

High dimensional regression with Gaussian mixtures and partially latent response variables

Antoine Deleforge, Florence Forbes, and Radu Horaud

Perception and Mistis teams
INRIA Grenoble Rhone-Alpes
first.last@inria.fr

Statistics and Computing, Springer, 2014
team.inria.fr/perception/research/high-dim-regression/

Dealing with high dimensional data

Find f between $\mathbf{y} \in \mathbb{R}^D$ and $\mathbf{x} \in \mathbb{R}^L$ with $D \gg L$

$$f : \mathbf{y} \in \mathbb{R}^D \longrightarrow \mathbf{x} \in \mathbb{R}^L$$

from a learning sample $\{(\mathbf{y}_n, \mathbf{x}_n), n = 1, \dots, N\}$

Difficulty : D large \implies curse of dimensionality

Solutions : via dimensionality reduction

- Reduce dimension of \mathbf{y} before regression: eg. PCA on the \mathbf{y}_n 's first

Risk: poor prediction of \mathbf{x}

- Take \mathbf{x} into account: PLS [Rosipal et al 06], SIR [Li 91], Kernel SIR [Wu 08], PC based methods [Cook 07, Adragani & Cook 09], etc.

\implies **two steps approaches not expressed as a single optimization problem**

Proposed: Inverse regression then forward prediction

Standard regression setting: Fully Observed Input and Output Variables

Learning (regression)

x

y

$$\begin{pmatrix} \bullet \\ \bullet \end{pmatrix} \Rightarrow \begin{pmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{pmatrix}$$

Testing (prediction)

y

x

$$\begin{pmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{pmatrix} \Rightarrow \begin{pmatrix} ? \\ ? \end{pmatrix}$$

Proposed Method: An inverse regression strategy

- $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^L$ low-dimensional space,
- $\mathbf{Y} \in \mathcal{Y} \subset \mathbb{R}^D$ high-dimensional space,
- (\mathbf{y}, \mathbf{x}) realization of $(\mathbf{Y}, \mathbf{X}) \sim p(\mathbf{Y}, \mathbf{X}; \boldsymbol{\theta})$, $\boldsymbol{\theta}$ parameters

- Inverse conditional density: $p(\mathbf{Y} | \mathbf{X}; \boldsymbol{\theta})$
 \mathbf{Y} is a noisy function of \mathbf{X}

Modeled via mixtures \rightarrow tractable $\boldsymbol{\theta}$ estimation

- Forward conditional density: $p(\mathbf{X} | \mathbf{Y}; \boldsymbol{\theta}^*)$, with $\boldsymbol{\theta}^* = g(\boldsymbol{\theta})$
 \rightarrow high-to-low prediction, eg. $\hat{\mathbf{x}} = E[\mathbf{X} | \mathbf{Y} = \mathbf{y}; \boldsymbol{\theta}^*]$

Gaussian Locally-linear Mapping (GLLiM)

- $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^L$ low-dimensional space,
- $\mathbf{Y} \in \mathcal{Y} \subset \mathbb{R}^D$ high-dimensional space,
- A piecewise affine model: Introduce a missing variable Z

$$p(\mathbf{Y} = \mathbf{y}, \mathbf{X} = \mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^K p(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}, Z = k; \boldsymbol{\theta}) p(\mathbf{X} = \mathbf{x} | Z = k; \boldsymbol{\theta}) p(Z = k; \boldsymbol{\theta})$$

$Z = k \Leftrightarrow \mathbf{Y}$ is the image of \mathbf{X} by an affine transformation τ_k

Hierarchical definition

$$\mathbf{Y} = \sum_{k=1}^K \mathbb{I}(Z = k)(\mathbf{A}_k \mathbf{X} + \mathbf{b}_k + \mathbf{E}_k)$$

\mathbb{I} Indicator function, \mathbf{A}_k $D \times L$ matrix, \mathbf{b}_k D -dim vector

\mathbf{E}_k : observation noise in \mathbb{R}^D and reconstruction error, Gaussian, centered, independent on \mathbf{X} , \mathbf{Y} , and Z

$$p(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}, Z = k; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}; \mathbf{A}_k \mathbf{x} + \mathbf{b}_k, \boldsymbol{\Sigma}_k)$$

- Affine transformations are local: mixture of K Gaussians

$$\begin{aligned} p(\mathbf{X} = \mathbf{x} | Z = k; \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{x}; \mathbf{c}_k, \boldsymbol{\Gamma}_k) \\ p(Z = k; \boldsymbol{\theta}) &= \pi_k \end{aligned}$$

- The set of all model parameters is:

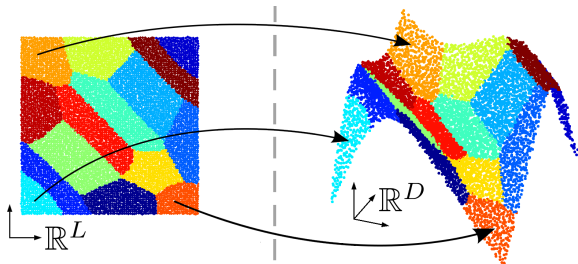
$$\boldsymbol{\theta} = \{\mathbf{c}_k, \boldsymbol{\Gamma}_k, \pi_k, \mathbf{A}_k, \mathbf{b}_k, \boldsymbol{\Sigma}_k, k = 1 \dots K\}$$

Usually $\{\boldsymbol{\Sigma}_k = \sigma \mathbf{I}_D, k = 1 \dots K\}$ (isotropic reconstruction error)

Geometric Interpretation

This model induces a **partition of \mathbb{R}^L** into K regions \mathcal{R}_k where the transformation τ_k is the most probable.

If $|\Gamma_1| = \dots = |\Gamma_K|$: $\{\mathcal{R}_k, k = 1 \dots K\}$ define a Voronoi diagram of centroids $\{\mathbf{c}_k, k = 1 \dots K\}$ (Mahalanobis distance $\|\cdot\|_{\Gamma}$).



$$L = 2, D = 3, K = 15.$$

Low-to-high (Inverse) Regression

If \mathbf{X} and \mathbf{Y} are both observed

- The parameter vector, $\boldsymbol{\theta}$, can be estimated in closed-form using an **EM inference procedure**
- This yields the *inverse conditional density* which is a Gaussian mixture:

$$p(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^K \frac{\pi_k \mathcal{N}(\mathbf{x}; \mathbf{c}_k, \boldsymbol{\Gamma}_k)}{\underbrace{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}; \mathbf{c}_j, \boldsymbol{\Gamma}_j)}_{\beta_k}} \mathcal{N}(\mathbf{y}; \underbrace{\mathbf{A}_k \mathbf{x} + \mathbf{b}_k}_{\boldsymbol{\mu}_k}, \boldsymbol{\Sigma}_k)$$

High-to-low (Forward) Regression

- The forward parameter vector, θ^* , has an analytic expression as a **function of θ**
- This yields the *forward conditional density* which is a Gaussian mixture as well:

$$p(\mathbf{X} = \mathbf{x} | \mathbf{Y} = \mathbf{y}; \theta^*) = \sum_{k=1}^K \frac{\pi_k^* \mathcal{N}(\mathbf{y}; \mathbf{c}_k^*, \mathbf{\Gamma}_k^*)}{\underbrace{\sum_{j=1}^K \pi_j^* \mathcal{N}(\mathbf{y}; \mathbf{c}_j^*, \mathbf{\Gamma}_j^*)}_{\beta_k^*}} \mathcal{N}(\mathbf{x}; \underbrace{\mathbf{A}_k^* \mathbf{y} + \mathbf{b}_k^*}_{\boldsymbol{\mu}_k^*}, \boldsymbol{\Sigma}_k^*)$$

The forward parameter vector $\boldsymbol{\theta}^*$ from $\boldsymbol{\theta}$

$$\mathbf{c}_k^* = \mathbf{A}_k \mathbf{c}_k + \mathbf{b}_k,$$

$$\boldsymbol{\Gamma}_k^* = \boldsymbol{\Sigma}_k + \mathbf{A}_k \boldsymbol{\Gamma}_k \mathbf{A}_k^\top,$$

$$\pi_k^* = \pi_k,$$

$$\mathbf{A}_k^* = \boldsymbol{\Sigma}_k^* \mathbf{A}_k^\top \boldsymbol{\Sigma}_k^{-1},$$

$$\mathbf{b}_k^* = \boldsymbol{\Sigma}_k^* (\boldsymbol{\Gamma}_k^{-1} \mathbf{c}_k - \mathbf{A}_k^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{b}_k),$$

$$\boldsymbol{\Sigma}_k^* = (\boldsymbol{\Gamma}_k^{-1} + \mathbf{A}_k^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{A}_k)^{-1}.$$

Regression functions

Both densities are Gaussian mixtures parameterized by θ .
Therefore, to obtain:

- A low-to-high *inverse* regression function:

$$\mathbb{E}[\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}; \theta] = \sum_{k=1}^K \frac{\pi_k \mathcal{N}(\mathbf{x}; \mathbf{c}_k, \mathbf{\Gamma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}; \mathbf{c}_j, \mathbf{\Gamma}_j)} (\mathbf{A}_k \mathbf{x} + \mathbf{b}_k),$$

- A high-to-low *forward* regression function:

$$\mathbb{E}[\mathbf{X} = \mathbf{x} | \mathbf{Y} = \mathbf{y}; \theta] = \sum_{k=1}^K \frac{\pi_k \mathcal{N}(\mathbf{y}; \mathbf{c}_k^*, \mathbf{\Gamma}_k^*)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{y}; \mathbf{c}_j^*, \mathbf{\Gamma}_j^*)} (\mathbf{A}_k^* \mathbf{y} + \mathbf{b}_k^*).$$

Low-to-High or High-to-Low?

If θ is unconstrained

GLLiM \Leftrightarrow Joint GMM on (\mathbf{X}, \mathbf{Y}) (JGMM)

- \mathbf{X} and \mathbf{Y} roles are symmetric
- Low-to-High or High-to-Low estimation are equivalent

Intractable for high D :

- JGMM requires inversion of K matrices of size $(D + L) \times (D + L)$

Low-to-High or High-to-Low?

Error vectors \mathbf{E}_k assumed isotropic Gaussians: $\forall k, \Sigma_k = \sigma \mathbf{I}_D$ (θ is constrained)

Example: $D = 1000$, $L = 2$, $K = 10$

- **Low-to-high** regression:
 $K(1 + L + DL + L(L + 1)/2 + D) = 30,060$ parameters.
- **High-to-low** regression:
 $K(1 + D + LD + D(D + 1)/2 + L) = 5,035,030$ parameters.
Requires inversion of 1000×1000 covariance matrices.

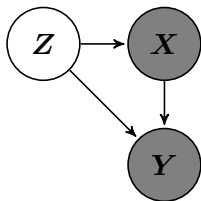
Therefore it is better to perform a low-dimensional-to-high-dimensional (inverse) regression, and then deduce the forward density.

Extension: the Hybrid GLLiM model

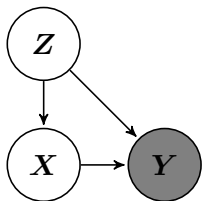
Incorporate a **latent component** into the **low-dimensional** variable:

$$\mathbf{X} = \begin{bmatrix} \mathbf{T} \\ \mathbf{W} \end{bmatrix}$$

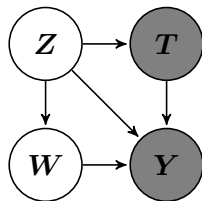
where $\mathbf{T} \in \mathbb{R}^{L_t}$ is observed and $\mathbf{W} \in \mathbb{R}^{L_w}$ is latent ($L = L_t + L_w$)



Supervised GLLiM (regression)



Unsupervised GLLiM (dim reduction)



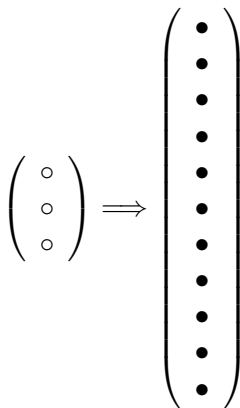
Hybrid GLLiM

Fully-latent Output Variable: Dimensionality reduction, eg. PPCA

Learning

x

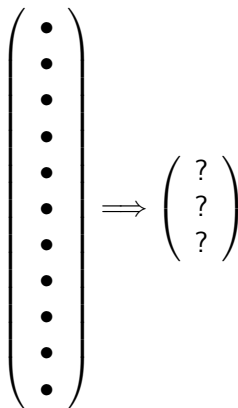
y



Testing

y

x

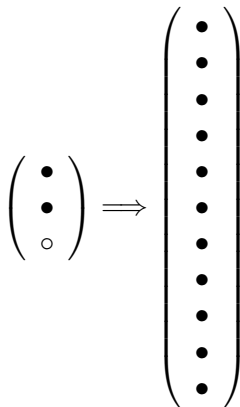


Partially-latent Output Variable : Hybrid GLLiM

Learning

x

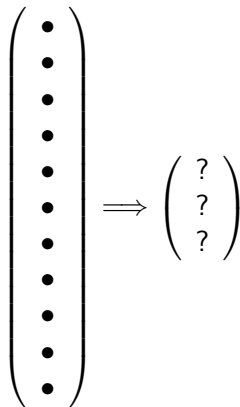
y



Testing

y

x



The hybrid GLLiM model

Hybrid between regression and dimensionality reduction:

$$\mathbf{X} = \begin{bmatrix} \mathbf{T} \\ \mathbf{W} \end{bmatrix}$$

- Observed pairs $\{(\mathbf{y}_n, \mathbf{t}_n), n = 1 \dots N\}$ ($\mathbf{T} \in \mathbb{R}^{L_t}$)
- Additional latent variable \mathbf{W} ($\mathbf{W} \in \mathbb{R}^{L_w}$)
- Assuming the independence of \mathbf{T} and \mathbf{W} given Z :

$$p(\mathbf{X} = (\mathbf{t}, \mathbf{w})^\top \mid Z = k) = \mathcal{N}_L((\mathbf{t}, \mathbf{w})^\top; \mathbf{c}_k, \mathbf{\Gamma}_k)$$

$$\text{with } \mathbf{c}_k = \begin{bmatrix} \mathbf{c}_k^t \\ \mathbf{c}_k^w \end{bmatrix}, \mathbf{\Gamma}_k = \begin{bmatrix} \mathbf{\Gamma}_k^t & \mathbf{0} \\ \mathbf{0} & \mathbf{\Gamma}_k^w \end{bmatrix}$$

The hybrid GLLiM model

With $\mathbf{A}_k = \begin{bmatrix} \mathbf{A}_k^t & \mathbf{A}_k^w \end{bmatrix}$,

$$\mathbf{Y} = \sum_{k=1}^K \mathbb{I}(Z = k) (\mathbf{A}_k^t \mathbf{T} + \mathbf{A}_k^w \mathbf{W} + \mathbf{b}_k + \mathbf{E}_k)$$

with $\mathbf{E}_k \sim \mathcal{N}_D(0, \boldsymbol{\Sigma}_k)$

rewrites

$$\mathbf{Y} = \sum_{k=1}^K \mathbb{I}(Z = k) (\mathbf{A}_k^t \mathbf{T} + \mathbf{b}_k + \mathbf{A}_k^w \mathbf{c}_k^w + \mathbf{E}'_k)$$

with $\mathbf{E}'_k \sim \mathcal{N}_D(0, \boldsymbol{\Sigma}_k + \mathbf{A}_k^w \boldsymbol{\Gamma}_k^w \mathbf{A}_k^{w\top})$

- Supervised GLLiM with *unconventional covariance* structure
- Diagonal $\boldsymbol{\Sigma}_k \rightarrow$ **Factor analysis** with L_w factors (at most)
- A compromise between full $O(D^2)$ and diagonal $O(D)$ covariances

Link to other models

Assuming $\Sigma_k = \sigma_k^2 \mathbf{I}_D$,

$$(\Sigma'_k = \Sigma_k + \mathbf{A}_k^w \mathbf{\Gamma}_k^w \mathbf{A}_k^{w\top})$$

- $L_w = 0$, Supervised case, $\Sigma'_k = \Sigma_k$:
Mixture of local linear experts (MLE) [Xu et al 95]
- $L_w = D$, Σ'_k general covariance matrix:
JGMM model [Qiao et al 09], the most general GLLiM model
Over-parameterized, intractable ($(D + L) \times (D + L)$ matrices)
- $0 < L_w < D$: a wide variety of models *between* MLE and JGMM.

Gaussian Process Latent Variable Model [Lawrence 05, Fusi & al 12]:

Regression with partially-latent *input*, but not with partially-latent *response*

GPLVM mapping not invertible (non-linear nature of the kernels used in practice)

Particular instances of the hybrid GLLiM model

First three rows: supervised GLLiM methods ($L_w = 0$)

Last six rows: unsupervised GLLiM methods ($L_t = 0$)

Model	\mathbf{c}_k	$\mathbf{\Gamma}_k$	π_k	\mathbf{A}_k	\mathbf{b}_k	$\mathbf{\Sigma}_k$	L_t	L_w	K
MLE [Xu et al 95]	-	-	-	-	-	diag	-	0	-
MLR [Jedidi et al 96]	$\mathbf{0}_L$	$\infty \mathbf{I}_L$	-	-	-	iso+eq	-	0	-
JGMM [Qiao et al 09]	-	-	-	-	-	-	-	0	-
PPAM [Deleforge et al 12]	-	eq	eq	-	-	diag+eq	-	0	-
GTM [Bishop et al 98]	fixed	$\mathbf{0}_L$	eq.	eq.	$\mathbf{0}_D$	iso+eq	0	-	-
PPCA [Tipping et al 99a]	$\mathbf{0}_L$	\mathbf{I}_L	-	-	-	iso	0	-	1
MPPCA [Tipping et al 99b]	$\mathbf{0}_L$	\mathbf{I}_L	-	-	-	iso	0	-	-
MFA [Ghahramani et al 96]	$\mathbf{0}_L$	\mathbf{I}_L	-	-	-	diag	0	-	-
PCCA [Bach et al 05]	$\mathbf{0}_L$	\mathbf{I}_L	-	-	-	block	0	-	1
RCA [Kalaitzis et al 11]	$\mathbf{0}_L$	\mathbf{I}_L	-	-	-	fixed	0	-	1

Expectation Maximization for Hybrid GLLiM

2 data augmentation schemes: Convergence speed/M-step tractability tradeoff

(other: Alternating ECM [Meng & Rubin 97] eg. for MFA [McLachlan et al 03])

- General hybrid GLLiM-EM: augmenting with both (Z, \mathbf{W})
 - Closed-form expressions for a wide range of $\{\mathbf{\Gamma}_k, \mathbf{\Sigma}_k, k = 1 \dots K\}$
- Marginal-hGLLiM: integrating out the \mathbf{W}
 - Less general, closed form only for distinct isotropic $\{\mathbf{\Sigma}_k, k = 1 \dots K\}$
 - Algorithmic insight: alternation of a regression & reduction step
 - Natural initialization strategy

Identifiability issues

As for latent variable models for dimensionality reduction (eg MPPCA, MFA):

$\{\mathbf{c}_k^w\}_{k=1}^K$ and $\{\mathbf{\Gamma}_k^w\}_{k=1}^K$ must be fixed , eg. $\mathbf{c}_k^w = \mathbf{0}$ and $\mathbf{\Gamma}_k^w = \mathbf{I}_{L_w}$

$$\mathbf{Y} = \sum_{k=1}^K \mathbb{I}(Z = k) (\mathbf{A}_k^t \mathbf{T} + \mathbf{A}_k^w \mathbf{W} + \mathbf{b}_k + \mathbf{E}_k)$$

$$\begin{aligned} (\mathbf{W} | Z = k) &\sim \mathcal{N}(\mathbf{c}_k^w, \mathbf{\Gamma}_k^w) \implies \\ (\mathbf{A}_k^w \mathbf{W} + \mathbf{b}_k | Z = k) &\sim \mathcal{N}(\mathbf{A}_k^w \mathbf{c}_k^w + \mathbf{b}_k, \mathbf{A}_k^w \mathbf{\Gamma}_k^w \mathbf{A}_k^{w\top}) \end{aligned}$$

The general Hybrid GLLiM EM algorithm

Observed $\{(\mathbf{Y}_n, \mathbf{T}_n), n = 1 : N\}$ and Missing variables $\{(Z_n, \mathbf{W}_n), n = 1 : N\}$

At iteration (i) , update:

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

E-step: compute posterior distributions: $\forall n, \forall k$

$$p(Z_n = k | \mathbf{t}_n, \mathbf{y}_n; \boldsymbol{\theta}^{(i)}) = \frac{\pi_k^{(i)} p(\mathbf{y}_n, \mathbf{t}_n | Z_n = k; \boldsymbol{\theta}^{(i)})}{\sum_{j=1}^K \pi_j^{(i)} p(\mathbf{y}_n, \mathbf{t}_n | Z_n = j; \boldsymbol{\theta}^{(i)})}$$

$$p(\mathbf{w}_n | Z_n = k, \mathbf{t}_n, \mathbf{y}_n; \boldsymbol{\theta}^{(i)}) \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_{nk}^{\mathbf{w}}, \tilde{\mathbf{S}}_k^{\mathbf{w}}) \quad (\text{Factor Analysis like})$$

The general Hybrid GLLiM EM algorithm

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

M-step: divides in two

- **Updating $\pi_k, \mathbf{c}_k^t, \boldsymbol{\Gamma}_k^t$:** standard Gaussian mixture step on $\{\mathbf{t}_n, n = 1 \dots N\}$
- **Updating the mapping parameters $\mathbf{A}_k, \mathbf{b}_k, \boldsymbol{\Sigma}_k$**
 - $L_w = 0$: \mathbf{A}_k is that of standard linear regression from $\{\mathbf{t}_n, n = 1 \dots N\}$ to $\{\mathbf{y}_n, n = 1 \dots N\}$ weighted by $\{\tilde{r}_{nk}, n = 1 \dots N\}$
 - $L_t = 0$: principal components update of PPCA

M-GMM step

With $\tilde{r}_k = \sum_{n=1}^N \tilde{r}_{nk}$

$$\begin{aligned}\tilde{\pi}_k &= \frac{\tilde{r}_k}{N}, \\ \tilde{\mathbf{c}}_k^t &= \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} \mathbf{t}_n, \\ \tilde{\mathbf{\Gamma}}_k^t &= \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} (\mathbf{t}_n - \tilde{\mathbf{c}}_k^t)(\mathbf{t}_n - \tilde{\mathbf{c}}_k^t)^\top.\end{aligned}$$

M-mapping step

$$\tilde{\mathbf{A}}_k = \tilde{\mathbf{Y}}_k \tilde{\mathbf{X}}_k^\top (\tilde{\mathbf{S}}_k^x + \tilde{\mathbf{X}}_k \tilde{\mathbf{X}}_k^\top)^{-1}$$

where:

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

$$\tilde{\mathbf{X}}_k = \frac{1}{\sqrt{\tilde{r}_k}} \left[\sqrt{\tilde{r}_{1k}}(\tilde{\mathbf{x}}_{1k} - \tilde{\mathbf{x}}_k), \quad \dots, \quad \sqrt{\tilde{r}_{Nk}}(\tilde{\mathbf{x}}_{Nk} - \tilde{\mathbf{x}}_k) \right],$$

$$\tilde{\mathbf{Y}}_k = \frac{1}{\sqrt{\tilde{r}_k}} \left[\sqrt{\tilde{r}_{1k}}(\mathbf{y}_1 - \tilde{\mathbf{y}}_k), \quad \dots, \quad \sqrt{\tilde{r}_{Nk}}(\mathbf{y}_N - \tilde{\mathbf{y}}_k) \right],$$

$$\tilde{\mathbf{x}}_{nk} = [\mathbf{t}_n; \tilde{\boldsymbol{\mu}}_{nk}^w] \in \mathbb{R}^L, \quad \tilde{\mathbf{x}}_k = \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} \tilde{\mathbf{x}}_{nk}, \quad \tilde{\mathbf{y}}_k = \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} \mathbf{y}_n.$$

And

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

Practical setting

Algorithm initialization

No straightforward way of choosing $r_{nk}^{(0)}$, $\boldsymbol{\mu}_{nk}^{w(0)}$, $\mathbf{S}_k^{w(0)}$ or a complete set $\boldsymbol{\theta}^{(0)}$ including all affine transformations

→ Use one iteration of Marginal hGLLiM EM to get $\boldsymbol{\theta}^{(0)}$

Latent dimension L_w estimation

$$BIC(\tilde{\boldsymbol{\theta}}, N) = -2\mathcal{L}(\tilde{\boldsymbol{\theta}}) + \mathcal{D}(\tilde{\boldsymbol{\theta}}) \log N,$$

\mathcal{L} : observed-data log-likelihood

$\mathcal{D}(\tilde{\boldsymbol{\theta}})$: dimension of the complete parameter vector

The Marginal Hybrid GLLiM-EM

$$\mathbf{c}_k^w = 0 \text{ and } \mathbf{\Gamma}_k^w = \mathbf{I}_{L_w} \implies$$

$$\mathbf{Y} = \sum_{k=1}^K \mathbb{I}(Z = k) (\mathbf{A}_k^t \mathbf{T} + \mathbf{b}_k + \mathbf{E}'_k)$$

with $\mathbf{E}'_k \sim \mathcal{N}_D(0, \mathbf{\Sigma}_k + \mathbf{A}_k^w \mathbf{A}_k^{w\top})$

- No E-W step (marginalized)
- Same E-Z step (\tilde{r}_{nk} initialized via eg. K-means)
- Same M-GMM step ($\pi_k, \mathbf{c}_k^t, \mathbf{\Gamma}_k^t$)
- M-regression step ($\mathbf{A}_k^t, \mathbf{b}_k$) : standard, does not involve noise variance
- M-residual step ($\mathbf{A}_k^w, \mathbf{\Sigma}_k$) : PPCA like on residuals $\mathbf{y}_n - \mathbf{A}_k^t \mathbf{t}_n - \mathbf{b}_k$ (time consuming)

The Marginal Hybrid GLLiM M-step

M-regression-step: Weighted affine regression from

$\{\mathbf{t}_n, n = 1 : N\}$ to $\{\mathbf{y}_n, n = 1 : N\}$ with weights \tilde{r}_{nk} ,

$$\tilde{\mathbf{A}}_k^t = \tilde{\mathbf{Y}}_k \tilde{\mathbf{T}}_k^\top (\tilde{\mathbf{T}}_k \tilde{\mathbf{T}}_k^\top)^{-1}, \quad \tilde{\mathbf{b}}_k = \sum_{n=1}^N \frac{\tilde{r}_{kn}}{\tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{A}}_k^t \mathbf{t}_n),$$

with

$$\tilde{\mathbf{T}}_k = \left[\sqrt{\tilde{r}_{1k}} (\mathbf{t}_1 - \tilde{\mathbf{t}}_k) \dots \sqrt{\tilde{r}_{Nk}} (\mathbf{t}_N - \tilde{\mathbf{t}}_k) \right] \sqrt{\tilde{r}_k}$$

and

$$\tilde{\mathbf{t}}_k = \sum_{n=1}^N (\tilde{r}_{kn} / \tilde{r}_k) \mathbf{t}_n$$

M-residual-step: Minimization of the following criterion:

$$Q_k(\boldsymbol{\Sigma}_k, \mathbf{A}_k^w) = -\frac{1}{2} \left(\log |\boldsymbol{\Sigma}_k + \mathbf{A}_k^w \mathbf{A}_k^{w\top}| + \sum_{n=1}^N \mathbf{u}_{kn}^\top (\boldsymbol{\Sigma}_k + \mathbf{A}_k^w \mathbf{A}_k^{w\top})^{-1} \mathbf{u}_{kn} \right),$$

where $\mathbf{u}_{kn} = \sqrt{\tilde{r}_{nk} / \tilde{r}_k} (\mathbf{y}_n - \tilde{\mathbf{A}}_k^t \mathbf{t}_n - \tilde{\mathbf{b}}_k)$.

Experiments and results

High dimensional function regression

$$\boldsymbol{\phi} = (\phi_1 \dots \phi_d \dots \phi_D)^\top$$

$$\boldsymbol{\phi} = \mathbf{f}, \mathbf{g}, \mathbf{h}$$

$$\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^D \text{ with } f_d(t, w_1) = \alpha_d \cos(\eta_d t/10 + \phi_d) + \gamma_d w_1^3$$

$$\mathbf{g} : \mathbb{R}^2 \rightarrow \mathbb{R}^D \text{ with } g_d(t, w_1) = \alpha_d \cos(\eta_d t/10 + \beta_d w_1 + \phi_d)$$

$$\mathbf{h} : \mathbb{R}^3 \rightarrow \mathbb{R}^D \text{ with}$$

$$h_d(t, w_1, w_2) = \alpha_d \cos(\eta_d t/10 + \beta_d w_1 + \phi_d) + \gamma_d w_2^3$$

$$\boldsymbol{\xi} = \{\alpha_d, \eta_d, \phi_d, \beta_d, \gamma_d\}_{d=1}^D \text{ in } [0, 2], [0, 4\pi], [0, 2\pi], [0, \pi], [0, 2]$$

High dimensional function regression

100 $\mathbf{f}, \mathbf{g}, \mathbf{h}$ functions generated using different random values for ξ

N training couples $\{(t_n, \mathbf{y}_n)\}_{n=1}^N$

N' test couples $\{(t'_n, \mathbf{y}'_n)\}_{n=1}^{N'}$

by randomly drawing $t \in [0, 10]$ and $\mathbf{w} \in [-1, 1]$ (\mathbf{f}, \mathbf{g}) or $\in [-1, 1]^2$ (\mathbf{h})

and adding some random isotropic Gaussian noise $\mathbf{y} = \phi(t, \mathbf{w}) + \mathbf{e}$.

Training couples: train the different regression algorithms tested (h-GLLiM, SIR, RVM, MLE, JGMM)

Task: Estimate \hat{t}'_n given a test observation $\mathbf{y}'_n = \phi(t'_n, \mathbf{w}'_n) + \mathbf{e}'_n$

High dimensional function regression $D = 50$

Average, standard deviation and % of extreme values of the absolute error $|\hat{t}'_n - t'_n|$. $N = 200, N' = 200, K = 5$ (MLE, JGMM, hGLLiM)

MLE: $L_w = 0$, JGMM : $L_w > D$

Method	f			g			h		
	Avg	Std	Ex	Avg	Std	Ex	Avg	Std	Ex
JGMM	1.78	2.21	19.5	2.45	2.76	28.4	2.26	2.87	22.4
SIR-1	1.28	1.07	5.92	1.73	1.39	14.9	1.64	1.31	13.0
SIR-2	0.60	0.69	1.02	1.02	1.02	4.20	1.03	1.06	4.91
RVM	0.59	0.53	0.30	0.86	0.68	0.52	0.93	0.75	1.00
MLE	0.36	0.53	0.50	0.36	0.34	0.04	0.61	0.69	0.99
hGLLiM-1	0.20	0.24	0.00	0.25	0.28	0.01	0.46	0.48	0.22
hGLLiM-2	0.23	0.24	0.00	0.25	0.25	0.00	0.36	0.38	0.04
hGLLiM-3	0.24	0.24	0.00	0.26	0.25	0.00	0.34	0.34	0.01
hGLLiM-4	0.23	0.23	0.01	0.28	0.27	0.00	0.35	0.34	0.01
hGLLiM-BIC	0.18	0.21	0.00	0.24	0.26	0.00	0.33	0.35	0.06

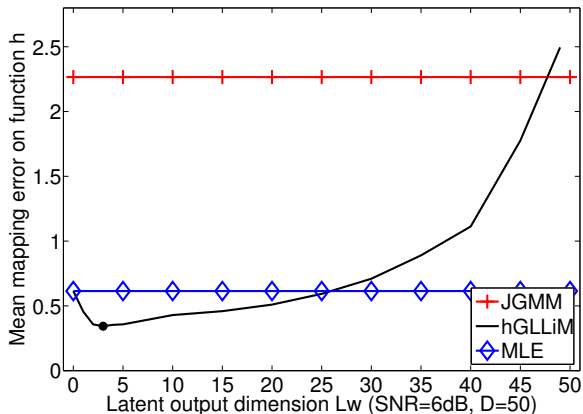
hGLLiM-BIC minimizes BIC for $0 < L_w < 10$: expected latent dimension L_w ($L_w = 2$ or $L_w = 1$) selected 72 times over 100 (non-linear effects could be modeled by higher L_w)

Influence of L_w

Influence of the parameter L_w of hGLLiM on the mean mapping error of \mathbf{h} .

Each point corresponds to an average error over 10,000 tests on 50 distinct functions

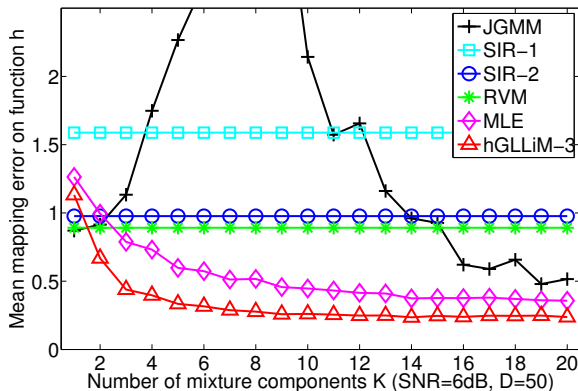
MLE: $L_w = 0$, JGMM : $L_w > D$



Influence of K

Influence of K in MLE, JGMM and hGLLiM-3 on the mean mapping error of synthetic function h .

Each point corresponds to an average error over 10,000 tests on 50 distinct functions

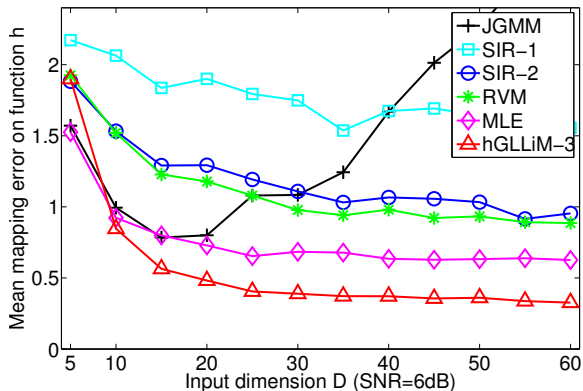


Errors generally decrease with K . Overfitting for $K > 10$ for JGMM

Influence of D

Influence of D on the mean mapping error of synthetic functions h

Each point corresponds to an average error over 10,000 tests on 50 distinct functions

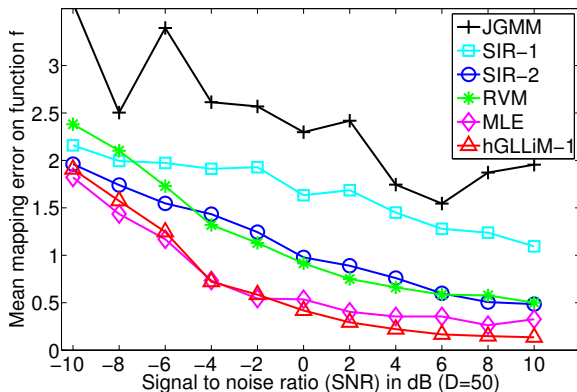


h-GLLiM performs better in high-dimension

Influence of the SNR

Influence of the signal-to-noise ratio (SNR) on the mean mapping error of synthetic functions f

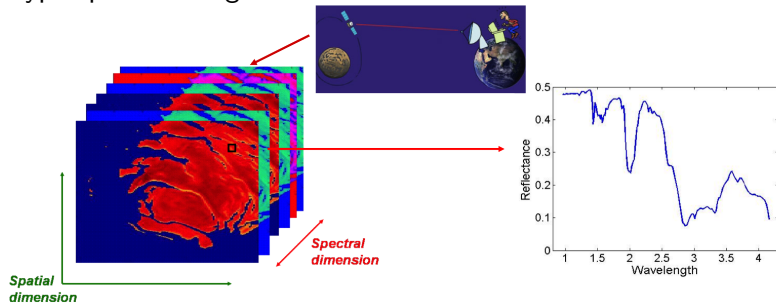
Each point corresponds to an average error over 10,000 tests on 50 distinct functions



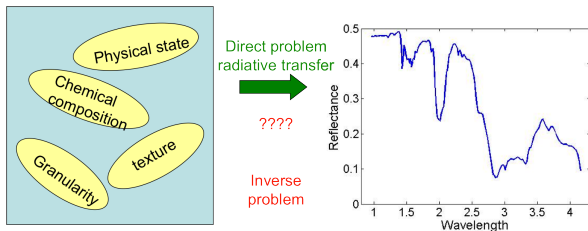
All methods perform similarly under extreme noise (SNR=-10dB) (except for JGMM)

Retrieval of Mars physical properties

Hyperspectral images



Radiative transfer model



Synthetic data

15,407 spectra ($D = 184$ wavelengths) and $L = 5$ real parameters
(proportion of water ice, of CO₂ ice, of dust, grain size of water ice, of CO₂ ice)

Proportion of water ice & grain size of CO₂ ice ignored from training

Method	Proportion of CO ₂ ice	Proportion of dust	Grain size of water ice
JGMM	0.83 ± 1.61	0.62 ± 1.00	0.79 ± 1.09
SIR-1	1.27 ± 2.09	1.03 ± 1.71	0.70 ± 0.94
SIR-2	0.96 ± 1.72	0.87 ± 1.45	0.63 ± 0.88
RVM	0.52 ± 0.99	0.40 ± 0.64	0.48 ± 0.64
MLE	0.54 ± 1.00	0.42 ± 0.70	0.61 ± 0.92
hGLLiM-1	0.36 ± 0.70	0.28 ± 0.49	0.45 ± 0.75
hGLLiM-2*†	0.34 ± 0.63	0.25 ± 0.44	0.39 ± 0.71
hGLLiM-3	0.35 ± 0.66	0.25 ± 0.44	0.39 ± 0.66
hGLLiM-4	0.38 ± 0.71	0.28 ± 0.49	0.38 ± 0.65
hGLLiM-5	0.43 ± 0.81	0.32 ± 0.56	0.41 ± 0.67
hGLLiM-20	0.51 ± 0.94	0.38 ± 0.65	0.47 ± 0.71
hGLLiM-BIC	0.34 ± 0.63	0.25 ± 0.44	0.39 ± 0.71

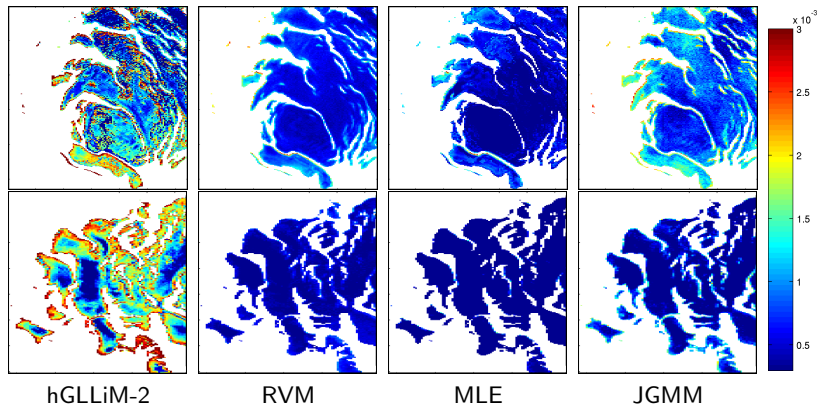
NRMSE for Mars surface physical properties recovered from synthetic spectra:
cross validation with 10,000 training pairs at random and 5,407 test pairs ($\times 20$)

$K = 50$ for MLE, LGMM, hGLLiM

Hyperspectral images of South polar cap of Mars

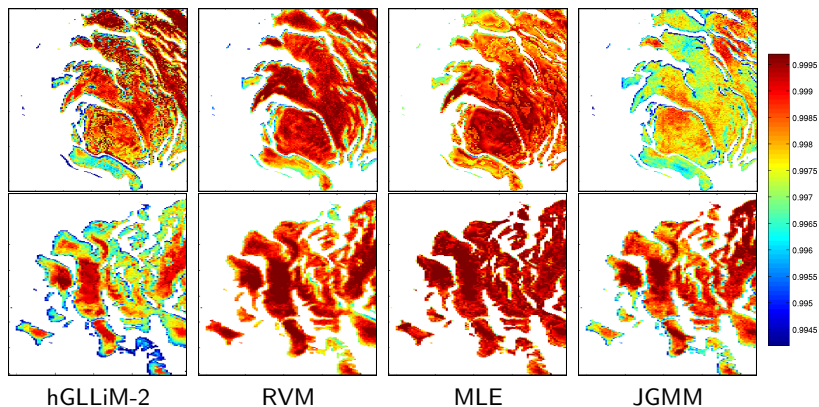
Omega instrument, Mars Express

Proportions of dust for the South polar cap of Mars: orbits 41 and 61



Hyperspectral images of South polar cap of Mars

Proportions of CO₂ ice for the South polar cap of Mars: orbits 41 and orbit 61



Conclusion/ Perspectives

- We propose a novel **inverse** approach to high-dimensional regression based on mixture- and latent-variable models.
- Latent component allows to capture behaviors that cannot be easily modeled

- Adaptive latent dimension L_w selection
- More complex dependencies between variables (eg. $(Z_1 \dots Z_N)$ is a MRF)
- More complex noise models, eg, Student for outliers accommodation and robustness

Matlab code available at: https://team.inria.fr/perception/gllim_toolbox/

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

A. Deleforge, F. Forbes and R. Horaud, Hyper-spectral Image Analysis with Partially-Latent Regression. EUSIPCO, Lisbon, Portugal, September 2014.

MRF modelling

