

ROBUST STREAMING TENSOR FACTORISATION VIA ONLINE VARIATIONAL BAYESIAN INFERENCE

Avinandan Bose Spandan Senapati Ketan Rejawat

Abstract—Streaming Tensor Factorization is a powerful tool for processing high-volume and multi-way temporal data. It finds use in Internet Networks, recommender systems, image/video data analysis and Estimated Time of Arrival(ETA) predictions. This work puts forward the use of time-series and tensor completion algorithms. Our contribution is extending the previously existing Robust Streaming Tensor Factorization methods and incorporating first order auto-regressive time series in it.

I. INTRODUCTION

Tensors are extensively used in various application domains including recommender systems, computer vision etc. Strategies such as CP factorizations, have been proposed to obtain low-rank representations of tensors. Streaming tensors appear sequentially in the time domain. Incorporating temporal relationships in tensor data analysis can give significant advantages in traffic prediction, recommender systems etc. Various Tensor Decomposition methods such as the Tucker and the CP decomposition have been proposed in the past few decades. They differ in the way they solve the numerical optimization problem. However these methods assume fixed rank and fail to capture outliers in the data. However streaming tensor factorisations follow bayesian inference that leads to an Automatic Rank Determination(ARD) and also enables one to capture the sparse outliers in the data.

We consider first low-rank robust subspace filtering approach for online tensor imputation and prediction. Different from the existing tensor completion formulations, we consider low-rank matrices whose underlying subspace evolves according to a first order Markovian Chain. As slices of tensor data matrix arrive sequentially over time, the low rank components as well as the state-space model are learned in an online fashion using the Variational Bayes formalism. In this paper we discuss the traffic estimation problem and show that the Variational Bayesian approach incorporated with a first order auto regressive model can be used to learn the model in online setting which estimates the missing entries(i.e speed) when only a fraction of other entries is available. Robust version of the Variational Bayesian Subspace Filtering algorithm is proposed for outlier removal and data cleansing in a dynamic setting.

II. PROBABILISTIC MODEL

A. Details about Notations

Given data \mathcal{Y}_{Ω_T} in the form of a tensor of N dimensions, we wish to find its representation in its CP Factorisation form. $\mathcal{Y}_{\Omega_T} = \langle A^{(1)}, A^{(2)}, \dots, A^{(N)}; b_T \rangle + S_{\Omega_T}$, where $A^{(1)}, \dots, A^{(N)}$ are the factor matrices which evolve over time, b_t is temporal factor which follows the first order autoregressive model $p(b_t | J, b_{t-1}) = \mathcal{N}(b_t | Jb_{t-1}, I_r)$ and S_{Ω_T} accounts for the sparse outliers. The matrix B is of dimensions $R \times T$, (R is the rank of the low-rank CP-decomposition and determined on its own, i.e. not present) and the t^{th} row of B is b_t which is like a weighing factor in the CP decomposition for the factor matrices. We use Variational Bayesian Inference to compute the optimal values if the parameters.

B. Variational Inference

Variational Inference forms one of the means of Approximate Inference i.e. there are other forms of approximate inference like Markov Chain Monte Carlo Stimulation(MCMC) etc. Solving Optimisation Problems by Bayesian Inference can prove to be a handy alternative in many situations especially when convex optimisation methods can't be adopted for the same. In this we assume prior distributions over latent variables (variables that aren't observable but affect the model predictions inherently) and look for an iterative scheme to reach convergence of a function termed the Evidence Lower Bound(ELBO).

$$\log p(x) = \underbrace{\mathcal{L}_q}_{\text{ELBO}} + \underbrace{\text{KL}(q(z) || p(z|x))}_{\text{KL Divergence serves as a measure of dissimilarity}}.$$

The optimal changes to the distribution q is obtained as,

$$q_j^* = \mathbb{E}_{i \neq j} \left[\ln p(X, z_1, z_2, \dots, z_n) \right]$$

where z_1, z_2, \dots, z_n denote the latent variables in the model and q_z denotes the assumed prior distribution on the latent variable z . The conditional distributions are obtained from Bayes Theorem as,

$$p(z|x) = \frac{p(z, x)}{p(x)} = \frac{p(x|z)p(z)}{\int p(x|z)p(z)dz}$$

C. Likelihood

$$\begin{aligned} \text{ELBO} &= p(Y_{\Omega_T}, \tilde{D}_{t, \Omega_t} | A^{(1)}, A^{(2)}, A^{(3)} \dots A^{(N)}, B, S_{\Omega_T}, \beta) \\ &= \mathcal{N}(y_{\omega_N} | < \hat{a}_{i_1}^{(1)}, \hat{a}_{i_2}^{(2)} \dots \hat{a}_{i_N}^{(N)}, b_T > + S_{\omega_N}, \beta^{-1}) \end{aligned}$$

where $\tilde{D}_t = Y_t - S_t$ Here \tilde{D}_t is the cleaned up data at time t after removing the sparse outliers and ω_N denotes the subscript i_1, i_2, \dots, i_N . We could hence use ω_N and the subscript Notation interchangeably.

D. Prior Distributions

$A^{(1)}, A^{(2)}, \dots, A^{(N)}$ are the factor matrices in the CP decomposition of input data which evolve over time.

$$\prod_{i=1}^N p(A^{(i)} | \lambda) = \prod_{k=1}^N \prod_{i_k}^{I_k} \mathcal{N}(\hat{a}_{i_k}^{(k)} | 0, \Lambda^{-1}) \quad (1)$$

where $\Lambda = \text{Diag}(\lambda)$

$$p(\lambda) = \prod_{i=1}^r \Gamma(\lambda | c_0, d_0) \quad (2)$$

B is a matrix whose rows are the temporal components. Its rows follow a first order auto-regressive model. The following 3 equations describe this.

$$p(B | J) = N(b_i; \mu_1, \zeta_1) \prod_{t=2}^T \mathcal{N}(b_t; Jb_{t-1}, I_r) \quad (3)$$

$$p(J | \nu) = \prod_{i=1}^r \prod_{j=1}^r \mathcal{N}(J_{ij} | 0, \nu_i^{-1} I) \quad (4)$$

$$p(\nu) = \prod_{i=1}^r \frac{1}{\nu_i} \quad (5)$$

S_{Ω_T} is a tensor which represents the sparse outliers in the input data.

$$p(S_{\Omega_T} | \gamma) = \prod_{i \in \Omega_T} \mathcal{N}(S_{i_1 i_2 \dots i_N} | 0, \gamma_{i_1 i_2 \dots i_N}) \quad (6)$$

$$p(\gamma) = \prod_{i \in \Omega_T} \Gamma(\gamma_{i_1, i_2 \dots i_N} | a_0^\gamma, b_0^\gamma) \quad (7)$$

β is the precision factor in our likelihood function.

$$p(\beta) = \Gamma(\beta | a_\beta, b_\beta) \quad (8)$$

III. VARIATIONAL BAYESIAN UPDATES

A. Factor Matrices Updates

$$q(A^{(n)}) = \prod_{i_n=1}^{I_n} \mathcal{N}(\hat{a}_{i_n}^{(n)} | \bar{a}_{i_n}^{(n)}, V_{i_n}^{(n)}) \quad (9)$$

$$V_{i_n}^{(n)} = (\mathbb{E}_q[\beta] \sum_{t=i}^T \mathbb{E}_q[A_{i_n}^{(n)T} A_{i_n}^{(n)}]_{\Omega_t} + \mathbb{E}_q[\Lambda])^{-1} \quad (10)$$

$$\begin{aligned} \bar{a}_{i_n}^{(n)} &= V_{i_n}^{(n)} (\mathbb{E}_q[A_{i_n}^{(n)T}]_{\Omega_T} \text{vec}(y_{\Omega_T} - \mathbb{E}[S_{\Omega_T}])) + \\ &= \sum_{t=i}^{T-1} \mathbb{E}_q[A_{i_n}^{(n)T}]_{\Omega_T} \text{vec}(\tilde{D}_{t, \Omega_t}) \end{aligned}$$

$$\mathbb{E}_q[A_{i_n}^{(n)}]_{\Omega_T} = (\mathbb{E}_q[\bigodot_{j \neq n} A^{(j)}])_{I_{i_n}} \quad (11)$$

The matrix $A_{i_n}^{(n)}$ is $\prod_{j \neq n} I_j R$ and the Indicator function I_{i_n} samples the row $(i_1, \dots, i_{n-1}, i_{n+1}, \dots, i_{N+1})$ if the entry $(i_1, \dots, i_{n-1}, i_{n+1}, \dots, i_{N+1})$ is in Ω_t . The Matrix $A^{(N+1)}$ is the temporal component and is equal to B , and we will look at its update separately in Section 3.4 which is different from the updates of the non-temporal components.

B. Posterior Distribution of λ

$$q(\lambda) = \prod_{r=1}^R \Gamma(\lambda_r | c_M^r, d_M^r) \quad (12)$$

$$c_M^r = c_0 + 1 + \frac{\sum_{i=1}^N I_i}{2} \quad (13)$$

$$d_M^r = d_0 + 1 + \frac{1}{2} \sum_{i=1}^{N+1} \mathbb{E}_q[(a_r^n)^T a_r^n] \quad (14)$$

$$\mathbb{E}_q[\Lambda] = \text{Diag}(\frac{c_M^1}{d_M^1}, \frac{c_M^2}{d_M^2}, \dots, \frac{c_M^R}{d_M^R}) \quad (15)$$

where R is the learned Rank and M denotes the number of iterations.

C. Updates for J

$$q_{j_i} = \mathcal{N}(j_i | \mu_i^J, \Xi_i^J) \quad (16)$$

$$\Xi_i^J = (\text{Diag}(\hat{\nu}) + \sum_{t=1}^{\tau} \Sigma_{t, t-1}^B)^{-1} \quad (17)$$

$$\mu_i^J = [\Xi_i^J \Sigma_{t, t-1}^B]_{.i} \quad (18)$$

where Ξ_i^J corresponds to the Covariance Matrix corresponding to the i th row of J and τ denotes current time.

$$\Sigma_{t, t-1}^B = \mu_t^B (\mu_{t-1}^B)^T + \Xi_{t, t-1}^B \quad (19)$$

$\Xi_{t,t-1}^B$ denotes the Covariance Matrix corresponding to the $t-1$ th Column of B which has been obtained already at time $t-1$

D. Updates for B

Define

$$\hat{\beta} = \frac{|\Omega_t|}{\sum_{t=1}^{\tau} \left(\sum_{i_1, i_2, \dots, i_N \in \Omega_t} \left(\Psi_t^1 + \Psi_t^2 \right) \right)} \quad (20)$$

$$\begin{aligned} \Psi_t^1 &= y_{i_1, i_2, \dots, i_N}^2 - 2(y_{i_1, i_2, \dots, i_N} - [\mu_{i_1, \dots, i_N}^S]_{t=t})(\mu_i^A)^T \mu_{\tau}^B \\ &\quad - 2y_{i_1, i_2, \dots, i_N} [\mu_{i_1, i_2, \dots, i_N}^S]_{t=t} \\ \Psi_t^2 &= ([\mu_{i_1, i_2, \dots, i_N}^S]_{t=t})^2 + \Xi_{i_1, i_2, \dots}^S + \text{Trace}(\Sigma_i^A \Sigma_{t,t}^B) \end{aligned}$$

$y_{i_1, i_2, \dots, i_N} - [\mu_{i_1, \dots, i_N}^S]_{t=t} = [D_{i_1, i_2, \dots, i_N}]_{t=t}$ for $t \leq \tau-1$ and for $t = \tau$ we express it otherwise. Also $\mu_{i_1, i_2, \dots, i_N}^S$ corresponds to the mean of the Univariate Normal corresponding to the i_1, i_2, \dots, i_n entry of S .

Now we describe how to define A a matrix of dimensions $\prod_{i=1}^N I_i \times R$, and subsequently defining the mean and covariance for each of its rows and this serves useful for the update of B.

$A = A^{(1)} \odot A^{(2)} \odot \dots A^{(N)}$ and \odot denotes the Khatri Rao product.

$$\mu^A = \left[\langle A^{(1)} \rangle \odot \langle A^{(2)} \rangle \odot \langle A^{(3)} \rangle \odot \dots \odot \langle A^{(N)} \rangle \right] \quad q(\gamma_{i_1, i_2, i_3, \dots, i_N}) = \Gamma \left(\gamma_{i_1, i_2, i_3, \dots, i_N} | a_0^\gamma + \frac{1}{2}, b^{\gamma_{i_1}, \dots, i_N} \right) \quad (26)$$

and μ_i^A denotes mean corresponding to i th row of A that is and the angular brackets denote the expectations of each Factor Matrix.

Define

$$\Xi_i^A = \prod_{k=1}^N \text{elementwise} (V_{i_k}^{(k)} + \bar{a}_{i_k}^k (\bar{a}_{i_k}^k)^T) - \prod_{k=1}^N \text{elementwise} \bar{a}_{i_k}^k (\bar{a}_{i_k}^k)^{\gamma_{i_1}, \dots, i_N} = b_0^\gamma + \frac{1}{2} \cdot \left((\mu_{i_1, i_2, \dots, i_N})^2 + \sigma_{i_1, i_2, \dots, i_N}^2 \right) \quad \text{where}$$

Ξ_i^A denotes to Covariance Matrix of i th row of A and i_1, i_2, \dots, i_N are in the same order as the rows are selected in the Khatri Rao Product that gives A. The next few lines define some terms which are used to reduce the clumsiness of the expressions.

$$\Sigma_i^A = \mu_i^A (\mu_i^A)^T + \Xi_i^A$$

$$\Sigma_{\tau, \tau}^B = \mu_{\tau}^B (\mu_{\tau}^B)^T + \Xi_{\tau, \tau}^B$$

where $\Sigma_{\tau, \tau}^B$ is an $r_t \times r_t$ matrix

$$\Xi_B = \hat{\beta} \cdot \text{Diag}(\Sigma_{(1)}^A, \Sigma_{(2)}^A \dots \Sigma_{(\tau)}^A) + \underbrace{\begin{bmatrix} \zeta_1^{-1} & -\hat{J} & \dots & \dots & 0 \\ -\hat{J} & I_r + \Sigma_J & -\hat{J} & \dots & \dots \\ \dots & \dots & \dots & -\hat{J} & I_r \end{bmatrix}}_{\mathcal{X}} \quad (21)$$

where \mathcal{X} is a triangular matrix and $\Sigma_{(t)}^A = \sum_{i \in \Omega_t} \Sigma_i^A$

$$\mu_B = \Xi_B \begin{bmatrix} \hat{\beta} \sum_{i \in \Omega_1} [y_{i_1, i_2, \dots, i_n}]_{t=1} \cdot \mu_i^A + \zeta_1^{-1} \mu_1 \\ \hat{\beta} \sum_{i \in \Omega_2} [y_{i_1, i_2, \dots, i_n}]_{t=2} \cdot \mu_i^A \\ \dots \\ \dots \\ \hat{\beta} \sum_{i \in \Omega_{\tau}} [y_{i_1, i_2, \dots, i_n}]_{t=\tau} \cdot \mu_i^A \end{bmatrix} \quad (22)$$

where $\hat{J} := \mathbb{E}_q[J|y_{\Omega}]$

E. Posterior Distribution of Sparse S

$$q(S_{\Omega_T}) = \prod_{i_1 i_2 \dots i_N \in \Omega_T} \mathcal{N}(S_{i_1 i_2 \dots i_N} | \mu_{i_1 i_2 \dots i_N}^S, \sigma_{i_1 i_2 \dots i_N}^2) \quad (23)$$

$$\mu_{i_1 i_2 \dots i_N}^S = \sigma_{i_1 i_2 \dots i_N}^2 \mathbb{E}_q[\beta] (\mathcal{Y}_{i_1 i_2 \dots i_N} - \mathbb{E}_q[\langle \hat{a}_{i_1}^{(1)}, \dots, \hat{a}_{i_N}^{(N)}; b_T \rangle]) \quad (24)$$

$$\sigma_{i_1 i_2 \dots i_N}^2 = (\mathbb{E}_q[\gamma_{i_1 i_2 \dots i_N}] + \mathbb{E}_q[\beta])^{-1} \quad (25)$$

F. Posterior Distribution of γ

G. Posterior Distribution of β

$$q^\beta = \Gamma(\beta | a_0, b_0) \quad (28)$$

$$a_0 = a_\beta + \frac{|\Omega_\tau|}{2} \quad (29)$$

$$b_0 = b_\beta + \frac{\mathbb{E}_q[||\mathcal{Y} - \langle \hat{a}_{i_1}^{(1)}, \hat{a}_{i_2}^{(2)}, \dots, \hat{a}_{i_N}^{(N)}; b_T \rangle - S_T ||_F^2]}{2} \quad (30)$$

H. Update for ν

$$\hat{\nu}_i = \frac{T}{\sum_{k \in \text{rows}} [\mu_k^J]^2 + [\Xi_k^J]_{ii}} \quad (31)$$

$$\nu = \Gamma(\nu | \Theta, \kappa) \quad (32)$$

such that

$$\frac{\Theta}{\kappa} = \hat{\nu}_i \quad (33)$$

where $[\mu_k^J]_i$ denotes the i th term of μ_k^J and $[\Xi_k^J]_{ii}$ denotes the ii th term of Ξ_k^J . $\hat{\nu}_i$ denotes the updated Mean of the Gamma Distribuion of ν .

REFERENCES

- [1] Cole Hawkins, Zheng Zhang, Variational Bayesian Inference for Robust Streaming Tensor Factorization and Completion.
- [2] Charul, Uttkarsha Bhatt, Pravesh Biyani, Ketan Rajawat, Traffic estimation and prediction via Online Variational Bayesian Subspace Filtering