$  indexed from  $\{1, \dots, K\}$ .

We test the three conjugacy approximation methods, “KL”, “Moment”, and “Moment-KL”, for computer model calibration on the Borehole function [12], a widely used example for illustrating various methods in computer experiments. This function models the flow rate of water through a borehole, and has the following form:

$$f(\mathbf{x}) = \log \left\{ \frac{2\pi x_1 x_6}{\log(x_5/x_2) \left[ 1 + \frac{2x_3 x_1}{\log(x_5/x_2) x_2^2 x_7} + \frac{x_1}{x_4} \right]} \right\}, \quad (37)$$

where  $\mathbf{x} = (x_1, \dots, x_7)^\top$ , and the ranges and units of inputs  $x_1 \sim x_7$  are given in Table 3. Inputs  $x_1$ – $x_5$  are the control variables of this system, and inputs  $x_6$  and  $x_7$  are the qualitative calibration parameters. Function (37) is used as the computer model, and the true physical system is specified as

$$\hat{\eta}(\mathbf{x}) = \log \left\{ \frac{2\pi x_1 \times 401}{\log(x_5/x_2) \left[ 1 + \frac{2x_3 x_1}{\log(x_5/x_2) x_2^2 \times 11000} + \frac{x_1}{x_4} \right]} \right\} + N(0, 1). \quad (38)$$

Table 3: Ranges of the inputs  $x_1 \sim x_7$  on the Borehole function.

| Variable | Range        | Unit     | Variable | Range      | Unit   |
|----------|--------------|----------|----------|------------|--------|
| $x_1$    | 63070-115600 | $m^2/yr$ | $x_5$    | 100-50000  | $M$    |
| $x_2$    | 0.05-0.15    | $M$      | $x_6$    | 170-410    | $M$    |
| $x_3$    | 1120-1680    | $M$      | $x_7$    | 9588-12045 | $m/yr$ |
| $x_4$    | 63.1-116     | $m^2/yr$ |          |            |        |

In our experiments, we compute  $\hat{y}_k^{n+1}$  in each step according to (35), where  $\hat{\eta}(\cdot)$  function is given by (38). We let  $x_6$  be a qualitative parameter with three equally spaced levels, and we consider two different numbers of levels for parameter  $x_7$ , 10 and 17, which gives 30 and 51 level combinations in total, respectively. For each qualitative level, we generate the design points of the control variables  $\mathbf{x}$  using a 5-dimensional Latin hypercube design with eight runs.



Figure 3: Average opportunity cost at each step over 1000 steps for each methods “KL”, “Moment”, and “Moment-KL” in the computer model calibration with borehole function (Section 5.3) with 30 and 51 qualitative levels over 500 replications.

Table 4: The mean and standard deviation of the final opportunity cost (at the 1000-th step) for methods “KL”, “Moment”, and “Moment-KL” in the computer model calibration with borehole function (Section 5.3) with various level combinations and number of samples for prior estimation.

| K  | # prior | KL        |        | Moment    |        | Moment-KL |        |
|----|---------|-----------|--------|-----------|--------|-----------|--------|
|    |         | Opp. cost | Error  | Opp. cost | Error  | Opp. cost | Error  |
| 30 | 20      | 0.0315    | 0.0003 | 0.0196    | 0.0002 | 0.0334    | 0.0004 |
|    | 50      | 0.0226    | 0.0002 | 0.0148    | 0.0001 | 0.0151    | 0.0001 |
| 51 | 20      | 0.0347    | 0.0004 | 0.0194    | 0.0002 | 0.0223    | 0.0003 |
|    | 50      | 0.0288    | 0.0002 | 0.0205    | 0.0002 | 0.0215    | 0.0002 |

Table 4 and Figure 3 show the performances of three different methods. Consistent with what we have observed in the multivariate normal case, method “Moment” performs better than the other two methods, especially in the case when the number of qualitative levels is small, and a small number of samples are used to estimate the prior distribution. In cases with larger number of qualitative levels, the performances of the proposed two new conjugacy approximation methods “Moment” and “Moment-KL” are competitive, and both significantly outperform method “KL”. We also see that, when a larger number of samples are used to estimate the prior, the performances of two “KL” based methods are significantly improved. However, this is not the case for method “Moment”. This shows that method “Moment” is less sensitive to the accuracy of the prior distribution.

## 6 Concluding Remarks

We have proposed two alternative conjugacy approximation methods for Bayesian ranking and selection. Unlike the distribution match conjugacy approximation in [14], our proposal is developed based on moment matching. We have conducted comprehensive numerical experiments, including the applications of the Bayesian ranking and selection on wind farm placement and computer model calibration. Our experiment results have shown the superiority of the proposed methods.

## References

- [1] S. Chick. Bayesian ideas and discrete event simulation: why, what and how. In L. Perrone, F. Wieland, J. Liu, B. Lawson, D. Nicol, and R. Fujimoto, editors, *Proceedings of the Winter Simulation Conference*, pages 96–105, 2006.
- [2] S. Chick and P. Frazier. Sequential sampling for selection with economics of selection procedures. *Management Science*, 58(3):550–569, 2012.
- [3] B.A. Cosgrove, D. Lohmann, K.E. Mitchell, P.R. Houser, E.F. Wood, J.C. Schaake, A. Robock, et al. Real-time and retrospective forcing in the north american land data assimilation system (nldas) project. *Journal of Geophysical Research*, 108(D22):8842–8853, 2003.
- [4] M.H. DeGroot. *Optimal statistical decisions*. John Wiley and Sons, 2004.
- [5] P.I. Frazier, W.B. Powell, and S. Dayanik. The knowledge-gradient policy for correlated normal rewards. *INFORMS Journal on Computing*, 21(4):599–613, 2009.
- [6] G.H. Golub and C.F. Van Loan. *Matrix Computations*. Johns Hopkins Studies in the Mathematical Science, 1996.
- [7] A.K. Gupta and D.K. Nagar. *Matrix variate distributions*. Chapman & Hall, 2000.
- [8] L.J. Hong and B.L. Nelson. A brief introduction to optimization via simulation. In M.D. Rosetti, R.R. Hill, B. Johansson, A. Dunkin, and R.G. Ingalls, editors, *Proceedings of the Winter Simulation Conference*, pages 75–85, 2009.
- [9] S.-H. Kim and B.L. Nelson. A fully sequential procedure for indifference-zone selection in simulation. *ACM Transactions on Modeling and Computer Simulation*, 11(3):251–273, 2001.
- [10] S.-H. Kim and B.L. Nelson. On the asymptotic validity of fully sequential selection procedures for steady-state simulation. *Operations Research*, 54(3):475–488, 2006.
- [11] S.-H. Kim and B.L. Nelson. Recent advances in ranking and selection. In S. G. Henderson, B. Biller, M.-H. Hsieh, J. Shortle, J. D. Tew, and R. R. Barton, editors, *Proceedings of the Winter Simulation Conference*, pages 162–172, 2007.
- [12] M.D. Morris, T.J. Mitchell, and D. Ylvisaker. Bayesian design and analysis of computer experiments: Use of derivatives in surface prediction. *Technometrics*, 35:243–255, 1993.
- [13] W.B. Powell and I.O. Ryzhov. *Optimal learning*. John Wiley and Sons, 2012.
- [14] H. Qu, I.O. Ryzhov, M.C. Fu, and Z. Ding. Sequential selection with unknown correlation structures. *Operations Research*, 63(4):931–948, 2015.

- [15] W. Scott, P. Frazier, and W.B. Powell. The correlated knowledge gradient for simulation optimization of continuous parameters using gaussian process regression. *SIAM Journal on Optimization*, 21(3):996–1026, 2011.
- [16] R. Tuo and J.C.F. Wu. Efficient calibration for imperfect computer models. *The Annals of Statistics*, 43(6):2331–2352, 2015.
- [17] J. Xie and P. Frazier. Sequential bayes-optimal policies for multiple comparisons with a known standard. *Operations Research*, 61(3):1174–1189, 2013.
- [18] Q. Zhang and Y. Song. Simulation selection for empirical model comparison. In *Proceedings of the 2015 Winter Simulation Conference*, 2015.

## A Proof of Proposition 1

*Proof. Proof of Proposition 1 (a)* Given  $\Sigma$  and  $\hat{y}_{k_{n+1}}$ , the density function of  $\mu$  is given by:

$$p^{n+1}(\mu|\Sigma, \hat{y}_{k_{n+1}}) \propto \exp \left\{ -\frac{q^n}{2}(\mu - \theta^n)^\top \Sigma^{-1}(\mu - \theta^n) - \frac{(\hat{y}_{k_{n+1}} - \mu_k)^2}{2\Sigma_{kk}} \right\}. \quad (39)$$

To show (39) is a multivariate normal distribution, we need to find  $\tilde{\theta}$  and  $\tilde{\Sigma}$  that satisfy

$$p^{n+1}(\mu|\Sigma, \hat{y}_{k_{n+1}}) \propto \exp \left\{ -\frac{q^{n+1}}{2}(\mu - \tilde{\theta})^\top \tilde{\Sigma}^{-1}(\mu - \tilde{\theta}) \right\}. \quad (40)$$

By comparing (39) and (40),  $\tilde{\theta}$  and  $(q^{n+1})^{-1}\tilde{\Sigma}$  should satisfy

$$q^{n+1}\tilde{\Sigma}^{-1}\tilde{\theta} = q^n\Sigma^{-1}\theta^n + \frac{\hat{y}_k^{n+1}}{\Sigma_{k,k}}e_k \quad (41)$$

and

$$q^{n+1}\tilde{\Sigma}^{-1} = q^n\Sigma^{-1} + \frac{1}{\Sigma_{kk}}e_k(e_k)^\top \quad (42)$$

where  $e_k$  is a  $K$ -dimensional vector whose  $k$ -th entry equals to 1 and other entries equal to 0. By applying the Sherman-Morrison-Woodbury matrix formula [6], we obtain the formula of  $\tilde{\theta}$  and  $\tilde{\Sigma}$  as in Proposition 1 (a).

**Proof of Proposition 1 (b)** Since  $A$ ,  $a$ ,  $\tilde{a}$ , and  $c$  are functions of  $\Sigma$ , we first derive the density function of  $\Sigma$ . According to the proof in Proposition 1 (a), we have that

$$\begin{aligned} p^{n+1}(\mu|\Sigma, \hat{y}_{k_{n+1}}) &\propto \exp \left\{ -\frac{q^n}{2}(\mu - \theta^n)^\top \Sigma^{-1}(\mu - \theta^n) - \frac{(\hat{y}_{k_{n+1}} - \mu_k)^2}{2\Sigma_{kk}} \right\} \\ &= \exp \left\{ -\frac{q^{n+1}}{2}(\mu - \tilde{\theta})^\top \tilde{\Sigma}^{-1}(\mu - \tilde{\theta}) - \frac{q^n(\hat{y}_{k_{n+1}} - \theta_k^n)^2}{2(q^n + 1)\Sigma_{kk}} \right\}. \end{aligned}$$

Thus,

$$p^{n+1}(\Sigma|\hat{y}_{k_{n+1}}) = \int p^{n+1}(\mu, \Sigma|\hat{y}_{k_{n+1}})d\mu$$



$$\begin{aligned}
&= |\Sigma|^{-\frac{b^n+K+2}{2}} \Sigma_{kk}^{-1/2} \cdot \exp \left\{ -\frac{q^n(\hat{y}_{k_{n+1}} - \theta_k^n)^2}{2(q^n+1)\Sigma_{kk}} - \frac{1}{2} \text{tr}(\mathbf{B}^n \Sigma^{-1}) \right\} \\
&\quad \cdot \int \exp \left\{ -\frac{q^{n+1}}{2} (\boldsymbol{\mu} - \tilde{\boldsymbol{\theta}})^\top \tilde{\Sigma}^{-1} (\boldsymbol{\mu} - \tilde{\boldsymbol{\theta}}) \right\} d\boldsymbol{\mu} \\
&\propto |\Sigma|^{-\frac{b^n+K+2}{2}} \Sigma_{kk}^{-1/2} \cdot \exp \left\{ -\frac{q^n(\hat{y}_{k_{n+1}} - \theta_k^n)^2}{2(q^n+1)\Sigma_{kk}} - \frac{1}{2} \text{tr}(\mathbf{B}^n \Sigma^{-1}) \right\} \cdot |\tilde{\Sigma}|^{1/2}.
\end{aligned} \tag{43}$$

We now transform the variables in (43) in terms of  $A$ ,  $a$ , and  $c$ . Since

$$A = \tilde{\Sigma}_{-k|k} = \frac{q^{n+1}}{q^n} \Sigma_{-k|k}, \tag{44}$$

$$a = \tilde{\Sigma}_{kk}^{-1} \tilde{\Sigma}_{-k,k} = \Sigma_{kk}^{-1} \Sigma_{-k,k}, \tag{45}$$

and

$$c = \tilde{\Sigma}_{kk} = \frac{q^{n+1}}{q^n+1} \Sigma_{kk}, \tag{46}$$

we express

$$|\tilde{\Sigma}| = |\tilde{\Sigma}_{-k|k}| \cdot \tilde{\Sigma}_{kk} = |A| \cdot c, \tag{47}$$

$$|\Sigma| = |\Sigma_{-k|k}| \cdot \Sigma_{kk} \propto |A| \cdot c, \tag{48}$$

and

$$\begin{aligned}
\text{tr}(\mathbf{B}^n \Sigma^{-1}) &= \frac{\mathbf{B}_{kk}^n}{\Sigma_{kk}} + \text{tr}(\mathbf{B}_{-k|k}^n \Sigma_{-k|k}^{-1}) + \mathbf{B}_{kk}^n \left( \frac{\Sigma_{-k,k}}{\Sigma_{kk}} - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right)^\top \Sigma_{-k|k}^{-1} \left( \frac{\Sigma_{-k,k}}{\Sigma_{kk}} - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right) \\
&= \frac{q^{n+1} \mathbf{B}_{kk}^n}{(q^n+1)c} + \frac{q^{n+1}}{q^n} \text{tr}(\mathbf{B}_{-k|k}^n A^{-1}) + \frac{q^{n+1}}{q^n} \mathbf{B}_{kk}^n \left( a - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right)^\top A^{-1} \left( a - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right).
\end{aligned} \tag{49}$$

The determinant of the Jacobin matrix for the transformations in (44)–(46) is  $c^{K-1}$ . Therefore, we express

$$\begin{aligned}
p^{n+1}(A, a, c | \hat{y}_{k_{n+1}}) &\propto |A|^{-\frac{b^n+K+1}{2}} \exp \left\{ -\frac{q^{n+1}}{2q^n} \text{tr}(\mathbf{B}_{-k|k}^n A^{-1}) \right\} \\
&\quad \cdot c^{-\frac{b^n-K+4}{2}} \exp \left\{ -\frac{q^{n+1} \mathbf{B}_{kk}^n}{2(q^n+1)c} - \frac{q^n q^{n+1} (\hat{y}_{k_{n+1}} - \theta_k^n)^2}{2(q^n+1)^2 c} \right\} \\
&\quad \cdot \exp \left\{ -\frac{\mathbf{B}_{kk}^n q^{n+1}}{2q^n} \left( a - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right)^\top A^{-1} \left( a - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right) \right\}.
\end{aligned} \tag{50}$$

According to (50), we have

$$p^{n+1}(a | A, c, \hat{y}_{k_{n+1}}) \propto \exp \left\{ -\frac{\mathbf{B}_{kk}^n q^{n+1}}{2q^n} \left( a - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right)^\top A^{-1} \left( a - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right) \right\},$$

which further leads to the multivariate normal distributions of  $a$  and  $\tilde{a}$ . By integrating over  $a$  in (50), we have

$$\begin{aligned}
p^{n+1}(A, c | \hat{y}_{k_{n+1}}) &\propto |A|^{-\frac{b^n+K}{2}} \exp \left\{ -\frac{q^{n+1}}{2q^n} \text{tr}(\mathbf{B}_{-k|k}^n A^{-1}) \right\} \\
&\quad \cdot c^{-\frac{b^n-K+4}{2}} \exp \left\{ -\frac{q^{n+1} \mathbf{B}_{kk}^n}{2(q^n+1)c} - \frac{q^n q^{n+1} (\hat{y}_{k_{n+1}} - \theta_k^n)^2}{2(q^n+1)^2 c} \right\},
\end{aligned}$$

which leads to the independent Inverse-Wishart distributions of  $A$  and  $c$ .  $\square$

## B Proof of Proposition 2

*Proof.* We update  $\boldsymbol{\theta}^{n+1}$  and  $\mathbf{B}^{n+1}$  by matching them with the posterior moments of  $\tilde{\boldsymbol{\theta}}$  and  $\tilde{\boldsymbol{\Sigma}}$  in Proposition 1, i.e.,

$$\boldsymbol{\theta}^{n+1} = \mathbb{E}(\tilde{\boldsymbol{\theta}}|\hat{y}_{k_{n+1}}) \quad (51)$$

and

$$\mathbf{B}^{n+1} = (b^{n+1} - K - 1)\mathbb{E}(\tilde{\boldsymbol{\Sigma}}|\hat{y}_{k_{n+1}}). \quad (52)$$

Therefore, the tasks in this proposition is to derive  $\mathbb{E}(\tilde{\boldsymbol{\theta}}|\hat{y}_{k_{n+1}})$  and  $\mathbb{E}(\tilde{\boldsymbol{\Sigma}}|\hat{y}_{k_{n+1}})$ .

We first derive  $\mathbb{E}(\tilde{\boldsymbol{\theta}}|\hat{y}_{k_{n+1}})$ . Recall that

$$\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^n + \frac{(\hat{y}_{k_{n+1}} - \theta_k^n)\boldsymbol{\Sigma}_{\cdot,k}}{(q^n + 1)\boldsymbol{\Sigma}_{kk}}.$$

Thus,

$$\mathbb{E}(\tilde{\boldsymbol{\theta}}|\hat{y}_{k_{n+1}}) = \boldsymbol{\theta}^n + \frac{(\hat{y}_{k_{n+1}} - \theta_k^n)}{(q^n + 1)} \mathbb{E} \frac{\boldsymbol{\Sigma}_{\cdot,k}}{\boldsymbol{\Sigma}_{kk}}.$$

According to the proof of Proposition 1 (b),  $\boldsymbol{\Sigma}_{\cdot,k}/\boldsymbol{\Sigma}_{kk}$  is a vector whose  $k$ -th component equals to 1, and other components equal to the entries in  $a$  defined in Proposition 1. We see from Proposition 1 that, given  $\hat{y}_{k_{n+1}}$  and  $A$ ,  $a$  follows a normal distribution with mean  $\mathbf{B}_{-k,k}^n/\mathbf{B}_{kk}^n$ . Thus, we obtain the expression of  $\boldsymbol{\theta}^{n+1}$ .

Now we derive  $\mathbb{E}(\tilde{\boldsymbol{\Sigma}}|\hat{y}_{k_{n+1}})$ . According to the definition of  $A$ ,  $a$ ,  $\tilde{a}$  and  $c$  in Proposition 1, we have

$$\begin{aligned} \tilde{\boldsymbol{\Sigma}}_{-k,-k} &= A + caa^\top, \\ \tilde{\boldsymbol{\Sigma}}_{-k,k} &= \tilde{a}, \end{aligned}$$

and

$$\tilde{\boldsymbol{\Sigma}}_{k,k} = c.$$

The distributions of  $A$ ,  $a$ ,  $\tilde{a}$  and  $c$  are given in Proposition 1 (b). The expectations of  $A$ ,  $\tilde{a}$ , and  $c$  can be directly given as

$$\mathbb{E}(A|\hat{y}_{k_{n+1}}) = \frac{q^{n+1}\mathbf{B}_{-k|k}^n}{q^n(b^n - K)}, \quad (53)$$

$$\mathbb{E}(c|\hat{y}_{k_{n+1}}) = \frac{q^{n+1}}{(q^n + 1)(b^n - K)} \left[ \mathbf{B}_{kk}^n + \frac{q^n}{q^n + 1} (\hat{y}_{k_{n+1}} - \theta_k^n)^2 \right], \quad (54)$$

and

$$\begin{aligned} \mathbb{E}(\tilde{a}|\hat{y}_{k_{n+1}}) &= \mathbb{E} \{ \mathbb{E}(\tilde{a}|A, c, \hat{y}_{k_{n+1}})|\hat{y}_{k_{n+1}} \} = \mathbb{E}(c|\hat{y}_{k_{n+1}}) \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \\ &= \frac{q^{n+1}}{(q^n + 1)(b^n - K)} \left[ 1 + \frac{q^n(\hat{y}_{k_{n+1}} - \theta_k^n)^2}{(q^n + 1)\mathbf{B}_{kk}^n} \right] \mathbf{B}_{-k,k}^n. \end{aligned} \quad (55)$$

Now we consider  $\mathbb{E}(caa^\top|\hat{y}_{k_{n+1}})$ . According to the proof of Proposition 1(b)

$$\begin{aligned} \mathbb{E}(caa^\top|\hat{y}_{k_{n+1}}) &= \mathbb{E} \{ \mathbb{E}(caa^\top|A, c, \hat{y}_{k_{n+1}})|\hat{y}_{k_{n+1}} \} \\ &= \mathbb{E} \{ \mathbb{E}(caa^\top|A, c, \hat{y}_{k_{n+1}})|\hat{y}_{k_{n+1}} \} \\ &= \mathbb{E} \left\{ c \left[ \text{Var}(a|A, c, \hat{y}_{k_{n+1}}) + \mathbb{E}(a|A, c, \hat{y}_{k_{n+1}})\mathbb{E}(a^\top|A, c, \hat{y}_{k_{n+1}}) \right] |\hat{y}_{k_{n+1}} \right\} \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left\{ c \left[ \frac{q^n A}{q^{n+1} \mathbf{B}_{k,k}^n} + \frac{\mathbf{B}_{-k,k}^n \mathbf{B}_{k,-k}^n}{(\mathbf{B}_{kk}^n)^2} \right] \middle| \hat{y}_{k_{n+1}} \right\} \\
&= \mathbb{E}(c | \hat{y}_{k_{n+1}}) \left[ \frac{q^n \mathbb{E}(A | \hat{y}_{k_{n+1}})}{q^{n+1} \mathbf{B}_{k,k}^n} + \frac{\mathbf{B}_{-k,k}^n \mathbf{B}_{k,-k}^n}{(\mathbf{B}_{kk}^n)^2} \right] \\
&= \frac{q^{n+1}}{(q^n + 1)(b^n - K)} \left[ 1 + \frac{q^n (\hat{y}_{k_{n+1}} - \boldsymbol{\theta}_k^n)^2}{(q^n + 1) \mathbf{B}_{kk}^n} \right] \left[ \frac{\mathbf{B}_{-k|k}^n}{b^n - K} + \frac{\mathbf{B}_{-k,k}^n \mathbf{B}_{k,-k}^n}{\mathbf{B}_{kk}^n} \right] \tag{56}
\end{aligned}$$

Combining the results in (52) and (44)–(56), we obtain the updating formulas for  $\mathbf{B}^{n+1}$ .  $\square$

## C Proof of Proposition 3

*Proof.* We decompose the density function of  $\tilde{\Sigma}$  to

$$\xi(\tilde{\Sigma}) \propto \xi^0 \xi(A) \xi(a|A) \xi(c), \tag{57}$$

where  $A$ ,  $a$  and  $c$  are defined in Proposition 1, and

$$\begin{aligned}
\xi^0 &= |\mathbf{B}|^{b^{n+1}/2} = |\mathbf{B}_{-k|k}|^{b^{n+1}/2} \cdot \mathbf{B}_{kk}^{b^{n+1}/2}, \\
\xi(a|A) &= \exp \left\{ -\frac{1}{2} \left( a - \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}} \right)^\top \left( \frac{A}{\mathbf{B}_{kk}} \right)^{-1} \left( a - \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}} \right) \right\}, \\
\xi(c) &= c^{-\frac{b^{n+1}+K+1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{B}_{kk} c^{-1} \right\},
\end{aligned}$$

and

$$\xi(A) = |A|^{-\frac{b^{n+1}+K+1}{2}} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{B}_{-k|k} A^{-1}) \right\}.$$

According to the properties of the Inverse-Wishart distribution, we have

$$a|A \sim N_{K-1}(\mathbf{B}_{-k,k}/\mathbf{B}_{kk}, A/\mathbf{B}_{kk}), \tag{58}$$

$$c \sim IW_1(\mathbf{B}_{kk}, b^{n+1} - K + 1), \tag{59}$$

and

$$A \sim IW_{K-1}(\mathbf{B}_{-k|k}, b^{n+1}). \tag{60}$$

According to Proposition 2, we have that the variance of  $\boldsymbol{\mu}|\Sigma, \hat{y}_k^{n+1}$  is  $\tilde{\Sigma}$ . According to the proof of Proposition 2, the density function of  $\tilde{\Sigma}|\hat{y}_k^{n+1}$  can be decomposed by

$$p^{n+1}(\tilde{\Sigma}|\hat{y}_k^{n+1}) \propto p(a|A, \hat{y}_k^{n+1}) p(A|\hat{y}_k^{n+1}) p(c|\hat{y}_k^{n+1}), \tag{61}$$

where

$$\begin{aligned}
p(a|A, \hat{y}_k^{n+1}) &= \exp \left\{ -\frac{q^{n+1} \mathbf{B}_{k,k}^n}{2q^n} \left( a - \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{k,k}^n} \right)^\top A^{-1} \left( a - \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{k,k}^n} \right) \right\}, \\
p(c|\hat{y}_k^{n+1}) &= c^{-\frac{b^n+K+2}{2}} \exp \left\{ -\frac{q^{n+1}}{2(q^n + 1)} \left[ \mathbf{B}_{kk}^n + \frac{q^n}{q^n + 1} (\hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n)^2 \right] c^{-1} \right\},
\end{aligned}$$

and

$$p(A|\hat{y}_k^{n+1}) = |A|^{-\frac{b^n+K+1}{2}} \exp \left\{ -\frac{q^{n+1}}{2q^n} \text{tr}(\mathbf{B}_{-k|k}^n A^{-1}) \right\}.$$

Notice that  $p(A|\hat{y}_k^{n+1})$ ,  $p(c|\hat{y}_k^{n+1})$ , and  $p(a|A, \hat{y}_k^{n+1})$  are not necessarily the density functions of  $A$ ,  $c$  and  $a|A$ .

Therefore, we have

$$D_{KL}(\mathbf{B}) = \log \xi^0 + \text{E} \log \frac{\xi(A)}{p(A|\hat{y}_{k_{n+1}})} + \text{E} \log \frac{\xi(c)}{p(c|\hat{y}_{k_{n+1}})} + \text{E} \log \frac{\xi(a|A)}{p(a|A, \hat{y}_{k_{n+1}})}. \quad (62)$$

We now derive the terms in (62) one by one.

First,

$$\log \xi^0 = \frac{b^{n+1}}{2} \log |\mathbf{B}_{-k|k}| + \frac{b^{n+1}}{2} \log \mathbf{B}_{kk}. \quad (63)$$

Second, according to the Inverse Wishart distribution of  $c$ , we have

$$\text{E} \log c^{-1} \propto \log \mathbf{B}_{kk}^{-1}$$

and

$$\text{E} c = (b^{n+1} - K + 1) \mathbf{B}_{kk}.$$

Thus, we obtain

$$\text{E} \log \frac{\xi(c)}{p(c|\hat{y}_{k_{n+1}})} \propto \frac{b^n - b^{n+1} + 1}{2} \log \mathbf{B}_{kk} + \frac{q^{n+1}(b^{n+1} - K + 1) \mathbf{B}_{kk}^n}{2(q^n + 1) \mathbf{B}_{kk}} \left[ 1 + \frac{q^n(\hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n)^2}{(q^n + 1) \mathbf{B}_{kk}^n} \right]. \quad (64)$$

Third, according to the Inverse Wishart distribution of  $A$ , we have

$$\text{E} \log |A| \propto \log |\mathbf{B}_{-k|k}|$$

and

$$\text{E} A^{-1} = b^{n+1} \mathbf{B}_{-k|k}^{-1}.$$

Thus, we obtain that

$$\text{E} \log \frac{\xi(A)}{p(A|\hat{y}_{k_{n+1}})} \propto \frac{b^n - b^{n+1}}{2} \log |\mathbf{B}_{-k|k}| + \frac{q^{n+1} b^{n+1}}{2q^n} \text{tr}(\mathbf{B}_{-k|k}^n \mathbf{B}_{-k|k}^{-1}). \quad (65)$$

Lastly, according to the multivariate normal distribution of  $a|A$ , we have

$$\text{E} \log \frac{\xi(a|A)}{p(a|A, \hat{y}_{k_{n+1}})} \propto \left( \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}} - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right)^\top \mathbf{B}_{-k|k}^{-1} \left( \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}} - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right). \quad (66)$$

Combine (63)–(66), the objective function can be expressed as

$$D_{KL}(\mathbf{B}) = \frac{b^n + 1}{2} \log \mathbf{B}_{kk} + \frac{q^{n+1}(b^{n+1} - K + 1) \mathbf{B}_{kk}^n}{2(q^n + 1) \mathbf{B}_{kk}} \left[ 1 + \frac{q^n(\hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n)^2}{(q^n + 1) \mathbf{B}_{kk}^n} \right] \quad (67a)$$

$$+ \frac{b^n}{2} \log |\mathbf{B}_{-k|k}| + \frac{q^{n+1} b^{n+1}}{2q^n} \text{tr}(\mathbf{B}_{-k|k}^n \mathbf{B}_{-k|k}^{-1}) \quad (67b)$$

$$+ \frac{q^{n+1} \mathbf{B}_{kk}^n}{2q^n} \left( \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}} - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right)^\top \mathbf{B}_{-k|k}^{-1} \left( \frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}} - \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n} \right). \quad (67c)$$

We next minimize  $D_{KL}(\mathbf{B})$  with respect to  $\mathbf{B}$ . It is clear that this can be done by minimizing  $D_{KL}(\mathbf{B})$  with respect to  $\mathbf{B}_{kk}$ ,  $\frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}}$ , and  $\mathbf{B}_{-k|k}$ . We first observe that only (67c) involves term  $\frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}}$ . For any fixed  $\mathbf{B}_{kk}$  and  $\mathbf{B}_{-k|k}$ , the minimizer of (67c) is given by  $(\frac{\mathbf{B}_{-k,k}}{\mathbf{B}_{kk}})^* = \frac{\mathbf{B}_{-k,k}^n}{\mathbf{B}_{kk}^n}$ , and the corresponding minimum of (67c) is 0. We then notice that (67a) only involves  $\mathbf{B}_{kk}$ , and (67b) only involves  $\mathbf{B}_{-k|k}$ , by optimizing (67a) and (67b) with respect to  $\mathbf{B}_{kk}$  and  $\mathbf{B}_{-k|k}$ , respectively, we get:

$$(\mathbf{B}_{-k|k})^* = \frac{b^{n+1}q^{n+1}}{b^n q^n} \mathbf{B}_{-k|k}^n,$$

$$(\mathbf{B}_{k,k})^* = \frac{q^{n+1}(b^{n+1} - K + 1) \left[ \mathbf{B}_{k,k}^n + \frac{q^n}{q^{n+1}-1} (\hat{y}_k^{n+1} - \boldsymbol{\theta}_k^n)^2 \right]}{(b^n + 1)(q^n + 1)}.$$

□