*Article*

# Variational Bayesian Approximation (VBA): A comparison between three optimization algorithms

**Azadeh Fallah[1] and Ali Mohammad-Djafari [2,3,\*] orcid number:0000-0003-0678-7759**

[1]    Ferdowsi University of Mashhad, Mashhad, Iran; azadeh.fallah@mail.um.ac.ir
[2]    International Science Consulting and Training (ISCT), 91440 Bures sur Yvette, France; djafari@ieee.org
[3]    Scientific leader of Shanfeng company, Shaoxing, China
       Correspondence: azadeh.fallah@mail.um.ac.ir

**Abstract:** In many Bayesian computations, first, we obtain the expression of the joint distribution of all the unknown variables given the observed data. In general, this expression is not separable in those variables. Thus, obtaining their marginals for each variable and computing the expectations are difficult and costly. This problem becomes even more difficult in high dimensional quandaries, which is an important issue in inverse problems. We may then try to propose a surrogate expression with which we can do approximate computations. Often a separable expression approximation can be useful enough. The Variational Bayesian Approximation (VBA) is a technique that approximates the joint distribution $p$ with an easier, for example separable, one $q$ by minimizing Kullback–Leibler Divergence $KL(q|p)$. When $q$ is separable in all the variables, the approximation is also called Mean Field Approximation (MFA) and so $q$ is the product of the approximated marginals. A first standard and general algorithm is alternate optimization of $KL(q|p)$ with respect to $q_i$. A second general approach is its optimization in the Riemannian manifold. However, in this paper, for practical reasons, we consider the case where $p$ is in the exponential family and so is $q$. For this case, $KL(q|p)$ becomes a function of the parameters $\boldsymbol{\theta}$ of the exponential family. Then, we can use any other optimization algorithm to obtain those parameters. In this paper, we compare three optimization algorithms: standard alternate optimization, a gradient-based algorithm and a natural gradient algorithm and study their relative performances on three examples.

**Keywords:** Variational Bayesian Approach (VBA); Kullback–Leibler Divergence; Mean Field Approximation (MFA); Optimization Algorithm

## 1. Introduction

In many applications, with direct or undirect observations, the use of the Bayesian computations starts with obtaining the expression of the joint distribution of all the unknown variables given the observed data. Then, we have to use it to do inference. In general, this expression is not separable in all the variables of the problem. So, the computations becomes hard and costly. For example, obtaining the marginals for each variable and computing the expectations are difficult and costly. This problem becomes even more crutial in high dimensional quandaries, which is an important issue in inverse problems. We may then need to propose a surrogate expression with which we can do approximate computations.

The Variational Bayesian Approximation (VBA) is a technique that approximates the joint distribution $p$ with an easier, for example a separable one $q$ by minimizing Kullback–Leibler Divergence $KL(q|p)$, which makes the marginal computations much easier. For example, in the case of two

32 variables, $p(x, y)$ is approximated by $q(x, y) = q_1(x)q_2(y)$ via minimizing $KL(q_1q_2|p)$. When $q$ is

33 separable in all the variables of $p$, the approximation is also called Mean Field Approximation (MFA).

34 To obtain the approximate marginals $q_1$ and $q_2$ we have to minimize $KL(q_1q_2|p)$. A first standard

35 and general algorithm is alternate optimization of $KL(q_1q_2|p)$ with respect to $q_1$ and $q_2$. Finding the

36 expression of the functional derivatives of $KL(q_1q_2|p)$ with respect to $q_1$ and $q_2$ and equating them

37 to zero alternatively, we obtain an iterative optimization algorithm. A second general approach is its

38 optimization in the Riemannian manifold. However, in this paper, for practical reasons, we consider

39 the case where $p$ is in the exponential family and so are $q_1$ and $q_2$. For this case, $KL(q_1q_2|p)$ becomes

40 a function of the parameters $\theta$ of the exponential family. Then, we can use any other optimization

41 algorithm to obtain those parameters.

42 In this paper, we compare three optimization algorithms: standard alternate optimization (Algorithm 1),

43 a gradient-based algorithm [1,2] (Algorithm 2) and a natural gradient algorithm [3–5] (Algorithm 3).

44 Between the main advantages of the VBA for inference problems, such as inverse problems and

45 machine learning, we can mention the following:

46 - First, VBA builds a sufficient model according to prior information and the final posterior distribution.

47 Especially in the Mean-Field Approximation (MFA), the result ends in an explicit form for each

48 unknown component using conjugate priors, and works well for small sample sizes, [6–8].

49 - The second benefit is, for example in machine learning, a robust way for classification based on the

50 predictive posterior distribution and diminishes parameter over-trained, [7].

51 - The third privilege is that the target structure has uncertainty in the VBA recursive processes. This

52 feature prevents further error propagation and increases the robustness of VBA, [9].

53 Besides all these preponderance, the VBA has some weaknesses, such as difficulty in the solution

54 of integrals and expectations to get a posterior distribution, and there is no evidence to find an

55 exact posterior, [6]. Its most significant drawback arises when there are strong dependencies

56 between unknown parameters, and the VBA ignores them. Then estimates, computed based on

57 this approximation, may be very far from the exact values. However, it works well when the amounts

58 of dependencies are low [8].

59 In this article, we work on three different estimating algorithms of the unknown parameters in a

60 model concerning prior information. The first iterative algorithm is standard alternate optimization

61 based on VBA, which start from some initial points. Sometimes, the points are estimated from an

62 available dataset, but most of the time, we do not have enough data on the parameters to make some

63 pre-estimation of them. To solve this obstacle, we can start the algorithm with some desired points,

64 and then by repeating the process, they approach the true values using the posterior distribution. The

65 second two algorithms are gradient-based and natural gradient algorithms, whose based function

66 are Kullback–Leibler divergence. First, the gradient of Kullback–Leibler for all unknown parameters

67 have to be found, then start from some points either estimated from data or desired choices. Then, we

68 repeat the iterative algorithm till it convergences to some points. If we denote the unknown parameter

69 space with $\theta$, then the recursive formula is $\tilde{\theta}^{(k+1)} = \tilde{\theta}^{(k)} - \Delta KL(\tilde{\theta}^{(k)})$ for gradient-based algorithm.

70 The formula for natural gradient is pretty similar and is $\tilde{\theta}^{(k+1)} = \tilde{\theta}^{(k)} - \frac{\Delta KL(\tilde{\theta}^{(k)})}{\|\Delta KL(\tilde{\theta}^{(k)})\|}$.

71 Also, we consider three examples, Normal-Inverse-Gamma, multivariate Normal, and linear

72 inverse problem for checking the performance and convergence speed of the algorithms.

73 We propose the following organization of this paper: In section 2, we present a brief

74 explanation of the basic VBA analytical part. In section 3, we explain our first example related

75 to Normal-Inverse-Gamma distribution analytically and, in practice, explain the outcomes of three

76 algorithms to estimate the unknown parameters. In section 4, we study a more complex example of a

77 multivariate Normal distribution whose means and variance-covariance matrix are unknown and have

78 Normal Inverse Wishart distribution. The aim of this section is to demonstrate marginal distributions

79 of $\tilde{\mu}$ and $\tilde{\Sigma}$ using a set of multivariate Normal observations using these mean and variance. In section

80 5, the example is more close to realistic situations and is a linear inverse problem. In section 6, we

present a summary of the work done in the article and compare the three recursive algorithms in three different examples.

## 2. Variational Bayesian Approach (VBA)

As we mentioned earlier, VBA uses Kullback–Leibler Divergence. Kullback–Leibler Divergence [10] $KL(q|p)$ is an information measure of discrepancy between two probability functions defined as follow. Let $p(x)$ and $q(x)$ be two density functions of a continuous random variable $x$ respect to support set $\mathbb{S}_X$. $KL(q|p)$ function is introduced as

$$KL(q|p) = \int_{x \in \mathbb{S}_X} q(x) \ln \frac{q(x)}{p(x)} \, dx. \tag{1}$$

For simplicity, we assume a bivariate case of distribution $p(x,y)$, and want to assess it via VBA, then we have:

$$KL(q|p) = -H(q_1) - H(q_2) - < \ln \ p(x,y) >_{q_1 q_2}, \tag{2}$$

where

$$H(q_1) = -\int_{x \in \mathbb{S}_X} q_1(x) \ln \ q_1(x) \, dx \quad \text{and} \quad H(q_2) = -\int_{y \in \mathbb{S}_Y} q_2(y) \ln \ q_2(y) \, dy$$

are, respectively, the Shannon entropies of $x$ and of $y$, and

$$< \ln p(x,y) >_{q_1 q_2} = \int \int_{(x,y) \in \mathbb{S}_{XY}} q_1(x) q_2(y) \ln \ p(x,y) \, dx \, dy.$$

Now, differentiating the equation (2) with respect to $q_1$, and then with respect to $q_2$ and equating them to zero, we obtain:

$$q_1(x) \propto \exp \left\{ < \ln \ p(x,y) >_{q_2(y)} \right\} \quad \text{and} \quad q_2(y) \propto \exp \left\{ < \ln \ p(x,y) >_{q_1(x)} \right\} \tag{3}$$

These results can be easily extended to more dimensions [11]. They do not have any closed form, because they depend on the expression of $p(x,y)$ and those of $q_1$ and $q_2$. An interesting case is the case of exponential families and conjugate priors, where writing

$$p(x,y) = p(x|y)p(y), \text{ and } p(y|x) = \frac{p(x,y)}{p(x)} = \frac{p(x|y)p(y)}{p(x)}, \tag{4}$$

we can consider $p(y)$ as prior, $p(x|y)$ as the likelihood, and $p(y|x)$ as the posterior distributions. Then, if $p(y)$ is a conjugate prior for the likelihood $p(x|y)$, then the posterior $p(y|x)$ will be in the same family as the prior $p(y)$. To illustrate all these properties, we give details of these expressions for a first simple example of Normal-Inverse-Gamma $p(x,y) = \mathcal{N}(x|\mu,y)\mathcal{IG}(y|\alpha,\beta)$ with $q_1(x) = \mathcal{N}(x|\mu,v)$ and $q_2 = \mathcal{IG}(y|\alpha,\beta)$. For this simple case, first we give the expression of $KL(q|p)$ with $q_1(x) = \mathcal{N}(x|\widetilde{\mu},\widetilde{v})$ and $q_2(y) = \mathcal{IG}(y|\widetilde{\alpha},\widetilde{\beta})$ as a function of the parameters $\theta = (\widetilde{\mu},\widetilde{v},\widetilde{\alpha},\widetilde{\beta})$ and then the expressions of the three above-mentioned algorithms and we study their convergence.

## 3. Normal-Inverse-Gamma Distribution Example

The purpose of this section is to explain in detail the process of performing calculations in VBA. For this we consider a simple case for which, we have all the necessary expressions. The objective here is to compare the three different algorithms mentioned earlier. Also, its practical application can be explained as follows:
We have a sensor which measures a quantity $X$, N times $x_1, ..., x_N$. We want to model these data. In a first step, we model it as $N(x|\mu,v)$ with fixed $\mu$ and $v$. Then, it is easy to estimate the parameters $(\mu,v)$ either by Maximal Likelihood or Bayesian strategy. If we assume that the model is Gaussian with unknown variance and call this variance $y$ and assign an $\mathcal{IG}$ prior to it, then we have a model

100   $\mathcal{NIG}$ for $p(x,y)$. We showed that the margins are $St$ and $\mathcal{IG}$. Working directly with $St$ is difficult. So,
101   we want to approximate it with a Gaussian $q_1(x)$. This is equivalent to approximating $p(x,y)$ with
102   $q_1(x)q_2(y)$. Now, we want to find the parameters $\mu$, $v$, $\alpha$, and $\beta$, which minimize $KL(q_1q_2|p)$. This
103   process is called VBA. Then, we want to compare three algorithms to obtain the parameters which
104   minimize $KL(\cdot|\cdot)$. $KL(\cdot|\cdot)$ is convex with respect to $q_1$ if $q_2$ is fixed and is convex with respect to $q_2$ if
105   $q_1$ is fixed. So, we hope that the iterative algorithm converges. However, $KL(\cdot|\cdot)$ may not be convex in
106   the space of parameters. So, we have to study the shape of this criterion concerning the parameters $\widetilde{v}$,
107   $\widetilde{\alpha}$, and $\widetilde{\beta}$.

The practical problem considered here is in the following: A sensors delivers a few samples
$x = \{x_1, x_2, \cdots, x_N\}$ of a physical quantity $X$. We want to find $p(x)$. For this process, we assume a
simple Gaussian model, but with unknown variance $y$. So that, the forward model can be written as
$p(x,y) = \mathcal{N}(x|\mu, y)\mathcal{IG}(y|\alpha, \beta)$. In this simple example, we know that $p(x)$ is a Student-t distribution
obtained by:

$$S(x|m, \alpha, \beta) = \int \mathcal{N}(x|\mu, y)\mathcal{IG}(y|\alpha, \beta)\, \mathrm{d}y \tag{5}$$

108   Our objective is to find the three parameters $\theta = (\mu, \alpha, \beta)$ from the data $x$ and an approximate marginal
109   $q(x)$ for $p(x)$.

The main idea is to find such $q_1(x)q_2(y)$ as an approximation of $p(x,y)$. Here, we show the
VBA, step by step. For this, we start by choosing the conjugate families $q_1(x) = \mathcal{N}(x|\widetilde{\mu}, \widetilde{v})$ and
$q_2(y) = \mathcal{IG}(y|\widetilde{\alpha}, \widetilde{\beta})$.
In the first step, we have to calculate $\ln p(x,y)$

$$\ln p(x,y) = c - \frac{1}{2}\ln y - \frac{1}{2y}(x - \widetilde{\mu})^2 - (\widetilde{\alpha} + \frac{1}{2})\ln y - \frac{\widetilde{\beta}}{y}. \tag{6}$$

where $c$ is a constant value term independent of $x$ and $y$. First of all, we have to start by finding $q_2$, so
the integration of $\ln p(x,y)$ is concerning $q_1$

$$< \ln\ p(x,y) >_{q_1} = c - \frac{1}{2y} < (x - \widetilde{\mu})^2 >_{q_1} - (\widetilde{\alpha} + 1)\ln\ y - \frac{\widetilde{\beta}}{y}. \tag{7}$$

Since the mean of $x$ is the same in prior and posterior distribution, $\widetilde{\mu} = \widetilde{\mu}'$, and then $< (x - \widetilde{\mu})^2 >_{q_1} = \widetilde{v}$.
Thus

$$q_2(y) \propto \exp\left[-(\widetilde{\alpha} + 1)\ln\ y - (\frac{\widetilde{v}}{2} + \widetilde{\beta})\frac{1}{y}\right]. \tag{8}$$

Thus, the function $q_2(y)$ is equivalent to an inverse gamma distribution $\mathcal{IG}(\widetilde{\alpha}, \frac{\widetilde{v}}{2} + \widetilde{\beta})$. Now, we have
to take integral of $\ln\ p(x,y)$ over $q_2$ to find $q_1$

$$< \ln\ p(x,y) >_{q_2} = c - (\widetilde{\alpha} + 1) < \ln\ y >_{q_2} - (\widetilde{\beta} + \frac{1}{2}(x - \widetilde{\mu})^2) < \frac{1}{y} >_{q_2} \tag{9}$$

Note that the first term does not depend on $x$ and in the second term we have $< \frac{1}{y} >_{q_2} = \frac{2\widetilde{\alpha}}{2\widetilde{\beta}+\widetilde{v}}$, so

$$q_1(x) \propto \exp\left[-\frac{2\widetilde{\alpha}}{2\widetilde{\beta} + \widetilde{v}}(\widetilde{\beta} + \frac{1}{2}(x - \widetilde{\mu})^2)\right] \propto \exp\left[-\frac{(x - \widetilde{\mu})^2}{2\frac{2\widetilde{\beta}+\widetilde{v}}{2\widetilde{\alpha}}}\right]. \tag{10}$$

110   We see that $q_1$ is again a Normal distribution but with updated parameters $\mathcal{N}(\widetilde{\mu}, \frac{2\widetilde{\beta}+\widetilde{v}}{2\widetilde{\alpha}})$, so $\widetilde{v} = \frac{2\widetilde{\beta}+\widetilde{v}}{2\widetilde{\alpha}}$.
111   Note that, we obtained the conjugacy property: If $p(x|y) = \mathcal{N}(x|\mu, y)$ and $p(y) = \mathcal{IG}(y|\alpha, \beta)$, then
112   $p(y|x) = \mathcal{IG}(y|\alpha', \beta')$ where $\mu'$, $\alpha'$ and $\beta'$ are $\mu' = \mu$, $\alpha' = \alpha$ , $\beta' = \beta + \frac{2\beta+v}{4\alpha}$. In this case, we also know
113   that $p(x|\alpha, \beta) = St(x|\mu', \alpha, \beta)$.

In standard alternate optimization (algorithm 1), there is no need for an iterative process for $\widetilde{\mu}$ and $\widetilde{\alpha}$, which are approximated by $\widetilde{\mu} = \mu_0$ and $\widetilde{\alpha} = \alpha_0$, respectively. The situation for $\widetilde{\beta}$ and $\widetilde{v}$ is different because there are circular dependencies among them. So, the approximation needs an iterative process, staring from $\widetilde{\mu}^{(1)} = \mu_0$, $\widetilde{v}^{(1)} = v_0$, $\alpha^{(1)} = \alpha_0$, and $\beta^{(1)} = \beta_0$. The main conclusion for this case is that the VBA algorithm becomes:

**Algorithm 1:**

$$
\begin{aligned}
\widetilde{\alpha}^{(k+1)} &= \widetilde{\alpha}^{(k)}, \\
\widetilde{\beta}^{(k+1)} &= \widetilde{\beta}^{(k)} + \frac{\widetilde{v}^{(k)}}{2}, \\
\widetilde{\mu}^{(k+1)} &= \widetilde{\mu}^{(k)}, \\
\widetilde{v}^{(k+1)} &= \frac{2\widetilde{\beta}^{(k)} + \widetilde{v}^{(k)}}{2\widetilde{\alpha}^{(k)}}.
\end{aligned}
$$

For the two other algorithms, gradient and natural gradient-based, it requires to find the expression of $KL(q_1 q_2 : p)$ as a function of the parameters $\boldsymbol{\theta} = (\alpha, \beta, \mu, v)$

$$
KL(\tilde{\boldsymbol{\theta}}) = 2\ln \Gamma(\tilde{\alpha} - \frac{1}{2}) - (2\tilde{\alpha} + \frac{3}{2})\psi_0(\tilde{\alpha} - \frac{1}{2}) + \frac{(2\tilde{\alpha} - 1)(\tilde{v} + 2\tilde{\beta})}{4\tilde{\beta}} + \tilde{\alpha} + \frac{5}{2}\ln \tilde{\beta} - \frac{1}{2}\ln \tilde{v} - 1, \quad (11)
$$

Then, we also need the gradient expression of $\nabla KL(\widetilde{\boldsymbol{\theta}})$ for $\widetilde{\boldsymbol{\theta}}$

$$
\nabla KL(\tilde{\boldsymbol{\theta}}) = \left( \frac{\widetilde{v}}{2\widetilde{\beta}} - (2\widetilde{\alpha} + \frac{3}{2})\psi_1(\widetilde{\alpha} - \frac{1}{2}) + 2, \quad -(2\widetilde{\alpha} - 1)\frac{\widetilde{v}}{4\widetilde{\beta}^2} + \frac{5}{2\widetilde{\beta}}, \quad 0, \quad \frac{2\widetilde{\alpha} - 1}{4\widetilde{\beta}} - \frac{1}{2\widetilde{v}} \right), \quad (12)
$$
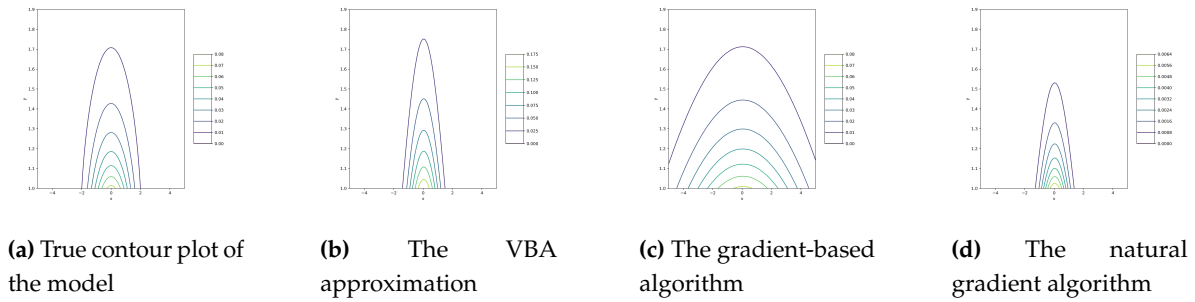
As we can see these expressions do not depend on $\widetilde{\mu}$, so their derivatives are zero. The means are preserved.

**Algorithm 2 and 3:**

$$
\begin{aligned}
\widetilde{\alpha}^{(k+1)} &= \widetilde{\alpha}^{(k)} - \gamma \left( \frac{\widetilde{v}^{(k)}}{2\widetilde{\beta}^{(k)}} - (2\widetilde{\alpha}^{(k)} + \frac{3}{2})\psi_1(\widetilde{\alpha}^{(k)} - \frac{1}{2}) + 2 \right), \\
\widetilde{\beta}^{(k+1)} &= \widetilde{\beta}^{(k)} - \gamma \left( -(2\widetilde{\alpha}^{(k)} - 1)\frac{\widetilde{v}^{(k)}}{4[\widetilde{\beta}^{(k)}]^2} + \frac{5}{2\widetilde{\beta}^{(k)}} \right), \\
\widetilde{\mu}^{(k+1)} &= \widetilde{\mu}^{(k)}, \\
\widetilde{v}^{(k+1)} &= \widetilde{v}^{(k+1)} - \gamma \left( \frac{2\widetilde{\alpha}^{(k)} - 1}{4\widetilde{\beta}^{(k)}} - \frac{1}{2\widetilde{v}^{(k)}} \right),
\end{aligned}
$$

where $\gamma$ is fixed for the gradient algorithm, and is proportional to $1/\|\nabla KL\|$ for the natural gradient algorithm. To start the numerical computations, we generate $n = 100$ samples from the model $p(x, y) = \mathcal{N}(x|1, y)\mathcal{IG}(y|3, 1)$. Thus, we know the exact values of the unknown parameters, just keeping in mind, not used in algorithms. $\tilde{\boldsymbol{\theta}}_1$, $\tilde{\boldsymbol{\theta}}_1$, and $\tilde{\boldsymbol{\theta}}_1$ are the estimated parameters using the alternative, gradient-based, and natural gradient algorithms along with their surface and contour plots in fig 1, respectively

$$
\tilde{\boldsymbol{\theta}}_1 = \begin{cases} \widetilde{\mu} = 0.044724, \\ \widetilde{v} = 0.590894, \\ \widetilde{\alpha} = 3.792337, \\ \widetilde{\beta} = 1.765735 \end{cases} \quad
\tilde{\boldsymbol{\theta}}_2 = \begin{cases} \widetilde{\mu} = 0.044724, \\ \widetilde{v} = 7.910504, \\ \widetilde{\alpha} = 4.991621, \\ \widetilde{\beta} = 3.002594 \end{cases} \quad
\tilde{\boldsymbol{\theta}}_3 = \begin{cases} \widetilde{\mu} = 0.044724, \\ \widetilde{v} = 0.423761, \\ \widetilde{\alpha} = 4.415239, \\ \widetilde{\beta} = 0.706049. \end{cases}
$$

**(a)** True contour plot of the model

**(b)** The VBA approximation

**(c)** The gradient-based algorithm

**(d)** The natural gradient algorithm

**Figure 1.** The true model has $\mathcal{N}(x|0,y)\mathcal{IG}(y|3,1)$. The numbers of convergence are diverse in the algorithms.

All three algorithms try to minimize the same criterion. So, the objectives are all the same, but the number of steps may differ. The requirements must reach the minimum $KL(\cdot)$. In this simple example of the Normal-Inverse-Gamma distribution, the convergence step numbers of VBA, gradient-based, and natural gradient are 1, 2, and 1. The overall performance of the standard alternate optimization (VBA) is more precise than any others. The poorest estimation is from a gradient-based algorithm. So, the algorithms are able to approximate the joint density function with a separable one, but with different accuracy. In the following section, we will tackle a more complex model.

## 4. Multivariate Normal Inverse Wishart Example

In previous section, we explain how to preform VBA to approximate a complicated joint distribution function by tractable marginal factorials over a simple case study. In this section, a multivariate Normal case $p(x) = \mathcal{N}(x|\widetilde{\mu}, \widetilde{\Sigma})$ is considered which is approximated by $q(x) = \prod_i \mathcal{N}(x_i|\widetilde{\mu}_i, \widetilde{v}_i)$ for different shapes for the covariance matrix $\widetilde{\Sigma}$.

We assume that the basic structure of an available data set is multivariate Normal with unknown mean vector $\widetilde{\mu}$ and variance-covariance matrix $\widetilde{\Sigma}$. Their joint prior distribution is a Normal Inverse Wishart distribution of $\mathcal{NIW}(\widetilde{\mu}, \widetilde{\Sigma}|\widetilde{\mu}_0, \tilde{\kappa}, \widetilde{\Psi}, \tilde{v})$. The posteriors are multivariate Normal for mean vector and Inverse Wishart for variance-covariance matrix. Since the Normal Inverse Wishart distribution is a conjugate prior distribution for multivariate Normal, the posterior distribution of $\widetilde{\mu}$ and $\widetilde{\Sigma}$ is again belong to the same family, and their corresponding margins are

$$\mathcal{MN}\left(\frac{\tilde{\kappa}\widetilde{\mu}_0 + n\bar{x}}{\tilde{\kappa} + n}, \frac{1}{\tilde{\kappa} + n}\tilde{\Lambda}\right), \quad \mathcal{IW}\left(\tilde{\Lambda} + \widetilde{\Psi} + \sum_{i=1}^{n}(x_i - \bar{x})(x_i - \bar{x})^T, \tilde{v} + n\right), \tag{13}$$

where $n$ is the sample size. To present the performance of the three algorithms, we work on a data set coming from $x \sim \mathcal{NIW}(x|\mu, \Sigma)$ whose parameters have the below density structure
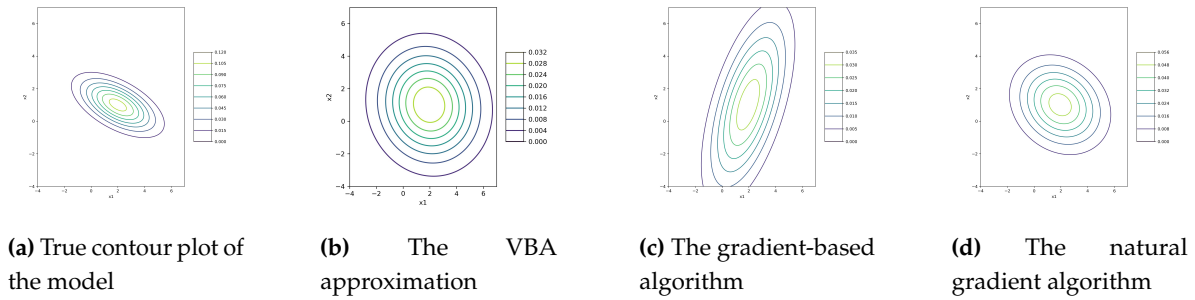
$$\mu \sim \mathcal{MN}\left(\mu\Big| \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \frac{1}{2}\begin{bmatrix} 3 & -1 \\ -1 & 1 \end{bmatrix}\right), \quad \Sigma \sim \mathcal{IW}\left(\Sigma\Big| \begin{bmatrix} 3 & -1 \\ -1 & 1 \end{bmatrix}, 6\right).$$

We use only the data of $x$ in the estimation processes. The results of algorithms are drowned in fig 2 along with the true contour plot of the model. The VBA estimation is the most separable distribution compared with gradient and natural gradient methods. The next best case is the natural gradient algorithm, but its weakness is transferring the dependency a little bit to the approximation. The result for the gradient-based is completely shown the dependency, and its inability to get a separable model.

## 5. Simple Linear Inverse Problem

Finally the third example is the case of linear inverse problems with $g = Hf + \epsilon$ with priors $p(\epsilon) = \mathcal{N}(\epsilon|0, v_\epsilon I)$ and $p(f) = \mathcal{N}(f|0, \text{diag}[v])$ with $f = [f_1, f_2, \cdots, f_N]$ and $v = [v_{f_1}, v_{f_2}, \cdots, v_{f_N}]$

**(a)** True contour plot of the model

**(b)** The VBA approximation

**(c)** The gradient-based algorithm

**(d)** The natural gradient algorithm

**Figure 2.** The number of steps to get the result for 2b, 2c, and 2d are 2, 280, and 4, respectively. In this example, the convergence numbers are different between three algorithms, in the gradient-based. Although the gradient-based is not able to approximate the joint density function with a separable one, the VBA and natural gradient algorithms estimate the distribution with separable ones in fewer steps.

for which we have $p(f, v|g) \propto p(g|f, v_\epsilon)p(f|v)p(v)$ with $p(g|f, v_\epsilon) = \mathcal{N}(g|Hf, v_\epsilon I)$, $p(f|v) = \mathcal{N}(f|0, \text{diag}[v])$ and $p(v|\alpha, \beta) = \prod_j \mathcal{IG}(v_j|\alpha, \beta)$, [12]. Thus, the joint distribution of $g$, $f$, and $v$ is estimated by VBA as the following relation

$$p(g, f, v) \propto p(g|f)p(f|v)p(v). \tag{14}$$

Although we can mathematically compute the margins in this particular example, we desire to approximate them via the iterative VBA algorithm $q(g, \widetilde{f}, \tilde{v}) = q_1(g|\widetilde{f})q_2(\widetilde{f})q_3(\tilde{v})$ compared them with gradient and natural gradient-based algorithms. The objective function is the estimation of $q_2(\widetilde{f})$, but in the recursive process, $q_1(g|\widetilde{f})$ and $q_3(\tilde{v})$ are updated, too. For simplicity, we suppose that the transposition matrix $H$ is an identical matrix $I$. The final outputs are as follows

$$\begin{aligned}
\widetilde{f} &\sim \mathcal{MN}\left( \frac{\widetilde{\mu}_{\widetilde{f}}}{1 + \frac{2\tilde{v}_\epsilon \check{\alpha}}{n(\check{v}+\widetilde{\mu}_{\widetilde{f}}^2+2\check{\beta})}}, \text{diag}\left[ \frac{\tilde{v}_\epsilon(\check{v}+\widetilde{\mu}_{\widetilde{f}}^2+2\check{\beta})}{n(\check{v}+\widetilde{\mu}_{\widetilde{f}}^2)+2n\check{\beta}+2\tilde{v}_\epsilon\check{\alpha}} \right] \right), \\
g &\sim \mathcal{N}(\widetilde{\mu}_{\widetilde{f}}, \tilde{v}_\epsilon I), \quad \text{and} \quad \tilde{v}_k \sim \mathcal{IG}\left( \widetilde{\alpha}_k, \frac{\tilde{v}_k + \widetilde{\mu}_{\widetilde{f}_k}^2}{2} + \widetilde{\beta}_k \right), \; k = 1, \cdots, p.
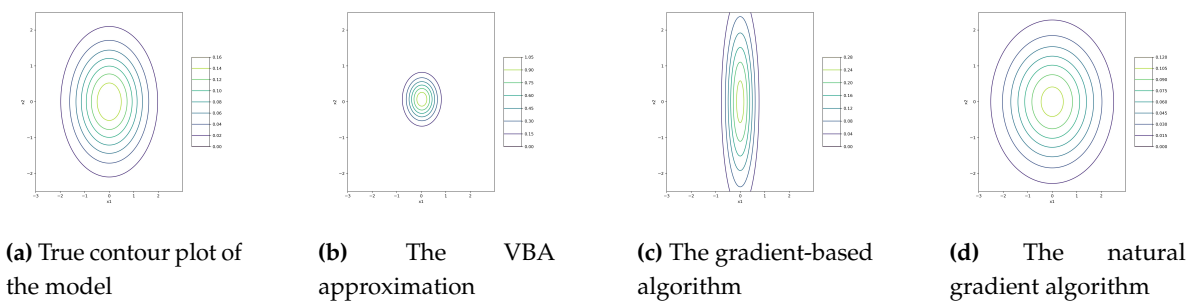\end{aligned} \tag{15}$$

We choose a model to see the performance of these margins and compare them with gradient and natural gradient algorithms. The selected model is $g = Hf + \epsilon$ with the following knowledge

$$H = I, \quad f \sim \mathcal{MN}(f|0, \text{diag}[v_1, v_2]), \quad v_1 \sim \mathcal{IG}(v_1|3, 2), \quad v_2 \sim \mathcal{IG}(v_2|4, 3), \quad \epsilon \sim \mathcal{MN}(\epsilon|0, I).$$

In the assessment procedure, we do not apply the above information. The output of algorithms are shown in fig 3, as well as the actual contour plot. In this example, the best diagnosis is from the natural gradient algorithm. The VBA by construction is separable and cannot be the same as the original.

## 6. Conclusion

This paper presents an approximation method of the unknown density functions for hidden variables called VBA. It is compared with gradient and natural gradient algorithms. We also consider three examples Normal Inverse Gamma, Normal Inverse Wishart, and linear inverse problem. We put the details of the first model here and give the details for two other examples in the whole paper. In all three models, the parameters are unexplored and need to be estimated by recursive algorithms. We try to approximate the joint complex distribution with a simpler version of the margin factorials that looks like independent cases. The VBA and natural gradient converge pretty soon. The major discrepancy in

**(a)** True contour plot of the model

**(b)** The VBA approximation

**(c)** The gradient-based algorithm

**(d)** The natural gradient algorithm

**Figure 3.** The true model is almost separable. The natural gradient algorithm works better in this example, and the poorest approximation is for the gradient-based algorithm.

algorithms comes from the accuracy of the results. They estimate the intricate joint distribution with separable ones. The overall performance of VBA is the best here.

## References

1. Acerbi, L. Variational bayesian monte carlo. *Advances in Neural Information Processing Systems* **2018**, *31*.

2. Kuusela, M.; Raiko, T.; Honkela, A.; Karhunen, J. A gradient-based algorithm competitive with variational Bayesian EM for mixture of Gaussians. 2009 International Joint Conference on Neural Networks. IEEE, 2009, pp. 1688–1695.

3. Gharsalli, L.; Duchêne, B.; Mohammad-Djafari, A.; Ayasso, H. A gradient-like variational Bayesian approach: Application to microwave imaging for breast tumor detection. 2014 IEEE International Conference on Image Processing (ICIP). IEEE, 2014, pp. 1708–1712.

4. Zhang, G.; Sun, S.; Duvenaud, D.; Grosse, R. Noisy natural gradient as variational inference. International Conference on Machine Learning. PMLR, 2018, pp. 5852–5861.

5. Lin, W.; Khan, M.E.; Schmidt, M. Fast and simple natural-gradient variational inference with mixture of exponential-family approximations. International Conference on Machine Learning. PMLR, 2019, pp. 3992–4002.

6. Ghahramani, Z.; Beal, M. Variational inference for Bayesian mixtures of factor analysers. *Advances in neural information processing systems* **1999**, *12*.

7. Watanabe, S.; Minami, Y.; Nakamura, A.; Ueda, N. Variational Bayesian estimation and clustering for speech recognition. *IEEE Transactions on Speech and Audio Processing* **2004**, *12*, 365–381.

8. Montesinos-López, O.A.; Montesinos-López, A.; Crossa, J.; Montesinos-López, J.C.; Luna-Vázquez, F.J.; Salinas-Ruiz, J.; Herrera-Morales, J.R.; Buenrostro-Mariscal, R. A variational Bayes genomic-enabled prediction model with genotype× environment interaction. *G3: Genes, Genomes, Genetics* **2017**, *7*, 1833–1853.

9. Babacan, S.D.; Molina, R.; Katsaggelos, A.K. Variational Bayesian super resolution. *IEEE Transactions on Image Processing* **2010**, *20*, 984–999.

10. Kullback, S.; Leibler, R.A. On information and sufficiency. *The annals of mathematical statistics* **1951**, *22*, 79–86.

11. Fox, C.W.; Roberts, S.J. A tutorial on variational Bayesian inference. *Artificial intelligence review* **2012**, *38*, 85–95.

12. Ayasso, H.; Mohammad-djafari, A. Joint image restoration and segmentation using Gauss-Markov-Potts prior models and variational Bayesian computation. 2009 16th IEEE International Conference on Image Processing (ICIP), 2009, pp. 1297–1300. doi:10.1109/ICIP.2009.5413589.

13. Gupta, M.; Srivastava, S. Parametric Bayesian estimation of differential entropy and relative entropy. *Entropy* **2010**, *12*, 818–843.

**182   Appendix A  Computation of $KL(\tilde{\boldsymbol{\theta}})$ and its Gradient**

The mathematical computations of Kullback–Leibler divergence $KL(\tilde{\boldsymbol{\theta}})$ in section 3 is based on (2) as follows

$$
KL(q_1 q_2 | p) = -H(q_1) - H(q_2) - \int\int q_1(x) q_2(y) \ln \left\{ \frac{1}{\sqrt{2\pi y}} \exp\{-\frac{(x-\tilde{\mu})^2}{2y}\} \right. \tag{A1}
$$

$$
\left. \frac{\widetilde{\beta}^{\tilde{\alpha}-\frac{1}{2}}}{\Gamma(\tilde{\alpha}-\frac{1}{2})} y^{-(\tilde{\alpha}+\frac{1}{2})} \exp\{-\frac{\tilde{\beta}}{y}\} \right\} dxdy
$$

$$
= -H(q_1) - H(q_2) - \int\int q_1(x) q_2(y) \left\{ -\frac{1}{2} \ln (2\pi) - \frac{(x-\tilde{\mu})^2}{2y} + (\tilde{\alpha}-\frac{1}{2}) \ln \tilde{\beta} \right.
$$

$$
\left. -\ln \Gamma(\tilde{\alpha}-\frac{1}{2}) - (\tilde{\alpha}+1) \ln y - \frac{\tilde{\beta}}{y} \right\} dxdy.
$$

Since, $H(q_1)$ and $H(q_2)$ are Shannon entropy, we have

$$
-H(q_1) = \int q_1(x) \ln \left( \frac{1}{\sqrt{2\pi\tilde{v}}} \exp\{-\frac{(x-\tilde{\mu})^2}{2\tilde{v}}\} \right) dx \tag{A2}
$$

$$
= -\frac{1}{2} \ln (2\pi\tilde{v}) - \frac{1}{2},
$$

and

$$
-H(q_2) = \int q_2(y) \ln \left( \frac{\widetilde{\beta}^{\tilde{\alpha}-\frac{1}{2}}}{\Gamma(\tilde{\alpha}-\frac{1}{2})} y^{-(\tilde{\alpha}+\frac{1}{2})} \exp\{-\frac{\tilde{\beta}}{y}\} \right) dy \tag{A3}
$$

$$
= \tilde{\alpha} - \frac{1}{2} + \ln (\tilde{\beta}\Gamma(\tilde{\alpha}-\frac{1}{2})) - (\tilde{\alpha}+\frac{1}{2})\psi_0(\tilde{\alpha}-\frac{1}{2}).
$$

Also, we need to know $- <\ln p(x,y)>_{q_1 q_2}$ as well

$$
- <\ln p(x,y)>_{q_1 q_2} = \int\int q_1(x) q_2(y) \left( \frac{1}{2} \ln (2\pi) + \frac{(x-\tilde{\mu})^2}{2y} - (\tilde{\alpha}-\frac{1}{2}) \ln \tilde{\beta} \right. \tag{A4}
$$

$$
\left. + \ln \Gamma(\tilde{\alpha}-\frac{1}{2}) + (\tilde{\alpha}+1) \ln y + \frac{\tilde{\beta}}{y} \right) dxdy
$$

$$
= \frac{1}{2} \ln (2\pi) - (\tilde{\alpha}-\frac{1}{2}) \ln \tilde{\beta} + \ln \Gamma(\tilde{\alpha}-\frac{1}{2})
$$

$$
+ \int\int q_1(x) q_2(y) \left( \frac{(x-\tilde{\mu})^2}{2y} + (\tilde{\alpha}+1) \ln y + \frac{\tilde{\beta}}{y} \right) dxdy
$$

$$
= \frac{1}{2} \ln (2\pi) - (\tilde{\alpha}-\frac{1}{2}) \ln \tilde{\beta} + \ln \Gamma(\tilde{\alpha}-\frac{1}{2})
$$

$$
+ (\tilde{\alpha}+1)(\ln \tilde{\beta} - \psi_0(\tilde{\alpha}-\frac{1}{2})) + \frac{(2\tilde{\alpha}-1)(\tilde{v}+2\tilde{\beta})}{4\tilde{\beta}}
$$

Thus, the desire function $KL(\tilde{\boldsymbol{\theta}})$ is

$$
KL(\tilde{\boldsymbol{\theta}}) = 2\ln \Gamma(\tilde{\alpha}-\frac{1}{2}) - (2\tilde{\alpha}+\frac{3}{2})\psi_0(\tilde{\alpha}-\frac{1}{2}) + \frac{(2\tilde{\alpha}-1)(\tilde{v}+2\tilde{\beta})}{4\tilde{\beta}} + \tilde{\alpha} + \frac{5}{2}\ln \tilde{\beta} - \frac{1}{2}\ln \tilde{v} - 1, \tag{A5}
$$

where, $\psi_0(\cdot)$ is the polygamma function of order 0, or called digamma function. The gradient expression respect to $\tilde{\boldsymbol{\theta}} = (\widetilde{\alpha}, \widetilde{\beta}, \widetilde{\mu}, \widetilde{v})$ is

$$
\nabla KL(\tilde{\boldsymbol{\theta}}) = \left( \frac{\widetilde{v}}{2\widetilde{\beta}} - (2\widetilde{\alpha} + \frac{3}{2})\psi_1(\widetilde{\alpha} - \frac{1}{2}) + 2, \quad -(2\widetilde{\alpha} - 1)\frac{\widetilde{v}}{4\widetilde{\beta}^2} + \frac{5}{2\widetilde{\beta}}, \quad 0, \quad \frac{2\widetilde{\alpha} - 1}{4\widetilde{\beta}} - \frac{1}{2\widetilde{v}} \right), \quad (A6)
$$

where, $\psi_1(\cdot)$ is the polygamma function of order 1. We substitute (A6) in the gradient-based and natural gradient formulas, so their algorithms are as follows, respectively

**Algorithm 2:**

$$
\begin{aligned}
\widetilde{\alpha}^{(k+1)} &= \widetilde{\alpha}^{(k)} + \gamma \frac{\partial KL}{\partial \widetilde{\alpha}}(\widetilde{\alpha}^{(k)}, \widetilde{\beta}^{(k)}, \widetilde{v}^{(k)}) \\
&= \widetilde{\alpha}^{(k)} - \gamma \left( \frac{\widetilde{v}^{(k)}}{2\widetilde{\beta}^{(k)}} - (2\widetilde{\alpha}^{(k)} + \frac{3}{2})\psi_1(\widetilde{\alpha}^{(k)} - \frac{1}{2}) + 2 \right), \\
\widetilde{\beta}^{(k+1)} &= \widetilde{\beta}^{(k)} - \gamma \frac{\partial KL}{\partial \widetilde{\beta}}(\widetilde{\alpha}^{(k)}, \widetilde{\beta}^{(k)}, \widetilde{v}^{(k)}) \\
&= \widetilde{\beta}^{(k)} - \gamma \left( -(2\widetilde{\alpha}^{(k)} - 1)\frac{\widetilde{v}^{(k)}}{4[\widetilde{\beta}^{(k)}]^2} + \frac{5}{2\widetilde{\beta}^{(k)}} \right), \\
\widetilde{\mu}^{(k+1)} &= \widetilde{\mu}^{(k)}, \\
\widetilde{v}^{(k+1)} &= \widetilde{v}^{(k)} - \gamma \frac{\partial KL}{\partial \widetilde{v}}(\widetilde{\alpha}^{(k)}, \widetilde{\beta}^{(k)}, \widetilde{v}^{(k)}) \\
&= \widetilde{v}^{(k+1)} - \gamma \left( \frac{2\widetilde{\alpha}^{(k)} - 1}{4\widetilde{\beta}^{(k)}} - \frac{1}{2\widetilde{v}^{(k)}} \right).
\end{aligned}
$$

**Algorithm 3:**

$$
\begin{aligned}
\widetilde{\alpha}^{(k+1)} &= \widetilde{\alpha}^{(k)} - \frac{1}{\|\Delta KL\|}\frac{\partial KL}{\partial \widetilde{\alpha}}(\widetilde{\alpha}^{(k)}, \widetilde{\beta}^{(k)}, \widetilde{v}^{(k)}) \\
&= \widetilde{\alpha}^{(k)} - \frac{1}{\|\Delta KL\|} \left( \frac{\widetilde{v}^{(k)}}{2\widetilde{\beta}^{(k)}} - (2\widetilde{\alpha}^{(k)} + \frac{3}{2})\psi_1(\widetilde{\alpha}^{(k)} - \frac{1}{2}) + 2 \right), \\
\widetilde{\beta}^{(k+1)} &= \widetilde{\beta}^{(k)} - \frac{1}{\|\Delta KL\|}\frac{\partial KL}{\partial \widetilde{\beta}}(\widetilde{\alpha}^{(k)}, \widetilde{\beta}^{(k)}, \widetilde{v}^{(k)}) \\
&= \widetilde{\beta}^{(k)} - \frac{1}{\|\Delta KL\|} \left( -(2\widetilde{\alpha}^{(k)} - 1)\frac{\widetilde{v}^{(k)}}{4[\widetilde{\beta}^{(k)}]^2} + \frac{5}{2\widetilde{\beta}^{(k)}} \right), \\
\widetilde{\mu}^{(k+1)} &= \widetilde{\mu}^{(k)}, \\
\widetilde{v}^{(k+1)} &= \widetilde{v}^{(k)} - \frac{1}{\|\Delta KL\|}\frac{\partial KL}{\partial \widetilde{v}}(\widetilde{\alpha}^{(k)}, \widetilde{\beta}^{(k)}, \widetilde{v}^{(k)}) \\
&= \widetilde{v}^{(k)} - \frac{1}{\|\Delta KL\|} \left( \frac{2\widetilde{\alpha}^{(k)} - 1}{4\widetilde{\beta}^{(k)}} - \frac{1}{2\widetilde{v}^{(k)}} \right).
\end{aligned}
$$

## Appendix B  Conjugate Posterior of Normal Inverse Wishart of $(\widetilde{\mu}, \widetilde{\Sigma})$

The Normal Inverse Wishart distribution function $\mathcal{NIW}(\widetilde{\mu}, \widetilde{\Sigma} | \widetilde{\mu}_0, \tilde{\kappa}, \tilde{\Psi}, \tilde{v})$ is

$$
p(\widetilde{\mu}, \widetilde{\Sigma} | \widetilde{\mu}_0, \tilde{\kappa}, \tilde{\Psi}, \tilde{v}) = \mathcal{N} \left( \widetilde{\mu} | \widetilde{\mu}_0, \frac{1}{\tilde{\kappa}}\widetilde{\Sigma} \right) \mathcal{IW}(\widetilde{\Sigma} | \tilde{\Psi}, \tilde{v}). \quad (A7)
$$

The explicit form of joint distribution $p(x, \widetilde{\mu}, \widetilde{\Sigma})$ is

$$p(x, \widetilde{\mu}, \widetilde{\Sigma}) = p(x|\widetilde{\mu}, \widetilde{\Sigma}) p(\widetilde{\mu}, \widetilde{\Sigma}|\widetilde{\mu}_0, \check{\kappa}, \check{\Psi}, \check{\nu}) \tag{A8}$$

$$= \frac{1}{(2\pi)^{\frac{np}{2}} |\widetilde{\Sigma}|^{\frac{n}{2}}} \exp\left\{-\frac{1}{2}\sum_{i=1}^{n}(x_i - \widetilde{\mu})^T \widetilde{\Sigma}^{-1}(x_i - \widetilde{\mu})\right\}$$

$$\frac{\sqrt{\check{\kappa}}}{(2\pi)^{\frac{p}{2}} |\widetilde{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{\check{\kappa}}{2}(\widetilde{\mu} - \widetilde{\mu}_0)^T \widetilde{\Sigma}^{-1}(\widetilde{\mu} - \widetilde{\mu}_0)\right\} \frac{|\check{\Psi}|^{\frac{\check{\nu}}{2}}}{2^{\frac{\check{\nu}p}{2}} \Gamma_p(\frac{\check{\nu}}{2})} |\widetilde{\Sigma}|^{-\frac{\check{\nu}+p+1}{2}} \exp\left\{-\frac{1}{2}\mathrm{Tr}\left[\check{\Psi}\widetilde{\Sigma}^{-1}\right]\right\}.$$

We rewrite the exponential expression of multivariate Normal Distribution as follows

$$\sum_{i=1}^{n}(x_i - \widetilde{\mu})^T \widetilde{\Sigma}^{-1}(x_i - \widetilde{\mu}) = n(\overline{x} - \widetilde{\mu})^T \widetilde{\Sigma}^{-1}(\overline{x} - \widetilde{\mu}) + \sum_{i=1}^{n}(x_i - \overline{x})^T \widetilde{\Sigma}^{-1}(x_i - \overline{x}). \tag{A9}$$

By substituting (A9) into (A8), we get

$$p(x, \widetilde{\mu}, \widetilde{\Sigma}) \propto |\widetilde{\Sigma}|^{-\frac{\check{\nu}+p+n+2}{2}} \exp\left\{-\frac{n}{2}(\overline{x} - \widetilde{\mu})^T \widetilde{\Sigma}^{-1}(\overline{x} - \widetilde{\mu}) - \frac{\check{\kappa}}{2}(\widetilde{\mu} - \widetilde{\mu}_0)^T \widetilde{\Sigma}^{-1}(\widetilde{\mu} - \widetilde{\mu}_0) \right. \tag{A10}$$

$$\left. -\frac{1}{2}\mathrm{Tr}\left[\left(\check{\Psi} + \sum_{i=1}^{n}(x_i - \overline{x})(x_i - \overline{x})^T\right)\widetilde{\Sigma}^{-1}\right]\right\}.$$

We have to make integration over $x$ and $\widetilde{\Sigma}$ to get the margin of $\widetilde{\mu}$

$$<\ln p(\widetilde{\mu}, \widetilde{\Sigma})>_{q_1 q_3} \propto -\frac{n}{2}<(\overline{x} - \widetilde{\mu})^T \widetilde{\Sigma}^{-1}(\overline{x} - \widetilde{\mu})>_{q_1 q_3} -\frac{\check{\kappa}}{2}<(\widetilde{\mu} - \widetilde{\mu}_0)^T \widetilde{\Sigma}^{-1}(\widetilde{\mu} - \widetilde{\mu}_0)>_{q_1 q_3} \tag{A11}$$

$$\propto -\frac{\check{\kappa}+n}{2}(\widetilde{\mu} - \frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n})^T \widetilde{\Lambda}^{-1}(\widetilde{\mu} - \frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n}),$$

where, $<\widetilde{\Sigma}^{-1}>_{q_3} = \widetilde{\Lambda}^{-1}$. Thus, $\widetilde{\mu}$ has a Normal distribution $\mathcal{N}\left(\frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n}, \frac{1}{\check{\kappa}+n}\widetilde{\Lambda}\right)$. Also, we have the same for the margin of $\widetilde{\Sigma}$

$$<\ln p(\widetilde{\mu}, \widetilde{\Sigma})>_{q_1 q_2} \propto -\frac{\check{\nu}+p+n+2}{2}\ln|\widetilde{\Sigma}| - \frac{n}{2}<(\overline{x} - \widetilde{\mu})^T \widetilde{\Sigma}^{-1}(\overline{x} - \widetilde{\mu})>_{q_1 q_2} \tag{A12}$$

$$-\frac{\check{\kappa}}{2}<(\widetilde{\mu} - \widetilde{\mu}_0)^T \widetilde{\Sigma}^{-1}(\widetilde{\mu} - \widetilde{\mu}_0)>_{q_1 q_2} -\frac{1}{2}\mathrm{Tr}\left[\left(\check{\Psi} + \sum_{i=1}^{n}(x_i - \overline{x})(x_i - \overline{x})^T\right)\widetilde{\Sigma}^{-1}\right]$$

$$\propto -\frac{\check{\nu}+p+n+2}{2}\ln|\widetilde{\Sigma}| - \frac{\check{\kappa}+n}{2}<(\widetilde{\mu} - \frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n})^T \widetilde{\Sigma}^{-1}(\widetilde{\mu} - \frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n})>_{q_1 q_2}$$

$$-\frac{1}{2}\mathrm{Tr}\left[\left(\check{\Psi} + \sum_{i=1}^{n}(x_i - \overline{x})(x_i - \overline{x})^T\right)\widetilde{\Sigma}^{-1}\right]$$

$$\propto -\frac{\check{\nu}+p+n+2}{2}\ln|\widetilde{\Sigma}| - \frac{1}{2}\mathrm{Tr}\left[\left((\check{\kappa}+n)<(\widetilde{\mu} - \frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n})(\widetilde{\mu} - \frac{\check{\kappa}\widetilde{\mu}_0 + n\overline{x}}{\check{\kappa}+n})^T>_{q_1 q_2}\right.\right.$$

$$\left.\left. +\check{\Psi} + \sum_{i=1}^{n}(x_i - \overline{x})(x_i - \overline{x})^T\right)\widetilde{\Sigma}^{-1}\right]$$

$$\propto -\frac{\check{\nu}+p+n+2}{2}\ln|\widetilde{\Sigma}| - \frac{1}{2}\mathrm{Tr}\left[\left(\widetilde{\Lambda} + \check{\Psi} + \sum_{i=1}^{n}(x_i - \overline{x})(x_i - \overline{x})^T\right)\widetilde{\Sigma}^{-1}\right]$$

So, $\widetilde{\boldsymbol{\Sigma}}$ has the structure of an inverse Wishart distribution with below parameters as following

$$\mathcal{IW}\left(\tilde{\boldsymbol{\Lambda}} + \check{\boldsymbol{\Psi}} + \sum_{i=1}^{n}(\boldsymbol{x}_i - \overline{\boldsymbol{x}})(\boldsymbol{x}_i - \overline{\boldsymbol{x}})^T, \tilde{v} + n\right). \tag{A13}$$

It is clear that $\boldsymbol{x}$ has again a multivariate Normal distribution whose mean and variance-covariance matrix are in the following expression

$$\mathcal{MN}\left(\frac{\tilde{\kappa}\widetilde{\boldsymbol{\mu}}_0 + n\overline{\boldsymbol{x}}}{\tilde{\kappa} + n}, \frac{1}{\tilde{v} + n - p - 1}\left[\tilde{\boldsymbol{\Lambda}} + \check{\boldsymbol{\Psi}} + \sum_{i=1}^{n}(\boldsymbol{x}_i - \overline{\boldsymbol{x}})(\boldsymbol{x}_i - \overline{\boldsymbol{x}})^T\right]\right). \tag{A14}$$

In the next step, we want to make an iterative algorithm of the above mentation marginal distributions. We need some pre-guesses for mean $\widetilde{\boldsymbol{\mu}}_0^{(1)} = \boldsymbol{\mu}_0$, variance $\tilde{\boldsymbol{\Lambda}}^{(1)} = \boldsymbol{\Lambda}_0$, and precision $\tilde{\kappa}^{(1)} = \kappa_0$ of $\widetilde{\boldsymbol{\mu}}$ and also for scale matrix $\check{\boldsymbol{\Psi}}^{(1)} = \boldsymbol{\Psi}_0$ and degree of freedom $\tilde{v}^{(1)} = v_0 + n$ of $\widetilde{\boldsymbol{\Sigma}}$. So, we can consider the alternative algorithm for Normal Inverse Wishart distribution as follows        **Algorithm 1:**

$$
\begin{aligned}
\tilde{\kappa}^{(k+1)} &= \tilde{\kappa}^{(k)}, \\
\widetilde{\boldsymbol{\mu}}_0^{(k+1)} &= \frac{\tilde{\kappa}^{(k)}\widetilde{\boldsymbol{\mu}}_0^{(k)} + n\overline{\boldsymbol{x}}}{\tilde{\kappa}^{(k)} + n}, \\
\tilde{\boldsymbol{\Lambda}}^{(k+1)} &= \frac{1}{\tilde{\kappa}^{(k)} + n}\tilde{\boldsymbol{\Lambda}}^{(k)}, \\
\check{\boldsymbol{\Psi}}^{(k+1)} &= [\tilde{\boldsymbol{\Lambda}}^{(k)}]^{-1} + \check{\boldsymbol{\Psi}}^{(k)} + \sum_{i=1}^{n}(\boldsymbol{x}_i - \overline{\boldsymbol{x}})(\boldsymbol{x}_i - \overline{\boldsymbol{x}})^T, \\
\tilde{v}^{(k+1)} &= \tilde{v}^{(k)}.
\end{aligned}
$$

## Appendix C  Kullback-Leibler Divergence of Normal Inverse Wishart

First of all, we define $\widetilde{\boldsymbol{\theta}} = (\widetilde{\boldsymbol{\mu}}_0, \tilde{\kappa}, \tilde{\boldsymbol{\Lambda}}, \check{\boldsymbol{\Psi}}, \tilde{v})$, and the corresponding Kullback-Leibler function

$$
\begin{aligned}
KL(\tilde{\boldsymbol{\theta}}) &= \int\int q_1(\widetilde{\boldsymbol{\mu}})q_2(\widetilde{\boldsymbol{\Sigma}})\ln\left(\frac{q_1(\widetilde{\boldsymbol{\mu}})q_2(\widetilde{\boldsymbol{\Sigma}})}{p(\widetilde{\boldsymbol{\mu}},\widetilde{\boldsymbol{\Sigma}})}\right)d\widetilde{\boldsymbol{\mu}}\,d\widetilde{\boldsymbol{\Sigma}} \\
&= \int\int q_1(\widetilde{\boldsymbol{\mu}})q_2(\widetilde{\boldsymbol{\Sigma}})\ln\left(\frac{q_2(\widetilde{\boldsymbol{\Sigma}})}{p(\widetilde{\boldsymbol{\mu}},\widetilde{\boldsymbol{\Sigma}})}\right)d\widetilde{\boldsymbol{\mu}}\,d\widetilde{\boldsymbol{\Sigma}} - H(q_1) \\
&= \int q_2(\widetilde{\boldsymbol{\Sigma}})\left(\ln q_2(\widetilde{\boldsymbol{\Sigma}}) - \int q_1(\widetilde{\boldsymbol{\mu}})\ln p(\widetilde{\boldsymbol{\mu}},\widetilde{\boldsymbol{\Sigma}})d\widetilde{\boldsymbol{\mu}}\right)d\widetilde{\boldsymbol{\Sigma}} - H(q_1).
\end{aligned}
\tag{A15}
$$

We need to compute the internal integral first

$$
\begin{aligned}
\int q_1(\widetilde{\boldsymbol{\mu}})\ln p(\widetilde{\boldsymbol{\mu}},\widetilde{\boldsymbol{\Sigma}})d\widetilde{\boldsymbol{\mu}} &= \int q_1(\widetilde{\boldsymbol{\mu}})\ln\left\{\frac{\sqrt{\tilde{\kappa}}|\check{\boldsymbol{\Psi}}|^{\frac{\tilde{v}}{2}}}{(2\pi)^{\frac{p}{2}}2^{\frac{\tilde{v}p}{2}}\Gamma_p(\frac{\tilde{v}}{2})}\right. \\
&\qquad\left. |\widetilde{\boldsymbol{\Sigma}}|^{-\frac{\tilde{v}+p+2}{2}}\exp\{-\frac{\tilde{\kappa}}{2}(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0)^T\widetilde{\boldsymbol{\Sigma}}^{-1}(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0) - \frac{1}{2}\mathrm{Tr}\left[\check{\boldsymbol{\Psi}}\widetilde{\boldsymbol{\Sigma}}^{-1}\right]\}\right\}d\widetilde{\boldsymbol{\mu}} \\
&= \ln\left\{\frac{\sqrt{\tilde{\kappa}}|\check{\boldsymbol{\Psi}}|^{\frac{\tilde{v}}{2}}}{(2\pi)^{\frac{p}{2}}2^{\frac{\tilde{v}p}{2}}\Gamma_p(\frac{\tilde{v}}{2})}|\widetilde{\boldsymbol{\Sigma}}|^{-\frac{\tilde{v}+p+2}{2}}\exp\{-\frac{1}{2}\mathrm{Tr}\left[\check{\boldsymbol{\Psi}}\widetilde{\boldsymbol{\Sigma}}^{-1}\right]\}\right\} \\
&\qquad - \frac{\tilde{\kappa}}{2}\int q_1(\widetilde{\boldsymbol{\mu}})\mathrm{Tr}\left[(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0)(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0)^T\widetilde{\boldsymbol{\Sigma}}^{-1}\right]d\widetilde{\boldsymbol{\mu}} \\
&= \ln\left\{\frac{\sqrt{\tilde{\kappa}}|\check{\boldsymbol{\Psi}}|^{\frac{\tilde{v}}{2}}}{(2\pi)^{\frac{p}{2}}2^{\frac{\tilde{v}p}{2}}\Gamma_p(\frac{\tilde{v}}{2})}|\widetilde{\boldsymbol{\Sigma}}|^{-\frac{\tilde{v}+p+2}{2}}\exp\{-\frac{1}{2}\mathrm{Tr}\left[(\tilde{\boldsymbol{\Lambda}} + \check{\boldsymbol{\Psi}})\widetilde{\boldsymbol{\Sigma}}^{-1}\right]\}\right\}.
\end{aligned}
\tag{A16}
$$

By substituting (A16) in side of (A15), we get

$$KL(\tilde{\boldsymbol{\theta}}) = \int q_2(\widetilde{\boldsymbol{\Sigma}}) \left( \ln q_2(\widetilde{\boldsymbol{\Sigma}}) - \ln \left\{ \frac{\sqrt{\tilde{\kappa}}|\tilde{\boldsymbol{\Psi}}|^{\frac{\tilde{\nu}}{2}}}{(2\pi)^{\frac{p}{2}} 2^{\frac{\tilde{\nu}p}{2}} \Gamma_p(\frac{\tilde{\nu}}{2})} |\widetilde{\boldsymbol{\Sigma}}|^{-\frac{\tilde{\nu}+p+2}{2}} \exp\left\{ -\frac{1}{2} \mathrm{Tr}\left[ (\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}) \widetilde{\boldsymbol{\Sigma}}^{-1} \right] \right\} \right\} \right) d\widetilde{\boldsymbol{\Sigma}} - H(q_1).$$
(A17)

Since $q_2(\widetilde{\boldsymbol{\Sigma}}) = \mathcal{IW}(\tilde{\boldsymbol{\Psi}}, \tilde{\nu})$, we can rewrite $KL(\tilde{\boldsymbol{\theta}})$ as following

$$KL(\tilde{\boldsymbol{\theta}}) = -\ln\left\{ \frac{\sqrt{\tilde{\kappa}}|\tilde{\boldsymbol{\Psi}}|^{\frac{\tilde{\nu}}{2}} \Gamma_p(\frac{\tilde{\nu}+1}{2})}{\pi^{\frac{p}{2}} \Gamma_p(\frac{\tilde{\nu}}{2})|\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}|^{\frac{\tilde{\nu}+1}{2}}} \right\} + \int q_2(\widetilde{\boldsymbol{\Sigma}}) \ln\left( \frac{\mathcal{IW}(\tilde{\boldsymbol{\Psi}}, \tilde{\nu})}{\mathcal{IW}(\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}, \tilde{\nu}+1)} \right) d\widetilde{\boldsymbol{\Sigma}} - H(q_1).$$
(A18)

The second term is again a Kullback-Leibler function respect to two Inverse Wishart distributions with different parameters calculated in [13]

$$\int q_2(\widetilde{\boldsymbol{\Sigma}}) \ln\left( \frac{\mathcal{IW}(\tilde{\boldsymbol{\Psi}}, \tilde{\nu})}{\mathcal{IW}(\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}, \tilde{\nu}+1)} \right) d\widetilde{\boldsymbol{\Sigma}} = \ln\left( \frac{\Gamma_p(\frac{\tilde{\nu}+1}{2})}{\Gamma_p(\frac{\tilde{\nu}}{2})} \right) + \frac{\tilde{\nu}}{2} \mathrm{Tr}\left[ \tilde{\boldsymbol{\Psi}}^{-1} \tilde{\boldsymbol{\Lambda}} + \mathbf{I} \right] - \frac{p\tilde{\nu}}{2}$$
(A19)
$$- \frac{\tilde{\nu}}{2} \ln|\tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I}| - \frac{1}{2}\sum_{i=1}^{p} \psi_0(\frac{\tilde{\nu}-p+i}{2}).$$

The last term of (A15) is the Shannon entropy of $\widetilde{\boldsymbol{\mu}}$

$$-H(q_1) = \int q_1(\widetilde{\boldsymbol{\mu}}) \ln q_1(\widetilde{\boldsymbol{\mu}}) d\widetilde{\boldsymbol{\mu}}$$
(A20)
$$= \int q_1(\widetilde{\boldsymbol{\mu}}) \ln\left( \frac{\tilde{\kappa}^{\frac{p}{2}}}{(2\pi)^{\frac{p}{2}}|\tilde{\boldsymbol{\Lambda}}|^{\frac{1}{2}}} \exp\left\{ -\frac{\tilde{\kappa}}{2}(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0)^T \tilde{\boldsymbol{\Lambda}}^{-1}(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0) \right\} \right) d\widetilde{\boldsymbol{\mu}}$$
$$= \frac{p}{2}\ln(\tilde{\kappa}) - \frac{p}{2}\ln(2\pi) - \frac{1}{2}\ln|\tilde{\boldsymbol{\Lambda}}| - \frac{\tilde{\kappa}}{2}\int q_1(\widetilde{\boldsymbol{\mu}})(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0)^T \tilde{\boldsymbol{\Lambda}}^{-1}(\widetilde{\boldsymbol{\mu}} - \widetilde{\boldsymbol{\mu}}_0) d\widetilde{\boldsymbol{\mu}}$$
$$= \frac{p}{2}\ln(\tilde{\kappa}) - \frac{p}{2}\ln(2\pi) - \frac{1}{2}\ln|\tilde{\boldsymbol{\Lambda}}| - \frac{1}{2}.$$

The final expression of $KL(\tilde{\boldsymbol{\theta}})$ is equivalent to the following after substitution of (A19) and (A20) into (A18)

$$KL(\tilde{\boldsymbol{\theta}}) = -\frac{1}{2}\ln\tilde{\kappa} - \frac{\tilde{\nu}}{2}\ln|\tilde{\boldsymbol{\Psi}}| + \frac{p}{2}\ln\pi - \ln\left( \frac{\Gamma_p(\frac{\tilde{\nu}+1}{2})}{\Gamma_p(\frac{\tilde{\nu}}{2})} \right) + \frac{\tilde{\nu}+1}{2}\ln|\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}|$$
(A21)
$$+ \ln\left( \frac{\Gamma_p(\frac{\tilde{\nu}+1}{2})}{\Gamma_p(\frac{\tilde{\nu}}{2})} \right) + \frac{\tilde{\nu}}{2}\mathrm{Tr}\left[ \tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I} \right] - \frac{p\tilde{\nu}}{2} - \frac{\tilde{\nu}}{2}\ln|\tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I}| - \frac{1}{2}\sum_{i=1}^{p}\psi_0(\frac{\tilde{\nu}-p+i}{2})$$
$$+ \frac{p}{2}\ln(\tilde{\kappa}) - \frac{p}{2}\ln(2\pi) - \frac{1}{2}\ln|\tilde{\boldsymbol{\Lambda}}| - \frac{1}{2}$$
$$\propto \frac{p-1}{2}\ln\tilde{\kappa} - \frac{\tilde{\nu}}{2}\ln\left( |\tilde{\boldsymbol{\Psi}}| \cdot |\tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I}| \right) + \frac{\tilde{\nu}+1}{2}\ln|\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}|$$
$$+ \frac{\tilde{\nu}}{2}\mathrm{Tr}\left[ \tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I} \right] - \frac{p\tilde{\nu}}{2} - \frac{1}{2}\sum_{i=1}^{p}\psi_0(\frac{\tilde{\nu}-p+i}{2}) - \frac{1}{2}\ln|\tilde{\boldsymbol{\Lambda}}|$$
$$\propto \frac{p-1}{2}\ln\tilde{\kappa} - \frac{p\tilde{\nu}}{2} + \frac{1}{2}\ln|\mathbf{I} + \tilde{\boldsymbol{\Lambda}}^{-1}\tilde{\boldsymbol{\Psi}}| + \frac{\tilde{\nu}}{2}\mathrm{Tr}\left[ \tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I} \right] - \frac{1}{2}\sum_{i=1}^{p}\psi_0(\frac{\tilde{\nu}-p+i}{2}).$$

The gradient of $KL(\tilde{\boldsymbol{\theta}})$ respect to $\tilde{\boldsymbol{\theta}} = (\tilde{\boldsymbol{\mu}}_0, \tilde{\kappa}, \tilde{\boldsymbol{\Lambda}}, \tilde{\boldsymbol{\Psi}}, \tilde{v})$ is

$$
\nabla KL(\tilde{\boldsymbol{\theta}}) = \Big( 0, \quad \frac{p-1}{2\tilde{\kappa}}, \quad \frac{1}{2}[\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}]^{-1} - \frac{1}{2}\tilde{\boldsymbol{\Lambda}}^{-1} + \frac{\tilde{v}}{2}\tilde{\boldsymbol{\Psi}}^{-1}, \tag{A22}
$$

$$
\frac{1}{2}[\tilde{\boldsymbol{\Lambda}} + \tilde{\boldsymbol{\Psi}}]^{-1} - \frac{\tilde{v}}{2}\tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}}\tilde{\boldsymbol{\Psi}}^{-1}, \quad -\frac{p}{2} + \frac{1}{2}\operatorname{Tr}\left[\tilde{\boldsymbol{\Psi}}^{-1}\tilde{\boldsymbol{\Lambda}} + \mathbf{I}\right] - \frac{1}{4}\sum_{i=1}^{p}\psi_1(\frac{\tilde{v}-p+i}{2}) \Big).
$$

So, the two second algorithms are

**Algorithm 2:**

$$
\begin{aligned}
\tilde{\boldsymbol{\mu}}_0^{(k+1)} &= \tilde{\boldsymbol{\mu}}_0^{(k)}, \\
\tilde{\kappa}^{(k+1)} &= \tilde{\kappa}^{(k)} - \frac{\gamma}{2}\left(\frac{p-1}{\tilde{\kappa}^{(k)}}\right), \\
\tilde{\boldsymbol{\Lambda}}^{(k+1)} &= \tilde{\boldsymbol{\Lambda}}^{(k)} - \frac{\gamma}{2}\left([\tilde{\boldsymbol{\Lambda}}^{(k)} + \tilde{\boldsymbol{\Psi}}^{(k)}]^{-1} - [\tilde{\boldsymbol{\Lambda}}^{(k)}]^{-1} + \tilde{v}^{(k)}[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\right), \\
\tilde{\boldsymbol{\Psi}}^{(k+1)} &= \tilde{\boldsymbol{\Psi}}^{(k)} - \frac{\gamma}{2}\left([\tilde{\boldsymbol{\Lambda}}^{(k)} + \tilde{\boldsymbol{\Psi}}^{(k)}]^{-1} - \tilde{v}^{(k)}[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\tilde{\boldsymbol{\Lambda}}^{(k)}[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\right), \\
\tilde{v}^{(k+1)} &= \tilde{v}^{(k)} - \frac{\gamma}{2}\left(\operatorname{Tr}\left[[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\tilde{\boldsymbol{\Lambda}}^{(k)} + \mathbf{I}\right] - \frac{1}{2}\sum_{i=1}^{p}\psi_1(\frac{\tilde{v}^{(k)}-p+i}{2}) - p\right),
\end{aligned}
$$

and

**Algorithm 3:**

$$
\begin{aligned}
\tilde{\boldsymbol{\mu}}_0^{(k+1)} &= \tilde{\boldsymbol{\mu}}_0^{(k)}, \\
\tilde{\kappa}^{(k+1)} &= \tilde{\kappa}^{(k)} - \frac{1}{2\|\Delta KL\|}\left(\frac{p-1}{2\tilde{\kappa}^{(k)}}\right), \\
\tilde{\boldsymbol{\Lambda}}^{(k+1)} &= \tilde{\boldsymbol{\Lambda}}^{(k)} - \frac{1}{2\|\Delta KL\|}\left([\tilde{\boldsymbol{\Lambda}}^{(k)} + \tilde{\boldsymbol{\Psi}}^{(k)}]^{-1} - [\tilde{\boldsymbol{\Lambda}}^{(k)}]^{-1} + \tilde{v}^{(k)}[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\right), \\
\tilde{\boldsymbol{\Psi}}^{(k+1)} &= \tilde{\boldsymbol{\Psi}}^{(k)} - \frac{1}{2\|\Delta KL\|}\left([\tilde{\boldsymbol{\Lambda}}^{(k)} + \tilde{\boldsymbol{\Psi}}^{(k)}]^{-1} - \tilde{v}^{(k)}[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\tilde{\boldsymbol{\Lambda}}^{(k)}[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\right), \\
\tilde{v}^{(k+1)} &= \tilde{v}^{(k)} - \frac{1}{2\|\Delta KL\|}\left(\operatorname{Tr}\left[[\tilde{\boldsymbol{\Psi}}^{(k)}]^{-1}\tilde{\boldsymbol{\Lambda}}^{(k)} + \mathbf{I}\right] - \frac{1}{2}\sum_{i=1}^{p}\psi_1(\frac{\tilde{v}^{(k)}-p+i}{2}) - p\right).
\end{aligned}
$$

## Appendix D   VBA Approximation the Linear Inverse Problem

In this example, the process is pretty the same. First, we have to rewrite $p(\boldsymbol{g}, \tilde{\boldsymbol{f}}, \tilde{v})$ and $\ln p(\boldsymbol{g}, \tilde{\boldsymbol{f}}, \tilde{v})$

$$
p(\boldsymbol{g}, \tilde{\boldsymbol{f}}, \tilde{v}) = \frac{1}{(2\pi)^{\frac{np}{2}}|\tilde{v}_\epsilon \boldsymbol{I}|^{\frac{n}{2}}}\exp\{-\frac{1}{2}\sum_{i=1}^{n}(\boldsymbol{g}_i - \tilde{\boldsymbol{f}})^T(\tilde{v}_\epsilon \boldsymbol{I})^{-1}(\boldsymbol{g}_i - \tilde{\boldsymbol{f}})\} \tag{A23}
$$

$$
\frac{1}{(2\pi)^{\frac{p}{2}}|\operatorname{diag}[\tilde{v}]|^{\frac{1}{2}}}\exp\{-\frac{1}{2}\tilde{\boldsymbol{f}}^T\operatorname{diag}[\tilde{v}]^{-1}\tilde{\boldsymbol{f}}\}\prod_{j=1}^{p}\frac{\tilde{\beta}_j^{\tilde{\alpha}_j}}{\Gamma(\tilde{\alpha}_j)}\tilde{v}_j^{-(\tilde{\alpha}_j+1)}\exp\{-\frac{\tilde{\beta}_j}{\tilde{v}_j}\},
$$

and

$$
\ln p(\boldsymbol{g}, \tilde{\boldsymbol{f}}, \tilde{v}) \propto -\frac{n}{2}\ln|\tilde{v}_\epsilon \boldsymbol{I}| - \frac{1}{2}\sum_{i=1}^{n}(\boldsymbol{g}_i - \tilde{\boldsymbol{f}})^T(\tilde{v}_\epsilon \boldsymbol{I})^{-1}(\boldsymbol{g}_i - \tilde{\boldsymbol{f}}) - \frac{1}{2}\ln|\operatorname{diag}[\tilde{v}]| \tag{A24}
$$

$$
-\frac{1}{2}\tilde{\boldsymbol{f}}^T\operatorname{diag}[\tilde{v}]^{-1}\tilde{\boldsymbol{f}} + \sum_{j=1}^{p}\tilde{\alpha}_j\ln\tilde{\beta}_j - \sum_{j=1}^{p}\ln\Gamma(\tilde{\alpha}_j) - \sum_{j=1}^{p}(\tilde{\alpha}_j+1)\ln\tilde{v}_j - \sum_{j=1}^{p}\frac{\tilde{\beta}_j}{\tilde{v}_j},
$$

where, $I$ is an identical matrix $p \times p$. According to the above expression, it is easy to find out all margins function starting from $\tilde{v}$

$$
< \ln p(g, \widetilde{f}, \tilde{v}) >_{g, \widetilde{f}} \propto -\frac{1}{2} \ln |\text{diag}\,[\tilde{v}]| - \frac{1}{2} < \widetilde{f}^T \text{diag}\,[\tilde{v}]^{-1} \widetilde{f} >_{\widetilde{f}} - \sum_{j=1}^{p} (\widetilde{\alpha}_j + 1) \ln \tilde{v}_j - \sum_{j=1}^{p} \frac{\widetilde{\beta}_j}{\tilde{v}_j} \quad \text{(A25)}
$$

$$
\propto -\frac{1}{2} \sum_{j=1}^{p} \ln \tilde{v}_i - \frac{1}{2} \sum_{j=1}^{p} \frac{< \tilde{f}_k^2 >_{\widetilde{f}}}{\tilde{v}_i} - \sum_{j=1}^{p} (\widetilde{\alpha}_j + 1) \ln \tilde{v}_j - \sum_{j=1}^{p} \frac{\widetilde{\beta}_j}{\tilde{v}_j}
$$

$$
\propto -(\widetilde{\alpha}_k + \frac{3}{2}) \ln \tilde{v}_k - (\frac{\tilde{v}_{\tilde{f}_k} + \widetilde{\mu}_{\tilde{f}_k}^2}{2} + \widetilde{\beta}_k) \frac{1}{\tilde{v}_k}.
$$

Thus, $\tilde{v}_k \sim \mathcal{IG}(\widetilde{\alpha}_k, \frac{\tilde{v}_k + \widetilde{\mu}_{\tilde{f}_k}^2}{2} + \widetilde{\beta}_k)$ for $k = 1, \cdots, p$. The process for computation of $g$ density function is similar

$$
< \ln p(g, \widetilde{f}, \tilde{v}) >_{\widetilde{f}, \tilde{v}} \propto -\frac{1}{2} \sum_{i=1}^{n} < (g_i - \widetilde{f})^T (\tilde{v}_\epsilon I)^{-1} (g_i - \widetilde{f}) >_{\widetilde{f}} \quad \text{(A26)}
$$

$$
= -\frac{1}{2\tilde{v}_\epsilon} \sum_{i=1}^{n} \sum_{j=1}^{p} < (g_{ij} - \tilde{f}_j)^2 >_{\widetilde{f}}
$$

$$
\propto -\frac{1}{2\tilde{v}_\epsilon} \sum_{i=1}^{n} \sum_{j=1}^{p} (g_{ij} - < \tilde{f}_j >_{\widetilde{f}})^2
$$

$$
= -\frac{1}{2} \sum_{i=1}^{n} (g_i - \widetilde{\mu}_{\widetilde{f}})^T (\tilde{v}_\epsilon I)^{-1} (g_i - \widetilde{\mu}_{\widetilde{f}}).
$$

So, $g$ has again a Normal distribution of $\mathcal{N}(\widetilde{\mu}_{\widetilde{f}}, \tilde{v}_\epsilon I)$. Since, we need the density of $\overline{g}$ in the proceeding, we can specify it now. Therefore, back to properties of Normal distribution, we know that $\overline{g}$ has also a Normal density shown by $\overline{g} \sim \mathcal{N}(\widetilde{\mu}_{\widetilde{f}}, \frac{1}{n}\tilde{v}_\epsilon I)$. Now, we redo all calculation on (A27) for the objective margin of $\widetilde{f}$, separately for each component denoted by $\tilde{f}_k$ for $k = 1, \cdots, p$

$$
< \ln p(g, \widetilde{f}, \tilde{v}) >_{g, \tilde{v}} \propto -\frac{1}{2} \sum_{i=1}^{n} < (g_i - \widetilde{f})^T (\tilde{v}_\epsilon I)^{-1} (g_i - \widetilde{f}) >_g -\frac{1}{2} < \widetilde{f}^T \text{diag}\,[\tilde{v}]^{-1} \widetilde{f} >_{\tilde{v}} \quad \text{(A27)}
$$

$$
\propto -\frac{1}{2} \left( n < (\overline{g} - \widetilde{f})^T (\tilde{v}_\epsilon I)^{-1} (\overline{g} - \widetilde{f}) >_g + < \widetilde{f}^T \text{diag}\,[\tilde{v}]^{-1} \widetilde{f} >_{\tilde{v}} \right)
$$

$$
= -\frac{1}{2} \left( \frac{n}{\tilde{v}_\epsilon} \sum_{j=1}^{p} < (\overline{g}_j - \tilde{f}_j)^2 >_g + \sum_{j=1}^{p} < \frac{1}{\tilde{v}_j} >_{\tilde{v}} \tilde{f}_j^2 \right)
$$

$$
\propto -\frac{1}{2} \sum_{j=1}^{p} (\frac{n}{\tilde{v}_\epsilon} + < \frac{1}{\tilde{v}_j} >_{\tilde{v}}) \left( \tilde{f}_j - \frac{< \overline{g}_j >_g}{1 + \frac{\tilde{v}_\epsilon}{n} < \frac{1}{\tilde{v}_j} >_{\tilde{v}}} \right)^2
$$

$$
= -\frac{1}{2} \sum_{j=1}^{p} (\frac{n}{\tilde{v}_\epsilon} + \frac{2\widetilde{\alpha}_j}{\tilde{v}_j + \widetilde{\mu}_{\tilde{f}_j}^2 + 2\widetilde{\beta}_j}) \left( \tilde{f}_j - \frac{\widetilde{\mu}_{g_i}}{1 + \frac{2\tilde{v}_\epsilon \widetilde{\alpha}_j}{n(\tilde{v}_j + \widetilde{\mu}_{\tilde{f}_j}^2 + 2\widetilde{\beta}_j)}} \right)^2.
$$

Thus, $\tilde{f}_k$ has a Normal distribution with these structures

$$\tilde{f}_k \sim \mathcal{N}\left(\frac{\widetilde{\mu}_{g_k}}{1+\frac{2\tilde{v}_\epsilon \widetilde{\alpha}_k}{n(\tilde{v}_k+\widetilde{\mu}^2_{\tilde{f}_k}+2\widetilde{\beta}_k)}}, \frac{\tilde{v}_\epsilon(\tilde{v}_k+\widetilde{\mu}^2_{\tilde{f}_k}+2\widetilde{\beta}_k)}{n(\tilde{v}_k+\widetilde{\mu}^2_{\tilde{f}_k})+2n\widetilde{\beta}_k+2\tilde{v}_\epsilon\widetilde{\alpha}_k}\right), \quad k=1,\cdots,p, \tag{A28}$$

so, So, $\widetilde{f}$ is separable due to the diag function. We can summarize all in a multivariate Normal distribution and $\widetilde{\mu}_g = \widetilde{\mu}_{\widetilde{f}}$

$$\widetilde{f} \sim \mathcal{MN}\left(\frac{\widetilde{\mu}_{\widetilde{f}}}{1+\frac{2\tilde{v}_\epsilon\widetilde{\boldsymbol{\alpha}}}{n(\tilde{\boldsymbol{v}}+\widetilde{\mu}^2_{\widetilde{f}}+2\widetilde{\boldsymbol{\beta}})}}, \mathrm{diag}\left[\frac{\tilde{v}_\epsilon(\tilde{v}+\widetilde{\mu}^2_{\widetilde{f}}+2\tilde{\beta})}{n(\tilde{v}+\widetilde{\mu}^2_{\widetilde{f}})+2n\tilde{\beta}+2\tilde{v}_\epsilon\tilde{\boldsymbol{\alpha}}}\right]\right). \tag{A29}$$

198  where all the mathematical operations on vectors are componentwise, and we get that $\widetilde{\mu}_g = \widetilde{\mu}_{\widetilde{f}}$. The
199  recursive algorithm is fixed for $\tilde{v}_\epsilon$ and $\tilde{\boldsymbol{\alpha}}$, which can be estimated from the data set. The alternative
200  algorithm is in the following for other parameters
201  **Algorithm 1:**

$$
\begin{aligned}
\tilde{v}_\epsilon^{(k+1)} &= \tilde{v}_\epsilon^{(k)}, \\
\tilde{\boldsymbol{\alpha}}^{(k+1)} &= \tilde{\boldsymbol{\alpha}}^{(k)}, \\
\tilde{\boldsymbol{\beta}}^{(k+1)} &= \frac{1}{2}(\tilde{v}^{(k)}+[\widetilde{\mu}_{\widetilde{f}}^{(k)}]^2)+\tilde{\boldsymbol{\beta}}^{(k)}, \\
\widetilde{\mu}_{\widetilde{f}}^{(k+1)} &= \widetilde{\mu}_{\widetilde{f}}^{(k)}\left(1+\frac{2\tilde{v}_\epsilon^{(k)}\tilde{\boldsymbol{\alpha}}^{(k)}}{n(\tilde{v}^{(k)}+[\widetilde{\mu}_{\widetilde{f}}^{(k)}]^2+2\tilde{\boldsymbol{\beta}}^{(k)})}\right)^{-1}, \\
\tilde{v}^{(k+1)} &= \mathrm{diag}\left[\frac{\tilde{v}_\epsilon^{(k)}(\tilde{v}^{(k)}+[\widetilde{\mu}_{\widetilde{f}}^{(k)}]^2+2\tilde{\boldsymbol{\beta}}^{(k)})}{n(\tilde{v}^{(k)}+[\widetilde{\mu}_{\widetilde{f}}^{(k)}]^2)+2n\tilde{\boldsymbol{\beta}}^{(k)}+2\tilde{v}_\epsilon^{(k)}\tilde{\boldsymbol{\alpha}}^{(k)}}\right].
\end{aligned}
$$

202  **Appendix E  $KL(\widetilde{\boldsymbol{\theta}})$ the Linear Inverse Problem**

We define $\widetilde{\boldsymbol{\theta}} = (\widetilde{\mu}_g, \tilde{v}_\epsilon, \tilde{v}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}})$ and its corresponding $KL(\widetilde{\boldsymbol{\theta}})$

$$KL(\widetilde{\boldsymbol{\theta}}) = -H(q_1)-H(q_2)-H(q_3)- <\ln p(\boldsymbol{g},\widetilde{f},\tilde{v})>_{q_1q_2q_3}. \tag{A30}$$

We separately decompose each term initializing from $-H(q_1)$

$$
\begin{aligned}
-H(q_1) &= \int q_1(\boldsymbol{g})\ln\left(\frac{1}{(2\pi)^{\frac{p}{2}}|\mathrm{diag}\,[\tilde{v}_\epsilon\mathcal{I}]|^{\frac{1}{2}}}\exp\{-\frac{1}{2}(\boldsymbol{g}-\widetilde{\mu}_g)^T\mathrm{diag}\,[\tilde{v}]^{-1}(\boldsymbol{g}-\widetilde{\mu}_g)\}\right)d\boldsymbol{g} \\
&= -\frac{p}{2}-\frac{p}{2}\ln(2\pi)-\frac{p}{2}\ln\tilde{v}_\epsilon,
\end{aligned}
\tag{A31}
$$

and

$$
\begin{aligned}
-H(q_2) &= \int q_2(\widetilde{f})\ln\left(\frac{1}{(2\pi)^{\frac{p}{2}}|\mathrm{diag}\,[\tilde{v}]|^{\frac{1}{2}}}\exp\{-\frac{1}{2}\widetilde{f}^T\mathrm{diag}\,[\tilde{v}]^{-1}\widetilde{f}\}\right)d\widetilde{f} \\
&= -\frac{p}{2}-\frac{p}{2}\ln(2\pi)-\frac{1}{2}\ln|\mathrm{diag}\,[\tilde{v}]|,
\end{aligned}
\tag{A32}
$$

and

$$
-H(q_3) = \int q_3(\tilde{v}) \ln \left( \prod_{j=1}^{p} \frac{\widetilde{\beta}_j^{\widetilde{\alpha}_j}}{\Gamma(\widetilde{\alpha}_j)} \tilde{v}_j^{-(\widetilde{\alpha}_j+1)} \exp\{-\frac{\widetilde{\beta}_j}{\tilde{v}_j}\} \right) d\tilde{v} \tag{A33}
$$

$$
= -\sum_{j=1}^{p} \widetilde{\alpha}_j - \sum_{j=1}^{p} \ln(\widetilde{\beta}_j \Gamma(\widetilde{\alpha}_j)) + \sum_{j=1}^{p} (\widetilde{\alpha}_j + 1)\psi_0(\widetilde{\alpha}_j),
$$

and

$$
- < \ln p(\boldsymbol{g}, \widetilde{\boldsymbol{f}}, \tilde{v}) >_{q_1 q_2 q_3} \propto \int \int \int q_1(\boldsymbol{g}) q_2(\widetilde{\boldsymbol{f}}) q_3(\tilde{v}) \left( \frac{n}{2} \ln |\tilde{v}_\epsilon \boldsymbol{I}| + \frac{1}{2} \sum_{i=1}^{n} (\boldsymbol{g}_i - \widetilde{\boldsymbol{f}})^T (\tilde{v}_\epsilon \boldsymbol{I})^{-1} (\boldsymbol{g}_i - \widetilde{\boldsymbol{f}}) \right.
$$

$$
\tag{A34}
$$

$$
+ \frac{1}{2} \ln |\text{diag}\,[\tilde{v}]| + \frac{1}{2} \widetilde{\boldsymbol{f}}^T \text{diag}\,[\tilde{v}]^{-1} \widetilde{\boldsymbol{f}} - \sum_{j=1}^{p} \widetilde{\alpha}_j \ln \widetilde{\beta}_j
$$

$$
\left. + \sum_{j=1}^{p} \ln \Gamma(\widetilde{\alpha}_j) + \sum_{j=1}^{p} (\widetilde{\alpha}_j + 1) \ln \tilde{v}_j + \sum_{j=1}^{p} \frac{\widetilde{\beta}_j}{\tilde{v}_j} \right) d\boldsymbol{g} d\widetilde{\boldsymbol{f}} d\tilde{v}
$$

$$
\propto \frac{np}{2} \ln(\tilde{v}_\epsilon) + \frac{1}{2\tilde{v}_\epsilon} \sum_{i=1}^{n} \sum_{j=1}^{p} \int \int q_1(\boldsymbol{g}) q_2(\widetilde{\boldsymbol{f}})(g_{ij}^2 + \tilde{f}_j^2) d\boldsymbol{g} d\widetilde{\boldsymbol{f}}
$$

$$
- \sum_{j=1}^{p} \widetilde{\alpha}_j \ln \widetilde{\beta}_j + \sum_{j=1}^{p} \ln \Gamma(\widetilde{\alpha}_j) + \int q_3(\tilde{v}) \left( \sum_{j=1}^{p} (\widetilde{\alpha}_j + \frac{3}{2}) \ln \tilde{v}_j + \sum_{j=1}^{p} \frac{\widetilde{\beta}_j}{\tilde{v}_j} \right) d\tilde{v}
$$

$$
\propto \frac{np}{2} \ln(\tilde{v}_\epsilon) + \frac{n}{2\tilde{v}_\epsilon} \sum_{j=1}^{p} (2\tilde{v}_j + \widetilde{\mu}_{g_j}^2) - \sum_{j=1}^{p} \widetilde{\alpha}_j \ln \widetilde{\beta}_j + \sum_{j=1}^{p} \ln \Gamma(\widetilde{\alpha}_j)
$$

$$
+ \sum_{j=1}^{p} (\widetilde{\alpha}_j + \frac{3}{2})(\ln \widetilde{\beta}_j - \psi_0(\widetilde{\alpha}_j)) + \sum_{j=1}^{p} \widetilde{\alpha}_j.
$$

Finally, we get

$$
KL(\widetilde{\boldsymbol{\theta}}) \propto - \frac{p}{2} \ln \tilde{v}_\epsilon - \frac{1}{2} \ln |\text{diag}\,[\tilde{v}]| - \sum_{j=1}^{p} \ln(\widetilde{\beta}_j \Gamma(\widetilde{\alpha}_j)) + \sum_{j=1}^{p} (\widetilde{\alpha}_j + 1)\psi_0(\widetilde{\alpha}_j) + \frac{np}{2} \ln(\tilde{v}_\epsilon) \tag{A35}
$$

$$
+ \frac{n}{2\tilde{v}_\epsilon} \sum_{j=1}^{p} (2\tilde{v}_j + \widetilde{\mu}_{g_j}^2) - \sum_{j=1}^{p} \widetilde{\alpha}_j \ln \widetilde{\beta}_j + \sum_{j=1}^{p} \ln \Gamma(\widetilde{\alpha}_j) + \sum_{j=1}^{p} (\widetilde{\alpha}_j + \frac{3}{2})(\ln \widetilde{\beta}_j - \psi_0(\widetilde{\alpha}_j)).
$$

$$
\propto \frac{p(n-1)}{2} \ln \tilde{v}_\epsilon + \frac{n}{2\tilde{v}_\epsilon} \sum_{j=1}^{p} (2\tilde{v}_j + \widetilde{\mu}_{g_j}^2) - \frac{1}{2} \sum_{j=1}^{p} \ln \tilde{v}_j - \frac{1}{2} \sum_{j=1}^{p} \psi_0(\widetilde{\alpha}_j) + \frac{1}{2} \sum_{j=1}^{p} \ln \widetilde{\beta}_j.
$$

To make the last two algorithms, we have to differentiate $KL(\widetilde{\boldsymbol{\theta}})$ respect to $\widetilde{\boldsymbol{\theta}} = (\widetilde{\boldsymbol{\mu}}_{\boldsymbol{g}}, \tilde{v}_\epsilon, \tilde{v}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}})$

$$
\nabla KL(\widetilde{\boldsymbol{\theta}}) = \left( \frac{n}{\tilde{v}_\epsilon} \widetilde{\boldsymbol{\mu}}_{\boldsymbol{g}}, \quad \frac{p(n-1)}{2\tilde{v}_\epsilon} - \frac{n}{2\tilde{v}_\epsilon^2} \sum_{j=1}^{p} (2\tilde{v}_j + \widetilde{\mu}_{g_j}^2), \quad \frac{n}{\tilde{v}_\epsilon} \boldsymbol{1} - \frac{1}{2} \tilde{v}^{-1}, \quad -\frac{1}{2} \psi_1(\tilde{\boldsymbol{\alpha}}), \quad \frac{1}{2} \tilde{\boldsymbol{\beta}}^{-1} \right), \tag{A36}
$$

where **1** is the all-ones vector, and also all operations are componentwise. The gradient and natural gradient algorithms are as follows, respectively

**Algorithm2** :

$$\widetilde{\mu}_g^{(k+1)} = \widetilde{\mu}_g^{(k)} - \gamma \left( \frac{n}{\tilde{v}_\epsilon} \widetilde{\mu}_g \right),$$

$$\tilde{v}_\epsilon^{(k+1)} = \tilde{v}_\epsilon^{(k)} - \gamma \left( \frac{p(n-1)}{2\tilde{v}_\epsilon} - \frac{n}{2\tilde{v}_\epsilon^2} \sum_{j=1}^p (2\tilde{v}_j + \widetilde{\mu}_{g_j}^2) \right),$$

$$\tilde{v}^{(k+1)} = \tilde{v}^{(k)} - \gamma \left( \frac{n}{\tilde{v}_\epsilon} \mathbf{1} - \frac{1}{2} \tilde{v}^{-1} \right),$$

$$\tilde{\alpha}^{(k+1)} = \tilde{\alpha}^{(k)} - \gamma \left( -\frac{1}{2} \psi_1(\tilde{\alpha}) \right),$$

$$\tilde{\beta}^{(k+1)} = \tilde{\beta}^{(k)} - \gamma \left( \frac{1}{2} \tilde{\beta}^{-1} \right),$$

**Algorithm3** :

$$\widetilde{\mu}_g^{(k+1)} = \widetilde{\mu}_g^{(k)} - \frac{1}{\|\Delta KL\|} \left( \frac{n}{\tilde{v}_\epsilon} \widetilde{\mu}_g \right),$$

$$\tilde{v}_\epsilon^{(k+1)} = \tilde{v}_\epsilon^{(k)} - \frac{1}{\|\Delta KL\|} \left( \frac{p(n-1)}{2\tilde{v}_\epsilon} - \frac{n}{2\tilde{v}_\epsilon^2} \sum_{j=1}^p (2\tilde{v}_j + \widetilde{\mu}_{g_j}^2) \right),$$

$$\tilde{v}^{(k+1)} = \tilde{v}^{(k)} - \frac{1}{\|\Delta KL\|} \left( \frac{n}{\tilde{v}_\epsilon} \mathbf{1} - \frac{1}{2} \tilde{v}^{-1} \right),$$

$$\tilde{\alpha}^{(k+1)} = \tilde{\alpha}^{(k)} - \frac{1}{\|\Delta KL\|} \left( -\frac{1}{2} \psi_1(\tilde{\alpha}) \right),$$

$$\tilde{\beta}^{(k+1)} = \tilde{\beta}^{(k)} - \frac{1}{\|\Delta KL\|} \left( \frac{1}{2} \tilde{\beta}^{-1} \right).$$

The objective of the inverse problem is to approximate the distribution of $f$. We use the distribution of $g$ in fig 3 to show the number of method accuracies. Here are the conjectures of the VBA, gradient-based, and natural gradient algorithms, respectively

$$\begin{cases} \mu_{\widetilde{f}} = \begin{bmatrix} 0.025537 \\ 0.069324 \end{bmatrix}, \\ \widetilde{\Sigma}_{\widetilde{f}} = \begin{bmatrix} 0.165003 & 0 \\ 0 & 0.147284 \end{bmatrix}, \end{cases} \quad \begin{cases} \mu_{\widetilde{f}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \\ \widetilde{\Sigma}_{\widetilde{f}} = \begin{bmatrix} 0.158464 & 0 \\ 0 & 2.410353 \end{bmatrix}, \end{cases} \quad \begin{cases} \mu_{\widetilde{f}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \\ \widetilde{\Sigma}_{\widetilde{f}} = \begin{bmatrix} 1.559747 & 0 \\ 0 & 1.291366 \end{bmatrix}. \end{cases}$$