

On the issue of using low-rank covariance matrix in mash

Dongyue Xie

1 Model

Assume

$$\begin{aligned}\mathbf{x}_i &\sim N_R(\boldsymbol{\mu}_i, \mathbf{I}), \\ \boldsymbol{\mu}_i &\sim \sum_k \pi_k N_R(\mathbf{0}, \mathbf{U}_k),\end{aligned}\tag{1}$$

where $N_R(\cdot, \cdot)$ denotes a R-dimensional multivariate normal distribution. When using the model for multiple testing, mash first estimates $\boldsymbol{\pi}, \mathbf{U}$ using empirical Baes procedure then obtains lfd, lfsr and the other posterior summaries.

Using rank-1 \mathbf{U}_k in above model breaks down one of the posterior summaries - lfsr. The issue generally appears when using rank-deficient covariance matrices. The assumption of rank-1 is not explicit specified in the mash model, but the ED algorithm preserves the rank of initializations. The initializations are generated using PCA hence most of them are rank-1 matrices. The TEEM algorithm does not preserve the rank but it truncates the eigenvalues and the estimated covaraince matrices could still be low rank.

The simplest solution is to use full-ranked matrices as intializations. For example, generate covaraince matrices from a Wishart distribution as intializations. But this generally increase the computation cost and the estimated covariance matrices might be very close to singular.

In the following sections we discuss two possible solutions.

A more systematic solution to this problem is adding a full ranked matrix to \mathbf{U}_k , and the simplest one perhaps is $\sigma_k^2 \mathbf{I}$. In the next section we will give a Bayesian interpretation of the choice and describe an empirical Bayes procedure for determining σ_k^2 .

2 Prior on covaraince matrix

Preferably, the parameter σ^2 should be data-dependent. Inspired by early works on empirical Bayes estimator of sample covariance matrix, we can employ the EB estimator of covariance matrix by introducing a prior on it. For more details, see Efron and Morris(1972), Efron and Morris(1976), and Haff(1980),

Define $\mathbf{T}_k = \mathbf{U}_k + \mathbf{I}$. When estimating \mathbf{U}_k using EM algorithm, the last iteration gives

$$\hat{\mathbf{T}}_k^{\text{mle}} = \frac{\sum_i \gamma_{ik} \mathbf{x}_i \mathbf{x}_i^T}{n_k},\tag{2}$$

where mle denotes the maximum likelihood estimator, γ_{ik} are posterior probabilities of \mathbf{x}_i coming from mixture component k and $n_k = \sum_i \gamma_{ik}$. (In practice one would use a modified $\hat{\mathbf{T}}_k^{\text{mle}}$ by truncating eigenvalues to ensure \mathbf{U}_k is a well-defined covariance matrix, and $\hat{\mathbf{U}}_k^{\text{mle}} = \hat{\mathbf{T}}_k^{\text{mle}} - \mathbf{I}$.) Most of the posterior probabilities γ_{ik} should be close to either 0 or 1 at the last iteration of EM. The reason is that assuming the model is correct, the posteriors should concentrate on the true mixture component assignments. The number n_k can be regarded as the effective sample size of mixture component k , i.e. the number of \mathbf{x}_i that are from the mixture k . Therefore, we have the following approximation

$$\hat{\mathbf{T}}_k^{\text{mle}} = \frac{\sum_i \gamma_{ik} \mathbf{x}_i \mathbf{x}_i^T}{n_k} \approx \frac{1}{n_k} \sum_{i \in \mathcal{K}} \mathbf{x}_i \mathbf{x}_i^T, \quad (3)$$

where \mathcal{K} is an index set of samples from mixture component k , so for $i \in \mathcal{K}$, $\mathbf{x}_i \sim N(\mathbf{0}, \mathbf{U}_k)$.

Let $\tilde{\mathbf{T}}_k = n_k \hat{\mathbf{T}}_k^{\text{mle}}$, then based on the approximation

$$\tilde{\mathbf{T}}_k \sim \mathcal{W}(\mathbf{T}_k, n_k), \quad (4)$$

where \mathcal{W} denotes Wishart distribution.

Assume \mathbf{T}_k has an inverse Wishart prior,

$$\mathbf{T}_k \sim \mathcal{W}^{-1}(\sigma_k^2 \mathbf{I}, \nu), \quad (5)$$

where ν is known integer and $\sigma_k^2 > 0$. The parameter ν is the degree of freedom and $\nu \geq R$. Larger ν gives smaller prior variance. In practice, we shall set $\nu = R$ or $\nu = R + 1$.

Haff(1980) derived an estimator of σ_k^2 by maximizing the marginal likelihood $f(\tilde{\mathbf{T}}_k; \sigma_k^2)$. Directly maximizing $f(\tilde{\mathbf{T}}_k; \sigma_k^2)$ is non-trivial and the estimator instead maximizes the first order approximation of it. The empirical Bayes estimator of σ_k^2 is

$$\hat{\sigma}_k^2 = \frac{R\nu}{n_k + \nu} \frac{1}{\text{Tr}(\tilde{\mathbf{T}}_k^{-1})}. \quad (6)$$

The estimator $\hat{\sigma}_k^2$ is strictly positive since $\frac{R\nu}{n_k + \nu} > 0$ and $\tilde{\mathbf{T}}_k$ is positive definite.

The posterior distribution of \mathbf{T}_k is

$$\mathbf{T}_k | \tilde{\mathbf{T}}_k \sim \mathcal{W}^{-1}(\tilde{\mathbf{T}}_k + \sigma_k^2 \mathbf{I}, n_k + \nu), \quad (7)$$

The posterior mean of \mathbf{T}_k is

$$\mathbb{E}(\mathbf{T}_k | \tilde{\mathbf{T}}_k, \sigma_k^2) = \frac{\tilde{\mathbf{T}}_k + \sigma_k^2 \mathbf{I}}{n_k + \nu - R - 1}. \quad (8)$$

The empirical Bayes estimator of \mathbf{T}_k is the posterior mean with estimated σ_k^2 ,

$$\hat{\mathbf{T}}_k^{\text{eb}} = \frac{\tilde{\mathbf{T}}_k + \hat{\sigma}_k^2 \mathbf{I}}{n_k + \nu - R - 1}. \quad (9)$$

To get the EB estimator of \mathbf{U}_k , we shall write

$$\begin{aligned}\hat{\mathbf{T}}_k^{\text{eb}} &= \hat{\mathbf{U}}_k^{\text{eb}} + \mathbf{I}, \\ \tilde{\mathbf{T}}_k &= n_k \hat{\mathbf{T}}_k^{\text{mle}} = n_k (\hat{\mathbf{U}}_k^{\text{mle}} + \mathbf{I}).\end{aligned}\tag{10}$$

Substituting (10) into (9), we have

$$\hat{\mathbf{U}}_k^{\text{eb}} = \frac{n_k}{n_k + \nu - R - 1} \hat{\mathbf{U}}_k^{\text{mle}} + \frac{R + 1 - \nu + \hat{\sigma}_k^2}{n_k + \nu - R - 1} \mathbf{I}.\tag{11}$$

When $\nu = R + 1$, we have

$$\hat{\mathbf{U}}_k^{\text{eb}} = \hat{\mathbf{U}}_k^{\text{mle}} + \frac{\hat{\sigma}_k^2}{n_k} \mathbf{I} = \hat{\mathbf{U}}_k^{\text{mle}} + \frac{R(R + 1)}{(n_k + R + 1) \text{Tr} \left((\hat{\mathbf{T}}_k^{\text{mle}})^{-1} \right)} \mathbf{I}\tag{12}$$

3 Lower bound of variance

We start with a simple one dimensional case. Assume $x_i \sim N(0, 1 + \sigma^2)$, our goal is to estimate σ^2 and possibly set $\hat{\sigma}^2$ to be the largest value of it consistent with data.

The log likelihood of $l(\sigma^2)$ is

$$l(\sigma^2) = -\frac{n}{2} \log(1 + \sigma^2) - \frac{ns^2}{2(1 + \sigma^2)},\tag{13}$$

where $s^2 = \sum_i x_i^2 / n$, the MLE of $(1 + \sigma^2)$.

Then we have

$$l(\sigma^2) - l(0) = -\frac{n}{2} \left(\log(1 + \sigma^2) - \frac{\sigma^2}{1 + \sigma^2} s^2 \right).\tag{14}$$

When $s^2 = 1$ (truncated) or equivalently $\hat{\sigma}^2 = 0$, and σ^2 is very small, using the approximation $\log(1 + \sigma^2) \approx \sigma^2$, we have

$$l(\sigma^2) - l(0) \approx -\frac{n}{2} \frac{(\sigma^2)^2}{1 + \sigma^2} \approx -\frac{n}{2} (\sigma^2)^2.\tag{15}$$

A rule of thumb is that 2 units drop in log likelihood gives 95% confidence intervals. For example, $l(\hat{\theta}) - 2 \approx l(\hat{\theta} + 2 * se(\hat{\theta}))$.

Solving for $l(\sigma^2) - l(0) = -2$, we have $\sigma^2 \approx \frac{2}{\sqrt{n}}$. This suggests that we can set the lower bound of σ^2 to be $\frac{2}{\sqrt{n}}$.

Now going back to the multivariate case. Let the eigendecomposition of \mathbf{U}_k be $\mathbf{Q}_k \mathbf{\Lambda}_k \mathbf{Q}_k^T$, and $\mathbf{U}_k + \mathbf{I} = \mathbf{Q}_k (\mathbf{\Lambda}_k + \mathbf{I}) \mathbf{Q}_k^T$. Then

$$\begin{aligned}\mathbf{x}_i &\sim N(\mathbf{0}, \mathbf{U}_k + \mathbf{I}), \\ \implies \mathbf{Q}_k^T \mathbf{x}_i &\sim N(\mathbf{0}, \mathbf{\Lambda}_k + \mathbf{I}).\end{aligned}\tag{16}$$

This suggests that we can set the eigenvalues of $\hat{\mathbf{U}}_k$ to be at least $\frac{2}{\sqrt{n_k}}$.