

High-Dimensional Regression with Gaussian Mixtures and Partially-Latent Response Variables (Supplementary Materials)

Antoine Deleforge · Florence Forbes · Radu Horaud

This document contains supplementary materials for the paper:

A. Deleforge, F. Forbes, and R. Horaud. High-Dimensional Regression with Gaussian Mixtures and Partially-Latent Response Variables. *Statistics and Computing*. 2014.

1 The general hybrid GLLiM-EM algorithm

Considering the complete data, with $(\mathbf{Y}, \mathbf{T})_{1:N}$ being the observed variables and $(Z, \mathbf{W})_{1:N}$ being the missing ones, the corresponding EM algorithm consists of estimating the parameter vector $\boldsymbol{\theta}^{(i+1)}$ that maximizes the expected complete-data log-likelihood, given the current parameter vector $\boldsymbol{\theta}^{(i)}$ and the observed data:

$$\boldsymbol{\theta}^{(i+1)} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}[\log p((\mathbf{y}, \mathbf{t}, \mathbf{W}, Z)_{1:N}; \boldsymbol{\theta}) | (\mathbf{y}, \mathbf{t})_{1:N}; \boldsymbol{\theta}^{(i)}]. \quad (1)$$

Using that $\mathbf{W}_{1:N}$ and $\mathbf{T}_{1:N}$ are independent conditionally on $Z_{1:N}$ and that $\{\mathbf{c}_k^w\}_{k=1}^K$ and $\{\mathbf{I}_k^w\}_{k=1}^K$ are fixed, maximizing (1) is then equivalent to maximizing the following expression:

$$\mathbb{E}_{r_Z^{(i+1)}} [\mathbb{E}_{r_{W|Z}^{(i+1)}} [\log p(\mathbf{y}_{1:N} | (\mathbf{t}, \mathbf{W}, Z)_{1:N}; \boldsymbol{\theta})] + \log p((\mathbf{t}, Z)_{1:N}; \boldsymbol{\theta})] \quad (2)$$

where $r_Z^{(i+1)}$ and $r_{W|Z}^{(i+1)}$ denote the posterior distributions

$$r_Z^{(i+1)} = p(\mathbf{Z}_{1:N} | (\mathbf{y}, \mathbf{t})_{1:N}; \boldsymbol{\theta}^{(i)}) \quad (3)$$

$$r_{W|Z}^{(i+1)} = p(\mathbf{W}_{1:N} | (\mathbf{y}, \mathbf{t}, Z)_{1:N}; \boldsymbol{\theta}^{(i)}). \quad (4)$$

It follows that the E-step splits into an **E-W** step and an **E-Z** step. The subsequent M-step can also be divided into two steps referred to as the **M-GMM** step and the **M-mapping** step. In what follows details are given for the **E-W** step and the **M-mapping** step as the **E-Z** and **M-GMM** steps are straightforward as explained in the main paper (equations (27) to (30). For the sake of readability, the current iteration superscript $(i+1)$ is replaced with a tilde. Hence, $\boldsymbol{\theta}^{(i+1)} = \tilde{\boldsymbol{\theta}}$ (the model parameter vector).

1.1 E-W-step

This step consists of identifying the conditional probability distribution $\tilde{r}_{W|Z}$ defined in (4). Given parameter estimates, it is fully defined by computing the distribution $p(\mathbf{w}_n|Z_n = k, \mathbf{t}_n, \mathbf{y}_n; \boldsymbol{\theta}^{(i)})$, denoted below by \tilde{r}_{nk}^w , for all n and all k . This density as a function of \mathbf{w}_n is proportional to

$$\tilde{r}_{nk}^w = p(\mathbf{w}_n|Z_n = k, \mathbf{t}_n, \mathbf{y}_n; \boldsymbol{\theta}^{(i)}) \propto p(\mathbf{y}_n|\mathbf{t}_n, \mathbf{w}_n, Z_n = k; \boldsymbol{\theta}^{(i)}) p(\mathbf{w}_n|Z_n = k; \boldsymbol{\theta}^{(i)}). \quad (5)$$

The second term in the right-hand side is equal to $\mathcal{N}(\mathbf{w}_n; \mathbf{c}_k^w, \boldsymbol{\Gamma}_k^w)$ by virtue of definitions (4) and (17) in the main paper, while from (18)

$$p(\mathbf{y}_n|\mathbf{t}_n, \mathbf{w}_n, Z_n = k; \boldsymbol{\theta}^{(i)}) = \mathcal{N}(\mathbf{y}_n; \mathbf{A}_k^{(i)}[\mathbf{t}_n; \mathbf{w}_n] + \mathbf{b}_k^{(i)}, \boldsymbol{\Sigma}_k^{(i)}). \quad (6)$$

Rewriting the expression in the exponential term of the right-hand side as a quadratic form in \mathbf{w}_n , it is easy to see that $p(\mathbf{y}_n|\mathbf{t}_n, \mathbf{w}_n, Z_n = k; \boldsymbol{\theta}^{(i)})$ induces a term over \mathbf{w}_n proportional to a Gaussian density,

$$\mathcal{N}(\mathbf{w}_n; \boldsymbol{\mu}_{nk}^{l(i)}, \boldsymbol{\Sigma}_k^{l(i)}),$$

with $\boldsymbol{\mu}_{nk}^{l(i)} = \boldsymbol{\Sigma}_k^{l(i)} \mathbf{A}_k^{w(i)\top} (\boldsymbol{\Sigma}_k^{(i)})^{-1} (\mathbf{y}_n - \mathbf{A}_k^{t(i)} \mathbf{t}_n - \mathbf{b}_k^{(i)})$, and $\boldsymbol{\Sigma}_k^{l(i)} = (\mathbf{A}_k^{w(i)\top} (\boldsymbol{\Sigma}_k^{(i)})^{-1} \mathbf{A}_k^{w(i)})^{-1}$.

It follows that \tilde{r}_{nk}^w is proportional to the product of two Gaussian distributions over \mathbf{w}_n which according to Lemma 1 in Section 3.1, implies that \tilde{r}_{nk}^w is Gaussian, with mean $\tilde{\boldsymbol{\mu}}_{nk}^w$ and covariance matrix $\tilde{\mathbf{S}}_k^w$ given by:

$$\begin{aligned} \tilde{\boldsymbol{\mu}}_{nk}^w &= \tilde{\mathbf{S}}_k^w (\mathbf{A}_k^{w(i)\top} (\boldsymbol{\Sigma}_k^{(i)})^{-1} (\mathbf{y}_n - \mathbf{A}_k^{t(i)} \mathbf{t}_n - \mathbf{b}_k^{(i)}) + (\boldsymbol{\Gamma}_k^{w(i)})^{-1} \mathbf{c}_k^{w(i)}) \\ \tilde{\mathbf{S}}_k^w &= ((\boldsymbol{\Gamma}_k^{w(i)})^{-1} + \mathbf{A}_k^{w(i)\top} (\boldsymbol{\Sigma}_k^{(i)})^{-1} \mathbf{A}_k^{w(i)})^{-1} \end{aligned} \quad (7)$$

1.2 M-mapping-step

This step consists of updating parameters $\{\mathbf{A}_k, \mathbf{b}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$. Omitting the terms in (2) that do not depend on $\{\mathbf{A}_k, \mathbf{b}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$, the M-mapping-step is equivalent to maximizing

$$\mathbb{E}_{r_Z^{(i+1)}} [\mathbb{E}_{r_{W|Z}^{(i+1)}} [\log p(\mathbf{y}_{1:N} | (\mathbf{t}, \mathbf{W}, Z)_{1:N}; \boldsymbol{\theta})]]. \quad (8)$$

Using $\tilde{r}_{nk} = p(Z_n = k|\mathbf{t}_n, \mathbf{y}_n; \boldsymbol{\theta}^{(i)})$ and $\tilde{r}_{nk}^w = p(\mathbf{w}_n|Z_n = k, \mathbf{t}_n, \mathbf{y}_n; \boldsymbol{\theta}^{(i)})$, this writes

$$\sum_{n=1}^N \sum_{k=1}^K \tilde{r}_{nk} \mathbb{E}_{\tilde{r}_{nk}^w} [\log p(\mathbf{y}_n | \mathbf{t}_n, \mathbf{W}_n, Z_n = k; \boldsymbol{\theta})]. \quad (9)$$

Using again the fact that $p(\mathbf{y}_n|\mathbf{t}_n, \mathbf{w}_n, Z_n = k; \boldsymbol{\theta})$ induces an expression over \mathbf{w}_n proportional to $\mathcal{N}(\mathbf{w}_n; \boldsymbol{\mu}_{nk}^l, \boldsymbol{\Sigma}_k^l)$, with $\boldsymbol{\mu}_{nk}^l = \boldsymbol{\Sigma}_k^l (\mathbf{A}_k^w)^\top (\boldsymbol{\Sigma}_k^{(i)})^{-1} (\mathbf{y}_n - \mathbf{A}_k^t \mathbf{t}_n - \mathbf{b}_k)$, and $\boldsymbol{\Sigma}_k^l = (\mathbf{A}_k^{w\top} (\boldsymbol{\Sigma}_k^{(i)})^{-1} \mathbf{A}_k^w)^{-1}$, it follows that

$$\log(\mathbf{y}_n|\mathbf{t}_n, \mathbf{w}_n, Z_n = k; \boldsymbol{\theta}) = \log C - \frac{1}{2} (\mathbf{w}_n - \boldsymbol{\mu}_{nk}^l)^T (\boldsymbol{\Sigma}_k^l)^{-1} (\mathbf{w}_n - \boldsymbol{\mu}_{nk}^l),$$

where C is a term that does not depend on \mathbf{w}_n . Then, using Lemma 2 in Section 3.2,

$$\begin{aligned} \mathbb{E}_{\tilde{r}_{nk}^w} [\log p(\mathbf{y}_n | \mathbf{t}_n, \mathbf{W}_n, Z_n = k; \boldsymbol{\theta})] &= \log C - \frac{1}{2} (\tilde{\boldsymbol{\mu}}_{nk}^w - \boldsymbol{\mu}_{nk}^l)^\top (\boldsymbol{\Sigma}_k^l)^{-1} (\tilde{\boldsymbol{\mu}}_{nk}^w - \boldsymbol{\mu}_{nk}^l) - \frac{1}{2} (\tilde{\mathbf{S}}_k^w \boldsymbol{\Sigma}_k^l) \\ &= \log p(\mathbf{y}_n | \mathbf{t}_n, \mathbf{W}_n = \tilde{\boldsymbol{\mu}}_{nk}^w, Z_n = k; \boldsymbol{\theta}) - \frac{1}{2} \text{tr}(\tilde{\mathbf{S}}_k^w \boldsymbol{\Sigma}_k^l). \end{aligned} \quad (10)$$

For the update of \mathbf{b}_k , it follows that the term to be maximized is only the first one and is the same as in a standard GMM where \mathbf{b}_k would be the mean of component k for observed data that would be $\mathbf{y}_n - \tilde{\mathbf{A}}_k \tilde{\mathbf{x}}_{nk}$, with $\tilde{\mathbf{x}}_{nk} = [\mathbf{t}_n; \tilde{\boldsymbol{\mu}}_{nk}^w] \in \mathbb{R}^L$. It comes the standard weighted mean formula,

$$\tilde{\mathbf{b}}_k = \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{A}}_k \tilde{\mathbf{x}}_{nk}). \quad (11)$$

For updating $\boldsymbol{\Sigma}_k$, using (10) and (6), as an expression over $\boldsymbol{\Sigma}_k$, the term in (9) to be maximized writes,

$$-\frac{1}{2} \sum_{n=1}^N \tilde{r}_{nk} (\log |\boldsymbol{\Sigma}_k| + \text{tr}(\mathbf{A}_k^w \tilde{\mathbf{S}}_k^w (\mathbf{A}_k^w)^\top \boldsymbol{\Sigma}_k^{-1}) + (\mathbf{y}_n - \mathbf{A}_k \tilde{\mathbf{x}}_{nk} - \mathbf{b}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{y}_n - \mathbf{A}_k \tilde{\mathbf{x}}_{nk} - \mathbf{b}_k)).$$

Denoting $\tilde{r}_k = \sum_{n=1}^N \tilde{r}_{nk}$, this is equivalent to minimize over $\boldsymbol{\Sigma}_k$

$$\tilde{r}_k \log |\boldsymbol{\Sigma}_k| + \text{tr}((\sum_{n=1}^N \tilde{r}_{nk} \mathbf{M}_{nk}) \boldsymbol{\Sigma}_k^{-1}),$$

where $\mathbf{M}_{nk} = \mathbf{A}_k^w \tilde{\mathbf{S}}_k^w (\mathbf{A}_k^w)^\top + (\mathbf{y}_n - \mathbf{A}_k \tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{b}}_k)(\mathbf{y}_n - \mathbf{A}_k \tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{b}}_k)^\top$ and $\sum_{n=1}^N \tilde{r}_{nk} \mathbf{M}_{nk}$ are positive definite matrices.

Using the result recalled in Lemma 3 in Section 3.3, it follows straightforwardly that,

$$\tilde{\boldsymbol{\Sigma}}_k = \tilde{\mathbf{A}}_k^w \tilde{\mathbf{S}}_k^w \tilde{\mathbf{A}}_k^{w\top} + \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{A}}_k \tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{b}}_k)(\mathbf{y}_n - \tilde{\mathbf{A}}_k \tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{b}}_k)^\top. \quad (12)$$

As regards \mathbf{A}_k , it comes from (9) and (10), omitting the terms independent of \mathbf{A}_k , that we have to minimize the following expression over \mathbf{A}_k ,

$$\text{tr}((\sum_{n=1}^N \tilde{r}_{nk} \tilde{\mathbf{M}}_{nk}) \boldsymbol{\Sigma}_k^{-1}), \quad (13)$$

where $\tilde{\mathbf{M}}_{nk}$ is the part of \mathbf{M}_{nk} that involves \mathbf{A}_k ,

$$\begin{aligned} \tilde{\mathbf{M}}_{nk} &= \mathbf{A}_k^w \tilde{\mathbf{S}}_k^w (\mathbf{A}_k^w)^\top + (\mathbf{A}_k \tilde{\mathbf{x}}_{nk})(\mathbf{A}_k \tilde{\mathbf{x}}_{nk})^\top - 2(\mathbf{y}_n - \mathbf{b}_k)(\mathbf{A}_k \tilde{\mathbf{x}}_{nk})^\top \\ &= \mathbf{A}_k \tilde{\mathbf{S}}_k^x \mathbf{A}_k^\top + \mathbf{A}_k \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top \mathbf{A}_k^\top - 2(\mathbf{y}_n - \mathbf{b}_k) \tilde{\mathbf{x}}_{nk}^\top \mathbf{A}_k^\top, \end{aligned}$$

where $\tilde{\mathbf{S}}_k^x$ is defined by

$$\tilde{\mathbf{S}}_k^x = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{S}}_k^w \end{bmatrix}, \quad (14)$$

so that $\mathbf{A}_k^w \tilde{\mathbf{S}}_k^w (\mathbf{A}_k^w)^\top = \mathbf{A}_k \tilde{\mathbf{S}}_k^x \mathbf{A}_k^\top$. It follows that,

$$\sum_{n=1}^N \tilde{r}_{nk} \tilde{\mathbf{M}}_{nk} = \tilde{r}_k \mathbf{A}_k \tilde{\mathbf{S}}_k^x \mathbf{A}_k^\top + \mathbf{A}_k (\sum_{n=1}^N \tilde{r}_{nk} \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top) \mathbf{A}_k^\top - 2(\sum_{n=1}^N \tilde{r}_{nk} (\mathbf{y}_n - \mathbf{b}_k) \tilde{\mathbf{x}}_{nk}^\top) \mathbf{A}_k^\top \quad (15)$$

and

$$\frac{\sum_{n=1}^N \tilde{r}_{nk} \tilde{\mathbf{M}}_{nk}}{\tilde{r}_k} = \mathbf{A}_k (\tilde{\mathbf{S}}_k^x + \sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top) \mathbf{A}_k^\top - 2 \left(\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\mathbf{y}_n - \mathbf{b}_k) \tilde{\mathbf{x}}_{nk}^\top \right) \mathbf{A}_k^\top \quad (16)$$

Deriving (see derivatives in Section 3.4) the trace in (13) with regard to \mathbf{A}_k leads to

$$2(\tilde{\mathbf{S}}_k^x + \sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top) \mathbf{A}_k^\top \boldsymbol{\Sigma}_k^{-1} - 2 \left(\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\mathbf{y}_n - \mathbf{b}_k) \tilde{\mathbf{x}}_{nk}^\top \right) \boldsymbol{\Sigma}_k^{-1}$$

Equating this to 0, leads to

$$\tilde{\mathbf{A}}_k = \left(\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{b}}_k) \tilde{\mathbf{x}}_{nk}^\top \right) \left(\tilde{\mathbf{S}}_k^x + \sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top \right)^{-1}$$

where $\tilde{\mathbf{b}}_k$ satisfies equation (11). Replacing $\tilde{\mathbf{b}}_k$ in the expression above further leads to a closed-form solution for \mathbf{A}_k ,

$$\tilde{\mathbf{A}}_k = \left(\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{y}}_k) \tilde{\mathbf{x}}_{nk}^\top \right) \left(\tilde{\mathbf{S}}_k^x + \sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top - \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top \right)^{-1}$$

where

$$\tilde{\mathbf{x}}_k = \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk}, \quad (17)$$

$$\tilde{\mathbf{y}}_k = \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} \mathbf{y}_n. \quad (18)$$

which can be also written,

$$\tilde{\mathbf{A}}_k = \left(\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{y}}_k) (\tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{x}}_k)^\top \right) \left(\tilde{\mathbf{S}}_k^x + \sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{x}}_k) (\tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{x}}_k)^\top \right)^{-1},$$

since

$$\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{y}}_k) \tilde{\mathbf{x}}_k^\top = 0$$

and

$$\sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} (\tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{x}}_k) (\tilde{\mathbf{x}}_{nk} - \tilde{\mathbf{x}}_k)^\top = \sum_{n=1}^N \frac{\tilde{r}_{nk}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk} \tilde{\mathbf{x}}_{nk}^\top - \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top.$$

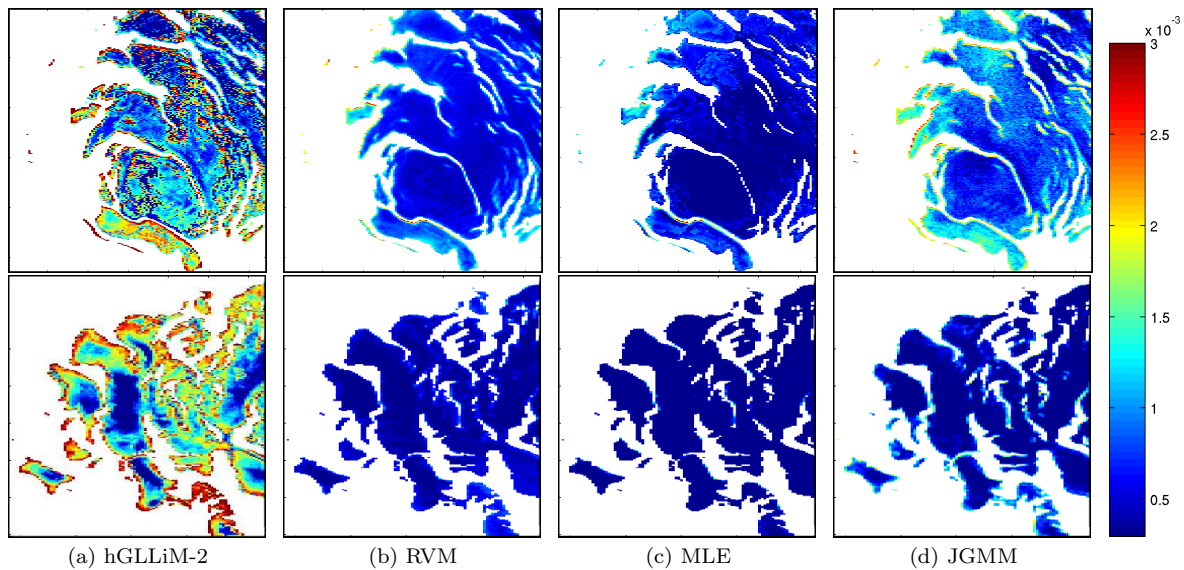


Fig. 1 Proportion of dust obtained with our method (hGLLiM-2) and with three other methods. The data correspond to hyperspectral images grabbed from two different viewpoints of the South polar cap of Mars. First row: orbit 41, second row: orbit 61. White areas correspond to unexamined regions, where the synthetic model does not apply.

2 Retrieval of Mars surface physical properties from hyperspectral images

Visible and near infrared imaging spectroscopy is a key remote sensing technique used to study and monitor planets. It records the visible and infrared light reflected from the planet in a given wavelength range and produces cubes of data where each observed surface location is associated with a spectrum. Physical properties of the planets' surface, such as chemical composition, granularity, texture, etc, are some of the most important parameters that characterize the morphology of spectra. In the case of Mars, radiative transfer models have been developed to numerically evaluate the link between these parameters and observable spectra. Such models allow to simulate spectra from a given set of parameter values. In practice, the goal is to scan the Mars ground from an orbit in order to observe gas and dust in the atmosphere and look for signs of specific materials such as silicates, carbonates and ice at the surface. We are thus interested in solving the associate inverse problem which is to deduce physical parameter values from the observed spectra.

The Hybrid GLLiM model with a 2-dimensional latent variable (hGLLiM-2) and three other methods, namely RVM, MLE and JGMM (see Section 6.3 in the main paper for details) were used to retrieve the physical properties of the South polar cap using two hyperspectral images of approximately the same area from different view points (orbit 41 and orbit 61). Since we are looking for proportions between 0 and 1, values smaller than 0 or higher than 1 are not acceptable and hence they were set to one of the bounds. As it can be seen in Fig. 1 and Fig. 2, hybrid GLLiM outputs proportion maps with similar characteristics for the two view points, which suggests good consistency. Such a consistency is not observed using the other tested methods. In addition, RVM and MLE output a much higher number of values falling outside the interval $[0, 1]$. Moreover, hGLLiM-2 is the only method featuring less dust at the South pole cap center and higher concentrations of dust at the boundaries of the CO₂ ice, which matches expected results from planetology [Douté et al., 2005]. Finally, note that the proportions of CO₂ ice and dust clearly seem to be complementary using hGLLiM-2, while this complementarity is less obvious using other methods.

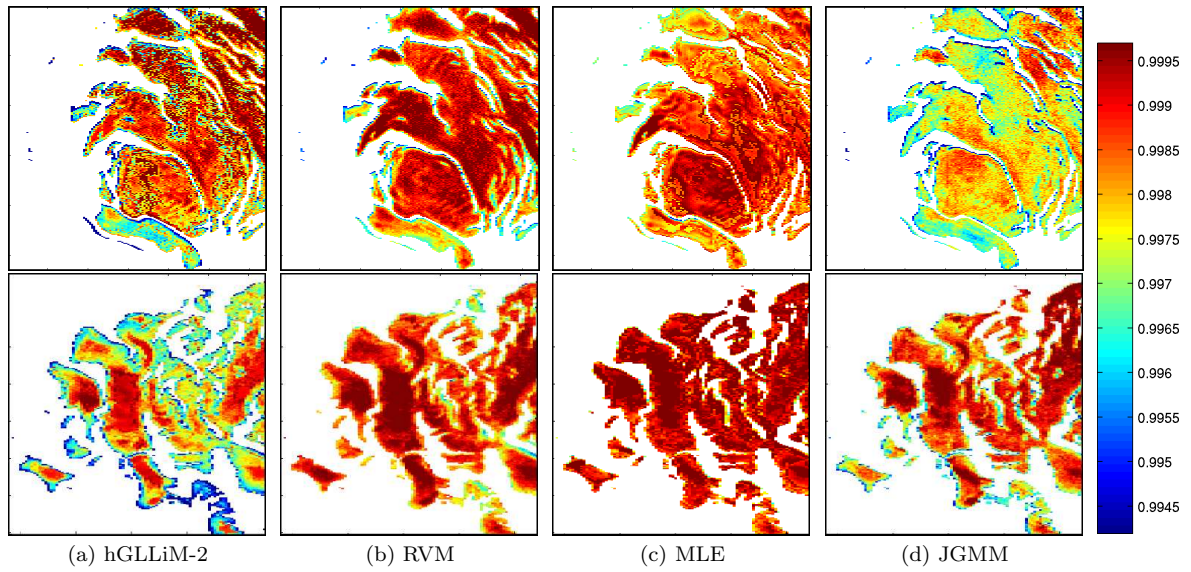


Fig. 2 Proportion of CO₂ ice obtained from hyperspectral images of two different viewpoints of the South polar cap of Mars. First row: orbit 41, second row: orbit 61. White areas correspond to unexamined regions, where the synthetic model does not apply.

3 Some properties of Gaussian distributions and useful results

3.1 Product of two multivariate Gaussian densities

Lemma 1 *The product of two multivariate Gaussian densities $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ is proportional to the multivariate Gaussian density with mean $\boldsymbol{\mu}_3$ and covariance matrix $\boldsymbol{\Sigma}_3$ given by:*

$$\begin{aligned} \boldsymbol{\mu}_3 &= \boldsymbol{\Sigma}_3 (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) \\ \boldsymbol{\Sigma}_3 &= (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1} \end{aligned} \quad (19)$$

Proof: The result follows easily by developing the terms in the sum of the two Gaussians quadratic forms and rearranging them into a single quadratic form over the variable \mathbf{x} . Note however that the product of two Gaussian densities is not normalized.

3.2 Mean of square forms

Lemma 2 *If \mathbf{X} is a variable with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, then,*

$$E[\mathbf{X}^\top \mathbf{M} \mathbf{X}] = \boldsymbol{\mu}^\top \mathbf{M} \boldsymbol{\mu} + \text{tr}(\mathbf{M} \boldsymbol{\Sigma})$$

Proof: The proof is omitted as this is a very standard result. See for instance [Schott, 1997] p.391

3.3 Maximum likelihood for a Gaussian sample

As shown in [Schott, 1997] p.347, when estimating using maximum likelihood principle, the mean and variance from a normally distributed sample, it comes a maximization problem similar to that in the following lemma.

Lemma 3 *The $(m \times m)$ matrix \mathbf{A} minimizing $\text{tr}(\mathbf{M}\mathbf{A}^{-1}) + \alpha \log |\mathbf{A}|$ where \mathbf{M} is a $m \times m$ symmetric definite positive matrix and α is a positive real number is $\mathbf{A} = \frac{\mathbf{M}}{\alpha}$.*

Proof: See [Schott, 1997] p.347 for a similar proof.

3.4 Derivatives of traces

It can be found for instance in [Schott, 1997] the following results regarding the derivatives of the following two expressions.

Lemma 4

$$\begin{aligned}\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{M}\mathbf{A}^\top) &= \mathbf{M} \\ \frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{A}\mathbf{B}\mathbf{A}^\top \mathbf{C}) &= \mathbf{C}\mathbf{A}\mathbf{B} + \mathbf{C}^\top \mathbf{A}\mathbf{B}^\top\end{aligned}$$

References

- S. Douté, B. Schmitt, J-P. Bibring, Y. Langevin, F. Altieri, G. Bellucci, B. Gondet, and the MEX OMEGA team. Nature and composition of the icy terrains of the south pole of Mars from MEX OMEGA observations. In *The 36th Lunar and Planetary Science Conference, (Lunar and Planetary Science XXXVI)*, March 2005.
- J. R. Schott. *Matrix Analysis for Statistics*. Wiley series in probability and statistics, 1997.