

Multivariate Modelling: Latent Variables & Dimensionality Reduction

Habib Ganjgahi

Department of Statistics
University of Oxford

November 6, 2018

Table of Contents

Neuroimaging

Latent Variable Modelling

Introduction to Neuroimaging

Go to Power point

Table of Contents

Neuroimaging

Latent Variable Modelling

Table of Contents

Neuroimaging

Latent Variable Modelling

Principal Component Analysis

Factor Analysis

Bayesian Latent Variable Modelling

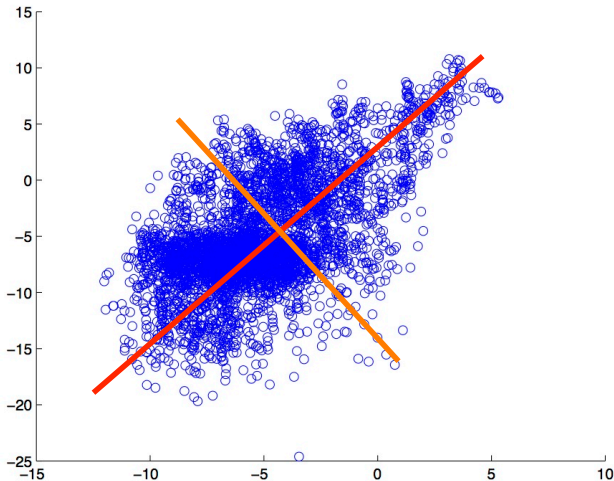
Multi-view Learning

Latent Variable Modelling

- ▶ Model high dimensional Y with a lower-dimensional approximation F that has N rows, but typically far fewer columns
- ▶ Latent variable modelling is used for:
 - ▶ Dimensionality reduction: replace Y with a lower-dimensional F
 - ▶ Outlier detection
 - ▶ Data visualization: display F in a scatterplot.
 - ▶ Factor discovering: discover important hidden **factors** underlying data.

PCA

Projects the data into a sub-space such that the projected data has maximum variance.



Principal Component Analysis (PCA)

- ▶ Projects the data into a sub-space such that the projected data has maximum variance.

$$Y \approx F\Lambda',$$

where $Y \in \mathbb{R}^{N \times P}$; $\Lambda \in \mathbb{R}^{P \times K}$ principal axis; $F \in \mathbb{R}^{N \times K}$ and $K \ll P$. PCA derivations:

- ▶ choose Λ minimizes construction error: $\operatorname{argmin} \|Y - F\Lambda'\|_F^2$
- ▶ choose Λ maximises the variance:
 $\operatorname{argmax} \sum_{i=1}^N \|\Lambda' Y_i\|^2 = \operatorname{Tr}(\Lambda' Y' Y \Lambda)$
- ▶ Λ is the eigenvector of sample covariance matrix $S = \frac{1}{N} Y' Y$

Table of Contents

Neuroimaging

Latent Variable Modelling

Principal Component Analysis

Factor Analysis

Bayesian Latent Variable Modelling

Multi-view Learning

Factor Analysis (FA)

Describes the correlation among variables using few latent variables.
Projects the data into a space that are not correlated

$$\mathbf{y}_i = \mathbf{f}_i \mathbf{\Lambda}' + \boldsymbol{\epsilon}_i, \quad (1)$$

where $k \ll N$ is; $\mathbf{f}_t \in \mathbb{R}^k$ is the k vector of latent variables (scores);
and $\mathbf{\Lambda} \in \mathbb{R}^{P \times k}$ is the loadings matrix. It is assumed that

- ▶ Latent variables are independent $\mathbf{f}_i \sim N(\mathbf{0}, \mathbf{I}_{k \times k})$
- ▶ Residuals are independent with $\boldsymbol{\epsilon}_t \sim N(\mathbf{0}, \mathbf{\Psi})$ where $\mathbf{\Psi} = \text{diag}(\sigma_1^2, \dots, \sigma_P^2)$
- ▶ The latent variables and residuals are mutually independent $\text{Cov}(\mathbf{f}_t, \boldsymbol{\epsilon}_s) = \mathbf{0}$ for all t and s .

Factor Analysis

The FA model:

$$y_i = f_i \Lambda' + \epsilon_i,$$

- ▶ $y_i | f_i \sim \mathcal{N}(f_i \Lambda', \Psi)$
- ▶ $y_i \sim \mathcal{N}(0, \Sigma), \Sigma = \Lambda \Lambda' + \Psi$
- ▶ log-likelihood function

$$l(\Lambda, \Psi) \propto -\frac{N}{2} \log |\Psi| - \frac{N}{2} \text{Tr}(S \Psi^{-1})$$

where $S = \frac{1}{N} \mathbf{Y}' \mathbf{Y}$ is the sample covariance matrix

- ▶ Parameter estimation: Maximise likelihood function
- ▶ Varimax rotation for interpretation

Dimensionality Reduction: Unified Approach

$$Y = F\Lambda' + E,$$

$$F \sim \mathcal{N}(\mathbf{0}, I), \quad E \sim \mathcal{N}(\mathbf{0}, \Psi),$$

$$Y \sim \mathcal{N}(\mathbf{0}, \Lambda\Lambda' + \Psi),$$

$$F|Y, \Lambda, \Psi \sim \mathcal{N}(m, V)$$

$$V = \left(I - \Lambda'\Psi^{-1}\Lambda\right)^{-1}, \quad m = V\Lambda'\Psi^{-1}Y$$

► PCA ¹

- $\Psi = \sigma^2 I$
- $\sigma^2 \rightarrow 0$ recovers classic PCA

► FA

- $\Psi = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$

Identifiability

The model (Eq., 1) is invariant under any orthogonal transformation:

$$\mathbf{y} = (\mathbf{F}\mathbf{P})(\mathbf{\Lambda}\mathbf{P})' + \boldsymbol{\epsilon} = \mathbf{F}^*\mathbf{\Lambda}^* + \boldsymbol{\epsilon}$$

where $\mathbf{P}\mathbf{P}' = \mathbf{I}$. That is mean that the rotated parameters $(\mathbf{F}^*, \mathbf{\Lambda}^*)$ and $(\mathbf{F}, \mathbf{\Lambda})$ are not distinguishable, hence further restrictions, so called identifiability condition, should be imposed to achieve unique solution.

Identifiability

- ▶ Constraint parameter space
 - ▶ Block of lower triangular matrix
 - ▶ Sparse mode
- ▶ Post process MCMC outputs or MAP estimates

Table of Contents

Neuroimaging

Latent Variable Modelling

Principal Component Analysis

Factor Analysis

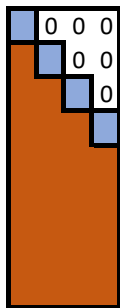
Bayesian Latent Variable Modelling

Multi-view Learning

Bayesian FA: Identifiability

Identifiability condition breaks the symmetry of likelihood function and guarantees unique mode of the posterior²:

- ▶ Λ is full column rank and has a square block of lower triangular with positive diagonal elements (LT).
- ▶ $\exists \mathbf{B}$ of size Λ and full column rank such that $\mathbf{B}'\Lambda$ has a block of lower triangular matrix.
- ▶ $\Lambda'\Lambda$ or $\Lambda'\Sigma^{-1}\Lambda$ is diagonal and diagonal elements are unique and arranged in descending order.



Inference

- ▶ Bayes rule:

$$\pi(\theta|Y) \propto \pi(\theta)\pi(Y|\theta),$$

where $\theta = (\Lambda, F, \Psi)$

- ▶ MCMC algorithm
- ▶ MAP estimates: Expectation Maximisation (EM)

$$\begin{aligned} Q(\theta) &= E_{F|Y}[\log(\pi(\theta|Y))] \\ \hat{\theta} &= \operatorname{argmax}(Q(\theta)) \end{aligned}$$

Identifiability: LT

▶ Caveats:

- ▶ Posterior may suffer from local maxima that affects the MCMC sampling schemes.
- ▶ Reduces the accuracy of parameter estimations ³.
- ▶ Ordering problem: LT constraint is invariant to order of variables
- ▶ Number of latent variables is influenced by order of variables ⁴.

Identifiability

- ▶ Identifiability condition that $\Lambda\Lambda'$ is diagonal with unique diagonal elements in descending order.
- ▶ Kaufmann and Schumacher⁵ used the following specification for the FA model in Eq.,(1)

$$\begin{aligned}y_t &= (f_t^*P)(\Lambda^*P)' + \epsilon, \\ &= f_t\Lambda' + \epsilon,\end{aligned}$$

where Λ^* is sparse; f_t^* is the latent variable vector corresponds to the sparse loading and P is an orthogonal matrix such that $\Lambda'\Lambda$ is diagonal with unique elements in descending order.

- ▶ P is derived from eigenvalue decomposition of $\Lambda^*\Lambda^*$.
- ▶ In their approach, a draw from Λ^* is used to get an identified solution, then latent variables correspond to the sparse loading is derived using the identified latent variable as $f_t^* = f_tP'$.

Identifiability: Sparse Modes

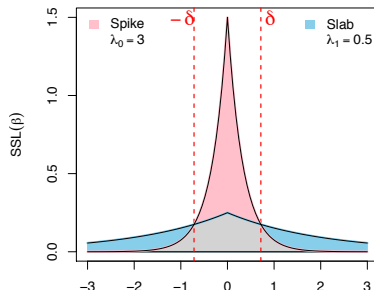
- ▶ Sparsity inducing priors can be used to achieve unique solution
- ▶ Helps interoperability of latent variables while only few variables enter to each latent variables.
- ▶ Improves predictive performances ⁶.
- ▶ Shrinkage prior:
 - ▶ Automatic relevance determination (ARD)
 - ▶ Horseshoe
- ▶ No exact zeros
- ▶ spike and slap prior

Spike and Slab prior

Spike and Slab (SSL) prior with Laplace components

$$\pi(\lambda|\gamma) = \gamma \text{Laplace}(\lambda|\sigma_1^2) + (1 - \gamma) \text{Laplace}(\lambda|\sigma_0^2)$$

- ▶ σ_1^2 **small**: to avoid over-shrinkage of large effects
- ▶ σ_0^2 **large**: to shrink ignorable coefficients to zero
- ▶ θ controls sparsity, $P(\gamma = 1|\theta) = \theta$
- ▶ Point-mass spike-and-slab is a limiting case as $\sigma_0^2 \rightarrow 0$



Latent variable: Spike and Slab prior

Hierarchical spike and slab prior⁷:

$$\pi(\lambda_{jk}|\gamma_{jk}) \sim \gamma \text{Laplace}(\lambda_{jk}|\sigma_1^2) + (1 - \gamma) \text{Laplace}(\lambda_{jk}|\sigma_0^2)$$

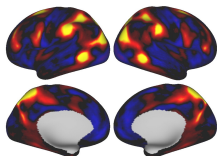
$$\gamma_{jk} \sim \text{Bernoulli}(\theta_k)$$

$$\theta_k = \prod_{l=1}^k \nu_l$$

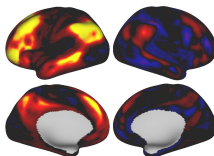
$$\nu_l \stackrel{iid}{\sim} \text{B}(\alpha, 1)$$

An Example: rsfMRI Modelling⁸

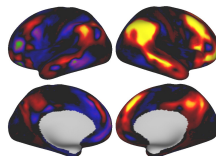
$$Y = F\Lambda' + E$$



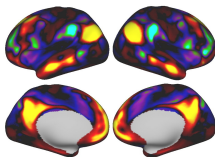
(a) PFM 13 (0.86). Dorsal attention (DA).



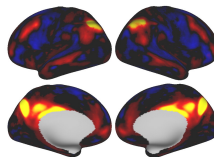
(b) PFM 18 (0.95). Left fronto-parietal control (IFPC).



(c) PFM 17 (0.77). Right fronto-parietal control (rFPC).



(d) PFM 15 (0.90). Default mode (DM).



(e) PFM 16 (0.98). Posteriomedial cortex, inferior parietal lobule (PMPL).



Table of Contents

Neuroimaging

Latent Variable Modelling

Principal Component Analysis

Factor Analysis

Bayesian Latent Variable Modelling

Multi-view Learning

Canonical Correlation Analysis

Linear combination such that maximises the correlation

$$\begin{array}{ccccc} \boxed{\begin{array}{c} Y_1 \\ (P_1 \times N) \end{array}} & = & \boxed{\begin{array}{c} \Lambda_1 \\ (P_1 \times K) \end{array}} & \times & \boxed{\begin{array}{c} F \\ (K \times N) \end{array}} + \boxed{\begin{array}{c} \epsilon_1 \end{array}} \\ \\ \boxed{\begin{array}{c} Y_2 \\ (P_2 \times N) \end{array}} & = & \boxed{\begin{array}{c} \Lambda_2 \\ (P_2 \times K) \end{array}} & \times & \boxed{\begin{array}{c} F \\ (K \times N) \end{array}} + \boxed{\begin{array}{c} \epsilon_2 \end{array}} \end{array}$$

CCA Application: UKBB

Go to power point

Canonical Correlation Analysis (CCA)

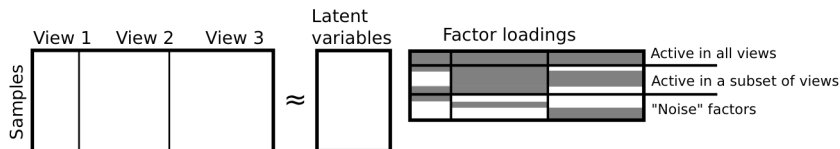
Bayesian CCA ⁹

$$\begin{aligned} \mathbf{y}_i^{(1)} &= \mathbf{A}^{(1)} \mathbf{f}_i^{(0)} + \mathbf{B}^{(1)} \mathbf{f}_i^{(1)} + \boldsymbol{\epsilon}_i^{(1)}, \\ \mathbf{y}_i^{(2)} &= \mathbf{A}^{(2)} \mathbf{f}_i^{(0)} + \mathbf{B}^{(2)} \mathbf{f}_i^{(2)} + \boldsymbol{\epsilon}_i^{(2)}, \end{aligned}$$

where $\mathbf{f}_i^{(0)} \in \mathbb{R}^{k_0 \times 1}$ is the shared latent variable; $\mathbf{f}_i^{(1)} \in \mathbb{R}^{k_1 \times 1}$ and $\mathbf{f}_i^{(2)} \in \mathbb{R}^{k_2 \times 1}$ are the view specific latent variables.

- ▶ $\boldsymbol{\epsilon}_i^{(w)} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi}^{(w)}), w = 1, 2$
- ▶ $\mathbf{f}_i^{(w)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ▶ Marginal distributions: $\text{cov}(\mathbf{y}_i)^{(w)} = \boldsymbol{\Sigma}^{(w)} = \mathbf{B}^{(w)} \mathbf{B}'^{(w)} + \boldsymbol{\Psi}^{(w)}$

Multi-view Learning



Generative model: Extend CCA to more than 2 views ¹⁰

$$\mathbf{y}_i^{(w)} = \mathbf{A}^{(w)} \mathbf{f}_i^{(0)} + \mathbf{B}^{(w)} \mathbf{f}_i^{(w)} + \boldsymbol{\epsilon}_i^{(w)}, \quad \text{for } w = 1, \dots, m$$

$$\mathbf{y}_i = \boldsymbol{\Lambda} \mathbf{f}_i + \boldsymbol{\epsilon}_i,$$

$$\boldsymbol{\Lambda} = \begin{bmatrix} \mathbf{A}^{(1)} & \mathbf{B}^{(1)} & \dots & \mathbf{0} \\ \mathbf{A}^{(2)} & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \mathbf{0} \\ \mathbf{A}^{(m)} & \mathbf{0} & \mathbf{0} & \mathbf{B}^{(m)} \end{bmatrix}$$

Multi-view Learning¹⁰

- ▶ Column-wise shrinkage: ARD prior on the loading columns

$$\begin{aligned}\lambda_{jh}^w &\sim \mathcal{N}\left(0, \left(\alpha_h^{(w)}\right)\right) \\ \alpha_h^{(w)} &\sim \text{Ga}(a_0, b_0)\end{aligned}$$

- ▶ Caveats:
 - ▶ Structured covariance matrix: Large m intractable
 - ▶ Can not capture covariance among arbitrary subset of data
 - ▶ ARD prior: Only induces column-wise sparsity

Multi-view Learning

$$\mathbf{y}_i = \Lambda \mathbf{f}_i + \boldsymbol{\epsilon}_i,$$

Structured sparsity prior on loadings ¹¹: sparse and dense factors

$$\begin{array}{ll} \text{Global} & \left\{ \begin{array}{l} \gamma^{(w)} \sim \text{Ga}(f, \nu) \\ \eta^{(w)} \sim \text{Ga}(e, \gamma^{(w)}) \end{array} \right. \\ \text{Factor — specific} & \left\{ \begin{array}{l} \tau_h^{(w)} \sim \text{Ga}(d, \eta^{(w)}) \\ \phi_h^{(w)} \sim \text{Ga}(c, \tau_h^{(w)}) \end{array} \right. \\ \text{Local} & \left\{ \begin{array}{l} \delta_{jh}^w \sim \text{Ga}(b, \phi_h^{(w)}) \\ \theta_{jh}^{(w)} \sim \pi^{(w)} \text{Ga}(a, \delta_{jh}^w) + (1 - \pi^{(w)}) \delta_{\phi_h^{(w)}(\cdot)} \end{array} \right. \\ & \lambda_{jh}^w \sim \mathcal{N}(0, \theta_{jh}^{(w)}) \end{array}$$

Bibliography I

- [1] Michael E. Tipping and Christopher M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622, 1999.
- [2] T. W. Anderson and Herman Rubin. Statistical inference in factor analysis. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 5: Contributions to Econometrics, Industrial Research, and Psychometry*, pages 111–150, Berkeley, Calif., 1956. University of California Press.
- [3] Joshua Chan, Roberto Leon-Gonzalez, and Rodney W. Strachan. Invariant inference and efficient computation in the static factor model. *Journal of the American Statistical Association*, 0(ja):0–0, 2017.

Bibliography II

- [4] Hedibert Freitas Lopes and Mike West. Bayesian model assessment in factor analysis. *Statistica Sinica*, 14(1):41–67, 2004.
- [5] Sylvia Kaufmann and Christian Schumacher. Bayesian estimation of sparse dynamic factor models with order-independent identification. Technical Report 13.04, April 2013.
- [6] David Knowles and Zoubin Ghahramani. Nonparametric bayesian sparse factor models with application to gene expression modeling. *Ann. Appl. Stat.*, 5(2B):1534–1552, 06 2011. doi: [10.1214/10-AOAS435](https://doi.org/10.1214/10-AOAS435).
- [7] Veronika Ročková and Edward I. George. Fast bayesian factor analysis via automatic rotations to sparsity. *Journal of the American Statistical Association*, 111(516):1608–1622, 2016.

Bibliography III

- [8] Samuel J. Harrison, Mark W. Woolrich, Emma C. Robinson, Matthew F. Glasser, Christian F. Beckmann, Mark Jenkinson, and Stephen M. Smith. Large-scale probabilistic functional modes from resting state fmri. *NeuroImage*, 109:217 – 231, 2015.
- [9] Arto Klami, Seppo Virtanen, and Samuel Kaski. Bayesian canonical correlation analysis. *Journal of Machine Learning Research*, 14(Apr):965–1003, 2013.
- [10] Arto Klami, Seppo Virtanen, Eemeli Leppäaho, and Samuel Kaski. Group factor analysis. *IEEE transactions on neural networks and learning systems*, 26(9):2136–2147, 2015.
- [11] Shiwen Zhao, Chuan Gao, Sayan Mukherjee, and Barbara E Engelhardt. Bayesian group factor analysis with structured sparsity. *The Journal of Machine Learning Research*, 17(1): 6868–6914, 2016.