



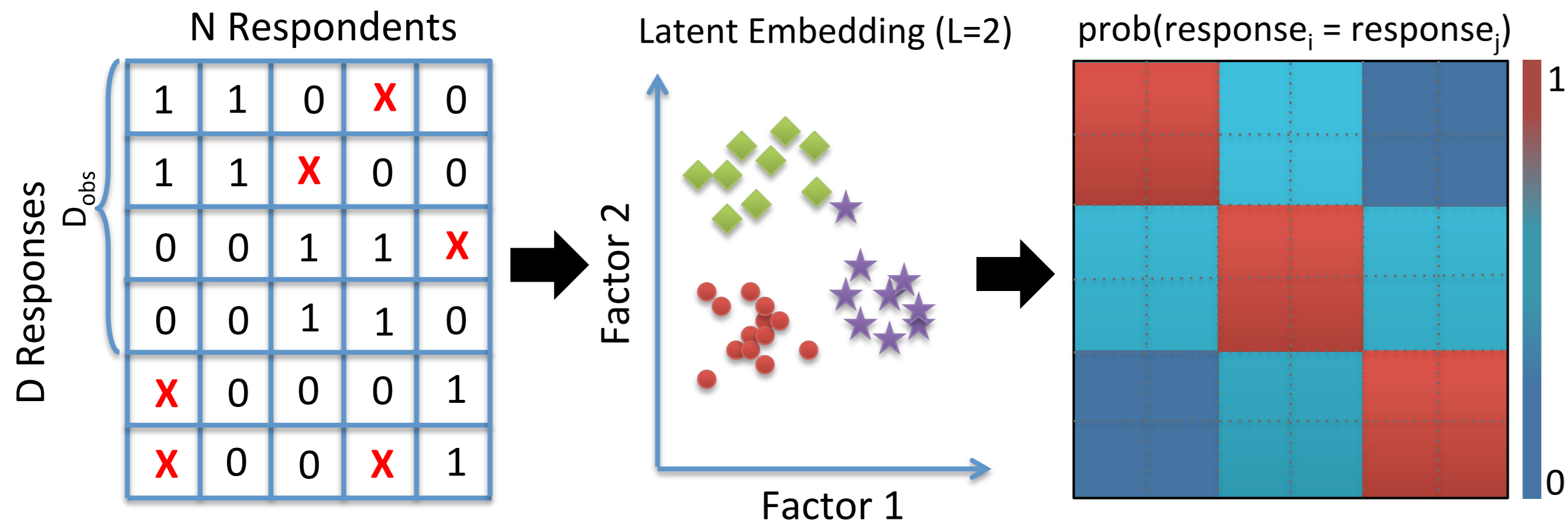
# Large-Scale Approximate Bayesian Learning for Non-Conjugate Latent Gaussian Models

Mohammad Emtyaz Khan, Shakir Mohamed and Kevin P. Murphy

University of British Columbia

## Introduction

**Motivation:** Analysis of high-dimensional data is essential in applications in information retrieval, econometrics, social sciences, and medical diagnostics. Such analysis can be carried out using latent Gaussian models, such as factor analysis and PCA.



### Contributions:

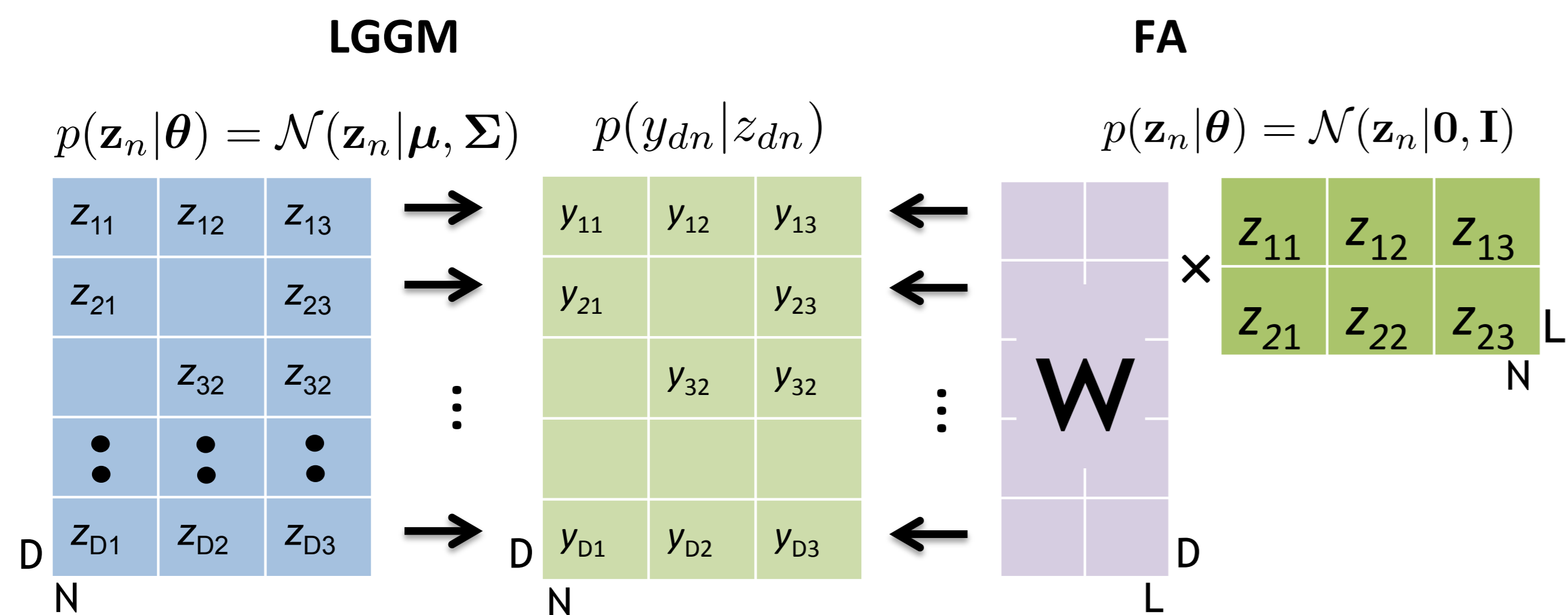
- We derive a variational EM algorithm where the lower bounds are obtained using **concave optimization**.
- For sparse data where  $D_{obs} \ll D < N$ , we **accelerate the E and M steps** resulting in a faster algorithm.

Non-conjugate Factor Analysis (FA) does not scale well to large data. For data with  $N$  observations,  $D$ -dimensional features and  $L$  latent factors, the computational cost is  $O(NL^3 + NDL^2)$  and the memory cost is  $O(NL^2)$ . For large data,  $L$  is usually close to  $D$ , bringing these costs for computation and memory closer to  $O(ND^3)$  and  $O(ND^2)$ , respectively.

- We use **latent Gaussian graphical models** as an alternative to factor analysis, for which we derive **fast inference and learning** algorithms.
- When the data has many missing values, we obtain a reduction in:
  - Computational cost from  $O(ND^3)$  to  $O(ND_{obs}^3)$
  - Memory cost from  $O(ND^2)$  to  $O(D^2)$

## Latent Gaussian Graphical Models

Latent Gaussian Graphical Models (LGGMs) use a correlated latent Gaussian vector of length equal to the data dimension  $D$ , unlike models such as FA that use a low-dimensional vector of length  $L \ll D$ . For the  $n$ th data vector, the generative process is: (1) Sample latent Gaussian vectors  $\mathbf{z}_n | \theta \in \mathbb{R}^D$ . (2) Draw data  $\mathbf{y}_n$  from an appropriate distribution given  $\mathbf{z}_n$ .



## Concave Lower Bounds

**Variational lower bound:** Computation of the marginal likelihood is intractable since the likelihood is not conjugate to the Gaussian prior. Using Jensen's inequality, we can obtain a lower bound to the marginal likelihood. We assume  $\mu = 0$ .

$$\mathcal{L}(\theta) = \sum_{n=1}^N \log \int \prod_{d \in D_{obs}} p(y_{dn} | \mathbf{z}, \theta) \mathcal{N}(\mathbf{z} | \mathbf{0}, \Sigma) d\mathbf{z} = \sum_{n=1}^N \log \int \frac{\prod_{d \in D_{obs}} p(y_{dn} | \mathbf{z}, \theta) \mathcal{N}(\mathbf{z} | \mathbf{0}, \Sigma)}{\mathcal{N}(\mathbf{z} | \mathbf{m}_n, \mathbf{V}_n)} \mathcal{N}(\mathbf{z} | \mathbf{m}_n, \mathbf{V}_n) d\mathbf{z}$$

$$\geq \sum_{n=1}^N \max_{\mathbf{m}_n, \mathbf{V}_n} -KL[\mathcal{N}(\mathbf{m}_n, \mathbf{V}_n) | \mathcal{N}(\mathbf{0}, \Sigma)] + \sum_{d \in D_{obs}} \int \log p(y_{dn} | \mathbf{z}, \theta) \mathcal{N}(\mathbf{z} | \mathbf{m}_n, \mathbf{V}_n) d\mathbf{z}$$

The expectation  $\mathbb{E}[\log p(y_{dn} | \mathbf{z})]$  is not available for most non-Gaussian distributions. We obtain a tractable bound using a local variational bound, such that  $\mathbb{E}[\log p(y_{dn} | \mathbf{z})] \geq f_b(y_{dn}, \mathbf{m}_{dn}, \mathbf{V}_{dn})$ .

$$\geq \sum_{n=1}^N \max_{\mathbf{m}_n, \mathbf{V}_n} \frac{1}{2} [\log |\mathbf{V}_n \Sigma^{-1}| - \text{tr}(\mathbf{V}_n \Sigma^{-1}) - \mathbf{m}_n^T \Sigma^{-1} \mathbf{m}_n + D] + \sum_{d \in D_{obs}} f_b(y_{dn}, \mathbf{m}_{dn}, \mathbf{V}_{dn}) \quad (1)$$

### Bound Optimization:

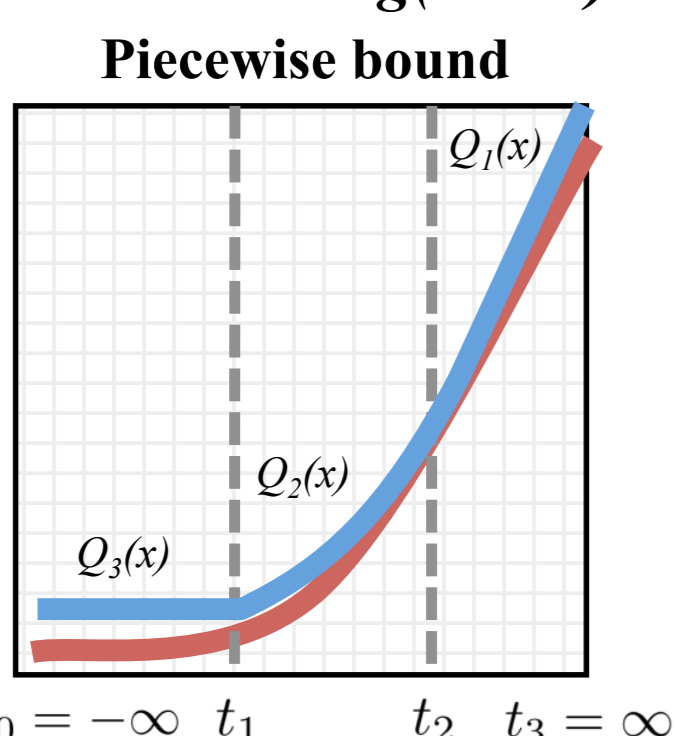
- Our variational bound is strictly concave when  $f_b$  is **jointly concave** with respect to  $\mathbf{m}_n, \mathbf{V}_n$ .
- Given  $\mathbf{V}_n$ , optimization w.r.t.  $\mathbf{m}_n$  is a **non-linear least-squares** function.
- Given  $\mathbf{m}_n$ , optimization w.r.t.  $\mathbf{V}_n$  is a form of **covariance selection** or graphical Lasso.

## Local Variational Bounds

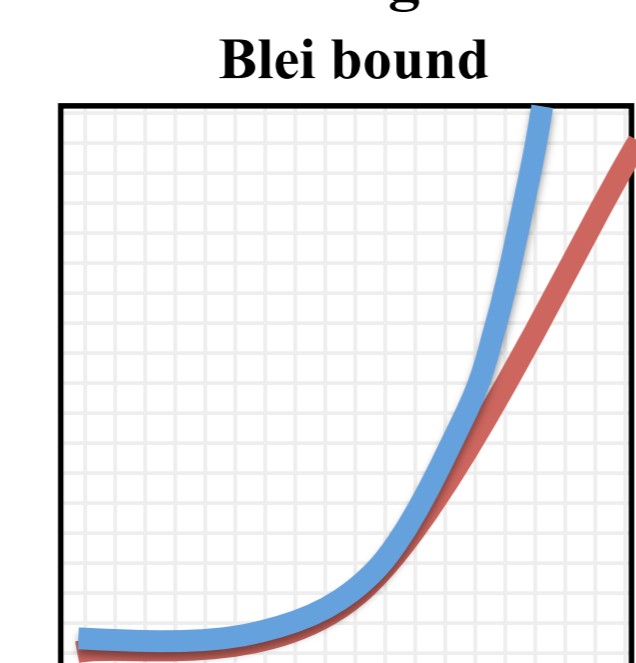
| Type        | Distribution      | $p(y \eta)$                                       | $\mathbb{E}[\log p(y \eta)]$  |
|-------------|-------------------|---|---|
| Count       | Poisson           | $p(y = k   \eta) = \frac{e^{-\eta} \eta^k}{k!}$   | $\eta m - \exp(m + \eta/2) - \log y!$   |
| Binary      | Bernoulli logit   | $p(y = 1   \eta) = \sigma(\eta)$                  | $\eta m - \mathbb{E}[\text{llp}(\eta)]$   |
| Categorical | Multinomial logit | $p(y = k   \eta) = e^{\eta_k - \text{lse}(\eta)}$ | $\mathbf{y}^T \mathbf{m} - \mathbb{E}[\text{lse}(\eta)]$                                  |
| Ordinal     | Cumulative logit  | $p(y \leq k   \eta) = \sigma(\phi_k - \eta)$      | $m - \mathbb{E}[\text{llp}(-\phi_y + \eta)] - \mathbb{E}[\text{llp}(-\phi_{y-1} + \eta)]$ |

The distributions above can be expressed in terms of  $\text{llp}(x) = \log(1 + \exp(x))$ , the logistic log-partition function (LLP), or the log-sum-exp (lse) function. Many convex bounds exist for these functions, ensuring concavity of the entire objective function.

### Bound for $\log(1 + e^x)$



### Bound for log-sum-exp



## Large-Scale Variational Inference

The number of variational parameters, i.e.  $\mathbf{m}_{1:N}$  and  $\mathbf{V}_{1:N}$ , is  $O(ND^2)$ , making optimization slow and inefficient. We **reparameterize** the objective function to reduce the number of parameters and the computational cost.

For sparse data, we only work with the observed entries. The set of observed entries will be indexed by  $O$ . We let  $[\mathbf{a}_O]$  be a vector of length  $D$ , with entries  $O$  equal to  $\mathbf{a}$  and the remaining entries being zero. Similarly,  $[\mathbf{A}_{OO}]$  is a  $D \times D$  matrix with zeros except the submatrix  $\mathbf{A}_{OO}$ . For simplicity, we drop the subscript  $n$ .

Using the derivative of the lower bound w.r.t.  $\mathbf{m}$  and  $\mathbf{V}$ , we can reparameterize the objective function in terms of new parameters  $\alpha$  and  $\lambda$ :

$$\text{Since } -\Sigma^{-1} \mathbf{m} + \left[ \frac{\partial f_b}{\partial \mathbf{m}} \right] = 0 \text{ allows us to reparameterize } \mathbf{m} = \Sigma_{\cdot, O} \alpha_O$$

$$\text{Similarly, } \mathbf{V}^{-1} - \Sigma^{-1} + \text{diag} \left( \left[ \frac{\partial f_b}{\partial \mathbf{V}} \right] \right) = 0 \text{ gives us } \mathbf{V} = (\Sigma^{-1} + \text{diag}([\lambda_O]))^{-1} = \Sigma - \Sigma_{\cdot, O} \underbrace{(\text{diag}(\lambda_O)^{-1} + \Sigma_{OO})^{-1}}_{\mathbf{B}} \Sigma_{O, \cdot}$$

Substituting this into equation 1, the objective functions for the  $n$ th data point is:

$$\mathcal{L}_R(\alpha_O, \lambda_O) = \frac{1}{2} (-\log(|\mathbf{B} \text{diag}(\lambda_O)|)) + \text{trace}(\mathbf{B}^{-1} \Sigma_{OO}) - \alpha_O^T \Sigma_{OO} \alpha_O + \sum_{d \in D_{obs}} f_b(y_d, m_d, V_{dd}) \quad (2)$$

We optimize w.r.t.  $\alpha$  and  $\lambda$ , bringing down the number of parameters to  $O(ND_{obs})$ . Computation of  $\mathbf{V}$  requires inversion of  $\mathbf{B}$ , making the computation  $O(ND_{obs}^3)$ .

## Large-Scale Variational Learning

Using the reparameterization, the update for  $\Sigma$  can be written in terms of  $\alpha_{On}$  and  $\mathbf{B}_n$  as:

$$\Sigma = \frac{1}{N} \sum_{n=1}^N \mathbf{V}_n + \mathbf{m}_n \mathbf{m}_n^T = \Sigma + \frac{1}{N} \Sigma \left( \sum_{n=1}^N [\alpha_{On} \alpha_{On}^T] - [\mathbf{B}_n^{-1}] \right) \Sigma$$

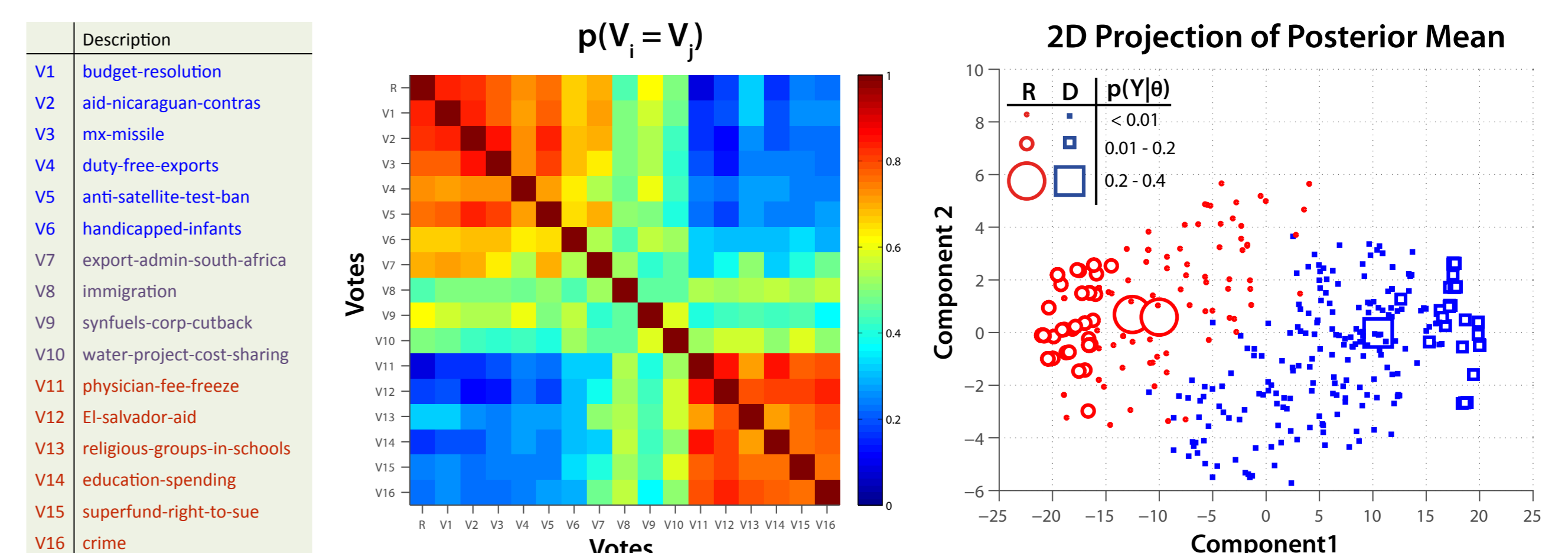
Thus, we never form the matrix  $\mathbf{V}_n$  and only work with  $\mathbf{B}_n$  by maintaining a  $D \times D$  matrix  $\mathbf{S}$  and filling parts of it with  $\mathbf{B}_n^{-1}$ . This reduces the memory cost to  $O(D^2)$ .

### Algorithm: EM for LGGM learning with missing entries

- Initialize  $\Sigma$ .
- Iterate the following until convergence.
  - Form  $\mathbf{S}$  of size  $D \times D$  with all entries set to 0.
  - For  $n = 1, \dots, N$ 
    - Get the observed data  $\mathbf{y}_{On}$  and corresponding subblock  $\Sigma_{OO}$ .
    - Compute  $\alpha_{On}$  and  $\lambda_{On}$  by maximizing equation 2.
    - $\mathbf{S}_{OO} \leftarrow \mathbf{S}_{OO} + [\alpha_{On} \alpha_{On}^T] - [\mathbf{B}_n^{-1}]$ .
  - Compute  $\Sigma \leftarrow \Sigma + \Sigma \mathbf{S} \Sigma / N$ .

## Results

**Illustrating LGGM on voting dataset** ( $D = 16, N = 435$ ). We plot the probability of two variables (votes) taking the same value. We also plot a 2-D projection of the posterior mean of the latent variables with size proportional to the log-likelihood.



**Large Synthetic Data** ( $D = 100, N = 240$ ) We generate a large binary synthetic dataset consisting of 12 patterns. Each pattern is replicated 20 times and 5% of the entries are flipped. We artificially create missing values in this dataset, such that only 35% of the dataset is observed. We impute the missing values using LGGM and FA.

