

---

# Variational Inference for Extreme Spatio-Temporal Matrix Completion

---

**Charul Paliwal**  
IIITD, India  
charuli@iiitd.ac.in

**Pravesh Biyani**  
IIITD, India  
praveshb@iiitd.ac.in

## Abstract

Missing data is a common problem in real-world sensor data collection. The performance of various approaches to impute data degrade rapidly in the extreme scenarios of low data sampling and noisy sampling, a case present in many real-world problems in the field of traffic sensing and environment monitoring, etc. However, jointly exploiting the spatiotemporal and periodic structure, which is generally not captured by classical matrix completion approaches, can improve the imputation performance of sensor data in such real-world conditions. We present a Bayesian approach toward spatiotemporal matrix completion wherein we estimate the underlying temporarily varying subspace using a Variational Bayesian technique. We jointly couple the low-rank matrix completion with the state space autoregressive framework along with a penalty function on the slowly varying subspace to model the temporal and periodic evolution in the data. We also propose a robust version of the above formulation, which improves the performance of imputation in the presence of outliers. Results demonstrate that the proposed method outperforms the recent state-of-the-art methods for real-world traffic and air pollution data. We demonstrate that fusing the subspace evolution over days can improve the imputation performance with even 15% of the data sampling.

## 1 Introduction

Data acquired from both static and moving sensors contains missing data due to sensor malfunction, irregularity in sensor measurements, etc. Additionally, moving sensors scheme uses relatively fewer sensors resulting in high data "gaps" in both spatial and temporal dimensions [1]. This motivates the problem of extreme matrix completion, where the percentage of data sampled may be as low as 15%. Thereby, a natural question to ask is: how to fill the high percentage missing entries within a reasonable error range? Can we leverage additional periodic information in the matrix completion framework to estimate the high percentage of missing entries. Also, in addition to the missing entries, sensor measurements can be contaminated with outliers emerging from the sensor malfunctioning, communication errors, or impulse noises. The occurrence of outliers in the measurements can further degrade the performance of data imputation. However, unlike the missing entries, the location and the value of outliers are unknown, which makes the problem more challenging. Therefore, the other question to ask is: how to estimate the missing data while detecting the noisy outliers? The answer to the above questions lies in exploiting the underlying structure available in the data. For instance, both the air pollution [2] and the traffic [3] data exhibit joint spatial and temporal correlation as well as periodicity in daily patterns, thereby generating redundancy that can be potentially exploited by performing an intelligent spatio-temporal extrapolation.

In this paper, we propose a Variational Bayesian Filtering with Subspace Information (VBFISI) to estimate the data in the case of extremely sparsely sampled data. We observe that even with a fraction of sampled observed data, we can estimate the remaining measurements with reasonable

accuracy. This work exploits the spatiotemporal and periodic pattern in the measurements and perform extreme matrix completion with reasonable accuracy. First, we enforce a low-rank structure to the spatiotemporal data. Second, we enforce the state-space model on the temporal embeddings to capture the temporal evolution in the data. Third, we enforce the subspace estimate of the matrix to be close to previously learned subspace distribution using the Mahalanobis distance. Exploiting the prior subspace over days results in a considerable reduction of the number of measurements required to estimate the matrix [4, 5] thereby boosting the performance in case of a low sampling rate. We use the Variational Bayes approach to update the parameters in an iterative fashion. In our work, the subspace distribution is chosen to allow automatic relevance determination, and unlike the matrix or tensor completion methods, the algorithm parameters such as rank, noise powers need not be specified or tuned. We compare the performance of the proposed algorithm with state-of-the-art algorithms on many real-world traffic and air quality datasets. The result shows that modeling the subspace evolution leads to improvement in performance even when a small percentage of random measurements are available for the purpose of imputation. A likely impact of our method is that cities with a low sensing budget can perform random drive-by sampling of the urban environment, and the suggested matrix completion framework can provide a reasonably accurate imputation leading to better decision making.

## 1.1 Related work

The traditional matrix completion framework is not applicable for time series data imputation, as it does not take into account the ordering among the temporal embeddings [6]. Autoregressive model can model temporal embeddings to capture the temporal evolution in the time series data [6–8]. However, these models fail to capture the prior subspace information that can be exploited to capture periodicity. Exploiting the 3-way pattern in the data using tensor completion-based frameworks can improve the imputation performance [9–11]. However, the temporal evolution and subspace evolution is not modeled in the traditional tensor completion frameworks. We propose to incorporate the periodicity in the matrix completion framework while modeling the temporal and subspace evolution over days. Variational Bayesian approaches are proposed for matrix/tensor completion and robust PCA [9, 10, 12–14]. A state-space model to capture the temporal evolution is also proposed in [7, 8]. However, these approaches do not explicitly model the evolution of the subspace to capture the periodicity in the data. Also, the proposed matrix and tensor completion methods are not evaluated for extreme matrix completion and extreme matrix completion in the presence of outliers.

## 2 Variational Bayesian Filtering with Subspace Information (VBFSI)

Let  $\mathbf{X} \in \mathbb{R}^{n \times t}$  be the data matrix for a day  $d$ , where  $n$  and  $t$  denote the number of spatial locations and time stamps respectively. The low rankness in the data can be imposed as

$$\mathcal{L}_1 = \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{P}_\Omega \odot (\mathbf{X} - \mathbf{UV}^T)\|_F \quad (1)$$

where  $\mathbf{U} \in \mathbb{R}^{n \times r}$  and  $\mathbf{V} \in \mathbb{R}^{t \times r}$  and  $r = \text{rank}(\mathbf{X}) \ll \min(n, t)$ . For sampling percentage of  $p$ , let  $\Omega$  denote the sampled data containing  $p \times n \times t$  samples.  $\mathbf{P}_\Omega$  is the indicator matrix where  $P_{ij} = 1 \forall (i, j) \in \Omega$ . Temporal evolution is captured by regularizing the columns of  $\mathbf{V}$  to follow an autoregressive model and regularize  $\mathbf{U}$  to capture the periodicity over days.

$$\mathcal{R}(\mathbf{V}) = \sum_{i=1}^t \|\mathbf{v}_i - \mathbf{F}\mathbf{v}_{i-1}\|, \quad \mathcal{R}(\mathbf{U}) = \eta \sum_{i=1}^n (\mathbf{u}_i - \mathbf{u}_i^{d-1})^T (\boldsymbol{\Xi}_i^{\mathbf{U}^{d-1}})^{-1} (\mathbf{u}_i - \mathbf{u}_i^{d-1}) \quad (2)$$

here  $\mathcal{R}(\mathbf{U})$  corresponds to the Mahalanobis distance between each row vector of  $\mathbf{U}^{d-1}$  (subspace estimate of previous day) and  $\mathbf{U}$  (current subspace estimate for a given day  $d$ ).  $\boldsymbol{\Xi}^{\mathbf{U}^{d-1}}$  denotes the covariance matrix of  $\mathbf{U}^{d-1}$ . Here  $\eta$  controls the effect of prior subspace ( $\mathbf{U}^{d-1}$ ,  $\boldsymbol{\Xi}_i^{d-1}$ ) in the estimation of  $\mathbf{U}$ .  $\mathbf{u}_i$  denotes the  $i$ -th row of  $\mathbf{U}$ . The entries of  $\mathbf{X}_\Omega$  can be defined as:

$$\mathbf{p}(\mathbf{X}_\Omega | \mathbf{U}, \mathbf{V}, \beta) = \prod_{(i,j) \in \Omega} \mathcal{N}(X_{ij} | \mathbf{u}_i \cdot \mathbf{v}_j^T, \beta^{-1}), \quad \mathbf{p}(\beta) = \beta^{-1} \quad (3)$$

here  $\beta$  is the noise precision. The prior on the noise is assumed to be non-informative Jeffrey’s prior.

Columns of  $\mathbf{U}$  are enforced with a sparsity profile using precision  $\gamma_i$  to automate the rank.

$$p(\mathbf{U} | \gamma) = \prod_{i=1}^r \mathcal{N}(\mathbf{u}_i | \mathbf{u}_i^{d-1}, \gamma_i^{-1} \mathbf{I}_n) \quad (4)$$

When  $\gamma_i$  are driven to a large value then the column mean will be  $\mathbf{u}_i^{d-1}$ , and we prune these columns and in turn, reducing the rank thereby modelling the low rank in the Bayesian framework, referred as Automatic Rank Determination [12]. Further, the autoregressive regularization can be modelled as

$$p(\mathbf{V} | \mathbf{F}) = \mathcal{N}(\mathbf{v}_1; \boldsymbol{\mu}_1, \boldsymbol{\Lambda}_1) \prod_{\tau=2}^t \mathcal{N}(\mathbf{v}_\tau | \mathbf{F}\mathbf{v}_{\tau-1}, \mathbf{I}_r) \quad (5)$$

$\mathbf{F}$  is assigned multivariate Gaussian priors with column-specific precisions  $v$ . Precision variables  $\gamma$  and  $v$  are selected to have non-informative Jeffrey's priors

$$p(\mathbf{F} | v) = \prod_{i=1}^r \mathcal{N}(\mathbf{f}_i | 0, v_i^{-1} \mathbf{I}_r), \quad p(\gamma_i) = \frac{1}{\gamma_i}, \quad p(v_i) = \frac{1}{v_i} \quad (6)$$

The overall joint distribution for spatio-temporal matrix completion can be expressed as

$$p(\mathbf{X}_\Omega, \mathbf{U}, \mathbf{V}, \mathbf{F}, \beta, \gamma, v) = p(\mathbf{X}_\Omega | \mathbf{U}, \mathbf{V}, \beta) p(\mathbf{U} | \gamma) p(\mathbf{V} | \mathbf{F}) p(\mathbf{F} | v) p(\beta) p(v) p(\gamma) \quad (7)$$

The Full Bayesian graphical model for spatio-temporal matrix completion is shown in Fig 1. The posterior distribution  $p(\mathcal{H} | \mathbf{x}_\Omega)$  is intractable as  $p(\mathbf{x}_\Omega)$  needs to be marginalized over all parameters ( $\mathcal{H}$ ). Therefore, we utilize the mean-field approximation and provide the Variational Bayesian update steps for all the parameters in the Appendix.

### 3 Robust VBFSI (RVBFSI)

RVBFSI estimates the missing data while detecting the noisy outliers. For robust matrix completion, we model  $\mathbf{X} = \mathbf{U}\mathbf{V}^T + \mathbf{E} + \mathbf{N}$ , where  $\mathbf{E}$  denotes the sparse outlier matrix and  $\mathbf{N}$  is the dense error matrix. The low rankness defined in the Eq. 1 is modified as

$$\mathcal{L}_1 = \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{P}_\Omega \odot (\mathbf{X} - \mathbf{U}\mathbf{V}^T - \mathbf{E})\|_F \quad (8)$$

The regularization on  $\mathbf{U}$  and  $\mathbf{V}$  follows the Eq. 2. Each entry of sparse outlier matrix  $E_{ij}$  is assigned a precision  $\alpha_{ij}$ , where the  $\alpha_{ij}$  have the non informative prior. This works similar to the ARD where instead of a column of the matrix, each entry of the matrix is assigned with precision. Whenever  $\alpha_{i,j}$  is driven to a large value  $E_{i,j} \rightarrow 0$  thereby enforcing sparsity.

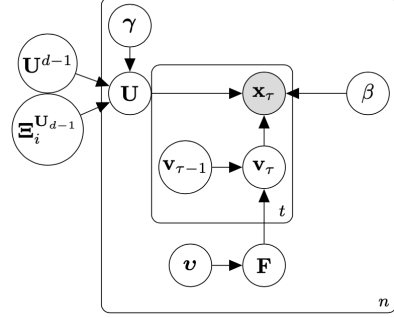
$$p(\mathbf{E} | \boldsymbol{\alpha}) = \prod_{i | ((i,j) \in \Omega)} \prod_{j | ((i,j) \in \Omega)} \mathcal{N}(E_{ij} | 0, \alpha_{ij}^{-1}), \quad p(\alpha_{ij}) = \alpha_{ij}^{-1} \quad (9)$$

## 4 Experimentation

**Dataset and Parameters:** We use traffic speed and air quality (PM 2.5) data for performance evaluation. Data (DT): Delhi traffic speed data [8]. Data (GT): Guangzhou urban traffic speed data [15]. Data (PT): Pems traffic speed data [16]. Data (CA): China Air Quality data [17]. Rank determination is automatic. We only tune the parameter  $\eta$ . We grid search the best  $\eta$  for different sampling percentages. Then, we fit exponential model for  $\eta$  vs.  $p$ . The initial subspace  $\mathbf{U}^0$  is calculated using the eight days average, then we run the algorithm in an online fashion for the next 30 days.

**Baseline Algorithms:** We compare VBFSI with state of the art matrix and tensor imputation methods. **VBFSI** [8], **VMC** [18], **TRMF** [6], **BTMF** [19], **BCPF** [10] and **TRLRF** [20]. We compare RVBFSI with the following Robust imputation methods. **RVBSF** [8], **RegL<sub>1</sub>** [21], **BRTF** [14].

Figure 1: VBFSI model



	p %	VBFSI	VBSF	VMC	BCPF	TRLRF	TRMF	BTMF
Data:DT	5%	<b>0.156 / 4.387</b>	0.782 / 22.03	0.998 / 28.18	0.164 / 4.6	0.901 / 25.45	0.183 / 5.146	0.157 / 4.394
	15%	<b>0.135 / 3.796</b>	0.162 / 4.552	0.97 / 27.39	0.147 / 4.137	0.682 / 19.25	0.151 / 4.24	<u>0.137 / 3.836</u>
	50%	<u>0.117 / 3.289</u>	0.119 / 3.344	0.131 / 3.687	<b>0.115 / 3.23</b>	0.171 / 4.785	0.121 / 3.409	0.119 / 3.342
	75%	<u>0.11 / 3.086</u>	0.11 / 3.099	0.117 / 3.28	<b>0.109 / 3.076</b>	0.13 / 3.642	0.117 / 3.262	0.115 / 3.224
Data:PT	5%	<b>0.144 / 8.608</b>	1 / 60.08	0.998 / 59.974	0.175 / 10.494	0.94 / 56.47	0.161 / 9.571	<u>0.151 / 9.084</u>
	15%	<b>0.111 / 6.625</b>	0.179 / 10.7	0.974 / 58.52	0.147 / 8.836	0.807 / 48.49	0.139 / 8.264	<u>0.118 / 7.06</u>
	50%	<b>0.084 / 5.056</b>	0.097 / 5.79	<u>0.087 / 5.213</u>	0.091 / 5.431	0.168 / 10.07	0.093 / 5.511	0.1 / 6.02
	75%	<u>0.081 / 4.841</u>	0.081 / 4.854	<b>0.069 / 4.135</b>	0.081 / 4.848	0.097 / 5.833	0.083 / 4.951	0.097 / 5.81
Data:GT	5%	0.159 / 6.384	1 / 40.31	0.993 / 40.03	<u>0.158 / 6.346</u>	0.863 / 34.76	0.184 / 6.614	<b>0.131 / 5.244</b>
	15%	<u>0.121 / 4.854</u>	0.148 / 5.91	0.382 / 15.33	0.138 / 5.541	0.492 / 19.8	0.162 / 5.845	<b>0.11 / 4.43</b>
	50%	<b>0.088 / 3.547</b>	0.112 / 4.501	<u>0.09 / 3.616</u>	0.097 / 3.902	0.112 / 4.503	0.128 / 4.624	0.095 / 3.801
	75%	<b>0.079 / 3.189</b>	0.1 / 4.027	<u>0.081 / 3.247</u>	0.088 / 3.515	0.091 / 3.652	0.12 / 4.303	0.092 / 3.712
Data:CA	5%	0.439 / 32.44	1 / 76.562	0.998 / 76.464	0.435 / 32.672	0.978 / 74.949	0.434 / 33.915	<b>0.414 / 31.431</b>
	15%	0.35 / 25.964	0.396 / 29.762	0.986 / 75.582	<b>0.341 / 25.471</b>	0.936 / 72.046	0.369 / 28.735	<u>0.344 / 25.94</u>
	50%	<u>0.222 / 16.466</u>	0.23 / 17.172	<b>0.213 / 15.892</b>	0.237 / 17.646	0.731 / 57.059	0.235 / 18.274	0.248 / 18.483
	75%	<u>0.198 / 14.648</u>	0.2 / 14.675	<b>0.171 / 12.636</b>	0.209 / 15.603	0.472 / 36.816	0.197 / 15.249	0.223 / 16.556

Table 1: MRE/RMSE scores for data imputation. The best two results are bold and underlined.

p %	o=5%				o=10%			
	10%	25%	50%	75%	10%	25%	50%	75%
RVBFSI	<b>0.167 / 4.672</b>	<b>0.14 / 3.91</b>	<b>0.126 / 3.544</b>	<b>0.119 / 3.337</b>	<b>0.17 / 4.78</b>	<b>0.14 / 3.925</b>	<b>0.128 / 3.573</b>	<b>0.118 / 3.314</b>
RVBSF	0.196 / 5.527	<u>0.154 / 4.313</u>	<u>0.132 / 3.696</u>	0.124 / 3.485	0.227 / 6.403	<u>0.16 / 4.492</u>	<u>0.132 / 3.728</u>	<u>0.124 / 3.487</u>
VBFSI	0.188 / 5.277	0.177 / 4.972	0.17 / 4.778	0.158 / 4.427	0.208 / 5.86	0.194 / 5.443	0.184 / 5.173	0.175 / 4.899
VBSF	0.305 / 8.583	0.201 / 5.656	0.177 / 4.98	0.166 / 4.65	0.391 / 11.01	0.237 / 6.65	0.198 / 5.561	0.184 / 5.16
VMC	0.995 / 28.05	0.798 / 22.4	0.368 / 10.35	0.306 / 8.605	0.996 / 28.08	0.937 / 26.41	0.562 / 15.8	0.44 / 12.38
BCPF	0.193 / 5.422	0.164 / 4.597	0.146 / 4.101	0.14 / 3.927	<u>0.205 / 5.763</u>	0.179 / 5.017	0.155 / 4.363	0.148 / 4.146
TRLRF	0.973 / 27.45	0.956 / 26.95	0.934 / 26.34	1.144 / 32.24	0.977 / 27.55	0.969 / 27.33	0.997 / 28.1	1.447 / 40.71
TRMF	0.382 / 10.76	0.419 / 11.78	0.265 / 7.426	0.217 / 6.05	0.565 / 15.91	0.56 / 15.75	0.338 / 9.492	0.291 / 8.139
BTMF	0.226 / 5.828	0.218 / 5.604	0.221 / 5.615	0.218 / 5.591	0.308 / 7.79	0.303 / 7.681	0.304 / 7.794	0.303 / 7.764
RegL <sub>1</sub>	0.521 / 14.66	0.489 / 13.78	0.192 / 5.429	0.132 / 3.718	0.699 / 19.69	0.498 / 14.01	0.242 / 6.825	0.155 / 4.387
BRTF	0.243 / 6.829	0.232 / 6.522	0.138 / 3.888	0.131 / 3.706	0.218 / 6.133	0.2 / 5.609	0.152 / 4.285	0.145 / 4.121

Table 2: MRE/RMSE for outlier corrupted data (DT), outlier percentage ( $o$ ) is 5% and 10%.

**Performance Comparison** The performance comparison of VBFSI with the current state of the art methods is shown in Table 1. VBFSI outperforms VBSF for all the datasets. The performance of VBSF is comparable to VBFSI for higher sampling. While for lower sampling, the performance of VBSF degrades. Incorporating even the noisy prior subspace information in the architecture can reduce the sampling complexity of the matrix by logarithmic factor [4]. Therefore, for low sampling VBFSI performs better than VBSF since we have incorporate the prior information using  $\eta$ . VMC experience a similar trend, where the performance is comparable for higher sampling and degrades for low sampling. VBFSI outperforms VMC for almost all the cases for traffic data (DT, PT, GT). However, for the air quality data (CA), the performance of VMC is better than VBFSI for a higher sampling percentage because VMC can capture the nonlinearity in the data. For low sampling percentage, VBFSI performance is comparable to BTMF in most of the cases. However, for higher sampling percentage, VBFSI outperforms BTMF. One of the disadvantage of BTMF and TRMF is that rank is not tuned automatically and uses gibbs sampling which is relatively slower than the Variational Bayesian approach for parameter estimation [22].

To compare the performance of VBFSI and RVBFSI in the case of outliers, we artificially add the outliers in the Data: DT. We randomly add 5% and 10% of the outliers in the total sampled data, i.e., the number of outliers is  $o \times p$  fraction of the overall data. The entries corrupted with outliers are uniformly distributed between  $[-\sigma, \sigma]$ , where  $\sigma$  is set as 100. Performance comparison of RVBFSI and VBFSI with other imputation methods are shown in Table 2. RVBFSI outperforms all imputation methods and robust imputation methods RVBSF, RegL<sub>1</sub> and BRTF significantly.

## 5 Conclusion

In this paper, we propose VBFSI for the imputation of spatiotemporal matrices that works even for extreme matrix completion. VBFSI simultaneously models the low rank, temporal evolution, and periodic evolution in one framework to capture the structure in the spatiotemporal data. We incorporate the prior subspace in our model to capture the periodic evolution in the data. We also propose a Robust VBFSI for missing data imputation in the presence of outliers. It is observed that RVBFSI performs better than the other imputation methods in the presence of the outliers.

## References

- [1] A. Anjomshoaa, F. Duarte, D. Rennings, T. J. Matarazzo, P. deSouza, and C. Ratti, “City scanner: Building and scheduling a mobile sensing platform for smart city services,” *IEEE Internet of Things Journal*, vol. 5, no. 6, pp. 4567–4579, 2018.
- [2] P. D. Sampson, A. A. Szpiro, L. Sheppard, J. Lindström, and J. D. Kaufman, “Pragmatic estimation of a spatio-temporal air quality model with irregular monitoring data,” *Atmospheric Environment*, vol. 45, no. 36, pp. 6593–6606, 2011.
- [3] W. Min and L. Wynter, “Real-time road traffic prediction with spatio-temporal correlations,” *Transportation Research Part C: Emerging Technologies*, vol. 19, no. 4, pp. 606–616, 2011.
- [4] A. Eftekhari, D. Yang, and M. B. Wakin, “Weighted matrix completion and recovery with prior subspace information,” *IEEE Transactions on Information Theory*, vol. 64, no. 6, pp. 4044–4071, 2018.
- [5] X. Zhang, W. Cui, and Y. Liu, “Matrix completion with prior subspace information via maximizing correlation,” *arXiv preprint arXiv:2001.01152*, 2020.
- [6] H.-F. Yu, N. Rao, and I. S. Dhillon, “Temporal regularized matrix factorization for high-dimensional time series prediction,” in *Advances in neural information processing systems*, 2016, pp. 847–855.
- [7] J. Luttinen, “Fast variational Bayesian linear state-space model,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2013, pp. 305–320.
- [8] C. Paliwal, U. Bhatt, P. Biyani, and K. Rajawat, “Traffic estimation and prediction via on-line variational bayesian subspace filtering,” *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–11, 2021.
- [9] M. T. Asif, N. Mitrovic, J. Dauwels, and P. Jaillet, “Matrix and tensor based methods for missing data estimation in large traffic networks,” *IEEE Trans. Intell. Transp. Syst.*, vol. 17, no. 7, pp. 1816–1825, 2016.
- [10] Q. Zhao, L. Zhang, and A. Cichocki, “Bayesian cp factorization of incomplete tensors with automatic rank determination,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1751–1763, 2015.
- [11] X. Chen, Z. He, and L. Sun, “A bayesian tensor decomposition approach for spatiotemporal traffic data imputation,” *Transportation research part C: emerging technologies*, vol. 98, pp. 73–84, 2019.
- [12] S. D. Babacan, M. Luessi, R. Molina, and A. K. Katsaggelos, “Sparse bayesian methods for low-rank matrix estimation,” *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 3964–3977, 2012.
- [13] L. Yang, J. Fang, H. Duan, H. Li, and B. Zeng, “Fast low-rank bayesian matrix completion with hierarchical gaussian prior models,” *IEEE Transactions on Signal Processing*, 2018.
- [14] Q. Zhao, G. Zhou, L. Zhang, A. Cichocki, and S.-I. Amari, “Bayesian robust tensor factorization for incomplete multiway data,” *IEEE transactions on neural networks and learning systems*, vol. 27, no. 4, pp. 736–748, 2015.
- [15] . Z. H. Xinyu Chen, Yixian Chen, “Urban traffic speed dataset of guangzhou, china [data set].zenodo.” 2018. [Online]. Available: <http://doi.org/10.5281/zenodo.1205229>
- [16] . X. C. Yixian Chen, “A large scale pems traffic speed dataset (version v1) [data set].zenodo.” 2020. [Online]. Available: <http://doi.org/10.5281/zenodo.3939793>
- [17] Y. Zheng, X. Yi, M. Li, R. Li, Z. Shan, E. Chang, and T. Li, “Forecasting fine-grained air quality based on big data,” in *Proceedings of the 21th SIGKDD conference on Knowledge Discovery and Data Mining*, August 2015. [Online]. Available: <http://research.microsoft.com/apps/pubs/?id=246398>
- [18] G. Ongie, R. Willett, R. D. Nowak, and L. Balzano, “Algebraic variety models for high-rank matrix completion,” in *International Conference on Machine Learning*, 2017, pp. 2691–2700.
- [19] X. Chen and L. Sun, “Bayesian temporal factorization for multidimensional time series prediction,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

- [20] L. Yuan, C. Li, D. Mandic, J. Cao, and Q. Zhao, “Tensor ring decomposition with rank minimization on latent space: An efficient approach for tensor completion,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 9151–9158.
- [21] Y. Zheng, G. Liu, S. Sugimoto, S. Yan, and M. Okutomi, “Practical low-rank matrix approximation under robust  $l_1$ -norm,” in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 1410–1417.
- [22] T. Salimans, D. Kingma, and M. Welling, “Markov chain monte carlo and variational inference: Bridging the gap,” in *International Conference on Machine Learning*, 2015, pp. 1218–1226.
- [23] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [24] M. E. Tipping, “Sparse bayesian learning and the relevance vector machine,” *Journal of machine learning research*, vol. 1, no. Jun, pp. 211–244, 2001.

## A Appendix

### A.1 Variational Bayesian Filtering with Subspace Information

#### A.1.1 Variational Bayesian Inference for Spatiotemporal Matrix Completion

We utilize the mean-field approximation, wherein the posterior distribution of parameters  $\theta := \{\mathbf{U}, \mathbf{V}, \mathbf{F}, \beta, \gamma, \mathbf{v}\}$  is factorized into a set of conditionally independent components. It is expressive as it captures the marginal density of the parameters. The main advantage of this assumption is that the optimization takes the form of coordinate ascent where the posterior distribution of each parameter can be found by taking expectation of all the other parameters in an iterative manner. The posterior distribution of parameters is factorized as:

$$p(\theta | \mathbf{x}_\Omega) = q_{\mathbf{U}}(\mathbf{U})q_{\mathbf{V}}(\mathbf{V})q_{\mathbf{F}}(\mathbf{F})q_{\mathbf{v}}(\mathbf{v})q_{\beta}(\beta)q_{\gamma}(\gamma). \quad (10)$$

The posterior distribution of all the parameters is determined by minimizing the Kullback–Leibler divergence of  $p(\theta|\mathbf{x}_\Omega)$  from  $q(\theta)$ , usually via an alternating minimization approach [23].

We use variational inference to estimate the posterior distribution of  $q_{\mathbf{U}}$ ,  $q_{\mathbf{V}}$ ,  $q_{\mathbf{F}}$ ,  $q_{\mathbf{v}}$ ,  $q_{\beta}$ , and  $q_{\gamma}$  for sampled data  $\mathbf{Z} = \mathbf{P}_\Omega(\mathbf{X})$ . The updates for posterior distribution of parameters are similar to the updates derived in [8, 12].

The posterior distribution for a row of  $\mathbf{U}$  is given by

$$q_{\mathbf{u}_i} = \mathcal{N}(\mathbf{u}_i | \boldsymbol{\mu}_i^{\mathbf{U}}, \boldsymbol{\Xi}_i^{\mathbf{U}}) \quad (11)$$

The updates of mean and covariance of  $\mathbf{U}$  are derived as

$$\begin{aligned} (\boldsymbol{\Xi}_i^{\mathbf{U}})^{-1} &= \hat{\gamma}_i \mathbf{I}_r + \hat{\beta} \sum_{\tau|(i,\tau) \in \Omega} (\boldsymbol{\mu}_\tau^{\mathbf{V}}(\boldsymbol{\mu}_\tau^{\mathbf{V}})^T + \boldsymbol{\Xi}_{\tau,\tau}^{\mathbf{V}}) \\ &\quad + \eta(\boldsymbol{\Xi}_i^{\mathbf{U}^{d-1}})^{-1} \end{aligned} \quad (12)$$

$$\boldsymbol{\mu}_i^{\mathbf{U}} = \boldsymbol{\Xi}_i^{\mathbf{U}}(\hat{\beta} \sum_{\tau|(i,\tau) \in \Omega} \boldsymbol{\mu}_\tau^{\mathbf{V}} Z_{i\tau} + \eta(\boldsymbol{\Xi}_i^{\mathbf{U}^{d-1}})^{-1} \boldsymbol{\mu}_i^{\mathbf{U}^{d-1}}) \quad (13)$$

The mean and covariance for the Posterior Distribution of  $\mathbf{V}$  are as follows:

$$q_{\mathbf{V}}(\mathbf{V}) = \mathcal{N}(\vec{\mathbf{V}} | \boldsymbol{\mu}^{\mathbf{V}}, \boldsymbol{\Xi}^{\mathbf{V}}) \quad (14)$$

$$\boldsymbol{\mu}^{\mathbf{V}} = \boldsymbol{\Xi}^{\mathbf{V}} \begin{bmatrix} \hat{\beta} \sum_{i|(i,1) \in \Omega} Z_{i,1} \boldsymbol{\mu}_i^{\mathbf{U}} + \boldsymbol{\Lambda}_1^{-1} \boldsymbol{\mu}_1 \\ \hat{\beta} \sum_{i|(i,2) \in \Omega} Z_{i,2} \boldsymbol{\mu}_i^{\mathbf{U}} \\ \vdots \\ \hat{\beta} \sum_{i|(i,t) \in \Omega} Z_{i,t} \boldsymbol{\mu}_i^{\mathbf{U}} \end{bmatrix} \quad (15)$$

$$\begin{aligned}
[\Xi^V]^{-1} &= \hat{\beta} \text{Diag}(\Xi_{(1)}^U, \dots, \Xi_{(t)}^U) + \\
&+ \begin{bmatrix} \Lambda_1^{-1} & -\hat{\mathbf{F}} & \dots & 0 \\ -\hat{\mathbf{F}} & \mathbf{I}_r + \Sigma^{\mathbf{F}} & -\hat{\mathbf{F}} & \dots \\ \vdots & \vdots & & \vdots \\ \dots & 0 & -\hat{\mathbf{F}} & \mathbf{I}_r \end{bmatrix}
\end{aligned} \tag{16}$$

The direct inversion of the dense matrix  $\Xi^V$  would be computationally costly. The block-tridiagonal structure  $(\Xi^V)^{-1}$  can be exploited to carry out the updates for  $\Xi^V$  in an efficient manner using LDL decomposition [7, 8].

The updates of the posterior distribution of  $\mathbf{F}$  are given by

$$q_{\mathbf{f}_i} = \mathcal{N}(\mathbf{f}_i \mid \boldsymbol{\mu}_i^{\mathbf{F}}, \Xi_i^{\mathbf{F}}) \tag{17}$$

$$\boldsymbol{\mu}_i^{\mathbf{F}} = [\Xi_i^{\mathbf{F}}(\boldsymbol{\mu}_\tau(\boldsymbol{\mu}_{\tau-1})^T + \Xi_{\tau, \tau-1}^{\mathbf{V}})]_i \tag{18a}$$

$$(\Xi_i^{\mathbf{F}})^{-1} = \text{Diag}(\hat{v}) + \sum_{\tau=1}^{t-1} (\boldsymbol{\mu}_\tau(\boldsymbol{\mu}_{\tau-1})^T + \Xi_{\tau, \tau-1}^{\mathbf{V}}) \tag{18b}$$

The posterior distribution for hyperparameters  $\{\beta, \gamma, \mathbf{v}\}$  are given by

$$q_\beta(\beta) = \text{Ga}(\beta; a^\beta, b^\beta) \tag{19a}$$

$$q_{\gamma_i}(\gamma_i) = \text{Ga}(\gamma_i; a_i^\gamma, b_i^\gamma) \tag{19b}$$

$$q_{v_i}(v_i) = \text{Ga}(v_i; a_i^v, b_i^v) \tag{19c}$$

where  $\text{Ga}(x, a, b)$  denotes the Gamma pdf with parameters  $a$  and  $b$ . The updates for  $\{\beta, \gamma, \mathbf{v}\}$  are given by

$$\hat{v}_i = \frac{r}{\sum_{k=1}^r ([\boldsymbol{\mu}_k^{\mathbf{F}}]_i^2 + [\Xi_k^{\mathbf{F}}]_{ii})} \tag{20a}$$

$$\hat{\gamma}_i = \frac{n+t}{\sum_{k=1}^n ([\boldsymbol{\mu}_k^{\mathbf{U}}]_i^2 + [\Sigma_k^{\mathbf{U}}]_{ii}) + \sum_{k=1}^t ([\boldsymbol{\mu}_k^{\mathbf{V}}]_i^2 + [\Sigma_k^{\mathbf{V}}]_{ii})} \tag{20b}$$

$$\hat{\beta} = \frac{pnt}{\|\mathbf{Z} - P_\Omega(\mathbf{U}\mathbf{V}^T)\|_F^2} \tag{20c}$$

We update the mean, covariance of  $\mathbf{U}$ ,  $\mathbf{V}$ ,  $\mathbf{F}$  and the hyperparameters  $\gamma, \mathbf{v}, \beta$  iteratively as shown in Algorithm 1

---

**Algorithm 1** VBSFI

---

**Input:**  $\Xi^{\mathbf{U}^{d-1}}, \boldsymbol{\mu}^{\mathbf{U}^{d-1}}, P_\Omega(\mathbf{X})$

**Initialization:**  $\gamma, \beta, \mathbf{v}, \Xi^{\mathbf{U}}, \boldsymbol{\mu}^{\mathbf{U}}, \Xi^{\mathbf{V}}, \boldsymbol{\mu}^{\mathbf{V}}, \Xi^{\mathbf{F}}, \boldsymbol{\mu}^{\mathbf{F}}, \mathbf{Z}$

- 1: **while**  $X_{conv} < 10^{-5}$  **do**
  - 2:    $\mathbf{X}_{old} = \hat{\mathbf{X}}$   
       Compute  $\mathbf{V}, \mathbf{F}, \mathbf{v}, \beta$  using (21, 22, 24, 26a, 26c)  
       Compute  $\mathbf{U}, \gamma, \beta$  using (18, 19, 26b, 26c)  
        $\hat{\mathbf{X}} = \boldsymbol{\mu}^{\mathbf{U}}(\boldsymbol{\mu}^{\mathbf{V}})^T$   
        $X_{conv} = \frac{\|\hat{\mathbf{X}} - \mathbf{X}_{old}\|_F}{\|\mathbf{X}_{old}\|_F}$
  - 3: **end while**
  - 4: **Output:**  $\hat{\mathbf{X}}$
-

## A.2 Robust Variational Bayesian Filtering with Subspace Information (RVBFSI)

### A.2.1 Bayesian Model

The conditional distribution of generating the entries of  $\mathbf{X}_\Omega$  can be defined as

$$p(\mathbf{X}_\Omega | \mathbf{U}, \mathbf{V}, \mathbf{E}, \beta) = \prod_{(i,j) \in \Omega} \mathcal{N}(X_{ij} | (\mathbf{u}_i \cdot \mathbf{v}_j^T + E_{ij}), \beta^{-1}) \quad (21)$$

Columns of  $\mathbf{U}$ ,  $\mathbf{V}$  and  $\mathbf{F}$  and precision variables  $\gamma, \beta, v$  follows same the prior distribution defined in (6-10).

Each entry of sparse outlier matrix  $E_{ij}$  is assigned a precision  $\alpha_{ij}$ .

$$p(\mathbf{E} | \alpha) = \prod_{i | ((i,j) \in \Omega)} \prod_{j | ((i,j) \in \Omega)} \mathcal{N}(E_{ij} | 0, \alpha_{ij}^{-1}) \quad (22)$$

where the  $\alpha_{ij}$  have the non informative prior

$$p(\alpha_{ij}) = \frac{1}{\alpha_{ij}} \quad (23)$$

This works similar to the ARD where instead of column of the matrix, each entry of the matrix is assigned with a precision. Whenever  $\alpha_{i,j}$  is driven to a large value, the  $E_{i,j} \rightarrow 0$  thereby enforcing sparsity. The overall joint distribution for Robust Spatio-Temporal Matrix Completion is expressed as

$$p(\mathbf{X}_\Omega, \mathbf{U}, \mathbf{V}, \mathbf{F}, \mathbf{E}, \beta, \gamma, v, \alpha) = p(\mathbf{X}_\Omega | \mathbf{U}, \mathbf{V}, \mathbf{E}, \beta) p(\mathbf{U} | \gamma) \\ \times p(\mathbf{V} | \mathbf{F}) p(\mathbf{F} | v) p(\mathbf{E} | \alpha) p(\beta) p(v) p(\gamma) \quad (24)$$

The full bayesian model for the Robust Spatio-Temporal Matrix Completion is depicted in Fig. 2.

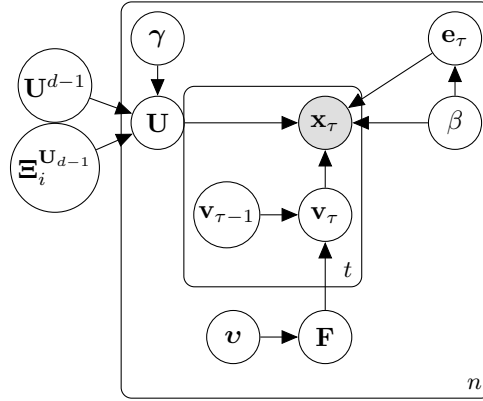


Figure 2: Robust Variational Bayesian Filtering with Subspace Information

### A.2.2 Variational Bayesian Inference

Approximate posterior distribution of parameters  $\theta_R := \{\mathbf{U}, \mathbf{V}, \mathbf{F}, \mathbf{E}, \beta, \gamma, v, \alpha\}$  are derived using Variational Inference.

We utilize the mean-field approximation, wherein the posterior distribution of parameters  $\theta_R$  is factorized as:

$$q_{\mathbf{U}}(\mathbf{U}) q_{\mathbf{V}}(\mathbf{V}) q_{\mathbf{F}}(\mathbf{F}) q_{\mathbf{E}}(\mathbf{E}) q_v(v) q_\beta(\beta) q_\gamma(\gamma) q_\alpha(\alpha).$$

The posterior distribution of  $q_{\mathbf{U}}$ ,  $q_{\mathbf{V}}$ ,  $q_{\mathbf{F}}$ ,  $q_v$ ,  $q_\beta$ , and  $q_\gamma$  takes the same form for  $\mathbf{Z} = \mathbf{P}_\Omega(\mathbf{X} - \mathbf{E})$  as shown in (17-26c). The posterior distribution for  $\mathbf{E}$  take the following form  $\forall (i, j) \in \Omega$ .

$$q(E_{ij}) = \mathcal{N}(E_{ij} | \mu_{ij}^{\mathbf{E}}, \Xi_{ij}^{\mathbf{E}}) \quad (25)$$



$$\Xi_{i,j}^{\mathbf{E}} = \frac{1}{\hat{\beta} + \alpha \hat{\alpha}_{i,j}} \quad (26)$$

$$\boldsymbol{\mu}_{i,j}^{\mathbf{E}} = \hat{\beta} \Xi_{i,j}^{\mathbf{E}} (X_{i,j} - \boldsymbol{\mu}_i^{\mathbf{A}} (\boldsymbol{\mu}_j^{\mathbf{B}})^T) \quad (27)$$

$$\hat{\alpha}_{i,j}^{new} = \frac{1 - \hat{\alpha}_{i,j}^{old} \Xi_{i,j}^{\mathbf{E}}}{(\boldsymbol{\mu}_{i,j}^{\mathbf{E}})^2} \quad (28)$$

$\hat{\alpha}_{i,j}^{new}$  is the fixed-point update for  $\alpha$ . This is used in the sparse bayesian learning that leads to much faster convergence and enhanced sparsity [12, 24]. For robust estimation of entries in the presence of outliers, we update the mean, covariance of  $\mathbf{U}$ ,  $\mathbf{V}$ ,  $\mathbf{F}$ ,  $\mathbf{E}$  and the hyperparameters  $\gamma$ ,  $\mathbf{v}$ ,  $\beta$ ,  $\alpha$  iteratively as shown in Algorithm 2.

---

**Algorithm 2** RVBSFI

---

**Input:**  $\Xi^{\mathbf{U}_{d-1}}, \boldsymbol{\mu}^{\mathbf{U}_{d-1}}, \mathbf{P}_{\Omega}(\mathbf{X})$

**Initialization:**  $\gamma, \beta, \mathbf{v}, \Xi^{\mathbf{U}}, \boldsymbol{\mu}^{\mathbf{U}}, \Xi^{\mathbf{V}}, \boldsymbol{\mu}^{\mathbf{V}}, \Xi^{\mathbf{F}}, \boldsymbol{\mu}^{\mathbf{F}}, \boldsymbol{\mu}^{\mathbf{E}}, \Xi^{\mathbf{E}}, \alpha, \mathbf{Z} = \mathbf{P}_{\Omega}(\mathbf{X} - \mathbf{E})$

- 1: **while**  $X_{conv} < 10^{-5}$  **do**
  - 2:    $\mathbf{X}_{old} = \hat{\mathbf{X}}$   
       Compute  $\mathbf{V}, \mathbf{F}, \mathbf{v}, \beta$  using (21, 22, 24, 26a, 26c)  
       Compute  $\mathbf{U}, \gamma, \mathbf{E}, \alpha, \beta$  (18, 19, 26b, 26c)  
       Compute  $\mathbf{E}, \alpha, \beta$  using (28-29, 26c)  
        $\hat{\mathbf{X}} = \boldsymbol{\mu}^{\mathbf{U}} (\boldsymbol{\mu}^{\mathbf{V}})^T$   
        $X_{conv} = \frac{\|\hat{\mathbf{X}} - \mathbf{X}_{old}\|_F}{\|\mathbf{X}_{old}\|_F}$
  - 3: **end while**
  - 4: **Output:**  $\hat{\mathbf{X}}$
-