

Gaussian mixtures, factor models, independent factor analysis, and mixtures of factor analyzers are all different restrictive types of Gaussian mixtures with density

$$p(\mathbf{y}_n) = \sum_{q=1}^Q \bar{w}_q \mathcal{N}(\mathbf{y}_n | \bar{\boldsymbol{\mu}}_q, \bar{\boldsymbol{\Sigma}}_q) \quad (1)$$

The mean and second moment of a Gaussian mixture is given by

$$\mathbb{E}[\mathbf{y}_n] = \sum_{q=1}^Q \bar{w}_q \bar{\boldsymbol{\mu}}_q \quad (2)$$

$$\mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top] = \sum_{q=1}^Q \bar{w}_q (\bar{\boldsymbol{\Sigma}}_q + \bar{\boldsymbol{\mu}}_q \bar{\boldsymbol{\mu}}_q^\top) \quad (3)$$

The covariance then follows from the standard definition,

$$\text{Cov}(\mathbf{y}_n) = \mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top] - \mathbb{E}[\mathbf{y}_n] \mathbb{E}[\mathbf{y}_n]^\top. \quad (4)$$

In the above we implicitly conditioned on a number of hyperparameters. If we assume that they are all independent and that $\mathbb{E}[\bar{w}_q] = 1/Q$ and $\mathbb{E}[\bar{\boldsymbol{\mu}}_q] = \mathbf{0}$ and $\mathbb{E}[\bar{\boldsymbol{\mu}}_q \bar{\boldsymbol{\mu}}_q^\top] = \boldsymbol{\Lambda}_{\bar{\mu}}$, then

$$\mathbb{E}[\mathbf{y}_n] = \mathbf{0} \quad (5)$$

$$\mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top] = \text{Cov}(\mathbf{y}_n) = \boldsymbol{\Lambda}_{\bar{\mu}} + \frac{1}{Q} \sum_{q=1}^Q \mathbb{E}[\bar{\boldsymbol{\Sigma}}_q] \quad (6)$$

We can now further assume that the covariance matrix has random low-rank structure of the form,

$$\bar{\boldsymbol{\Sigma}}_q = \mathbf{A} \boldsymbol{\Sigma}_q \mathbf{A}^\top + \boldsymbol{\Lambda}. \quad (7)$$

which leads to the intractable second moment,

$$\mathbb{E}[\mathbf{y}_n \mathbf{y}_n^\top] = \text{Cov}(\mathbf{y}_n) = \boldsymbol{\Lambda}_{\bar{\mu}} + \boldsymbol{\Lambda} + \frac{1}{Q} \sum_{q=1}^Q \mathbb{E}[\mathbf{A} \boldsymbol{\Sigma}_q \mathbf{A}^\top]. \quad (8)$$

To render it tractable, we can calculate the trace of the matrix, which is simply

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = \text{tr} \left(\mathbf{A}^\top \mathbf{A} \left(\sum_{q=1}^Q \bar{w}_q \boldsymbol{\Sigma}_q \right) \right) + \text{tr}(\mathbb{E}[\boldsymbol{\Lambda}]) + \text{tr}(\boldsymbol{\Lambda}_{\bar{\mu}}). \quad (9)$$

The total covariance is also amenable to marginalization over the hyperparameters yielding a,

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = \text{tr} \left(\mathbb{E}[\mathbf{A}^\top \mathbf{A}] \left(\frac{1}{Q} \sum_{q=1}^Q \mathbb{E}[\boldsymbol{\Sigma}_q] \right) \right) + \text{tr}(\mathbb{E}[\boldsymbol{\Lambda}]) + \text{tr}(\boldsymbol{\Lambda}_{\bar{\mu}}). \quad (10)$$

If we assume that \mathbf{A} is a standard Gaussian ensemble with $[\mathbf{A}]_{ij} \sim \mathcal{N}(0, 1)$ then $\mathbf{A}^\top \mathbf{A}$ follows a Wishart distribution with mean $d\mathbf{I}$ where d is the dimensionality of \mathbf{y}_n . The total covariance then simplifies to

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = \frac{d}{Q} \sum_{q=1}^Q \text{tr}(\mathbb{E}[\mathbf{\Sigma}_q]) + \text{tr}(\mathbb{E}[\mathbf{\Lambda}]) + \text{tr}(\mathbf{\Lambda}_{\bar{\mu}}). \quad (11)$$

If we assume that $\mathbf{\Lambda}_{\bar{\mu}} = \sigma_0^2 \mathbf{I}$ and $\mathbb{E}[\mathbf{\Lambda}] = m_{\lambda^2} \mathbf{I}$

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = \frac{d}{Q} \sum_{q=1}^Q \text{tr}(\mathbb{E}[\mathbf{\Sigma}_q]) + d\sigma_0^2 + dm_{\lambda^2}. \quad (12)$$

1 Centered Independent Factor Analysis

For a centered independent factor analysis model, there is no mean component and $\bar{\mathbf{\Sigma}}_q$ is diagonal and has inverse Gamma distributions on its diagonal elements. Take the distribution to have mean m_{σ^2} . The resulting total covariance is

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = dKm_{\sigma^2} + dm_{\lambda^2}. \quad (13)$$

If we set the total covariance to be equal to d and let m_{λ^2} range freely, this imposes the constraint

$$m_{\sigma^2} = \frac{1 - m_{\lambda^2}}{K}, \quad (14)$$

essentially distributing the remaining variance evenly across the source dimensions. To tune the variances, see section 3.

To encourage non-Gaussian patterns, we can impose a different inverse Gamma prior on one (or more of the mixture components). If we let the $q = 1$ cluster have prior mean $\rho\sigma^2$ instead, with $\rho \in (0, 1]$, then the total covariance becomes

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = dKm_{\sigma^2} - dKC^{-1}(1 - \rho)m_{\sigma^2} + dm_{\lambda^2}. \quad (15)$$

with resulting solution,

$$m_{\sigma^2} = \frac{1 - m_{\lambda^2}}{K(1 - C^{-1}(1 - \rho))}. \quad (16)$$

2 Projected Mixture

We define projected mixtures to be mixtures on a low-rank space that are then projected up into the observation space using a factor matrix \mathbf{A} as with a factor model. This corresponds to the above model structure except for the fact that we need to take the mean to be low-rank $\bar{\boldsymbol{\mu}}_q = \mathbf{A}\bar{\boldsymbol{\mu}}'_q$ where $\bar{\boldsymbol{\mu}}'_q \sim \mathcal{N}(\mathbf{0}, \alpha_0 \mathbf{I})$ so that

$$\text{tr}(\text{Cov}(\bar{\boldsymbol{\mu}}_q)) = \sigma_0^2 \text{tr}(\mathbb{E}[\mathbf{A}^\top \mathbf{A}]) = dK\sigma_0^2. \quad (17)$$

If we simultaneously let $\bar{\Sigma}_q$ follow an inverse Wishart which has mean

$$\mathbb{E}[\bar{\Sigma}_q] = \frac{\rho}{\nu - K - 1} \mathbf{I} \quad (18)$$

for scale matrix $\rho \mathbf{I}$ and degrees of freedom ν we get the total covariance,

$$\text{tr}(\text{Cov}(\mathbf{y}_n)) = dK\sigma_0^2 + \frac{dK\alpha}{\nu - K - 1} + dm_{\lambda^2}. \quad (19)$$

To balance the terms, we match the variance contribution of the mean and the components,

$$\sigma_0^2 = \frac{\rho}{\nu - K - 1}, \quad (20)$$

which leads to the following constraint on α and ν ,

$$\frac{\rho}{\nu - K - 1} = \frac{1 - m_{\lambda^2}}{2K} \quad (21)$$

We note that $[\bar{\Sigma}_q]_{ii}$ is distributed as an inverse gamma,

$$[\bar{\Sigma}_q]_{ii} \sim \text{IG}\left(\frac{\nu - K + 1}{2}, \frac{\rho}{2}\right) \quad (22)$$

with mean $\frac{\rho}{\nu - K - 1}$, so we propose tuning an inverse gamma $\text{IG}(\alpha, \beta)$ to have mean $\frac{1 - m_{\lambda^2}}{2K}$ and appropriate variance following section 3, and then setting

$$\nu = K + 2\alpha - 1, \quad \rho = 2\beta. \quad (23)$$

3 Tuning the Inverse Gamma Variances

For an inverse Gamma distribution $\text{InvGamma}(\alpha, \beta)$ the mean and variance does not exist unless $\alpha > 2$, so we consider a random variable $X \sim \text{InvGamma}(2 + \alpha, \beta)$. X then has mean and variance equal to

$$\mathbb{E}[X] = \frac{\beta}{\alpha + 1} \quad (24)$$

$$\text{Var}(X) = \frac{\beta^2}{\alpha(\alpha + 1)^2} \quad (25)$$

If we set $\mathbb{E}[X] = m$ and $\text{Var}(X) = v$, then we can isolate the alpha and beta parameters in terms of the mean and variance as,

$$\beta = m \left(1 + \frac{m^2}{v}\right), \quad \alpha = \frac{m^2}{v} \quad (26)$$

Employing Markov's inequality, we can bound the tail probability using both the mean and the variance as

$$\mathbb{P}[X \geq a] \leq \frac{\mathbb{E}[X]}{a}, \quad \mathbb{P}[X \geq a] \leq \frac{\text{Var}(X)}{(a - \mathbb{E}[X])^2}. \quad (27)$$

Assume that the mean is set equal to m . Then if we want less than t mass in the tails, we can select the variance to be

$$\text{Var}(X) = t(a - m)^2 \quad (28)$$

For a noise variable, we are likely to want a low mean variance like $m = 10^{-1}$, and we likely want to contain (at least) $1 - t = 0.95$ of the probability mass to $[0, 1]$ requiring $a = 1$. This results in a proposed variance of

$$\text{Var}(X) = \frac{1}{20}(1 - 10^{-1})^2 = 0.0405 \quad (29)$$

this roughly holds if $\beta = 0.1247$ and $\alpha = 0.2469$ (or $\alpha = 2.2469$ in the original parameterization).

3.1 Tuning the Projected Mixture

For the projected mixture, we need to fix four effective parameters. We make the assumption that,

$$K\alpha_0 = \frac{K\alpha}{\nu - K - 1} \quad (30)$$

which corresponds to matching the total covariance contribution of the component means and the observations. Assuming the total covariance is d (corresponding to an identity covariance matrix) then we have,

$$1 - m_{\lambda^2} = \frac{2K\alpha}{\nu - K - 1}. \quad (31)$$

Using this, we can isolate α ,

$$\alpha = (\nu - K - 1) \frac{1 - m_{\lambda^2}}{2K}. \quad (32)$$

For the inverse Wishart the diagonal elements $\bar{\Sigma}_{ii}$ have variance,

$$\text{Var}(\bar{\Sigma}_{ii}) = \frac{2\alpha^2}{(\nu - K - 1)^2(\nu - K - 3)}. \quad (33)$$

If we fix the variance, we can insert the above value for α and isolate ν as

$$\nu = K + 3 + 2 \frac{\left(\frac{1 - m_{\lambda^2}}{2K}\right)^2}{\text{Var}(\bar{\Sigma}_{ii})} \quad (34)$$

Adding it back in we get the following formula for α ,

$$\alpha = \left(1 + \frac{\left(\frac{1 - m_{\lambda^2}}{2K}\right)^2}{\text{Var}(\bar{\Sigma}_{ii})}\right) \frac{1 - m_{\lambda^2}}{K} \quad (35)$$

With Wishart priors, we start with

$$\alpha = \frac{1 - m_{\lambda^2}}{2K\nu}, \quad (36)$$

and using the variance expression,

$$\text{Var}(\bar{\Sigma}_{ii}) = 2\nu\alpha^2, \quad (37)$$

we can isolate ν ,

$$\nu = \frac{(1 - m_{\lambda^2})^2}{2K^2 \text{Var}(\bar{\Sigma}_{ii})} \quad (38)$$

and then find α ,

$$\alpha = \frac{K \text{Var}(\bar{\Sigma}_{ii})}{1 - m_{\lambda^2}}. \quad (39)$$

3.1.1 Alternative Calibration of the Inverse Wishart

Alternatively we can use that if $\bar{\Sigma} \sim \mathcal{W}^{-1}(\alpha \mathbf{I}, \nu)$ then,

$$\bar{\Sigma}_{ii} \sim \Gamma^{-1}\left(\frac{\alpha}{2}, \frac{\nu - K - 1}{2}\right) \quad (40)$$

which means we can calibrate the Wishart as an Inverse Gamma. We should then set

$$\alpha_0 = \frac{\alpha}{\nu - K - 1} \quad (41)$$

to balance

3.2 Low-rank Gaussian density

The Gaussian density function is given by,

$$\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (42)$$

If Σ is low-rank,

$$\Sigma = \mathbf{A} \mathbf{S} \mathbf{A}^\top + \mathbf{D} \quad (43)$$

then we can employ tricks to calculate the inverse more efficiently. The Woodbury matrix identity yields,

$$(\mathbf{A} \mathbf{S} \mathbf{A}^\top + \mathbf{D})^{-1} = \mathbf{D}^{-1} - \mathbf{D}^{-1} \mathbf{A} (\mathbf{S}^{-1} + \mathbf{A}^\top \mathbf{D}^{-1} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{D}^{-1}, \quad (44)$$

which can be simplified further if the Cholesky factorization $\mathbf{L} \mathbf{L}^\top = \mathbf{S}$ is known, as we can pull the factors outside,

$$(\mathbf{A} \mathbf{S} \mathbf{A}^\top + \mathbf{D})^{-1} = \mathbf{D}^{-1} - \mathbf{D}^{-1} \mathbf{A} \mathbf{L} (\mathbf{I} + \mathbf{L}^\top \mathbf{A}^\top \mathbf{D}^{-1} \mathbf{A} \mathbf{L})^{-1} \mathbf{L}^\top \mathbf{A}^\top \mathbf{D}^{-1}. \quad (45)$$

If we define $\mathbf{B} = \mathbf{D}^{-1}\mathbf{A}\mathbf{L}$ then this can be written fairly succinctly as,

$$(\mathbf{A}\mathbf{S}\mathbf{A}^\top + \mathbf{D})^{-1} = \mathbf{D}^{-1} - \mathbf{B}(\mathbf{I} + \mathbf{B}^\top \mathbf{D}\mathbf{B})^{-1}\mathbf{B}^\top. \quad (46)$$

which only requires inversion of diagonal matrices and matrices of the same shape as \mathbf{S} , which can be smaller than $\mathbf{\Sigma}$ by design.

We can similarly apply the matrix determinant lemma to calculate the determinant as

$$|\mathbf{A}\mathbf{S}\mathbf{A}^\top + \mathbf{D}| = |\mathbf{D}||\mathbf{I} + \mathbf{L}^\top \mathbf{A}^\top \mathbf{D}^{-1}\mathbf{A}\mathbf{L}| = |\mathbf{D}||\mathbf{I} + \mathbf{B}^\top \mathbf{D}\mathbf{B}| \quad (47)$$