

Bayesian Heteroscedastic Matrix Factorization

Leonara Alves Cesario da Silva^a and Ralph dos Santos Silva^b

^a *Credsystem, Brazil*

^b *Instituto de Matemática, Universidade Federal de Rio de Janeiro, Brazil*

ABSTRACT

This paper proposes a Bayesian heteroscedastic matrix factorization to deal with sparse and high dimensional data. This methodology uses techniques of latent factors, considered the state-of-the-art in recommender systems based on collaborative filtering models. The proposal is to include variations among users to accommodate divergent opinions on various items, which may or may not be more generous in their critiques of the products. For this reason, tailored priors are presented for the parameters to obtain greater scalability and a specific variation for each user. We compare the results, obtained with the probabilistic matrix factorization (PMF) and our Bayesian heteroscedastic matrix factorization (BHMF), in the Netflix and the MovieLens datasets, based on the root mean square error (RMSE).

KEYWORDS

Collaborative Filtering, Latent Factors, Recommender Systems, Sparse Datas

1. Introduction

In a reality where there is a growing increase in online sales, companies that leverage the digital age and adapt to these new types of consumers to profit are likely to sell more. With the popularization of this consumption, there is a noticeable increase in interaction between consumers and companies to determine if purchasing a particular product is worthwhile - through the assignment of certain reviews that identify whether the item would be a good buy or not. The result of these opinions could offer all kinds of companies an opportunity to see business prospects, as they specify unique characteristics of each customer regarding their experiences with the company.

If a company has specific information about all its customers' past purchases, one could conjecture the range of possibilities it would have if it could anticipate future purchases by its customers. The consequence of this is what will be referred to as recommender systems. These systems aim to synthesize information from each user and attempt to anticipate each customer's evaluations of a given product. Based on this, the process could preemptively recommend products that the user might be interested in before the customer even realizes they want that item.

With the increasing ease of obtaining user review data, it is important to develop specific analytical procedures that reasonably capture the obtained information and can help companies gain insights about their customers. Based on this, recommender

CONTACT Author^b. Email: ralph@im.ufrj.br

Article History

Received : 8 January 2025; Revised : 7 February 2025; Accepted : 16 February 2025; Published : 28 June 2025

To cite this paper

Leonara Alves Cesario da Silva & Ralph dos Santos Silva (2025). Bayesian Heteroscedastic Matrix Factorization. *Journal of Econometrics and Statistics*. 5(2), 129-148.

systems can be classified into various categories: collaborative filtering, content-based, knowledge-based, demographic, and hybrid. Collaborative filtering methods use information from similar users or items to make recommendations, while content-based methods use knowledge about past consumed products (such as genre, description, actors, among others).

Despite the various strategies for recommending products to users, collaborative filtering models are more commonly used because they feature characteristics similar to those applied in machine learning in the context of classification. This fact justifies the use of established models from the literature for making recommendations, such as neural networks, support vector machines, decision trees, Bayesian models, etc. Notwithstanding, in the era of big data, there is increasing concern about the amount of available information, and in this case, commonly found models may not be efficient.

In the literature, there are several techniques proposed for analyzing review data. Nevertheless, most of these data do not have high dimensionality - millions of users rating thousands of products. Canny [2] developed a new collaborative filtering procedure that protects user privacy while also applying probabilistic factor analysis to handle missing data. The author used a linear approach that generally extends singular value decomposition (SVD) methods and linear regression. For his applications, the author employed the iterative method known as the Expectation-Maximization (EM) algorithm, which is suited for handling sparse data and simple recursive definitions that can be combined with his main idea: privacy.

Xian et al. [11] applied an extension of SVD, called SVD++. This method is characterized by the addition of implicit feedback information from users. They used this methodology with a focus on respecting user privacy, employing a terminology called differential privacy, which ensures that other people cannot discover personal characteristics of individuals in the database.

Gemulla et al. [4] constructed an algorithm to approximate large matrices—with millions of user ratings for millions of products, amounting to billions of non-empty elements. This new algorithm is a distributed extension of the method using stochastic gradient descent, a stochastic optimization algorithm. Thus, it has been shown that it is possible to handle web-scale matrices with rapid convergence and scalability.

A study by Devooght, Kourtellis and Mantrach [3] explained that since it is commonly assumed that the distribution of missing ratings is the same as that of observed ratings, they introduced a dynamic matrix factorization structure allowing for explicit priors for these unknown values. The uniqueness of their work lies in the fact that for new users or items, the factorization can be updated regardless of the size of the data, allowing for rapid recommendations for new users.

Kabbur, Ning and Karypis [6] highlighted that the effectiveness of recommendation methods declines for the top- k nearest users or items with increased data sparsity. Based on this, they proposed an item-based method to generate top- k recommendations that learns the similarity matrix based on latent factors. Results from different datasets with varying degrees of sparsity indicated that their method performs better compared to methodologies seen in the literature.

Shen et al. [10] constructed a new recommendation system that utilizes specific characteristics of the locations users have frequented in the past to make recommendations. The technique is called landmark recommendation and attempts to identify, through a unified classifier, preference styles based on domain adaptation, leveraging photos from sites in the source domain and reference images in the target domain. The detected styles are then used to learn users' best preferences and make recommendations. They highlight that personalized methodologies for each user are a highly

effective approach in real-world travel data applications.

The overall objective of this paper is to present a procedure for predicting user ratings of specific items using real, high-dimensional data, particularly when conventional mechanisms employed in the literature are not suitable. We develop a methodology for estimating user ratings for a predefined set of items using matrix factorization, specifically probabilistic matrix decomposition, Bayesian approaches to probabilistic matrix decomposition, and propose generalizations that include heteroscedasticity among items or users. Moreover, we assess the ability of the proposed method to more accurately predict user ratings using appropriately selected training and test data. Finally, we compare, using the same dataset, the performance of the probabilistic matrix factorization and our Bayesian heteroscedastic matrix factorization.

The remainder of this paper is structured as follows: Section 2 briefly reviews latent factor models, including the probabilistic matrix factorization and the Bayesian matrix factorization models. Moreover, it introduces the new Bayesian heteroscedastic matrix factorization model. Section 3 brings two real data applications which are based on the Netflix and the MovieLens datasets. Section 4 concludes.

2. Latent factor models

With the development of technology, there is a growing interaction between users of a website on the internet and their opinions about the products consumed. In such cases, the seller may be interested in increasing sales based on these customer responses by employing techniques to evaluate the relevance of certain items to users. To this end, the merchant can leverage various data characteristics to make recommendations, such as similarities between users, similarities between items, and similarities between users and items. Here, “user” refers to the individual to whom the recommendation is provided, and “item” refers to the product that will be recommended to the user.

Let $\mathcal{R} = [r_{ij}]$ be the $m \times n$ matrix of user-item ratings, where m is the number of users, n is the number of items, and r_{ij} is the rating given by the i -th user for item j . It is important to note that there may be missing data in any of the rows or columns of this matrix.

In the basic matrix factorization model, a way to decompose the rating matrix \mathcal{R} is by approximating it as the product of a matrix \mathbf{U} ($m \times p$) and a matrix \mathbf{V} ($n \times p$) such that:

$$\mathcal{R} \approx \mathbf{U}\mathbf{V}^\top. \quad (1)$$

Each column of \mathbf{U} (or \mathbf{V}) can be referred to as a latent vector or component, and each row of \mathbf{U} (or \mathbf{V}) is classified as a latent factor. The approximation error is given by:

$$\|\mathcal{R} - \mathbf{U}\mathbf{V}^\top\|^2,$$

where $\|\cdot\|$ denotes the Frobenius norm (squared). Therefore, the objective function corresponds to the sum of the squares of the entries in the resulting residual matrix ($\mathcal{R} - \mathbf{U}\mathbf{V}^\top$):

$$e_{ij} = r_{ij} - \mathbf{v}_i^\top \mathbf{v}_j, \quad i = 1, 2, \dots, m, \quad \text{and} \quad j = 1, 2, \dots, n,$$

such that the smaller the value of this function, the better the previously indicated factorization ($\mathcal{R} \approx \mathbf{U}\mathbf{V}^\top$) will fit. The i -th row, $\mathbf{v}_i = (v_{i1}, \dots, v_{ip})^\top$, of \mathbf{U} is called the “user factor” and contains p entries corresponding to the affinities of user i with the p concepts (factors) in the rating matrix. Similarly, each column, $\mathbf{v}_j = (v_{j1}, \dots, v_{jp})^\top$, of \mathbf{V} is referred to as the “item factor” and represents the affinity of item j with respect to the p concepts.

From Equation (1), it is observed that each rating r_{ij} in \mathcal{R} can be approximately represented as the dot product of the i -th user factor and the j -th item factor, described by:

$$r_{ij} \approx \mathbf{v}_i^\top \mathbf{v}_j = \sum_{s=1}^p v_{is} \times v_{js}. \quad (2)$$

Singular value decomposition (SVD) is a matrix factorization method where the columns of \mathbf{U} and \mathbf{V} are constrained to be mutually orthogonal [1]. This approach has the advantage that the concepts can be completely uncorrelated. Consider the case where a fully specified matrix is available. The rating matrix \mathcal{R} can be factorized using a truncated SVD with rank $p \ll \min\{m, n\}$ (assuming that the matrix \mathcal{R} is of full rank) as follows:

$$\mathcal{R} = \mathbf{H}\mathbf{\Delta}\mathbf{K}^\top,$$

where \mathbf{H} , $\mathbf{\Delta}$, and \mathbf{K} are matrices of size $m \times p$, $p \times p$, and $n \times p$, respectively.

The matrices \mathbf{H} and \mathbf{K} contain the p largest eigenvectors of $\mathcal{R}\mathcal{R}^\top$ and $\mathcal{R}^\top\mathcal{R}$, respectively, and the diagonal matrix $\mathbf{\Delta}$ contains the p eigenvalues. The eigenvectors provide information about the directions of item-item (or user-user) correlations in the ratings, and consequently, enable the representation of each user (or item) in a reduced number of dimensions [1].

The diagonal matrix $\mathbf{\Delta}$ can be incorporated into the user factors \mathbf{H} or the item factors \mathbf{K} . By convention,

$$\mathbf{U} = \mathbf{H}\mathbf{\Delta} \quad \text{and} \quad \mathbf{V} = \mathbf{K}^\top.$$

Let the sum of squared residuals be defined by

$$S(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} e_{ij}^2 = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^\top \mathbf{v}_j)^2,$$

where \mathbb{I}_{ij} is the indicator variable that equals 1 if user i rated item j and 0 otherwise. Thus, the objective function can be formulated as an optimization problem over the matrices \mathbf{U} and \mathbf{V} , that is, $\min_{(\mathbf{U}, \mathbf{V})} S(\mathbf{U}, \mathbf{V})$ such that subject to columns of \mathbf{U} are mutually orthogonal and columns of \mathbf{V} are also mutually orthogonal.

Since the previously described scenario deals with fully observed matrices, applying an iterative process is computationally expensive, with a complexity of $O(m^2n + mn^2 + n^3)$ for an $m \times n$ matrix [5]. Therefore, this approach can be modified to an optimization problem with regularization [1, 9]. This treatment has the advantage of preventing overfitting by adding a bias to the model, favoring simplicity over complexity. Thus,

the incorporation of regularization is given by

$$\min_{(\mathbf{U}, \mathbf{V})} \left[S(\mathbf{U}, \mathbf{V}) + \frac{\lambda_{\mathbf{v}}}{2} \sum_{i=1}^m \|\mathbf{v}_i\|^2 + \frac{\lambda_{\mathbf{v}}}{2} \sum_{j=1}^n \|\mathbf{v}_j\|^2 \right], \quad (3)$$

with the same restrictions over \mathbf{U} and \mathbf{V} as before.

2.1. Probabilistic matrix factorization

The probabilistic matrix factorization (PMF), proposed by Salakhutdinov and Mnih [8], is a linear probabilistic model with normally distributed observed errors, where the distribution of \mathcal{R} conditional on \mathbf{U} and \mathbf{V} (likelihood function) and the prior distributions on \mathbf{U} and \mathbf{V} are given by

$$\begin{aligned} p(\mathcal{R}|\mathbf{U}, \mathbf{V}, \alpha) &= \prod_{i=1}^m \prod_{j=1}^n [\phi_1(r_{ij}|\mathbf{v}_i^\top \mathbf{v}_j, \alpha^{-1})]^{\mathbb{I}_{ij}} \\ p(\mathbf{U}|\alpha_{\mathbf{v}}) &= \prod_{i=1}^m \phi_p(\mathbf{v}_i|0, \alpha_{\mathbf{v}}^{-1}\mathbf{I}) \\ p(\mathbf{V}|\alpha_{\mathbf{v}}) &= \prod_{j=1}^n \phi_p(\mathbf{v}_j|0, \alpha_{\mathbf{v}}^{-1}\mathbf{I}), \end{aligned} \quad (4)$$

where $\phi_p(\cdot|\boldsymbol{\mu}, \boldsymbol{\Omega}^{-1})$ denotes the density function of a p -variate normal distribution with mean vector $\boldsymbol{\mu}$ and precision matrix $\boldsymbol{\Omega}$, for $p = 1, 2, \dots, \min\{m, n\}$.

Note that the structure of the model assumes that all components of \mathbf{U} , and similarly those of \mathbf{V} , are a priori independent with the same precision structure controlled by $\alpha_{\mathbf{v}}^{-1}$ (and $\alpha_{\mathbf{v}}^{-1}$). Thus, learning in this model is achieved by maximizing the log-posterior conditional on the user and item factors with fixed hyperparameters ($\alpha, \alpha_{\mathbf{v}}, \alpha_{\mathbf{v}}$):

$$\begin{aligned} \ln p(\mathbf{U}, \mathbf{V}|\mathcal{R}, \alpha, \alpha_{\mathbf{v}}, \alpha_{\mathbf{v}}) &= \ln p(\mathcal{R}|\mathbf{U}, \mathbf{V}, \alpha) + \ln p(\mathbf{U}|\alpha_{\mathbf{v}}) + \ln p(\mathbf{V}|\alpha_{\mathbf{v}}) \\ &= \eta - \frac{\alpha}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^\top \mathbf{v}_j)^2 \\ &\quad - \frac{\alpha_{\mathbf{v}}}{2} \sum_{i=1}^m \mathbf{v}_i^\top \mathbf{v}_i - \frac{\alpha_{\mathbf{v}}}{2} \sum_{j=1}^n \mathbf{v}_j^\top \mathbf{v}_j \\ &\quad - \frac{1}{2} \left[\left(\sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} \right) \ln \alpha^{-1} + mp \ln \alpha_{\mathbf{v}}^{-1} + np \ln \alpha_{\mathbf{v}}^{-1} \right], \end{aligned}$$

where η is a constant as well as the last term of the expression.

According to Salakhutdinov and Mnih [8] and Salakhutdinov and Mnih [7], maximizing this expression is equivalent to minimizing the sum of squared errors in the objective function with quadratic regularization terms given in Equation (3), where $\lambda_{\mathbf{v}} = \alpha_{\mathbf{v}}/\alpha$ and $\lambda_{\mathbf{v}} = \alpha_{\mathbf{v}}/\alpha$.

2.2. Bayesian matrix factorization

The PMF model has the drawback of not automatically tuning the parameters $\lambda_{\mathbf{v}}$, λ_{ν} , and α . In this case, it is useful to employ another methodology capable of automatically adjusting the complexity of the model using a fully Bayesian technique. This method is called Bayesian matrix factorization (BMF), as presented by Salakhutdinov and Mnih [7], which uses the same likelihood function proposed by PMF in Equation (4), but with normal prior distributions on the user factor matrix \mathbf{U} and the item factor matrix \mathbf{V} , expressed as:

$$p(\mathbf{U}|\boldsymbol{\mu}_{\mathbf{v}}, \boldsymbol{\Lambda}_{\mathbf{v}}) = \prod_{i=1}^m \phi_p(\mathbf{v}_i|\boldsymbol{\mu}_{\mathbf{v}}, \boldsymbol{\Lambda}_{\mathbf{v}}^{-1}) \quad \text{and} \quad p(\mathbf{V}|\boldsymbol{\mu}_{\nu}, \boldsymbol{\Lambda}_{\nu}) = \prod_{j=1}^n \phi_p(\nu_j|\boldsymbol{\mu}_{\nu}, \boldsymbol{\Lambda}_{\nu}^{-1}).$$

The most notable distinction of this proposal from the previous PMF model is the generalization, allowing the components of the vector \mathbf{v}_i (and also ν_j) to be correlated through a full precision matrix (full covariance matrix) and to have prior means different from zero. Nonetheless, it is worth noting that the covariance structure is the same between the columns of \mathbf{U} and \mathbf{V} , respectively. That is, it is assumed that a priori the users (and also the items) can be represented with fixed latent covariance structures.

Furthermore, a normal-Wishart prior distributions have been assigned to the parameters for users $\boldsymbol{\Theta}_{\mathbf{v}} = \{\boldsymbol{\mu}_{\mathbf{v}}, \boldsymbol{\Lambda}_{\mathbf{v}}\}$ and for items $\boldsymbol{\Theta}_{\nu} = \{\boldsymbol{\mu}_{\nu}, \boldsymbol{\Lambda}_{\nu}\}$, as follows:

$$\begin{aligned} p(\boldsymbol{\Theta}_{\mathbf{v}}|\boldsymbol{\Theta}_0) &= p(\boldsymbol{\mu}_{\mathbf{v}}|\boldsymbol{\Lambda}_{\mathbf{v}})p(\boldsymbol{\Lambda}_{\mathbf{v}}) = \phi_p(\boldsymbol{\mu}_{\mathbf{v}}|\mathbf{a}_{\mathbf{v}}, (b_{\mathbf{v}}\boldsymbol{\Lambda}_{\mathbf{v}})^{-1})\omega(\boldsymbol{\Lambda}_{\mathbf{v}}|\mathbf{C}_{\mathbf{v}}, d_{\mathbf{v}}) \\ p(\boldsymbol{\Theta}_{\nu}|\boldsymbol{\Theta}_0) &= p(\boldsymbol{\mu}_{\nu}|\boldsymbol{\Lambda}_{\nu})p(\boldsymbol{\Lambda}_{\nu}) = \phi_p(\boldsymbol{\mu}_{\nu}|\mathbf{a}_{\nu}, (b_{\nu}\boldsymbol{\Lambda}_{\nu})^{-1})\omega(\boldsymbol{\Lambda}_{\nu}|\mathbf{C}_{\nu}, d_{\nu}), \end{aligned}$$

where ω represents the density function of a random matrix with a Wishart distribution with ψ_0 degrees of freedom and a scale matrix \mathbf{W}_0 ($p \times p$). It can be expressed as $(\boldsymbol{\Lambda}|\mathbf{W}_0, \psi_0) \sim \mathcal{W}(\mathbf{W}_0, \psi_0)$. In this paper, $\boldsymbol{\Theta}_0$ represents the set of hyperparameters.

Moreover, the predictive distribution of the rating $\mathcal{R}^* = [r_{ij}^*]$ for the i -th user rating the j -th item is obtained by marginalizing the model parameters:

$$p(r_{ij}^*|\mathcal{R}, \boldsymbol{\Theta}_0) = \int p(r_{ij}^*|\mathbf{v}_i, \nu_j)p(\mathbf{U}, \mathbf{V}|\mathcal{R}, \boldsymbol{\Theta}_{\mathbf{v}}, \boldsymbol{\Theta}_{\nu})p(\boldsymbol{\Theta}_{\mathbf{v}}, \boldsymbol{\Theta}_{\nu}|\boldsymbol{\Theta}_0)d\{\mathbf{U}, \mathbf{V}, \boldsymbol{\Theta}_{\mathbf{v}}, \boldsymbol{\Theta}_{\nu}\}. \quad (5)$$

Since the exact value of this predictive distribution is analytically intractable due to the complexity of the posterior, inferential approximations are needed. As a result, Salakhutdinov and Mnih [7] used Markov chain Monte Carlo (MCMC) methods, specifically Gibbs sampler. Basically, MCMC methods use Monte Carlo approximation for the predictive distribution in Equation (5), expressed as:

$$p(r_{ij}^*|\mathcal{R}, \boldsymbol{\Theta}_0) \approx \frac{1}{K} \sum_{k=1}^K p(r_{ij}^*|\mathbf{v}_i^{(k)}, \nu_j^{(k)}),$$

where samples $\{\mathbf{v}_i^{(k)}, \nu_j^{(k)}\}$ are generated from the full conditional distributions of the current values of all variables (Gibbs sampler).

The choice of this approach is related to the use of normal distributions for the user and item factors, as their respective full conditional distributions, given the hyperparameter values $\boldsymbol{\Theta}_0$ and the matrix \mathcal{R} , are also normal (see Appendix A for our general case):

- Full conditional: $(\mathbf{v}_i | \mathcal{R}, \mathbf{V}, \Theta_{\mathbf{v}}, \alpha) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{v}_i}^*, [\boldsymbol{\Lambda}_{\mathbf{v}_i}^*]^{-1})$, where

$$\boldsymbol{\Lambda}_{\mathbf{v}_i}^* = \boldsymbol{\Lambda}_{\mathbf{v}} + \alpha \sum_{j=1}^n \mathbb{I}_{ij} [\boldsymbol{\nu}_j \boldsymbol{\nu}_j^\top] \quad \text{and} \quad \boldsymbol{\mu}_{\mathbf{v}_i}^* = [\boldsymbol{\Lambda}_{\mathbf{v}_i}^*]^{-1} \left(\boldsymbol{\Lambda}_{\mathbf{v}} \boldsymbol{\mu}_{\mathbf{v}} + \alpha \sum_{j=1}^n \mathbb{I}_{ij} [\boldsymbol{\nu}_j r_{ij}] \right);$$

- Full conditional: $(\boldsymbol{\nu}_j | \mathcal{R}, \mathbf{U}, \Theta_{\boldsymbol{\nu}}, \alpha) \sim \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\nu}_j}^*, [\boldsymbol{\Lambda}_{\boldsymbol{\nu}_j}^*]^{-1})$, where

$$\boldsymbol{\Lambda}_{\boldsymbol{\nu}_j}^* = \boldsymbol{\Lambda}_{\boldsymbol{\nu}} + \alpha \sum_{i=1}^m \mathbb{I}_{ij} [\mathbf{v}_i \mathbf{v}_i^\top] \quad \text{and} \quad \boldsymbol{\mu}_{\boldsymbol{\nu}_j}^* = [\boldsymbol{\Lambda}_{\boldsymbol{\nu}_j}^*]^{-1} \left(\boldsymbol{\Lambda}_{\boldsymbol{\nu}} \boldsymbol{\mu}_{\boldsymbol{\nu}} + \alpha \sum_{i=1}^m \mathbb{I}_{ij} [\mathbf{v}_i r_{ij}] \right);$$

- Full conditional: $(\boldsymbol{\mu}_{\mathbf{v}}, \boldsymbol{\Lambda}_{\mathbf{v}} | \mathbf{U}, \Theta_0) \sim \mathcal{NW}(\mathbf{a}_{\mathbf{v}}^*, [b_{\mathbf{v}}^* \boldsymbol{\Lambda}_{\mathbf{v}}]^{-1}, \mathbf{C}_{\mathbf{v}}^*, d_{\mathbf{v}}^*)$, where

$$b_{\mathbf{v}}^* = b_{\mathbf{v}} + m, \quad d_{\mathbf{v}}^* = d_{\mathbf{v}} + m, \quad \mathbf{a}_{\mathbf{v}}^* = (m\bar{\mathbf{v}} + b_{\mathbf{v}} \mathbf{a}_{\mathbf{v}}) / b_{\mathbf{v}}^*, \quad \bar{\mathbf{Q}}_{\mathbf{v}} = \sum_{i=1}^m (\mathbf{v}_i - \bar{\mathbf{v}})(\mathbf{v}_i - \bar{\mathbf{v}})^\top,$$

$$\bar{\mathbf{Q}}_{\mathbf{a}} = m b_{\mathbf{v}} (\mathbf{a}_{\mathbf{v}} - \bar{\mathbf{v}})(\mathbf{a}_{\mathbf{v}} - \bar{\mathbf{v}})^\top / b_{\mathbf{v}}^*, \quad [\mathbf{C}_{\mathbf{v}}^*]^{-1} = \mathbf{C}_{\mathbf{v}}^{-1} + \bar{\mathbf{Q}}_{\mathbf{v}} + \bar{\mathbf{Q}}_{\mathbf{a}}, \quad \bar{\mathbf{v}} = (1/m) \sum_{i=1}^m \mathbf{v}_i;$$

- Full conditional for: $p(\boldsymbol{\mu}_{\boldsymbol{\nu}}, \boldsymbol{\Lambda}_{\boldsymbol{\nu}} | \mathbf{V}, \Theta_0) \sim \mathcal{NW}(\mathbf{a}_{\boldsymbol{\nu}}^*, [b_{\boldsymbol{\nu}}^* \boldsymbol{\Lambda}_{\boldsymbol{\nu}}]^{-1}, \mathbf{C}_{\boldsymbol{\nu}}^*, d_{\boldsymbol{\nu}}^*)$, where

$$b_{\boldsymbol{\nu}}^* = b_{\boldsymbol{\nu}} + n, \quad d_{\boldsymbol{\nu}}^* = d_{\boldsymbol{\nu}} + n, \quad \mathbf{a}_{\boldsymbol{\nu}}^* = (n\bar{\boldsymbol{\nu}} + b_{\boldsymbol{\nu}} \mathbf{a}_{\boldsymbol{\nu}}) / b_{\boldsymbol{\nu}}^*, \quad \bar{\mathbf{Q}}_{\boldsymbol{\nu}} = \sum_{j=1}^n (\boldsymbol{\nu}_j - \bar{\boldsymbol{\nu}})(\boldsymbol{\nu}_j - \bar{\boldsymbol{\nu}})^\top,$$

$$\bar{\mathbf{Q}}_{\mathbf{b}} = n b_{\boldsymbol{\nu}} (\mathbf{a}_{\boldsymbol{\nu}} - \bar{\boldsymbol{\nu}})(\mathbf{a}_{\boldsymbol{\nu}} - \bar{\boldsymbol{\nu}})^\top / b_{\boldsymbol{\nu}}^*, \quad [\mathbf{C}_{\boldsymbol{\nu}}^*]^{-1} = \mathbf{C}_{\boldsymbol{\nu}}^{-1} + \bar{\mathbf{Q}}_{\boldsymbol{\nu}} + \bar{\mathbf{Q}}_{\mathbf{b}}, \quad \bar{\boldsymbol{\nu}} = (1/n) \sum_{j=1}^n \boldsymbol{\nu}_j.$$

Hyperparameter α

In the BMF methodology, the parameter α is treated as the precision error of the observations and it is fixed at the value of 2. Nevertheless, a prior can be applied to this parameter and included in the Gibbs sampler for estimation. It is assigned as $(\alpha | a_\alpha, b_\alpha) \sim \mathcal{G}(a_\alpha, b_\alpha)$, that is, a gamma distribution with mean a_α / b_α . Thus,

- Full conditional for: $(\alpha | \mathbf{U}, \mathbf{V}, \mathcal{R}) \sim \mathcal{G}(a_\alpha^*, b_\alpha^*)$, where $a_\alpha^* = \frac{N}{2} + a_\alpha$ with $N = \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij}$, and $b_\alpha^* = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^\top \boldsymbol{\nu}_j)^2 + b_\alpha$.

2.3. Bayesian heteroscedastic matrix factorization

In the BMF procedure, the factors for users (and items) can be correlated, but the variance for all users (and all items) is the same. This suggests that there are no discrepancies considered between user and item ratings. For example, it is not considered that one user might rate within a scale of 1 to 2, while a second user might rate within a scale of 2 to 5. Similarly, for items, a highly rated item might be classified only from 4 onwards, while a poorly rated item might be limited to a scale from 1 to 3.

To address this more generally, we propose a Bayesian heteroscedastic matrix factorization model (BHMF), which extends the BMF by adding a parameter $\lambda_{\mathbf{v}_i}$ for

each user and λ_{ν_j} for each item. This new model is novel and it can be presented as follows:

$$\begin{aligned} p(\mathcal{R}|\mathbf{U}, \mathbf{V}, \alpha) &= \prod_{i=1}^m \prod_{j=1}^n [\phi_1(r_{ij}|\mathbf{v}_i^\top \boldsymbol{\nu}_j, \alpha^{-1})]^{\mathbb{I}_{ij}}, \\ p(\mathbf{U}|\boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu, \boldsymbol{\lambda}_\nu) &= \prod_{i=1}^m \phi_p(\mathbf{v}_i|\boldsymbol{\mu}_\nu, [\boldsymbol{\lambda}_{\nu_i} \boldsymbol{\Lambda}_\nu]^{-1}), \\ p(\mathbf{V}|\boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu, \boldsymbol{\lambda}_\nu) &= \prod_{j=1}^n \phi_p(\boldsymbol{\nu}_j|\boldsymbol{\mu}_\nu, [\boldsymbol{\lambda}_{\nu_j} \boldsymbol{\Lambda}_\nu]^{-1}), \\ p(\alpha|a_\alpha, b_\alpha) &= \gamma(\alpha|a_\alpha, b_\alpha), \end{aligned}$$

where $\boldsymbol{\lambda}_\nu = (\lambda_{\nu_1}, \dots, \lambda_{\nu_m})$ and $\boldsymbol{\lambda}_\nu = (\lambda_{\nu_1}, \dots, \lambda_{\nu_n})$, whilst $\gamma(\cdot|a_\alpha, b_\alpha)$ denotes the probability density function of a gamma random variable with mean a_α/b_α .

Additionally, a gamma prior distribution is assigned to the user parameters, λ_{ν_i} , and item parameters, λ_{ν_j} , as well as to the parameters κ_ν and κ_ν , which refer to the rate and shape parameters of the gamma distribution for λ_{ν_i} and λ_{ν_j} , respectively:

$$\begin{aligned} p(\boldsymbol{\lambda}_\nu|\kappa_\nu) &= \prod_{i=1}^m p(\lambda_{\nu_i}|\kappa_\nu) = \prod_{i=1}^m \left[\gamma(\lambda_{\nu_i}|\frac{\kappa_\nu}{2}, \frac{\kappa_\nu}{2}) \right], \quad p(\kappa_\nu) = \gamma(\kappa_\nu|a_{\kappa_\nu}, b_{\kappa_\nu}), \\ p(\boldsymbol{\lambda}_\nu|\kappa_\nu) &= \prod_{j=1}^n p(\lambda_{\nu_j}|\kappa_\nu) = \prod_{j=1}^n \left[\gamma(\lambda_{\nu_j}|\frac{\kappa_\nu}{2}, \frac{\kappa_\nu}{2}) \right], \quad p(\kappa_\nu) = \gamma(\kappa_\nu|a_{\kappa_\nu}, b_{\kappa_\nu}). \end{aligned}$$

An important highlight is that modeling in this way considers λ_{ν_i} and λ_{ν_j} as auxiliary variables in the scale mixture of the normal distribution with the gamma distribution to obtain the marginal Student's t -distribution.

The normal-Wishart prior distributions for the user parameters $\boldsymbol{\Theta}_\nu = \{\boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu\}$ and item parameters $\boldsymbol{\Theta}_\nu = \{\boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu\}$ have been updated to

$$\begin{aligned} p(\boldsymbol{\Theta}_\nu|\boldsymbol{\Theta}_0) &= p(\boldsymbol{\mu}_\nu|\boldsymbol{\Lambda}_\nu)p(\boldsymbol{\Lambda}_\nu) = \phi_p(\boldsymbol{\mu}_\nu|\mathbf{a}_\nu, [b_\nu \boldsymbol{\Lambda}_\nu]^{-1})\omega(\boldsymbol{\Lambda}_\nu|\mathbf{C}_\nu, d_\nu) \\ p(\boldsymbol{\Theta}_\nu|\boldsymbol{\Theta}_0) &= p(\boldsymbol{\mu}_\nu|\boldsymbol{\Lambda}_\nu)p(\boldsymbol{\Lambda}_\nu) = \phi_p(\boldsymbol{\mu}_\nu|\mathbf{a}_\nu, [b_\nu \boldsymbol{\Lambda}_\nu]^{-1})\omega(\boldsymbol{\Lambda}_\nu|\mathbf{C}_\nu, d_\nu). \end{aligned}$$

Moreover, the predictive distribution of the rating $\mathcal{R}^* = [r_{ij}^*]$ for the i -th user who evaluated item j is obtained by marginalizing the model parameters:

$$\begin{aligned} p(r_{ij}^*|\mathcal{R}, \boldsymbol{\Theta}_0) &= \int p(r_{ij}^*|\mathbf{v}_i, \boldsymbol{\nu}_j, \alpha)p(\mathbf{U}, \mathbf{V}|\mathcal{R}, \boldsymbol{\Theta}_\nu, \boldsymbol{\Theta}_\nu, \boldsymbol{\lambda}_\nu, \boldsymbol{\lambda}_\nu) \\ &\times p(\boldsymbol{\Theta}_\nu, \boldsymbol{\Theta}_\nu|\boldsymbol{\Theta}_0)p(\alpha)p(\boldsymbol{\lambda}_\nu, \boldsymbol{\lambda}_\nu)d\{\mathbf{U}, \mathbf{V}, \boldsymbol{\Theta}_\nu, \boldsymbol{\Theta}_\nu, \alpha, \boldsymbol{\lambda}_\nu, \boldsymbol{\lambda}_\nu\}. \end{aligned}$$

Finally, the full conditionals for the BHMF model are described below (see Appendix A):

- Full conditional: $(\alpha|\mathcal{R}, \mathbf{U}, \mathbf{V}, a_\alpha, b_\alpha) \sim \mathcal{G}(a_\alpha^*, b_\alpha^*)$, where $a_\alpha^* = \frac{N}{2} + a_\alpha$ with $N = \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij}$, and $b_\alpha^* = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^\top \boldsymbol{\nu}_j)^2 + b_\alpha$;
- Full conditional: $(\mathbf{v}_i|\mathcal{R}, \mathbf{V}, \boldsymbol{\Theta}_\nu, \alpha) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{v}_i}^*, [\boldsymbol{\Lambda}_{\mathbf{v}_i}^*]^{-1})$, where

$$\mathbf{\Lambda}_{\mathbf{v}_i}^* = \lambda_{\mathbf{v}_i} \mathbf{\Lambda}_{\mathbf{v}} + \alpha \sum_{j=1}^n \mathbb{I}_{ij} [\boldsymbol{\nu}_j \boldsymbol{\nu}_j^\top] \quad \text{and} \quad \boldsymbol{\mu}_{\mathbf{v}_i}^* = [\mathbf{\Lambda}_{\mathbf{v}_i}^*]^{-1} \left(\lambda_{\mathbf{v}_i} \mathbf{\Lambda}_{\mathbf{v}} \boldsymbol{\mu}_{\mathbf{v}} + \alpha \sum_{j=1}^n \mathbb{I}_{ij} [r_{ij} \boldsymbol{\nu}_j] \right);$$

- Full conditional: $(\boldsymbol{\nu}_j | \mathcal{R}, \mathbf{U}, \boldsymbol{\Theta}_{\mathbf{v}}, \alpha) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{v}_j}^*, [\mathbf{\Lambda}_{\mathbf{v}_j}^*]^{-1})$, where

$$\mathbf{\Lambda}_{\mathbf{v}_j}^* = \lambda_{\mathbf{v}_j} \mathbf{\Lambda}_{\mathbf{v}} + \alpha \sum_{i=1}^m \mathbb{I}_{ij} [\mathbf{v}_i \mathbf{v}_i^\top], \quad \text{and} \quad \boldsymbol{\mu}_{\mathbf{v}_j}^* = [\mathbf{\Lambda}_{\mathbf{v}_j}^*]^{-1} \left(\lambda_{\mathbf{v}_j} \mathbf{\Lambda}_{\mathbf{v}} \boldsymbol{\mu}_{\mathbf{v}} + \alpha \sum_{i=1}^m \mathbb{I}_{ij} [r_{ij} \mathbf{v}_i] \right);$$

- Full conditional: $(\boldsymbol{\mu}_{\mathbf{v}}, \mathbf{\Lambda}_{\mathbf{v}} | \mathbf{U}, \boldsymbol{\Theta}_0) \sim \mathcal{NW}(\mathbf{a}_{\mathbf{v}}^*, [b_{\mathbf{v}}^* \mathbf{\Lambda}_{\mathbf{v}}]^{-1}, \mathbf{C}_{\mathbf{v}}^*, d_{\mathbf{v}}^*)$, where

$$b_{\mathbf{v}}^* = b_{\mathbf{v}} + m \bar{\lambda}_{\mathbf{v}}, \quad d_{\mathbf{v}}^* = d_{\mathbf{v}} + m, \quad \mathbf{a}_{\mathbf{v}}^* = \frac{1}{b_{\mathbf{v}}^*} (m \bar{\lambda}_{\mathbf{v}} \tilde{\mathbf{v}} + b_{\mathbf{v}} \mathbf{a}_{\mathbf{v}}),$$

$$\tilde{\mathbf{Q}}_{\mathbf{v}} = \sum_{i=1}^m \lambda_{\mathbf{v}_i} (\mathbf{v}_i - \tilde{\mathbf{v}})(\mathbf{v}_i - \tilde{\mathbf{v}})^\top, \quad \tilde{\mathbf{Q}}_{\mathbf{a}_{\mathbf{v}}} = \frac{m \bar{\lambda}_{\mathbf{v}} b_{\mathbf{v}}}{b_{\mathbf{v}}^*} (\mathbf{a}_{\mathbf{v}} - \tilde{\mathbf{v}})(\mathbf{a}_{\mathbf{v}} - \tilde{\mathbf{v}})^\top,$$

$$[\mathbf{C}_{\mathbf{v}}^*]^{-1} = \mathbf{C}_{\mathbf{v}}^{-1} + \tilde{\mathbf{Q}}_{\mathbf{v}} + \tilde{\mathbf{Q}}_{\mathbf{a}_{\mathbf{v}}}, \quad \bar{\lambda}_{\mathbf{v}} = \frac{1}{m} \sum_{i=1}^m \lambda_{\mathbf{v}_i}, \quad \tilde{\mathbf{v}} = \frac{1}{m \bar{\lambda}_{\mathbf{v}}} \sum_{i=1}^m \lambda_{\mathbf{v}_i} \mathbf{v}_i.$$

- Full conditional: $(\boldsymbol{\mu}_{\mathbf{v}}, \mathbf{\Lambda}_{\mathbf{v}} | \mathbf{V}, \boldsymbol{\Theta}_0) \sim \mathcal{NW}(\mathbf{a}_{\mathbf{v}}^*, [b_{\mathbf{v}}^* \mathbf{\Lambda}_{\mathbf{v}}]^{-1}, \mathbf{C}_{\mathbf{v}}^*, d_{\mathbf{v}}^*)$, where

$$b_{\mathbf{v}}^* = b_{\mathbf{v}} + n \bar{\lambda}_{\mathbf{v}}, \quad d_{\mathbf{v}}^* = d_{\mathbf{v}} + n, \quad \mathbf{a}_{\mathbf{v}}^* = \frac{1}{b_{\mathbf{v}}^*} (n \bar{\lambda}_{\mathbf{v}} \tilde{\boldsymbol{\nu}} + b_{\mathbf{v}} \mathbf{a}_{\mathbf{v}}),$$

$$\tilde{\mathbf{Q}}_{\mathbf{v}} = \sum_{j=1}^n \lambda_{\mathbf{v}_j} (\boldsymbol{\nu}_j - \tilde{\boldsymbol{\nu}})(\boldsymbol{\nu}_j - \tilde{\boldsymbol{\nu}})^\top, \quad \tilde{\mathbf{Q}}_{\mathbf{a}_{\mathbf{v}}} = \frac{n \bar{\lambda}_{\mathbf{v}} b_{\mathbf{v}}}{b_{\mathbf{v}}^*} (\mathbf{a}_{\mathbf{v}} - \tilde{\boldsymbol{\nu}})(\mathbf{a}_{\mathbf{v}} - \tilde{\boldsymbol{\nu}})^\top,$$

$$[\mathbf{C}_{\mathbf{v}}^*]^{-1} = \mathbf{C}_{\mathbf{v}}^{-1} + \tilde{\mathbf{Q}}_{\mathbf{v}} + \tilde{\mathbf{Q}}_{\mathbf{a}_{\mathbf{v}}}, \quad \bar{\lambda}_{\mathbf{v}} = \frac{1}{n} \sum_{j=1}^n \lambda_{\mathbf{v}_j}, \quad \tilde{\boldsymbol{\nu}} = \frac{1}{n \bar{\lambda}_{\mathbf{v}}} \sum_{j=1}^n \lambda_{\mathbf{v}_j} \boldsymbol{\nu}_j.$$

Just as with the prior distributions, the full conditional distributions of $\lambda_{\mathbf{v}_i}$ and $\lambda_{\mathbf{v}_j}$ (conditional on the set \mathcal{F} comprising everything else) are also gamma distributions (see Appendix A):

- Full conditional: $(\lambda_{\mathbf{v}_i} | \mathcal{F}) \sim \mathcal{G} \left(\frac{D_{\mathbf{v}} + \kappa_{\mathbf{v}}}{2}, \frac{1}{2} [\kappa_{\mathbf{v}} + (\mathbf{v}_i - \boldsymbol{\mu}_{\mathbf{v}})^\top \mathbf{\Lambda}_{\mathbf{v}} (\mathbf{v}_i - \boldsymbol{\mu}_{\mathbf{v}})] \right)$, where $D_{\mathbf{v}}$ is the dimension of the matrix $\mathbf{C}_{\mathbf{v}}$.
- Full conditional: $(\lambda_{\mathbf{v}_j} | \mathcal{F}) \sim \mathcal{G} \left(\frac{D_{\mathbf{v}} + \kappa_{\mathbf{v}}}{2}, \frac{1}{2} [\kappa_{\mathbf{v}} + (\boldsymbol{\nu}_j - \boldsymbol{\mu}_{\mathbf{v}})^\top \mathbf{\Lambda}_{\mathbf{v}} (\boldsymbol{\nu}_j - \boldsymbol{\mu}_{\mathbf{v}})] \right)$, where $D_{\mathbf{v}}$ is the dimension of the matrix $\mathbf{C}_{\mathbf{v}}$.

Finally, we have the full conditional distributions of $\kappa_{\mathbf{v}}$ and $\kappa_{\mathbf{v}}$. In this specific case, it is not possible to identify any known distribution (see Appendix A). Therefore,

- Full conditional:

$$p(\kappa_{\mathbf{v}} | \mathcal{F}) \propto \exp \left\{ -\frac{\kappa_{\mathbf{v}}}{2} \left[\sum_{i=1}^m (\lambda_{\mathbf{v}_i} - \ln(\lambda_{\mathbf{v}_i})) - m \ln \left(\frac{\kappa_{\mathbf{v}}}{2} \right) \right] - m \ln \Gamma \left(\frac{\kappa_{\mathbf{v}}}{2} \right) \right\} p(\kappa_{\mathbf{v}});$$

- Full conditional:

$$p(\kappa_{\nu}|\mathcal{F}) \propto \exp \left\{ -\frac{\kappa_{\nu}}{2} \left[\sum_{j=1}^n (\lambda_{\nu,j} - \ln(\lambda_{\nu,j})) - n \ln \left(\frac{\kappa_{\nu}}{2} \right) \right] - n \ln \Gamma \left(\frac{\kappa_{\nu}}{2} \right) \right\} p(\kappa_{\nu}).$$

Gibbs sampler, more generally, includes modeling the degrees of freedom κ_{ν} and κ_{ν} . The complete conditionals do not have a known form but are easy to sample using other numerical methods, such as Metropolis-Hastings, which should be used to generate values for these conditionals. Nonetheless, the goal is simply to model heteroscedasticity and obtain a lower mean squared prediction error without the algorithm taking too long. Thus, at this stage, the degrees of freedom were fixed at a number between 3 and 10, which will be specified in the applications.

3. Applications

In this section, analyses for two datasets are presented: the Netflix and the MovieLens, due to their importance in the field. The standard values for the parameter p , the latent dimension, tested across all procedures are 5, 10, 20, and 30. Generally, a higher value of p would lead to better predictions at the cost of increased computational expense. Additionally, all models were trained for 200 iterations or epochs, even though convergence may have been reached earlier in some cases. The programming language used was Python¹. The maximum processing time for Netflix data was 8 hours for PMF (for all values of p) and approximately 36 hours for each value of p for BMF and BHMF. For MovieLens data, the processing time was 4.5 hours for PMF and around 16 hours for the total processing time of all tested situations for BMF and BHMF.

3.1. *Netflix*

The Netflix is a service for streaming videos of movies, TV series, etc., via the internet. The training data, collected from Netflix, consists of a collection of over 100 million ratings, ranging on a scale from 1 to 5, from a random sample of 480,189 anonymous users and 17,700 movies between October 1998 and December 2005. Additionally, approximately 1.5 million ratings were provided as supplemental material for model preparation. The goal is to predict new ratings for about 3 million users on the same set of movies based on this database. The provided training sample is subdivided into two categories: training data and test data (75% for training and 25% for testing) for evaluating model efficiency. The challenge lies in obtaining good results with this data, as it is extremely sparse, with approximately 99% of the entries consisting of zeros.

In PMF, we used the same approach as Salakhutdinov and Mnih [8]. In BMF, we set $\mathbf{a}_{\nu} = \mathbf{a}_{\nu} = \mathbf{0}$, $b_{\nu} = b_{\nu} = 2$, $d_{\nu} = d_{\nu} = p$, $\mathbf{C}_{\nu} = \mathbf{C}_{\nu}$ as the identity matrix, and $\alpha = 2$. In our BHMF, we set $\mathbf{a}_{\nu} = \mathbf{a}_{\nu} = \mathbf{0}$, $b_{\nu} = b_{\nu} = 2$, $d_{\nu} = d_{\nu} = p$, and $\mathbf{C}_{\nu} = \mathbf{C}_{\nu}$ set as the identity matrix. The hyperparameters for the prior distribution of α are defined as $a_{\alpha} = b_{\alpha} = 0.1$. Moreover, we fix $\kappa_{\nu} = \kappa_{\nu} = 5$.

Table 1 presents the root mean squared error (RMSE) results of the PMF, BMF, and our BHMF models for the Netflix data.

¹<https://www.python.org>

Table 1. The RMSE obtained from PMF, BMF and our BHMf.

p	RMSE					
	Train data			Test data		
	PMF	BMF	BHMf	PMF	BMF	BHMf
5	0.8644	0.8344	0.8348	0.8864	0.8630	0.8630
10	0.8446	0.8026	0.7999	0.8755	0.8442	0.8431
20	0.8335	0.7774	0.7797	0.8732	0.8329	0.8356
30	0.8251	0.7688	0.7466	0.8719	0.8297	0.8315

Table 1 shows that increasing p reduces the mean squared error for all three models (PMF, BMF, and BHMf), confirming that a higher p leads to a lower RMSE. The BMF and BHMf models produced very similar results, with differences only in the third decimal place of the RMSE. Both models are thus considered suitable for Netflix data.

For training data, BHMf performed similarly to BMF across all p values, with the best result at $p = 30$. This was expected because BHMf is more robust (due to its use of the Student t -distribution with 5 degrees of freedom, which can better handle data discrepancies) compared to BMF. However, BHMf’s theoretical variance is higher than BMF’s, so BMF was expected to have a slightly better RMSE for test data, though the differences may be minor. Thus, BHMf might offer better RMSE adjustments at the possible cost of increased computational time.

The BMF and BHMf models outperformed PMF, which requires extensive hyperparameter tuning and showed less improvement with higher p values. Consequently, BMF and BHMf at $p = 30$ provided the best results for both training and test data, while PMF at $p = 5$ had the worst performance, with RMSE not falling below 0.86.

Figure 1 displays the learning curves of the PMF, BMF and BHMf models on Netflix data with two p values: 10 and 30. The PMF model effectively reduced RMSE for both training and test data, though it took over 25 epochs to start decreasing

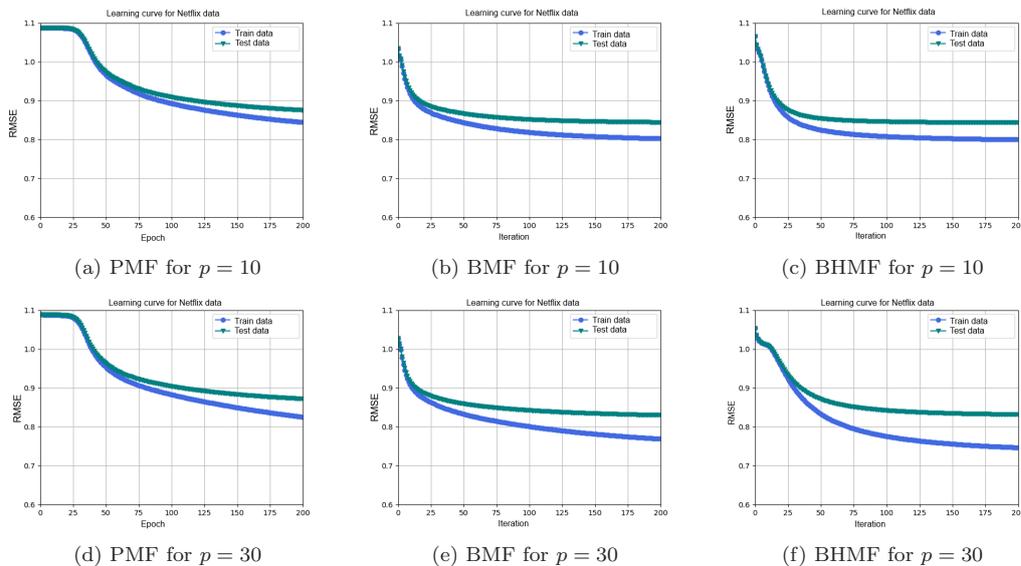


Figure 1. Learning curves of PMF, BMF, and BHMf for Netflix data at different values of p .

and remained above 0.80 even after 200 epochs. Extended training could potentially improve results. BMF also reduced RMSE for both datasets, with a rapid initial decrease in RMSE during the first 25 iterations. However, further reductions in test data were less significant, indicating that additional training might not be necessary. The BHMF model reduced RMSE similarly to BMF, with a high initial reduction but a slower decline as p increased. Further training had minimal impact on test data, suggesting limited benefit from extending iterations.

3.2. MovieLens

The GroupLens² research group collected and provided MovieLens datasets to aid film recommendation research. The data, which include over 10 million ratings on a 1 to 5 scale for 10,681 movies by 71,567 users, are very sparse (about 98.7% of the training set entries are zeros). The same training-test split provided by MovieLens was used.

For PMF, the best results were achieved with a mini-batch size of 100,000, a learning rate of $\epsilon = 5$, a gradient moment of 0.90, and a regularization parameter of $L_2 = 0.005$. For BMF, we used the following settings: $\mathbf{a}_v = \mathbf{a}_\nu = \mathbf{0}$, $b_v = b_\nu = 2$, $d_v = d_\nu = p$, $\mathbf{C}_v = \mathbf{C}_\nu$ as the identity matrix, and $\alpha = 2$. In BHMF, we again set $\mathbf{a}_v = \mathbf{a}_\nu = \mathbf{0}$, $b_v = b_\nu = 2$, $d_v = d_\nu = p$, and $\mathbf{C}_v = \mathbf{C}_\nu$ as the identity matrix. The prior distribution hyperparameters for α were set as $a_\alpha = b_\alpha = 0.1$, and we fixed $\kappa_v = \kappa_\nu = 5$. Initial values for \mathbf{U} and \mathbf{V} were set to 0.1, multiplied by random samples from a uniform distribution in $(0, 1)$.

Table 2 indicates that increasing p reduces the root mean squared error (RMSE) for all models (PMF, BMF, BHMF). Among these, BHMF achieved the best results on MovieLens data, especially with $p = 30$, showing RMSE below 0.80 for training and well below 0.90 for test data. This suggests BHMF is more effective at capturing data variability compared to PMF and BMF. On the other hand, PMF did not reach RMSE values below 0.79 (train data) despite various hyperparameter tests, and increasing p did not justify the computational cost. BMF performed between PMF and BHMF, making it a good candidate for further method comparisons.

Table 2. The RMSE obtained from PMF, BMF and our BHMF.

p	RMSE					
	Train data			Test data		
	PMF	BMF	BHMF	PMF	BMF	BHMF
5	0.8221	0.8045	0.7757	0.9137	0.8839	0.8680
10	0.8112	0.7986	0.7413	0.9068	0.8796	0.8550
20	0.7999	0.7945	0.7201	0.9016	0.8785	0.8508
30	0.7922	0.7916	0.7087	0.8976	0.8743	0.8494

Figure 2 shows the learning curve of the PMF, BMF and BHMF models on MovieLens data with different p values. For PMF, it illustrates that increasing p generally reduces RMSE for both training and test data, with a consistent decline after the 25th epoch and a high initial drop in RMSE, unlike Netflix data. RMSE continues to decrease even after 200 epochs, suggesting potential for further improvement with more iterations. Similar to PMF, BMF reduces RMSE across all p values, with a rapid decline in the first 25 iterations, followed by a slower decrease. RMSE continues to drop after 200 iterations, indicating that additional iterations might improve

²<https://grouplens.org>

results. The BHMF model also shows a reduction in RMSE for all p values, with a significant drop in the first 25 iterations. However, as p increases, the rate of RMSE reduction slows. Even after 200 iterations, RMSE continues to decrease, albeit more gradually, suggesting that extending training might be unnecessary beyond this point. This pattern is consistent across both MovieLens and Netflix data.

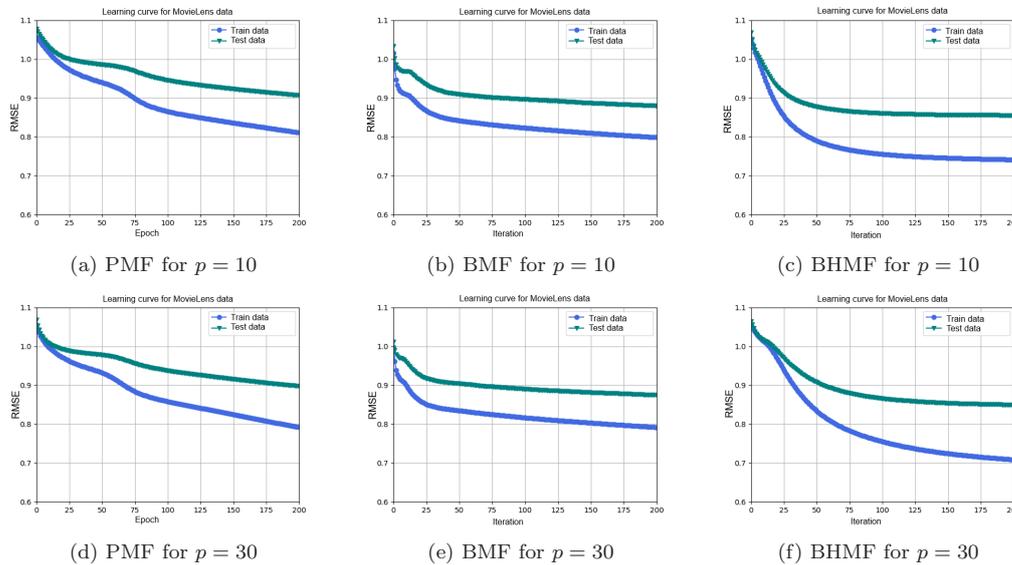


Figure 2. Learning curves of PMF, BMF, and BHMF for MovieLens data at different values of p .

4. Concluding remarks

This paper developed a method for estimating user ratings of movies with millions of observations. It used latent factor techniques, collaborative filtering, and a Bayesian approach in recommendation systems. The method was applied to real data from Netflix and MovieLens to test its effectiveness. Model performance was evaluated based on RMSE for both training and test data. The data were divided into 75% for training and 25% for testing for Netflix, and MovieLens data were split into two samples with exactly 10 test ratings per user.

The Bayesian heteroscedastic matrix factorization (BHMF) model extended the Bayesian matrix factorization (BMF) to account for individual user (or item) variation. Probabilistic matrix factorization (PMF) and BMF generally performed worse in this study, with BHMF sometimes yielding similar or slightly worse results compared to BMF, especially for $p = 20$ and $p = 30$ on Netflix data. For test data, PMF did not achieve RMSE values below 0.87 for Netflix or 0.89 for MovieLens in any configuration, unlike BMF and BHMF. A significant limitation was the high computational cost, as increasing p increases the computation time due to the larger data volume.

The proposed methodology demonstrated high accuracy in estimating user ratings for both Netflix and MovieLens, despite the computational cost. Future work could explore improvements to BHMF, use numerical methods for parameter sampling, investigate alternative distributions, apply different priors, and test other datasets and iteration counts.

References

- [1] Aggarwal, C., 2016. Recommender Systems: The Textbook, Springer. doi: 10.1007/978-3-319-29659-3
- [2] Canny, J., 2002. Collaborative filtering with privacy via factor analysis. SIGIR '02 Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval, 238–245. doi: 10.1145/564376.564419
- [3] Devooght, R., Kourtellis, N., Mantrach, A. (2015). Dynamic matrix factorization with priors on unknown values. KDD '15 Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 189–198. doi: 10.1145/2783258.2783346
- [4] Gemulla, R., Nijkamp, E., Haas, P. J., Sismanis, Y. 2011. Large-scale matrix factorization with distributed stochastic gradient descent. KDD '11 Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, 69–77. doi: 10.1145/2020408.2020426
- [5] Huffel, S., Vandewalle, J., Haegemans, A., 1987. An efficient and reliable algorithm for computing the singular subspace of a matrix, associated with its smallest singular values. Journal of Computational and Applied Mathematics, 19 (3), 313–330. doi: 10.1016/0377-0427(87)90201-9
- [6] Kabbur, S., Ning, X., Karypis, G. 2013. FISM: Factored item similarity models for top-N recommender systems. KDD '13 Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining, 659–667. doi: 10.1145/2487575.2487589
- [7] Salakhutdinov, R. and Mnih, A., 2008. Bayesian probabilistic matrix factorization using Markov chain Monte Carlo. ICML '08 Proceedings of the 25th International Conference on Machine Learning, 25, 880–887. doi: 10.1145/1390156.1390267
- [8] Salakhutdinov, R. and Mnih, A., 2008. Probabilistic matrix factorization. Advances in Neural Information Processing Systems 20, Canada, 2008. <http://papers.nips.cc/paper/3208-probabilistic-matrix-factorization.pdf>
- [9] Shen, H., Huang, J. Z., 2008. Sparse principal component analysis via regularized low rank matrix approximation. Journal of Multivariate Analysis, 99 (6), 1015–1034. doi: 10.1016/j.jmva.2007.06.007
- [10] Shen, J., Cheng, Z., Yang, M., Han, B., Li, S., 2019. Style-oriented personalized landmark recommendation. IEEE Transactions on Industrial Electronics, 66 (12), 9768–9776. doi: 10.1109/TIE.2019.2910043
- [11] Xian, Z., Li, Q., Li, G., Li, L. 2017. New collaborative filtering algorithms based on SVD++ and differential privacy. Mathematical Problems in Engineering, 2017 (1), 1–14. doi: 10.1155/2017/1975719.

Appendix A. The BHMF model and its full conditional distributions

The BHMF model, with $i = 1, 2, \dots, m$ and $j = 1, 2, \dots, n$, is given by:

$$\begin{aligned}
 (r_{ij} | \mathbf{v}_i, \boldsymbol{\nu}_j, \alpha) &\sim \mathcal{N}(\mathbf{v}_i^T \boldsymbol{\nu}_j, \alpha^{-1}), \\
 (\mathbf{v}_i | \boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_i}) &\sim \mathcal{N}(\boldsymbol{\mu}_v, [\lambda_{v_i} \boldsymbol{\Lambda}_v]^{-1}), \quad (\lambda_{v_i} | \kappa_v) \sim \mathcal{G}\left(\frac{\kappa_v}{2}, \frac{\kappa_v}{2}\right), \quad \kappa_v \sim \mathcal{G}(a_{\kappa_v}, b_{\kappa_v}), \\
 (\boldsymbol{\nu}_j | \boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu, \lambda_{\nu_j}) &\sim \mathcal{N}(\boldsymbol{\mu}_\nu, [\lambda_{\nu_j} \boldsymbol{\Lambda}_\nu]^{-1}), \quad (\lambda_{\nu_j} | \kappa_\nu) \sim \mathcal{G}\left(\frac{\kappa_\nu}{2}, \frac{\kappa_\nu}{2}\right), \quad \kappa_\nu \sim \mathcal{G}(a_{\kappa_\nu}, b_{\kappa_\nu}), \\
 (\boldsymbol{\mu}_v | \boldsymbol{\Lambda}_v) &\sim \mathcal{N}(\mathbf{a}_v, [b_v \boldsymbol{\Lambda}_v]^{-1}), \quad \boldsymbol{\Lambda}_v \sim \mathcal{W}(\mathbf{C}_v, d_v), \\
 (\boldsymbol{\mu}_\nu | \boldsymbol{\Lambda}_\nu) &\sim \mathcal{N}(\mathbf{a}_\nu, [b_\nu \boldsymbol{\Lambda}_\nu]^{-1}), \quad \boldsymbol{\Lambda}_\nu \sim \mathcal{W}(\mathbf{C}_\nu, d_\nu), \quad \text{and} \\
 \alpha &\sim \mathcal{G}(a_\alpha, b_\alpha).
 \end{aligned}$$

It follows that

$$\begin{aligned}
p(\mathcal{R}|U, V, \alpha) &= \prod_{i=1}^m \prod_{j=1}^n [\phi_1(r_{ij}|\mathbf{v}_i^T \boldsymbol{\nu}_j, \alpha^{-1})]^{\mathbb{I}_{ij}} \\
p(U|\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \boldsymbol{\lambda}_v) &= \prod_{i=1}^m \phi_p(\mathbf{v}_i|\boldsymbol{\mu}_v, [\lambda_{v_i} \boldsymbol{\Lambda}_v]^{-1}) \\
p(\boldsymbol{\lambda}_v|\kappa_v) &= \prod_{i=1}^m p(\lambda_{v_i}|\kappa_v) \\
p(V|\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \boldsymbol{\lambda}_v) &= \prod_{j=1}^n \phi_p(\boldsymbol{\nu}_j|\boldsymbol{\mu}_v, [\lambda_{\nu_j} \boldsymbol{\Lambda}_v]^{-1}) \\
p(\boldsymbol{\lambda}_v|\kappa_v) &= \prod_{j=1}^n p(\lambda_{\nu_j}|\kappa_v).
\end{aligned}$$

Moreover,

$$\begin{aligned}
p(\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v) &= p(\boldsymbol{\mu}_v|\boldsymbol{\Lambda}_v)p(\boldsymbol{\Lambda}_v) = \phi_p(\boldsymbol{\mu}_v|\mathbf{a}_v, [b_v \boldsymbol{\Lambda}_v]^{-1})\omega(\boldsymbol{\Lambda}_v|\mathbf{C}_v, d_v) \\
p(\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v) &= p(\boldsymbol{\mu}_v|\boldsymbol{\Lambda}_v)p(\boldsymbol{\Lambda}_v) = \phi_p(\boldsymbol{\mu}_v|\mathbf{a}_v, [b_v \boldsymbol{\Lambda}_v]^{-1})\omega(\boldsymbol{\Lambda}_v|\mathbf{C}_v, d_v).
\end{aligned}$$

From now on, consider $N = \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij}$ to be the effective sample size (counting only the observed values).

A.1. Full conditional distribution of α

$$\begin{aligned}
p(\alpha|\mathcal{F}) &\propto p(\mathcal{R}|U, V, \alpha)p(\alpha|a_\alpha, b_\alpha) \\
&\propto \alpha^{\frac{N}{2} + a_\alpha - 1} \exp \left\{ -\frac{\alpha}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^T \boldsymbol{\nu}_j)^2 - b_\alpha \alpha \right\} \\
&\propto \alpha^{(\frac{N}{2} + a_\alpha) - 1} \exp \left\{ -\alpha \left[\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^T \boldsymbol{\nu}_j)^2 + b_\alpha \right] \right\},
\end{aligned}$$

with $\mathcal{F} = (U, V, \mathcal{R})$. Thus,

$$(\alpha|\mathcal{F}) \sim \mathcal{G} \left(\frac{N}{2} + a_\alpha, \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^T \boldsymbol{\nu}_j)^2 + b_\alpha \right).$$

A.2. Full conditional distribution of \mathbf{v}_i

$$\begin{aligned}
p(\mathbf{v}_i|\mathcal{F}) &\propto p(\mathcal{R}|U, V, \alpha)p(\mathbf{v}_i|\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_i}) \\
&\propto \prod_{j=1}^n [\phi_1(r_{ij}|\mathbf{v}_i^T \boldsymbol{\nu}_j, \alpha^{-1})]^{\mathbb{I}_{ij}} p(\mathbf{v}_i|\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_i}) \\
&\propto \exp \left\{ -\frac{\alpha}{2} \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^T \boldsymbol{\nu}_j)^2 - \frac{1}{2} (\mathbf{v}_i - \boldsymbol{\mu}_v)^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] (\mathbf{v}_i - \boldsymbol{\mu}_v) \right\},
\end{aligned}$$

with $\mathcal{F} = (\mathbf{U}, \mathbf{V}, \mathcal{R}, \boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_i})$. Now, we have that

$$\begin{aligned}
SQ_{v_i} &= \alpha \sum_{j=1}^n \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^T \boldsymbol{\nu}_j)^2 + (\mathbf{v}_i - \boldsymbol{\mu}_v)^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] (\mathbf{v}_i - \boldsymbol{\mu}_v) \\
&= \alpha \sum_{j=1}^n \mathbb{I}_{ij} r_{ij}^2 - 2\alpha \sum_{j=1}^n \mathbb{I}_{ij} r_{ij} \mathbf{v}_i^T \boldsymbol{\nu}_j + \alpha \sum_{j=1}^n \mathbb{I}_{ij} (\mathbf{v}_i^T \boldsymbol{\nu}_j)^2 \\
&\quad + \mathbf{v}_i^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] \mathbf{v}_i - 2\mathbf{v}_i [\lambda_{v_i} \boldsymbol{\Lambda}_v] \boldsymbol{\mu}_v + \boldsymbol{\mu}_v^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] \boldsymbol{\mu}_v \\
&= c_1 + \mathbf{v}_i^T \left[\lambda_{v_i} \boldsymbol{\Lambda}_v + \alpha \sum_{j=1}^n \mathbb{I}_{ij} \boldsymbol{\nu}_j \boldsymbol{\nu}_j^T \right] \mathbf{v}_i \\
&\quad - 2\mathbf{v}_i^T \left[\lambda_{v_i} \boldsymbol{\Lambda}_v \boldsymbol{\mu}_v + \alpha \sum_{j=1}^n \mathbb{I}_{ij} r_{ij} \boldsymbol{\nu}_j \right] \\
&= c_1 + \mathbf{v}_i^T \boldsymbol{\Lambda}_{v_i}^* \mathbf{v}_i - 2\mathbf{v}_i^T \boldsymbol{\Lambda}_{v_i}^* \boldsymbol{\mu}_v^* \\
&= c_2 + \mathbf{v}_i^T \boldsymbol{\Lambda}_{v_i}^* \mathbf{v}_i - 2\mathbf{v}_i^T \boldsymbol{\Lambda}_{v_i}^* \boldsymbol{\mu}_{v_i}^* + \boldsymbol{\mu}_{v_i}^{*T} \boldsymbol{\Lambda}_{v_i}^* \boldsymbol{\mu}_{v_i}^* \\
&= c_2 + (\mathbf{v}_i - \boldsymbol{\mu}_{v_i}^*)^T \boldsymbol{\Lambda}_{v_i}^* (\mathbf{v}_i - \boldsymbol{\mu}_{v_i}^*),
\end{aligned}$$

where c_1 and c_2 are constants, whilst

$$\boldsymbol{\Lambda}_{v_i}^* = \lambda_{v_i} \boldsymbol{\Lambda}_v + \alpha \sum_{j=1}^n \mathbb{I}_{ij} \boldsymbol{\nu}_j \boldsymbol{\nu}_j^T \quad \text{and} \quad \boldsymbol{\mu}_{v_i}^* = [\boldsymbol{\Lambda}_{v_i}^*]^{-1} \left[\lambda_{v_i} \boldsymbol{\Lambda}_v \boldsymbol{\mu}_v + \alpha \sum_{j=1}^n \mathbb{I}_{ij} r_{ij} \boldsymbol{\nu}_j \right].$$

Moreover, note that

$$\begin{aligned}
\sum_{j=1}^n \mathbb{I}_{ij} (\mathbf{v}_i^T \boldsymbol{\nu}_j)^2 &= \sum_{j=1}^n \mathbb{I}_{ij} (\mathbf{v}_i^T \boldsymbol{\nu}_j) (\mathbf{v}_i^T \boldsymbol{\nu}_j) = \sum_{j=1}^n \mathbb{I}_{ij} (\mathbf{v}_i^T \boldsymbol{\nu}_j) (\mathbf{v}_i^T \boldsymbol{\nu}_j)^T \\
&= \sum_{j=1}^n \mathbb{I}_{ij} \mathbf{v}_i^T \boldsymbol{\nu}_j \boldsymbol{\nu}_j^T \mathbf{v}_i = \mathbf{v}_i^T \left[\sum_{j=1}^n \mathbb{I}_{ij} \boldsymbol{\nu}_j \boldsymbol{\nu}_j^T \right] \mathbf{v}_i.
\end{aligned}$$

It follows that $(\mathbf{v}_i | \mathcal{F}) \sim \mathcal{N}(\boldsymbol{\mu}_{v_i}^*, [\boldsymbol{\Lambda}_{v_i}^*]^{-1})$.

A.3. Full conditional distribution of $\boldsymbol{\nu}_j$

$$\begin{aligned}
p(\boldsymbol{\nu}_j | \mathcal{F}) &\propto p(\mathcal{R} | \mathbf{U}, \mathbf{V}, \alpha) p(\boldsymbol{\nu}_j | \boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_j}) \\
&\propto \prod_{i=1}^m [\phi_1(r_{ij} | \mathbf{v}_i^T \boldsymbol{\nu}_j, \alpha^{-1})]^{\mathbb{I}_{ij}} p(\boldsymbol{\nu}_j | \boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_j}) \\
&\propto \exp \left\{ -\frac{\alpha}{2} \sum_{i=1}^m \mathbb{I}_{ij} (r_{ij} - \mathbf{v}_i^T \boldsymbol{\nu}_j)^2 - \frac{1}{2} (\boldsymbol{\nu}_j - \boldsymbol{\mu}_v)^T [\lambda_{v_j} \boldsymbol{\Lambda}_v] (\boldsymbol{\nu}_j - \boldsymbol{\mu}_v) \right\},
\end{aligned}$$

with $\mathcal{F} = (\mathbf{U}, \mathbf{V}, \mathcal{R}, \boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \lambda_{v_j})$. By developing the expression above in the same manner as the full conditional distribution of \mathbf{v}_i , we obtain

$$\boldsymbol{\Lambda}_{v_j}^* = \lambda_{v_j} \boldsymbol{\Lambda}_v + \alpha \sum_{i=1}^m \mathbb{I}_{ij} \mathbf{v}_i \mathbf{v}_i^T \quad \text{and} \quad \boldsymbol{\mu}_{v_j}^* = [\boldsymbol{\Lambda}_{v_j}^*]^{-1} \left[\lambda_{v_j} \boldsymbol{\Lambda}_v \boldsymbol{\mu}_v + \alpha \sum_{i=1}^m \mathbb{I}_{ij} r_{ij} \mathbf{v}_i \right],$$

such that $(\boldsymbol{\nu}_j | \mathcal{F}) \sim \mathcal{N}(\boldsymbol{\mu}_{v_j}^*, [\boldsymbol{\Lambda}_{v_j}^*]^{-1})$.

A.4. Full conditional distribution of $(\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v)$

$$\begin{aligned}
p(\boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v | \mathcal{F}) &\propto p(\mathbf{U} | \boldsymbol{\mu}_v, \boldsymbol{\Lambda}_v, \boldsymbol{\lambda}_v) p(\boldsymbol{\mu}_v | \boldsymbol{\Lambda}_v) p(\boldsymbol{\Lambda}_v) \\
&\propto \prod_{i=1}^m \left[|\boldsymbol{\Lambda}_v|^{1/2} \exp \left\{ -\frac{1}{2} (\mathbf{v}_i - \boldsymbol{\mu}_v)^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] (\mathbf{v}_i - \boldsymbol{\mu}_v) \right\} \right] \\
&\times |\boldsymbol{\Lambda}_v|^{1/2} \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu}_v - \mathbf{a}_v)^T (b_v \boldsymbol{\Lambda}_v) (\boldsymbol{\mu}_v - \mathbf{a}_v) \right\} \\
&\times |\boldsymbol{\Lambda}_v|^{(d_v - D_v - 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{C}_v^{-1} \boldsymbol{\Lambda}_v) \right\} \\
&= |\boldsymbol{\Lambda}_v|^{m/2 + 1/2 + (d_v - D_v - 1)/2} \\
&\times \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{v}_i - \boldsymbol{\mu}_v)^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] (\mathbf{v}_i - \boldsymbol{\mu}_v) \right\} \\
&\times \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu}_v - \mathbf{a}_v)^T (b_v \boldsymbol{\Lambda}_v) (\boldsymbol{\mu}_v - \mathbf{a}_v) - \frac{1}{2} \text{tr}(\mathbf{C}_v^{-1} \boldsymbol{\Lambda}_v) \right\},
\end{aligned}$$

with $\mathcal{F} = (\mathbf{U}, \boldsymbol{\lambda}_v)$. Next, we split a few sums of squares and then combine them to determine the conditional distribution.

Part I

Let $\bar{\lambda}_v = \frac{1}{m} \sum_{i=1}^m \lambda_{v_i}$ and $\tilde{\mathbf{v}} = \frac{\sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i}{\sum_{i=1}^m \lambda_{v_i}} = \frac{1}{m \bar{\lambda}_v} \sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i$. Then,

$$\begin{aligned}
\text{SQ}_{\boldsymbol{\mu}, 1} &= \sum_{i=1}^m (\mathbf{v}_i - \boldsymbol{\mu}_v)^T [\lambda_{v_i} \boldsymbol{\Lambda}_v] (\mathbf{v}_i - \boldsymbol{\mu}_v) \\
&= \text{tr} \left(\sum_{i=1}^m \lambda_{v_i} (\mathbf{v}_i \mathbf{v}_i^T - \mathbf{v}_i \boldsymbol{\mu}_v^T - \boldsymbol{\mu}_v \mathbf{v}_i^T + \boldsymbol{\mu}_v \boldsymbol{\mu}_v^T) \boldsymbol{\Lambda}_v \right) \\
&= \text{tr} \left(\left[m \bar{\lambda}_v \boldsymbol{\mu}_v \boldsymbol{\mu}_v^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \boldsymbol{\mu}_v^T - m \bar{\lambda}_v \boldsymbol{\mu}_v \tilde{\mathbf{v}}^T + \sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T \right] \boldsymbol{\Lambda}_v \right) \\
&= \text{tr} \left(\left[m \bar{\lambda}_v \boldsymbol{\mu}_v \boldsymbol{\mu}_v^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \boldsymbol{\mu}_v^T - m \bar{\lambda}_v \boldsymbol{\mu}_v \tilde{\mathbf{v}}^T + m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T \right. \right. \\
&\quad \left. \left. - m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T + \sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T \right] \boldsymbol{\Lambda}_v \right) \\
&= \text{tr} \left(\left[m \bar{\lambda}_v (\boldsymbol{\mu}_v - \tilde{\mathbf{v}}) (\boldsymbol{\mu}_v - \tilde{\mathbf{v}})^T + \sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T \right] \boldsymbol{\Lambda}_v \right) \\
&= \text{tr} (m \bar{\lambda}_v (\boldsymbol{\mu}_v - \tilde{\mathbf{v}}) (\boldsymbol{\mu}_v - \tilde{\mathbf{v}})^T \boldsymbol{\Lambda}_v) + \text{tr} \left(\left[\sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T \right] \boldsymbol{\Lambda}_v \right) \\
&= (\boldsymbol{\mu}_v - \tilde{\mathbf{v}})^T [m \bar{\lambda}_v \boldsymbol{\Lambda}_v] (\boldsymbol{\mu}_v - \tilde{\mathbf{v}}) + \text{tr}(\tilde{\mathbf{S}}_v \boldsymbol{\Lambda}_v),
\end{aligned}$$

where

$$\begin{aligned}
\tilde{\mathbf{S}}_v &= \sum_{i=1}^m \lambda_{v_i} (\mathbf{v}_i - \tilde{\mathbf{v}}) (\mathbf{v}_i - \tilde{\mathbf{v}})^T \\
&= \sum_{i=1}^m (\lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T - \lambda_{v_i} \mathbf{v}_i \tilde{\mathbf{v}}^T - \lambda_{v_i} \tilde{\mathbf{v}} \mathbf{v}_i^T + \lambda_{v_i} \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T) \\
&= \sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T + m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T \\
&= \sum_{i=1}^m \lambda_{v_i} \mathbf{v}_i \mathbf{v}_i^T - m \bar{\lambda}_v \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T.
\end{aligned}$$

Part II

$$\begin{aligned}
\text{SQ}_{\mu,2} &= (\boldsymbol{\mu}_v - \tilde{\mathbf{v}})^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) (\boldsymbol{\mu}_v - \tilde{\mathbf{v}}) + (\boldsymbol{\mu}_v - \mathbf{a}_v)^T (b_v \boldsymbol{\Lambda}_v) (\boldsymbol{\mu}_v - \mathbf{a}_v) \\
&= \boldsymbol{\mu}_v^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \boldsymbol{\mu}_v - 2\boldsymbol{\mu}_v^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} + \tilde{\mathbf{v}}^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} \\
&\quad + \boldsymbol{\mu}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v - 2\boldsymbol{\mu}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v + \mathbf{a}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v \\
&= \boldsymbol{\mu}_v^T ((m\bar{\lambda}_v + b_v) \boldsymbol{\Lambda}_v) \boldsymbol{\mu}_v - 2\boldsymbol{\mu}_v^T \boldsymbol{\Lambda}_v (m\bar{\lambda}_v \tilde{\mathbf{v}} + b_v \mathbf{a}_v) \\
&\quad + \tilde{\mathbf{v}}^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} + \mathbf{a}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v \\
&= \boldsymbol{\mu}_v^T (b_v^* \boldsymbol{\Lambda}_v) \boldsymbol{\mu}_v - 2\boldsymbol{\mu}_v^T (b_v^* \boldsymbol{\Lambda}_v) \mathbf{a}_v^* + \mathbf{a}_v^{*T} (b_v^* \boldsymbol{\Lambda}_v) \mathbf{a}_v^* \\
&\quad + \tilde{\mathbf{v}}^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} + \mathbf{a}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v - \mathbf{a}_v^{*T} (b_v^* \boldsymbol{\Lambda}_v) \mathbf{a}_v^* \\
&= (\boldsymbol{\mu}_v - \mathbf{a}_v^*)^T (b_v^* \boldsymbol{\Lambda}_v) (\boldsymbol{\mu}_v - \mathbf{a}_v^*) + \tilde{\mathbf{v}}^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} \\
&\quad + \mathbf{a}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v - \mathbf{a}_v^{*T} (b_v^* \boldsymbol{\Lambda}_v) \mathbf{a}_v^*,
\end{aligned}$$

where

$$b_v^* = b_v + m\bar{\lambda}_v \quad \text{and} \quad \mathbf{a}_v^* = \frac{1}{b_v^*} (m\bar{\lambda}_v \tilde{\mathbf{v}} + b_v \mathbf{a}_v).$$

Thus, $(\boldsymbol{\mu}_v | \boldsymbol{\Lambda}_v, \mathcal{F}) \sim \mathcal{N}(\mathbf{a}_v^*, [b_v^* \boldsymbol{\Lambda}_v]^{-1})$.

Part III

$$\begin{aligned}
\text{SQ}_{\Lambda} &= \tilde{\mathbf{v}}^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} + \mathbf{a}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v - \mathbf{a}_v^{*T} (b_v^* \boldsymbol{\Lambda}_v) \mathbf{a}_v^* \\
&= \tilde{\mathbf{v}}^T (m\bar{\lambda}_v \boldsymbol{\Lambda}_v) \tilde{\mathbf{v}} + \mathbf{a}_v^T (b_v \boldsymbol{\Lambda}_v) \mathbf{a}_v - \frac{(m\bar{\lambda}_v \tilde{\mathbf{v}} + b_v \mathbf{a}_v)^T}{b_v^*} (b_v^* \boldsymbol{\Lambda}_v) \frac{(m\bar{\lambda}_v \tilde{\mathbf{v}} + b_v \mathbf{a}_v)^T}{b_v^*} \\
&= \frac{1}{b_v^*} [\tilde{\mathbf{v}}^T \boldsymbol{\Lambda}_v \tilde{\mathbf{v}} m\bar{\lambda}_v b_v^* + \mathbf{a}_v^T \boldsymbol{\Lambda}_v \mathbf{a}_v b_v b_v^* - (m\bar{\lambda}_v \tilde{\mathbf{v}} + b_v \mathbf{a}_v)^T \boldsymbol{\Lambda}_v (m\bar{\lambda}_v \tilde{\mathbf{v}} + b_v \mathbf{a}_v)] \\
&= \frac{1}{b_v^*} [\tilde{\mathbf{v}}^T \boldsymbol{\Lambda}_v \tilde{\mathbf{v}} (m^2 \bar{\lambda}_v^2 + m\bar{\lambda}_v b_v) + \mathbf{a}_v^T \boldsymbol{\Lambda}_v \mathbf{a}_v (m\bar{\lambda}_v b_v + b_v^2) \\
&\quad - m^2 \bar{\lambda}_v^2 \tilde{\mathbf{v}}^T \boldsymbol{\Lambda}_v \tilde{\mathbf{v}} - 2m\bar{\lambda}_v b_v \mathbf{a}_v^T \boldsymbol{\Lambda}_v \tilde{\mathbf{v}} - b_v^2 \mathbf{a}_v^T \boldsymbol{\Lambda}_v \mathbf{a}_v] \\
&= \frac{m\bar{\lambda}_v b_v}{b_v^*} [\mathbf{a}_v^T \boldsymbol{\Lambda}_v \mathbf{a}_v - 2\mathbf{a}_v^T \boldsymbol{\Lambda}_v \tilde{\mathbf{v}} + \tilde{\mathbf{v}}^T \boldsymbol{\Lambda}_v \tilde{\mathbf{v}}] \\
&= \frac{m\bar{\lambda}_v b_v}{b_v^*} (\mathbf{a}_v - \tilde{\mathbf{v}})^T \boldsymbol{\Lambda}_v (\mathbf{a}_v - \tilde{\mathbf{v}}) \\
&= \text{tr} \left(\frac{m\bar{\lambda}_v b_v}{b_v^*} (\mathbf{a}_v - \tilde{\mathbf{v}}) (\mathbf{a}_v - \tilde{\mathbf{v}})^T \boldsymbol{\Lambda}_v \right) = \text{tr} \left(\tilde{\mathbf{S}}_a \boldsymbol{\Lambda}_v \right),
\end{aligned}$$

where

$$\tilde{\mathbf{S}}_a = \frac{m\bar{\lambda}_v b_v}{b_v^*} (\mathbf{a}_v - \tilde{\mathbf{v}}) (\mathbf{a}_v - \tilde{\mathbf{v}})^T.$$

Part IV

Now, we have that

$$\begin{aligned} p(\mathbf{\Lambda}_v | \mathcal{F}) &\propto |\mathbf{\Lambda}_v|^{(d_v+m-D_v-1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{C}_v^{-1} \mathbf{\Lambda}_v) - \frac{1}{2} \text{tr}(\tilde{\mathbf{S}}_v \mathbf{\Lambda}_v) - \frac{1}{2} \text{tr}(\tilde{\mathbf{S}}_a \mathbf{\Lambda}_v) \right\} \\ &= |\mathbf{\Lambda}_v|^{(d_v^*-D_v-1)/2} \exp \left\{ -\frac{1}{2} \text{tr}([\mathbf{C}_v^*]^{-1} \mathbf{\Lambda}_v) \right\}, \end{aligned}$$

where

$$d_v^* = d_v + m \quad \text{and} \quad [\mathbf{C}_v^*]^{-1} = \mathbf{C}_v^{-1} + \tilde{\mathbf{S}}_v + \tilde{\mathbf{S}}_a.$$

Thus, $(\mathbf{\Lambda}_v | \mathcal{F}) \sim \mathcal{W}(\mathbf{C}_v^*, d_v^*)$.

A.5. Full conditional distribution of $(\boldsymbol{\mu}_v, \mathbf{\Lambda}_v)$

Based on the calculations from the previous full conditional distribution, we have that

$$\begin{aligned} (\boldsymbol{\mu}_v | \mathbf{\Lambda}_v, \mathbf{V}, \boldsymbol{\lambda}_v) &\sim \mathcal{N}(\mathbf{a}_v^*, [b_v^* \mathbf{\Lambda}_v]^{-1}); \\ (\mathbf{\Lambda}_v | \mathbf{V}, \boldsymbol{\lambda}_v) &\sim \mathcal{W}(\mathbf{C}_v^*, d_v^*), \end{aligned}$$

where

$$\begin{aligned} b_v^* &= b_v + n \bar{\lambda}_v, \quad d_v^* = d_v + n, \quad \mathbf{a}_v^* = \frac{1}{b_v^*} (n \bar{\lambda}_v \tilde{\boldsymbol{\nu}} + b_v \mathbf{a}_v), \\ \tilde{\mathbf{Q}}_v &= \sum_{j=1}^n \lambda_{v_j} (\mathbf{v}_j - \tilde{\boldsymbol{\nu}})(\mathbf{v}_j - \tilde{\boldsymbol{\nu}})^T, \\ \tilde{\mathbf{Q}}_a &= \frac{n \bar{\lambda}_v b_v}{b_v^*} (\mathbf{a}_v - \tilde{\boldsymbol{\nu}})(\mathbf{a}_v - \tilde{\boldsymbol{\nu}})^T, \\ [\mathbf{C}_v^*]^{-1} &= \mathbf{C}_v^{-1} + \tilde{\mathbf{Q}}_v + \tilde{\mathbf{Q}}_a. \end{aligned}$$

A.6. Full conditional distribution of λ_{v_i}

$$\begin{aligned} p(\lambda_{v_i} | \mathcal{F}) &\propto p(\mathbf{v}_i | \boldsymbol{\mu}_v, \mathbf{\Lambda}_v, \lambda_{v_i}) p(\lambda_{v_i} | \kappa_v) \\ &\propto \lambda_{v_i}^{(D_v + \kappa_v)/2 - 1} \exp \left\{ -\frac{\lambda_{v_i}}{2} [\kappa_v + (\mathbf{v}_i - \boldsymbol{\mu}_v)^T \mathbf{\Lambda}_v (\mathbf{v}_i - \boldsymbol{\mu}_v)] \right\}, \end{aligned}$$

with $\mathcal{F} = (\mathbf{U}, \boldsymbol{\mu}_v, \mathbf{\Lambda}_v, \kappa_v)$. Thus,

$$(\lambda_{v_i} | \mathcal{F}) \sim \mathcal{G} \left(\frac{D_v + \kappa_v}{2}, \frac{1}{2} [\kappa_v + (\mathbf{v}_i - \boldsymbol{\mu}_v)^T \mathbf{\Lambda}_v (\mathbf{v}_i - \boldsymbol{\mu}_v)] \right).$$

A.7. Full conditional distribution of λ_{ν_j}

$$\begin{aligned} p(\lambda_{\nu_j}|\mathcal{F}) &\propto p(\boldsymbol{\nu}_j|\boldsymbol{\mu}_{\nu}, \boldsymbol{\Lambda}_{\nu}, \lambda_{\nu_j})p(\lambda_{\nu_j}|\kappa_{\nu}) \\ &\propto \lambda_{\nu_j}^{(D_{\nu}+\kappa_{\nu})/2-1} \exp\left\{-\frac{\lambda_{\nu_j}}{2} [\kappa_{\nu} + (\boldsymbol{\nu}_j - \boldsymbol{\mu}_{\nu})^T \boldsymbol{\Lambda}_{\nu}(\boldsymbol{\nu}_j - \boldsymbol{\mu}_{\nu})]\right\}. \end{aligned}$$

with $\mathcal{F} = (\mathbf{V}, \boldsymbol{\mu}_{\nu}, \boldsymbol{\Lambda}_{\nu}, \kappa_{\nu})$. Thus,

$$(\lambda_{\nu_j}|\mathcal{F}) \sim \mathcal{G}\left(\frac{D_{\nu} + \kappa_{\nu}}{2}, \frac{1}{2} [\kappa_{\nu} + (\boldsymbol{\nu}_j - \boldsymbol{\mu}_{\nu})^T \boldsymbol{\Lambda}_{\nu}(\boldsymbol{\nu}_j - \boldsymbol{\mu}_{\nu})]\right).$$

A.8. Full conditional distribution of κ_{ν}

$$\begin{aligned} p(\kappa_{\nu}|\boldsymbol{\lambda}_{\nu}) &\propto p(\boldsymbol{\lambda}_{\nu}|\kappa_{\nu})p(\kappa_{\nu}) = \left[\prod_{i=1}^m p(\lambda_{\nu_i}|\kappa_{\nu})\right] p(\kappa_{\nu}) \\ &\propto \frac{[\kappa_{\nu}/2]^{m\kappa_{\nu}/2}}{[\Gamma(\kappa_{\nu}/2)]^m} \left[\prod_{i=1}^m \lambda_{\nu_i}^{\kappa_{\nu}/2}\right] \exp\left\{-\frac{\kappa_{\nu}}{2} m\bar{\lambda}_{\nu}\right\} p(\kappa_{\nu}) \\ &= \exp\left\{-\frac{\kappa_{\nu}}{2} \left[\sum_{i=1}^m (\lambda_{\nu_i} - \ln(\lambda_{\nu_i})) - m \ln\left(\frac{\kappa_{\nu}}{2}\right)\right] - m \ln \Gamma\left(\frac{\kappa_{\nu}}{2}\right)\right\} p(\kappa_{\nu}) \end{aligned}$$

A.9. Full conditional distribution of κ_{ν}

$$\begin{aligned} p(\kappa_{\nu}|\boldsymbol{\lambda}_{\nu}) &\propto p(\boldsymbol{\lambda}_{\nu}|\kappa_{\nu})p(\kappa_{\nu}) = \left[\prod_{j=1}^n p(\lambda_{\nu_j}|\kappa_{\nu})\right] p(\kappa_{\nu}) \\ &\propto \frac{[\kappa_{\nu}/2]^{n\kappa_{\nu}/2}}{[\Gamma(\kappa_{\nu}/2)]^n} \left[\prod_{j=1}^n \lambda_{\nu_j}^{\kappa_{\nu}/2}\right] \exp\left\{-\frac{\kappa_{\nu}}{2} n\bar{\lambda}_{\nu}\right\} p(\kappa_{\nu}) \\ &= \exp\left\{-\frac{\kappa_{\nu}}{2} \left[\sum_{j=1}^n (\lambda_{\nu_j} - \ln(\lambda_{\nu_j})) - n \ln\left(\frac{\kappa_{\nu}}{2}\right)\right] - n \ln \Gamma\left(\frac{\kappa_{\nu}}{2}\right)\right\} p(\kappa_{\nu}). \end{aligned}$$