

Bayesian Matrix Completion with Hierarchical Gaussian Prior Models

Zebediah

2024 年 7 月 21 日

论文: Bayesian Matrix Completion with Hierarchical Gaussian Prior Models, 作者为 Linxiao Yang, Jun Fang, Huiping Duan, Hongbin Li and Bing Zeng.

代码: <https://www.junfang-uestc.net/codes/LRMC.rar>.

目录

1	引言	2
2	贝叶斯模型	2
3	变分贝叶斯推断	5
3.1	变分贝叶斯方法回顾	5
3.2	提出的算法	5
4	VB-GAMP	7
4.1	使用 GAMP 求解 (4.1)	8
4.2	VB-GAMP 算法	9

1 引言

低秩矩阵补全问题的规范形式可以表示为

$$\begin{aligned} \min_{\mathbf{X}} \quad & \text{rank}(\mathbf{X}), \\ \text{s.t.} \quad & \mathbf{Y} = \mathbf{\Omega} * \mathbf{X}, \end{aligned} \quad (1.1)$$

其中 $\mathbf{X} \in \mathbb{R}^{M \times N}$ 是一个未知的低秩矩阵, $\mathbf{\Omega} \in \{0, 1\}^{M \times N}$ 是一个二元矩阵, 用于表示观察到 \mathbf{X} 的哪些条目, $*$ 表示 Hadamard 乘积, $\mathbf{Y} \in \mathbb{R}^{M \times N}$ 是观察到的矩阵. 然而, 最小化矩阵的秩是 NP-hard 的. 最流行的替代方案是核范数 (矩阵的奇异值之和). 用核范数替换秩函数会产生以下凸优化问题:

$$\begin{aligned} \min_{\mathbf{X}} \quad & \|\mathbf{X}\|_*, \\ \text{s.t.} \quad & \mathbf{Y} = \mathbf{\Omega} * \mathbf{X}. \end{aligned} \quad (1.2)$$

本文开发了一种新的贝叶斯方法来解决低秩矩阵补全问题. 具体地, 假设低秩矩阵 \mathbf{X} 的列为相互独立并服从一个具有零均值和精度矩阵的高斯分布. 精度矩阵为一个随机参数, 并在其上指定了一个 Wishart 分布作为超先验. 本文使用了 GAMP 技术, 并将其嵌入变分贝叶斯 (VB) 推断中, 从而得到了一个用于矩阵补全的高效 VB-GAMP 算法. 但由于先验分布的不可因式分解形式, GAMP 技术不能直接使用. 为了解决这个问题, 我们构建了一个替代问题, 其后验分布正是 VB 推断所需的. 同时, 这个替代问题具有可因式分解的先验和噪声分布, 因此可以直接应用 GAMP 来获得近似后验分布. 这种技巧有助于大幅降低计算复杂度.

2 贝叶斯模型

在有噪声的情况下, 矩阵补全问题的标准形式可以表述为:

$$\begin{aligned} \min_{\mathbf{X}} \quad & \text{rank}(\mathbf{X}), \\ \text{s.t.} \quad & \mathbf{Y} = \mathbf{\Omega} * (\mathbf{X} + \mathbf{E}), \end{aligned} \quad (2.1)$$

其中 \mathbf{E} 表示加性噪声, $\mathbf{\Omega} \in \{0, 1\}^{M \times N}$. 不失一般性, 我们假设 $M \leq N$. 如前所述, 该问题是 NP-hard 的.

本文考虑在贝叶斯框架内对矩阵补全问题进行建模. 我们假设 \mathbf{E} 的条目是独立同分布 (i.i.d.) 的随机变量, 服从均值为零、方差为 γ^{-1} 的高斯分布. 为了学习 γ , 我们使用了一个 Gamma 超先验, 即

$$p(\gamma) = \text{Gamma}(\gamma \mid a, b) = \Gamma(a)^{-1} b^a \gamma^{a-1} e^{-b\gamma}, \quad (2.2)$$

其中 $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$ 是 Gamma 函数. 参数 a 和 b 被设置为较小的值, 例如 10^{-8} , 这使得 Gamma 分布成为一个无信息先验.

为了促进 \mathbf{X} 的低秩性, 我们提出了一个两层的层次高斯先验模型 (见图 1). 具体来说, 在第一层中, 假设 \mathbf{X} 的各列是相互独立的, 并且服从一个共同的高斯分布:

$$p(\mathbf{X} \mid \Sigma) = \prod_{n=1}^N p(\mathbf{x}_n \mid \Sigma) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n \mid \mathbf{0}, \Sigma^{-1}), \quad (2.3)$$

其中 \mathbf{x}_n 表示 \mathbf{X} 的第 i 列, $\boldsymbol{\Sigma} \in \mathbb{R}^{M \times M}$ 是精度矩阵. 第二层指定 Wishart 分布作为精度矩阵 $\boldsymbol{\Sigma}$ 的超先验:

$$p(\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{\frac{\nu-M-1}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{W}^{-1}\boldsymbol{\Sigma})\right), \quad (2.4)$$

其中 ν 和 $\mathbf{W} \in \mathbb{R}^{M \times M}$ 分别表示 Wishart 分布的自由度和尺度矩阵.

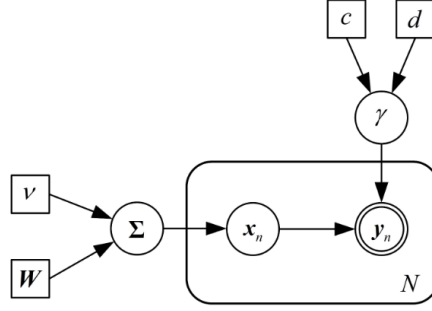


图 1: 分层高斯先验模型

我们将精度矩阵 $\boldsymbol{\Sigma}$ 积分掉得到 \mathbf{X} 的边缘分布:

$$\begin{aligned} p(\mathbf{X}) &= \int \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\Sigma}) p(\boldsymbol{\Sigma}) d\boldsymbol{\Sigma} \\ &\propto |\mathbf{W}^{-1} + \mathbf{X}\mathbf{X}^T|^{-\frac{\nu+N}{2}}. \end{aligned} \quad (2.5)$$

(2.5) 的推导细节如下:

$$\begin{aligned} p(\mathbf{X}) &= \int \prod_{i=1}^N p(\mathbf{x}_i | \boldsymbol{\Sigma}) p(\boldsymbol{\Sigma}) d\boldsymbol{\Sigma} \\ &\propto \int \left(\frac{|\boldsymbol{\Sigma}|}{(2\pi)^M} \right)^{\frac{N}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{X}^T \boldsymbol{\Sigma} \mathbf{X})\right) |\boldsymbol{\Sigma}|^{\frac{\nu-M-1}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{W}^{-1}\boldsymbol{\Sigma})\right) d\boldsymbol{\Sigma} \\ &\propto 2^{\frac{\nu M}{2}} \pi^{-\frac{MN}{2}} \Gamma_M\left(\frac{\nu+N}{2}\right) |\mathbf{W}^{-1} + \mathbf{X}\mathbf{X}^T|^{-\frac{\nu+N}{2}} \int \frac{|\boldsymbol{\Sigma}|^{\frac{\nu+N-M-1}{2}}}{2^{\frac{(\nu+N)M}{2}} |(\mathbf{W}^{-1} + \mathbf{X}\mathbf{X}^T)^{-1}|^{\frac{\nu+N}{2}} \Gamma_M\left(\frac{\nu+N}{2}\right)} d\boldsymbol{\Sigma}, \end{aligned}$$

其中

$$\Gamma_M(x) = \pi^{\frac{M(M-1)}{4}} \prod_{j=1}^M \Gamma\left(x + \frac{1-j}{2}\right).$$

注意到 $p(\mathbf{X})$ 右侧的积分项是一个标准 Wishart 分布, 自由度为 $\nu+N$, 尺度矩阵为 $(\mathbf{W} + \mathbf{X}\mathbf{X}^T)^{-1}$. 从而我们得到

$$\begin{aligned} p(\mathbf{X}) &\propto 2^{\frac{\nu M}{2}} \pi^{-\frac{MN}{2}} \Gamma_M\left(\frac{\nu+N}{2}\right) |\mathbf{W}^{-1} + \mathbf{X}\mathbf{X}^T|^{-\frac{\nu+N}{2}} \\ &\propto |\mathbf{W}^{-1} + \mathbf{X}\mathbf{X}^T|^{-\frac{\nu+N}{2}}. \end{aligned}$$

由 (2.5), 我们有

$$\log p(\mathbf{X}) \propto -\log |\mathbf{X}\mathbf{X}^T + \mathbf{W}^{-1}|. \quad (2.6)$$

如果我们选择 $\mathbf{W} = \epsilon^{-1}\mathbf{I}$, 并且令 ϵ 为一个小的正值, 那么

$$\begin{aligned} \log p(\mathbf{X}) &\propto -\log |\mathbf{X}\mathbf{X}^T + \epsilon\mathbf{I}| \\ &= -\sum_{m=1}^M \log(\lambda_m + \epsilon) \end{aligned} \quad (2.7)$$

其中 λ_m 表示 $\mathbf{X}\mathbf{X}^T$ 的第 m 个特征值. 显然, 在这种情况下, 先验 $p(\mathbf{X})$ 鼓励一个低秩解 \mathbf{X} . 这是因为最大化先验分布 $p(\mathbf{X})$ 等价于对 $\{\lambda_m\}$ 最小化 $\sum_{m=1}^M \log(\lambda_m + \epsilon)$. 而对数和函数 $\sum_{m=1}^M \log(\lambda_m + \epsilon)$ 是一种有效的促进稀疏的泛函.

除了 $\mathbf{W} = \epsilon^{-1}\mathbf{I}$, 参数 \mathbf{W} 还可以根据对 \mathbf{X} 的额外先验知识进行设计. 例如, 在一些应用中 (如图像修复), x_n 的相邻系数之间存在空间相关性. 为了捕获相邻系数之间的平滑性, \mathbf{W} 可以设置为

$$\mathbf{W} = \mathbf{F}^T \mathbf{F}, \quad (2.8)$$

其中 $\mathbf{F} \in \mathbb{R}^{M \times M}$ 是一个二阶差分算子, 其 (i, j) 项给出如下:

$$f_{i,j} = \begin{cases} -2, & i = j, \\ 1, & |i - j| = 1, \\ 0, & \text{else.} \end{cases} \quad (2.9)$$

另一种促进平滑解的 \mathbf{W} 选择是拉普拉斯矩阵, i.e.

$$\mathbf{W} = \mathbf{D} - \mathbf{A} + \hat{\epsilon}\mathbf{I}, \quad (2.10)$$

其中 \mathbf{A} 是图的邻接矩阵, 其元素为

$$a_{ij} = \exp\left(-\frac{|i - j|^2}{\theta^2}\right). \quad (2.11)$$

\mathbf{D} 称为度矩阵, 是一个对角元素为 $d_{ii} = \sum_j a_{ij}$ 的对角矩阵, $\hat{\epsilon}$ 是一个小的正值, 用于确保 \mathbf{W} 是满秩的. 可以证明, (2.8) 和 (2.10) 中定义的 \mathbf{W} 促进了 \mathbf{X} 的低秩性和平滑性. 为说明这一点, 我们首先介绍以下引理:

引理 1 对于一个正定矩阵 $\mathbf{W} \in \mathbb{R}^{M \times M}$, 对任意 $\mathbf{X} \in \mathbb{R}^{M \times N}$, 成立

$$\log |\mathbf{X}\mathbf{X}^T + \mathbf{W}^{-1}| = \log |\mathbf{W}^{-1}| + \log |\mathbf{I} + \mathbf{X}^T \mathbf{W} \mathbf{X}|. \quad (2.12)$$

由引理 1, 我们有

$$\begin{aligned} \log p(\mathbf{X}) &\propto -\log |\mathbf{X}\mathbf{X}^T + \mathbf{W}^{-1}| \\ &\propto -\log |\mathbf{I} + \mathbf{X}^T \mathbf{W} \mathbf{X}| \\ &= -\sum_{n=1}^N \log(\tilde{\lambda}_n + 1), \end{aligned} \quad (2.13)$$

其中 $\tilde{\lambda}_n$ 是 $\mathbf{X}^T \mathbf{W} \mathbf{X}$ 的第 n 个特征值. 我们看到, 最大化先验分布等价于对 $\{\tilde{\lambda}_n\}$ 最小化 $\sum_{n=1}^N \log(\tilde{\lambda}_n + 1)$.

从而得到稀疏的 $\{\lambda_n\}$, 即 $\mathbf{X}^T \mathbf{W} \mathbf{X}$ 是低秩的. 因为 \mathbf{W} 是满秩的, 因此 \mathbf{X} 是低秩的. 另一方面, 注意到 $\text{tr}(\mathbf{X}^T \mathbf{W} \mathbf{X})$ 是 $\log |\mathbf{I} + \mathbf{X}^T \mathbf{W} \mathbf{X}|$ 的一阶近似, 因此最小化 $\log |\mathbf{I} + \mathbf{X}^T \mathbf{W} \mathbf{X}|$ 会减少 $\text{tr}(\mathbf{X}^T \mathbf{W} \mathbf{X})$ 的值. 显然, 对于 (2.8) 和 (2.10) 中定义的 \mathbf{W} , 一个更平滑的解会导致 $\text{tr}(\mathbf{X}^T \mathbf{W} \mathbf{X})$ 的值更小. 因此, 当 \mathbf{W} 选取为 (2.8) 或 (2.10) 时, 所得的先验分布 $p(\mathbf{X})$ 可以促进低秩和平滑解.

3 变分贝叶斯推断

3.1 变分贝叶斯方法回顾

我们首先简要回顾一下变分贝叶斯 (VB) 方法, 在概率模型中, 令 \mathbf{y} 和 $\boldsymbol{\theta}$ 分别表示观测数据和隐藏变量. 很容易证明, 观测数据的边际概率可以分解为两个项:

$$\ln p(\mathbf{y}) = L(q) + \text{KL}(q\|p), \quad (3.1)$$

其中

$$L(q) = \int q(\boldsymbol{\theta}) \ln \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}, \quad \text{KL}(q\|p) = - \int q(\boldsymbol{\theta}) \ln \frac{p(\boldsymbol{\theta} | \mathbf{y})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}, \quad (3.2)$$

这里 $q(\boldsymbol{\theta})$ 是任何概率密度函数, $\text{KL}(q\|p)$ 是 $p(\boldsymbol{\theta} | \mathbf{y})$ 和 $q(\boldsymbol{\theta})$ 之间的 Kullback-Leibler 散度. 由于 $\text{KL}(q\|p) \geq 0$, 因此 $L(q)$ 是 $\ln p(\mathbf{y})$ 的严格下界. 因此最大化 $L(q)$ 等价于最小化 $\text{KL}(q\|p)$, 因此后验分布 $p(\boldsymbol{\theta} | \mathbf{y})$ 可以通过最大化 $L(q)$ 来近似为 $q(\boldsymbol{\theta})$.

最大化过程可以对每个潜在变量进行交替进行:

$$q_i(\theta_i) = \frac{e^{\langle \ln p(\mathbf{y}, \boldsymbol{\theta}) \rangle_{k \neq i}}}{\int e^{\langle \ln p(\mathbf{y}, \boldsymbol{\theta}) \rangle_{k \neq i}} d\theta_i}, \quad (3.3)$$

其中 $\langle \cdot \rangle_{k \neq i}$ 表示对所有 $k \neq i$ 的分布 $q_k(\theta_k)$ 求期望. 两边取对数可得

$$\ln q_i(\theta_i) = \langle \ln p(\mathbf{y}, \boldsymbol{\theta}) \rangle_{k \neq i} + \text{constant}. \quad (3.4)$$

3.2 提出的算法

现在我们进行对提出的层次模型的变分贝叶斯推理. 令 $\boldsymbol{\theta} \triangleq \{\mathbf{X}, \boldsymbol{\Sigma}, \gamma\}$ 表示所有隐藏变量, 我们的目标是找到后验分布 $p(\boldsymbol{\theta} | \mathbf{y})$. 由于计算 $p(\boldsymbol{\theta} | \mathbf{y})$ 通常不可行, 我们将 $p(\boldsymbol{\theta} | \mathbf{y})$ 近似为 $q(\mathbf{X}, \boldsymbol{\Sigma}, \gamma)$, 其因子化形式为隐藏变量 $\{\mathbf{X}, \boldsymbol{\Sigma}, \gamma\}$ 的因子化形式, i.e.

$$q(\mathbf{X}, \boldsymbol{\Sigma}, \gamma) = q_x(\mathbf{X}) q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) q_{\gamma}(\gamma). \quad (3.5)$$

如前面所述, $L(q)$ 的最大化可以对每个潜在变量交替进行, 因此我们有

$$\begin{aligned} \ln q_x(\mathbf{X}) &= \langle \ln p(\boldsymbol{\Sigma}, \gamma) \rangle_{q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) q_{\gamma}(\gamma)} + \text{constant}, \\ \ln q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) &= \langle \ln p(\mathbf{X}, \gamma) \rangle_{q_x(\mathbf{X}) q_{\gamma}(\gamma)} + \text{constant}, \\ \ln q_{\gamma}(\gamma) &= \langle \ln p(\mathbf{X}, \boldsymbol{\Sigma}) \rangle_{q_x(\mathbf{X}) q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma})} + \text{constant}, \end{aligned} \quad (3.6)$$

其中 $\langle \cdot \rangle_{q_1(\cdot) \dots q_K(\cdot)}$ 表示对分布 $\{q_k(\cdot)\}_{k=1}^K$ 的期望. 接下来将介绍算法的详细内容.

(1) $q_x(\mathbf{X})$ 的更新: $q_x(\mathbf{X})$ 的计算可以分解为一组独立任务, 每个任务计算 \mathbf{X} 每列的后验分布近似 $q_x(\mathbf{x}_n)$. 我们有

$$\begin{aligned} \ln q_x(\mathbf{x}_n) &\propto \langle \ln [p(\mathbf{y}_n | \mathbf{x}_n) p(\mathbf{x}_n | \boldsymbol{\Sigma})] \rangle_{q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) q_{\gamma}(\gamma)} \\ &\propto \left\langle -\gamma (\mathbf{y}_n - \mathbf{x}_n)^T \mathbf{O}_n (\mathbf{y}_n - \mathbf{x}_n) - \mathbf{x}_n^T \boldsymbol{\Sigma} \mathbf{x}_n \right\rangle \\ &\propto -\mathbf{x}_n^T (\langle \gamma \rangle \mathbf{O}_n + \langle \boldsymbol{\Sigma} \rangle) \mathbf{x}_n + 2 \langle \gamma \rangle \mathbf{x}_n^T \mathbf{O}_n \mathbf{y}_n, \end{aligned} \quad (3.7)$$

其中 \mathbf{y}_n 是 \mathbf{Y} 的第 n 列, $\mathbf{O}_n \triangleq \text{diag}(\mathbf{o}_n)$, 这里 \mathbf{o}_n 是 $\boldsymbol{\Omega}$ 的第 n 列. 由 (3.7), \mathbf{x}_n 服从高斯分布

$$q_x(\mathbf{x}_n) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, \mathbf{Q}_n), \quad (3.8)$$

其中

$$\boldsymbol{\mu}_n = \langle \gamma \rangle \mathbf{Q}_n \mathbf{O}_n \mathbf{y}_n, \quad \mathbf{Q}_n = (\langle \gamma \rangle \mathbf{O}_n + \langle \boldsymbol{\Sigma} \rangle)^{-1}. \quad (3.9)$$

可以看出, 为了计算 $q_x(\mathbf{x}_n)$, 需要对一个 $M \times M$ 的矩阵求逆, 计算复杂度为 $\mathcal{O}(M^3)$.

(2) $q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$ 的更新: 可以得到近似后验 $q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$ 为

$$\begin{aligned} \ln q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) &\propto \left\langle \ln \left[\prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\Sigma}) p(\boldsymbol{\Sigma}) \right] \right\rangle_{q_x(\mathbf{X})} \\ &\propto \left\langle \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr}(\mathbf{X}^T \boldsymbol{\Sigma} \mathbf{X}) + \frac{\nu - M - 1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr}(\mathbf{W}^{-1} \boldsymbol{\Sigma}) \right\rangle \\ &\propto \frac{\nu + N - M - 1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr}((\mathbf{W}^{-1} + \langle \mathbf{X} \mathbf{X}^T \rangle) \boldsymbol{\Sigma}). \end{aligned} \quad (3.10)$$

由 (3.10) 可知, $\boldsymbol{\Sigma}$ 服从 Wishart 分布, i.e.

$$q_{\boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) = \text{Wishart}(\boldsymbol{\Sigma}; \hat{\mathbf{W}}, \hat{\nu}), \quad (3.11)$$

其中

$$\hat{\mathbf{W}} = (\mathbf{W}^{-1} + \langle \mathbf{X} \mathbf{X}^T \rangle)^{-1}, \quad \hat{\nu} = \nu + N. \quad (3.12)$$

(3) $q_{\gamma}(\gamma)$ 的更新: 通过对 $q_{\gamma}(\gamma)$ 的变分优化得到

$$\begin{aligned} \ln q_{\gamma}(\gamma) &\propto \langle \ln p(\mathbf{Y} | \mathbf{X}, \gamma) p(\gamma) \rangle_{q_x(\mathbf{X})} \\ &\propto \left\langle \ln \prod_{(m,n) \in \mathbb{S}} p(y_{mn} | x_{mn}, \gamma) p(\gamma) \right\rangle \\ &\propto \left\langle \frac{L}{2} \ln \gamma - \frac{\gamma}{2} \sum_{(m,n) \in \mathbb{S}} (y_{mn} - x_{mn})^2 + (c-1) \ln \gamma - d\gamma \right\rangle \\ &= \left(\frac{L}{2} + c - 1 \right) \ln \gamma - \left(\frac{1}{2} \sum_{(m,n) \in \mathbb{S}} \langle (y_{mn} - x_{mn})^2 \rangle + d \right), \end{aligned} \quad (3.13)$$

其中 x_{mn} 和 y_{mn} 表示 \mathbf{X} 和 \mathbf{Y} 的第 (m, n) 项, $\mathbb{S} \triangleq \{(m, n) | \Omega_{mn} = 1\}$ 是由观测值的索引组成的集合,

$L \triangleq |\mathbb{S}|$ 是集合 \mathbb{S} 的基数, Ω_{mn} 表示 Ω 的第 (m, n) 项. 可以验证 $q_\gamma(\gamma)$ 服从 Gamma 分布

$$q_\gamma(\gamma) = \text{Gamma}(\gamma \mid \tilde{c}, \tilde{d}), \quad (3.14)$$

其中

$$\tilde{c} = \frac{L}{2} + c, \quad \tilde{d} = \frac{1}{2} \sum_{(m,n) \in \mathbb{S}} \langle (y_{mn} - x_{mn})^2 \rangle + d, \quad (3.15)$$

这里

$$\langle (y_{mn} - x_{mn})^2 \rangle = y_{mn}^2 - 2y_{mn} \langle x_{mn} \rangle + \langle x_{mn}^2 \rangle. \quad (3.16)$$

在更新过程中使用的一些期望和矩如下:

$$\begin{aligned} \langle \Sigma \rangle &= \hat{\mathbf{W}} \hat{\nu}, \\ \langle \mathbf{X} \mathbf{X}^T \rangle &= \langle \mathbf{X} \rangle \langle \mathbf{X} \rangle^T + \sum_{n=1}^N \mathbf{Q}_n, \\ \langle x_{mn}^2 \rangle &= \langle x_{mn} \rangle^2 + Q_n(m, m), \end{aligned} \quad (3.17)$$

其中 $Q_n(m, m)$ 表示 \mathbf{Q}_n 的第 m 个对角元. 上述过程总结如算法所示.

Algorithm 1 VB Algorithm for Matrix Completion

- 1: **Input:** \mathbf{Y}, Ω, ν and \mathbf{W} .
 - 2: **Initialization:** $\langle \Sigma \rangle$ and $\langle \gamma \rangle$.
 - 3: **while** not converge **do**
 - 4: **for** $n = 1$ to N **do**
 - 5: Update $q_x(\mathbf{x}_n)$ via (3.8), with $q_\Sigma(\Sigma)$ and $q_\gamma(\gamma)$ fixed;
 - 6: **end for**
 - 7: Update $q_\Sigma(\Sigma)$ via (3.11), with $q_x(\mathbf{X})$ and $q_\gamma(\gamma)$ fixed;
 - 8: Update $q_\gamma(\gamma)$ via (3.14);
 - 9: **end while**
 - 10: **Output:** $q_x(\mathbf{X}), q_\Sigma(\Sigma), q_\gamma(\gamma)$.
-

4 VB-GAMP

GAMP 算法要求先验分布和噪声分布都具有因子化形式. 然而在我们的模型中, 先验分布 $p(\mathbf{x}_n \mid \Sigma)$ 具有非因子化形式. 为了解决这个问题, 我们构建了一个替代问题, 该问题旨在从 $\mathbf{b} \in \mathbb{R}^M$ 中恢复 $\mathbf{x} \in \mathbb{R}^M$:

$$\mathbf{b} = \mathbf{U}^T \mathbf{x} + \mathbf{e}, \quad (4.1)$$

其中 $\mathbf{U} \in \mathbb{C}^{M \times M}$ 是通过将 $\langle \Sigma \rangle = \mathbf{U} \mathbf{S} \mathbf{U}^T$ 进行奇异值分解得到的, \mathbf{e} 表示加性高斯噪声, 均值为零, 协方差矩阵为 \mathbf{S}^{-1} . 我们假设 \mathbf{x} 的各个条目是相互独立的, 并服从以下分布:

$$p(x_m) = \begin{cases} \mathcal{N}(\kappa_m, \xi^{-1}), & \text{if } \pi_m = 1, \\ C, & \text{if } \pi_m = 0, \end{cases} \quad (4.2)$$

其中 π_m, x_m, κ_m 分别表示 $\boldsymbol{\pi}, \mathbf{x}, \boldsymbol{\kappa}$ 的第 m 项, C 是一个常数, $\boldsymbol{\pi}, \boldsymbol{\kappa} \in \mathbb{R}^{M \times 1}$, ξ 是已知参数. 考虑问题 (4.1), \mathbf{x} 的后验分布可以计算为:

$$\begin{aligned} p(\mathbf{x} | \mathbf{b}) &\propto p(\mathbf{b} | \mathbf{x})p(\mathbf{x}) \\ &\propto p(\mathbf{b} | \mathbf{x}) \prod_{m \in S} p(x_m) \\ &= \mathcal{N}(\mathbf{U}^T \mathbf{x}, \mathbf{S}^{-1}) \prod_{m \in S} \mathcal{N}(\kappa_m, \xi^{-1}), \end{aligned} \quad (4.3)$$

其中 $S \triangleq \{m | \pi_m = 1\}$.

对 $p(\mathbf{x} | \mathbf{b})$ 取对数有

$$\begin{aligned} \ln p(\mathbf{x} | \mathbf{b}) &\propto -\frac{1}{2} (\mathbf{b} - \mathbf{U}^T \mathbf{x})^T \mathbf{S} (\mathbf{b} - \mathbf{U}^T \mathbf{x}) - \frac{1}{2} \xi \sum_{m \in S} (x_m - \kappa_m)^2 \\ &= -\frac{1}{2} (\mathbf{b} - \mathbf{U}^T \mathbf{x})^T \mathbf{S} (\mathbf{b} - \mathbf{U}^T \mathbf{x}) - \frac{1}{2} \xi (\mathbf{x} - \boldsymbol{\kappa})^T \boldsymbol{\Pi} (\mathbf{x} - \boldsymbol{\kappa}) \\ &\propto -\frac{1}{2} \mathbf{x}^T (\mathbf{U} \mathbf{S} \mathbf{U}^T + \xi \boldsymbol{\Pi}) \mathbf{x} + (\mathbf{b}^T \mathbf{S} \mathbf{U}^T + \xi \boldsymbol{\kappa}^T \boldsymbol{\Pi}) \mathbf{x}, \end{aligned} \quad (4.4)$$

其中 $\boldsymbol{\Pi}$ 是一个对角矩阵, 其第 m 个对角元素等于 π_m . 显然, $p(\mathbf{x} | \mathbf{b})$ 服从高斯分布, 其均值 $\boldsymbol{\mu}$ 和协方差矩阵 \mathbf{Q} 为

$$\boldsymbol{\mu} = \mathbf{Q}(\mathbf{U} \mathbf{S} \mathbf{b} + \xi \boldsymbol{\Pi} \boldsymbol{\kappa}) \mathbf{Q} = (\mathbf{U} \mathbf{S} \mathbf{U}^T + \xi \boldsymbol{\Pi})^{-1} = (\langle \boldsymbol{\Sigma} \rangle + \xi \boldsymbol{\Pi})^{-1}. \quad (4.5)$$

比较 (3.9) 和 (4.5), 我们可以验证, 当 $\mathbf{b} = \mathbf{0}, \boldsymbol{\kappa} = \mathbf{y}_n, \boldsymbol{\pi} = \mathbf{o}_n$ (i.e. $\boldsymbol{\Pi} = \mathbf{O}_n$), 且 $\xi = \langle \gamma \rangle$ 时, $p(\mathbf{x} | \mathbf{b})$ 正好是所需的后验分布 $q_x(\mathbf{x}_n)$. 同时对于替代问题 (4.1), 先验分布和噪声分布都是可因子化的. 因此, 可以直接应用 GAMP 算法到 (4.1) 以找到后验分布 $p(\mathbf{x} | \mathbf{b})$ 的近似. 通过设置 $\mathbf{b} = \mathbf{0}, \boldsymbol{\kappa} = \mathbf{y}_n, \boldsymbol{\pi} = \mathbf{o}_n, \xi = \langle \gamma \rangle$ 可以有效地获得 (3.8) 中 $q_x(\mathbf{x}_n)$ 的近似.

4.1 使用 GAMP 求解 (4.1)

首先, GAMP 通过以下公式近似真实的边际后验分布 $p(x_m | \mathbf{b})$:

$$\hat{p}(x_m | \mathbf{b}, \hat{r}_m, \tau_m^r) = \frac{p(x_m) \mathcal{N}(x_m | \hat{r}_m, \tau_m^r)}{\int_x p(x_m) \mathcal{N}(x_m | \hat{r}_m, \tau_m^r) dx}, \quad (4.6)$$

其中 \hat{r}_m 和 τ_m^r 是在 GAMP 算法的迭代过程中迭代更新的量. 这里为了简化, 省略了对迭代次数 k 的显式依赖 (下同). 在 $\pi_m = 1$ 的情况下, 将先验分布 (4.2) 代入 (4.6) 中, 可以验证近似后验分布 $\hat{p}(x_m | \mathbf{b}, \hat{r}_m, \tau_m^r)$ 服从高斯分布, 其均值和方差分别为:

$$\mu_m^x = \phi_m^x (\xi \kappa_m + \hat{r}_m / \tau_m^r), \quad \phi_m^x = \frac{\tau_m^r}{1 + \xi \tau_m^r}. \quad (4.7)$$

同样, 在 $\pi_m = 0$ 的情况下, 近似后验分布 $\hat{p}(x_m | \mathbf{b}, \hat{r}_m, \tau_m^r)$ 也服从高斯分布, 其均值和方差分别为:

$$\mu_m^x = \hat{r}_m, \quad \phi_m^x = \tau_m^r. \quad (4.8)$$

另一种近似是对无噪声输出 $z_i \triangleq \mathbf{u}_i^T \mathbf{x}$ 进行, 其中 \mathbf{u}_i^T 表示 \mathbf{U}^T 的第 i 行. GAMP 通过以下公式近似真

实的边际后验分布 $p(z_i | \mathbf{b})$:

$$\hat{p}(z_i | \mathbf{b}, \hat{p}_i, \tau_i^p) = \frac{p(b_i | z_i) \mathcal{N}(z_i | \hat{p}_i, \tau_i^p)}{\int_z p(b_i | z_i) \mathcal{N}(z_i | \hat{p}_i, \tau_i^p) dz}, \quad (4.9)$$

其中 \hat{p}_i 和 τ_i^p 是在 GAMP 算法的迭代过程中迭代更新的量. 在加性白高斯噪声假设下, 我们有 $p(b_i | z_i) = \mathcal{N}(b_i | z_i, s_i^{-1})$, 其中 s_i 为 S 的第 i 个对角元素. 故 $\hat{p}(z_i | \mathbf{b}, \hat{p}_i, \tau_i^p)$ 也服从高斯分布, 其均值和方差分别为

$$\mu_i^z = \frac{\tau_i^p s_i b_i + \hat{p}_i}{1 + s_i \tau_i^p}, \quad \phi_i^z = \frac{\tau_i^p}{1 + s_i \tau_i^p}. \quad (4.10)$$

有了上述近似, 我们现在可以定义 GAMP 算法中使用的两个标量函数: $g_{\text{in}}(\cdot)$ 和 $g_{\text{out}}(\cdot)$. 输入标量函数 $g_{\text{in}}(\cdot)$ 简单地定义为后验均值 μ_m^x , i.e.

$$g_{\text{in}}(\hat{r}_m, \tau_m^r) = \mu_m^x = \begin{cases} \phi_m^x (\xi \kappa_m + \hat{r}_m / \tau_m^r), & \text{if } \pi_m = 1, \\ \hat{r}_m, & \text{if } \pi_m = 0. \end{cases} \quad (4.11)$$

$\tau_m^r g_{\text{in}}(\hat{r}_m, \tau_m^r)$ 对 \hat{r}_m 的偏导是后验方差 ϕ_m^x , i.e.

$$\tau_m^r \frac{\partial}{\partial \hat{r}_m} g_{\text{in}}(\hat{r}_m, \tau_m^r) = \phi_m^x = \begin{cases} \frac{\tau_m^r}{1 + \xi \tau_m^r}, & \text{if } \pi_m = 1, \\ \tau_m^r, & \text{if } \pi_m = 0. \end{cases} \quad (4.12)$$

输出标量函数 $g_{\text{out}}(\cdot)$ 与后验均值 μ_i^z 的关系如下:

$$g_{\text{out}}(\hat{p}_i, \tau_i^p) = \frac{1}{\tau_i^p} (\mu_i^z - \hat{p}_i) = \frac{s_i (b_i - \hat{p}_i)}{1 + s_i \tau_i^p}. \quad (4.13)$$

标量函数 $g_{\text{out}}(\hat{p}_i, \tau_i^p)$ 的偏导数与后验方差 $\phi_{i,n}^z$ 之间的关系如下:

$$\frac{\partial}{\partial \hat{p}_i} g_{\text{out}}(\hat{p}_i, \tau_i^p) = \frac{\phi_i^z - \tau_i^p}{(\tau_i^p)^2} = \frac{-s_i}{(1 + s_i \tau_i^p)} \quad (4.14)$$

综上所述, 可以总结适用于问题 (4.1) 的 GAMP 算法 (见算法 2), 其中 $u_{i,m}$ 表示 \mathbf{U}^T 的第 (i, m) 项.

4.2 VB-GAMP 算法

我们已经推导出一个有效的算法来获得问题 (4.1) 中 \mathbf{x} 的近似后验分布. 具体来说, x_m 的真实边际后验分布近似为一个高斯分布 $\hat{p}(x_m | \mathbf{b}, \hat{r}_m, \tau_m^r)$, 其均值和方差分别由 (4.7) 或 (4.8) 给出. 联合后验分布 $p(\mathbf{x} | \mathbf{b})$ 可以近似为近似边际后验分布的乘积:

$$p(\mathbf{x} | \mathbf{b}) \approx \hat{p}(\mathbf{x} | \mathbf{b}) = \prod_{m=1}^M \hat{p}(x_m | \mathbf{b}, \hat{r}_m, \tau_m^r). \quad (4.15)$$

如前所述, 通过设置 $b = 0, \kappa = y_n, \pi = o_n, \xi = \langle \gamma \rangle$ 可以使用 GAMP 算法得到的后验分布 $\hat{p}(x | b)$ 来近似 (3.8) 中的 $q_x(\mathbf{x}_n)$. VB-GAMP 算法如算法 3 所示.

Algorithm 2 GAMP Algorithm

-
- 1: **Input:** κ, π, b , and ξ .
 - 2: **Initialization:** Set $\hat{\psi}_i = 0, \forall i \in \{1, \dots, M\}$; $\{\mu_m^x\}_{m=1}^M$ are initialized as the mean variance of the prior distribution, and $\{\phi_m^x\}_{m=1}^M$ are set to small values, say 10^{-5} .
 - 3: **while** not converge **do**
 - 4: **Step 1.** $\forall i \in \{1, \dots, M\}$:

$$\hat{z}_i = \sum_m u_{i,m} \mu_m^x, \quad \tau_i^p = \sum_m u_{i,m}^2 \phi_m^x \hat{p}_i = \hat{z}_i - \tau_i^p \hat{\psi}_i.$$

- 5: **Step 2.** $\forall i \in \{1, \dots, M\}$:

$$\hat{\psi}_i = g_{\text{out}}(\hat{p}_i, \tau_i^p), \quad \tau_i^s = -\frac{\partial}{\partial \hat{p}_i} g_{\text{out}}(\hat{p}_i, \tau_i^p).$$

- 6: **Step 3.** $\forall m \in \{1, \dots, M\}$:

$$\tau_m^r = \left(\sum_i u_{i,m}^2 \tau_i^s \right)^{-1}, \quad \hat{r}_m = \mu_m^x + \tau_m^r \sum_i u_{i,m} \hat{\psi}_i.$$

- 7: **Step 4.** $\forall m \in \{1, \dots, M\}$:

$$\mu_m^x = g_{\text{in}}(\hat{r}_m, \tau_m^r), \quad \phi_m^x = \tau_m^r \frac{\partial}{\partial \hat{r}_m} g_{\text{in}}(\hat{r}_m, \tau_m^r).$$

- 8: **end while**
 - 9: **Output:** $\{\hat{r}_m, \tau_m^r\}, \{\hat{p}_i, \tau_i^p\}$, and $\{\mu_m^x, \phi_m^x\}$.
-

Algorithm 3 VB-GAMP Algorithm for Matrix Completion

-
- 1: **Input:** $\mathbf{Y}, \mathbf{\Omega}, \nu$ and \mathbf{W} .
 - 2: **Initialization:** $\langle \mathbf{X} \rangle, \langle \mathbf{\Sigma} \rangle$.
 - 3: **while** not converge **do**
 - 4: Calculate singular value decomposition of $\langle \mathbf{\Sigma} \rangle$;
 - 5: **for** $n = 1$ to N **do**
 - 6: Obtain an approximation of $q_x(\mathbf{x}_n)$ via Algorithm 2;
 - 7: **end for** Update $q_{\Sigma}(\mathbf{\Sigma})$ via (3.11); Update $q_{\gamma}(\gamma)$ via (3.14);
 - 8: **end while**
 - 9: **Output:** $q_x(\mathbf{X}), q_{\Sigma}(\mathbf{\Sigma})$, and $q_{\gamma}(\gamma)$.
-