

# Notes: MPLN model

## 1 Multivariate Poisson Lognormal (MPLN) model

The MPLN model (?) relates some  $p$ -dimensional observation vectors  $\mathbf{Y}_i$  to some  $p$ -dimensional vectors of Gaussian latent variables  $\mathbf{Z}_i$  as follows

$$\begin{aligned}\mathbf{Y}_i | \mathbf{Z}_i &\sim \text{Poisson}(e^{\mathbf{Z}_i}) \quad , \quad Y_{ij} | Z_{ij} \text{ indep.} \\ \mathbf{Z}_i &\sim \mathcal{N}(\mathbf{X}_i^\top \mathbf{B}_{d \times p}, \Sigma_{p \times p})\end{aligned}$$

where

- $Y_{ij}$  represents the  $j$ th code of  $i$ th patient with  $i = 1, \dots, n$  and  $j = 1, \dots, p$ ,  $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_n)^\top \in \mathbb{R}^{n \times p}$ ,
- $\mathbf{Z}_i$  is the patient-level embedding of  $i$ th patient,  $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)^\top \in \mathbb{R}^{n \times p}$ ,
- $\mathbf{X}_i$  is the covariates of  $i$ th patient,  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^\top \in \mathbb{R}^{n \times d}$

## 2 Leveraging external information for $\Sigma$

In the MPLN model,  $\Sigma$  describes the underlying residual relationships between the  $p$  codes. We may decompose  $\Sigma = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ , where  $\mathbf{V} \in \mathbb{R}^{p \times q}$  is orthonormal and the rank  $q \leq p$ ,  $\mathbf{\Lambda} \in \mathbb{R}^{q \times q}$  is diagonal. Assuming that across different institution/dataset, the orientation/rotation ( $\mathbf{V}$ ) is the same but the shape/scaling  $\mathbf{\Lambda}$  can be different, we may borrow  $\mathbf{V}$  that is obtained from external sources for model estimation. In that case, instead of estimating the full covariance matrix  $\Sigma$ , we only need to estimate the diagonal matrix  $\mathbf{\Lambda}$ .

An alternative way to write the model:

$$\begin{aligned}\mathbf{Y}_i | \mathbf{Z}_i &\sim \text{Poisson}(e^{\mathbf{Z}_i}) \quad , \quad Y_{ij} | Z_{ij} \text{ indep.} \\ \mathbf{Z}_i &= \mathbf{X}_i^\top \mathbf{B} + \mathbf{V}\mathbf{\Lambda}^{1/2}\mathbf{W}_i \quad , \\ \mathbf{W}_i &\sim \mathcal{N}(0_q, \mathbf{I}_{q \times q}) \quad iid\end{aligned}$$

### 2.1 Likelihood

In matrix notation,  $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{W}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top$ . The observation matrix  $\mathbf{Y}$  only depends on  $\mathbf{Z}$  through  $\mathbf{B}$ ,  $\mathbf{\Lambda}$  and  $\mathbf{W}$ , and the complete log-likelihood is therefore

$$\begin{aligned}\log p(\mathbf{Y}, \mathbf{W}; \mathbf{B}, \mathbf{\Lambda}) &= \sum_{i=1}^n \log p(\mathbf{Y}_i | \mathbf{W}_i; \mathbf{B}, \mathbf{\Lambda}) + \log p(\mathbf{W}_i) \\ &= \mathbf{1}_n^\top \left\{ \mathbf{Y} \odot (\mathbf{X}\mathbf{B} + \mathbf{W}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top) - \exp(\mathbf{X}\mathbf{B} + \mathbf{W}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top) \right\} \mathbf{1}_p \\ &\quad - \frac{\|\mathbf{W}\|_F^2}{2} - \frac{nq}{2} \log(2\pi) - \mathbf{1}_n^\top \log(\mathbf{Y}!) \mathbf{1}_p\end{aligned}$$

## 2.2 Variational bound of the likelihood

Maximize the evidence lower bound:

$$\begin{aligned}\mathcal{J}(\mathbf{B}, \mathbf{\Lambda}, q) &= \log p(\mathbf{Y}; \mathbf{B}, \mathbf{\Lambda}) - KL[q(\mathbf{W}), p(\mathbf{W} | \mathbf{Y}; \mathbf{B}, \mathbf{\Lambda})] \\ &= E_q[\log p(\mathbf{Y}, \mathbf{W}; \mathbf{B}, \mathbf{\Lambda})] - E_q[\log q(\mathbf{W})]\end{aligned}$$

Choose the set  $\mathcal{Q}$  of product distribution of  $q$ -dimensional multivariate Gaussian with diagonal covariance matrices:

$$\mathcal{Q} = \left\{ q : q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i), q_i(\mathbf{Z}_i) = \mathcal{N}(\mathbf{Z}_i; \mathbf{m}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i)), (\mathbf{m}_i, \mathbf{s}_i) \in \mathbb{R}^q \times \mathbb{R}_+^q \right\}$$

Let  $\mathbf{M} = (\mathbf{m}_1, \dots, \mathbf{m}_n)^\top \in \mathbb{R}^{n \times q}$  and  $\mathbf{S} = (\mathbf{s}_1, \dots, \mathbf{s}_n)^\top \in \mathbb{R}^{n \times q}$ . Results on first- and second-order moments of multivariate Gaussian show that

$$\begin{aligned}\mathcal{J}(\mathbf{B}, \mathbf{\Lambda}, \mathbf{M}, \mathbf{S}) &= \mathbf{1}_n^\top \left[ \mathbf{Y} \odot (\mathbf{X}\mathbf{B} + \mathbf{M}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top) - \mathbb{E}_q \left\{ \exp(\mathbf{X}\mathbf{B} + \mathbf{W}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top) \right\} \right] \mathbf{1}_p \\ &\quad - \frac{1}{2} \mathbf{1}_n^\top \{ \mathbf{M} \odot \mathbf{M} + \mathbf{S} \odot \mathbf{S} - 2 \log(\mathbf{S}) - \mathbf{1}_{n \times q} \} \mathbf{1}_q - \mathbf{1}_n^\top \log(\mathbf{Y}!) \mathbf{1}_p\end{aligned}$$

where

$$\mathbb{E}_q \left\{ \exp(\mathbf{X}\mathbf{B} + \mathbf{W}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top) \right\} = \exp \left\{ \mathbf{X}\mathbf{B} + \mathbf{M}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top + \frac{1}{2}(\mathbf{S} \odot \mathbf{S})(\mathbf{\Lambda}^{1/2}\mathbf{V}^\top \odot \mathbf{\Lambda}^{1/2}\mathbf{V}^\top) \right\} =: \mathbf{A}$$

## 2.3 Blockwise gradient

Find blockwise gradients of  $\mathcal{J}(\mathbf{B}, \mathbf{\Lambda}, \mathbf{M}, \mathbf{S})$ :

$$\begin{aligned}\frac{\partial \mathcal{J}}{\partial \mathbf{B}} &= (\mathbf{Y} - \mathbf{A})^\top \mathbf{X} \\ \frac{\partial \mathcal{J}}{\partial \mathbf{\Lambda}^{1/2}} &= \mathbf{V}^\top \left[ (\mathbf{Y} - \mathbf{A})^\top \mathbf{M} - \{ \mathbf{A}^\top (\mathbf{S} \odot \mathbf{S}) \} \odot \mathbf{V} \mathbf{\Lambda}^{1/2} \right] \in \mathbb{R}^{q \times q} \\ \frac{\partial \mathcal{J}}{\partial \mathbf{M}} &= (\mathbf{Y} - \mathbf{A}) \mathbf{V} \mathbf{\Lambda}^{1/2} - \mathbf{M} \\ \frac{\partial \mathcal{J}}{\partial \mathbf{S}} &= \mathbf{S}^\odot - \mathbf{S} - 2\mathbf{A}(\mathbf{V} \mathbf{\Lambda}^{1/2} \odot \mathbf{V} \mathbf{\Lambda}^{1/2}) \odot \mathbf{S}\end{aligned} \tag{1}$$

where  $\mathbf{S}^\odot$  is the element-wise inverse of  $\mathbf{S}$ ,  $S_{ij}^\odot = S_{ij}^{-1}$ . Unlike the standard MPLN, there is no closed form solution for  $\mathbf{B}$  and  $\mathbf{\Lambda}^{1/2}$  with  $\mathbf{M}$  and  $\mathbf{S}$  fixed. However, since the blockwise gradients can be easily derived, we may still use gradient-based local optimization algorithms for maximizing the ELBO.

In  $\mathbf{\Lambda}^{1/2}$ , we only have  $q$  parameters instead of  $q \times q$ .

Current implementation: optimization using NLOPT library in C++, similar to ?.

## 2.4 Simulation (Test code...)

Let  $n = 5000$ ,  $p = 50$ . Consider two scenarios: true covariance has full rank and has low rank.

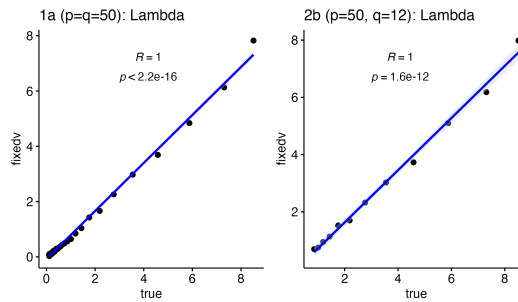
- **Scenario 1 (true covariance with full rank):**

1. Generate  $\Sigma = AR_1(0.8)$ . Find its eigen-decomposition  $\Sigma = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ .
2. Generate  $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \Sigma)$  and  $\mathbf{Y} \sim \text{Poisson}(e^{\mathbf{Z}})$ .
3. Fit model with  $\mathbf{V}$ ,
  - a. (correct) assuming  $\Sigma$  has full rank,  $\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$ , we use the full  $\mathbf{V}$  here.
  - b. (low rank approximation) assuming  $\Sigma$  has lower rank,  $\mathbf{\Lambda} \in \mathbb{R}^{q \times q}$  with  $q = 20, 10, 5$ , we only use  $\mathbf{V}_{[:,1:q]}$  here.

- **Scenario 2 (true covariance with low rank,  $q = 12$ ):**

1. Generate  $\Sigma = AR_1(0.8)$ . Find its eigen-decomposition  $\Sigma = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ .
2. Let  $\tilde{\Sigma} = \mathbf{V}_{[:,1:q]}\mathbf{\Lambda}_{[1:q,1:q]}\mathbf{V}_{[:,1:q]}^\top$ . Find its eigen-decomposition  $\tilde{\Sigma} = \tilde{\mathbf{V}}\tilde{\mathbf{\Lambda}}\tilde{\mathbf{V}}^\top$ . This will be our true covariance with known  $\text{rank} = q$ . (In reality we observe  $\tilde{\Sigma}$  from external sources)
3. Generate  $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma})$  and  $\mathbf{Y} \sim \text{Poisson}(e^{\mathbf{Z}})$ .
4. Fit model with  $\tilde{\mathbf{V}}$ ,
  - a. (neglecting low-rankness) assuming  $\Sigma$  has full rank,  $\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$ , we use the full  $\tilde{\mathbf{V}}$  here.
  - b. (correct) assuming  $\Sigma$  has lower rank,  $\mathbf{\Lambda} \in \mathbb{R}^{q \times q}$  with  $q = 20, 10, 5$ , we only use  $\tilde{\mathbf{V}}_{[:,1:q]}$  here. (In reality we don't know the true rank, so would also need to implement a procedure to select for the best  $q$ .)

- Compare  $\hat{\mathbf{\Lambda}}$  with the  $\mathbf{\Lambda}$  (for scenario 1) or  $\tilde{\mathbf{\Lambda}}$  (for scenario 2). Here I'm only showing 1a and 2b (with correct  $q = 12$ ), since the dimension of  $\hat{\mathbf{\Lambda}}$  doesn't match with the truth in 1b and 2a.



- Compare model-generated embeddings v.s. the true  $\mathbf{Z}$ . *fixedV* is our implementation and *full* is the original implementation (estimating the full  $\Sigma$  based on data). Results on next page.

Since *fixedV* has much less parameters even if  $p = q$ , the variance should be smaller compared with *full*?

- Look at the estimated mean and covariance: **problem: estimated mean is biased? looking into this**

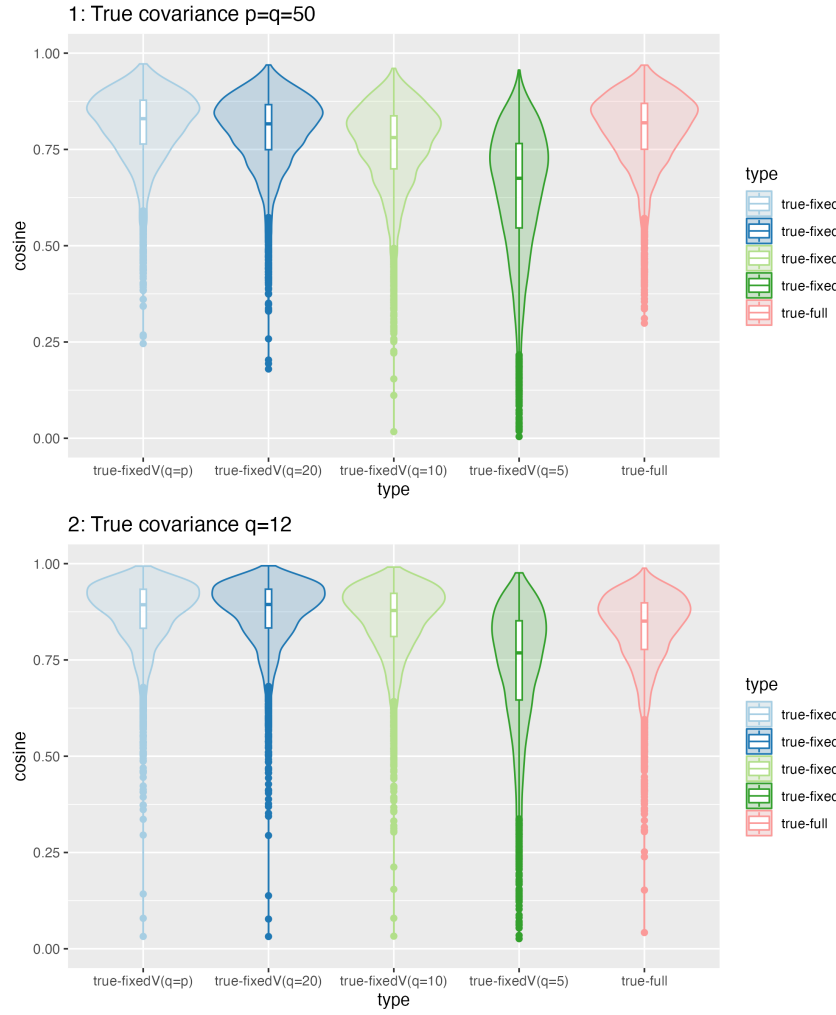


Figure 1: Cosine similarities between the latent variable  $\mathbf{Z}$  obtained with free  $\Sigma$  and fixed eigenvectors  $\mathbf{V}$ .

## Notes 4/5/2023

So far...

- Issue found in simulation: Parameter estimates are biased with `PLNModel` package even with small  $p$  and very large  $n$ .
- What causes the bias?
  - Based on ?, Variational Gaussian Approximation (VGA) for our model (iid and full rank) is inconsistent.
- Potential remedies?
  - (*known  $\Sigma$ , iid, full rank*) If we know the covariance parameters, try to re-parameterize the model and leverage independence in the VGA alg
  - (*unknown  $\Sigma$ , iid, low rank*) If we use the current VGA alg, need to allow  $p$  to grow with  $n$  while constraining rank  $q$  to be relatively small, so that parameter estimation is consistent in iid setting.
  - (*unknown  $\Sigma$ , iid, full rank*) Explore other methods that might help reduce the bias, e.g. black box variational inference (BBVI) (?).

### Why biased?

*Theory of Gaussian Variational Approximation for a Poisson Mixed Model. ?*

- They considered repeated measures data:
  - Data:  $(X_{ij}, Y_{ij})$ ,  $1 \leq i \leq m$  corresponds to  $m$  patients,  $1 \leq j \leq n$  corresponds to  $n$  repeated measures on those patients.
  - Model:  $Y_{ij} \mid (X_{ij}, U_i)$  indep Poisson with mean  $\exp(\beta_0 + \beta_1 X_{ij} + U_i)$ ,  $U_i$  independent  $N(0, \sigma^2)$
  - Results: Estimators of the model parameters based on Gaussian variational approximation is consistent at rate  $m^{-1/2} + n^{-1}$ .
    - \* Number of model parameters:  $d + 1$
    - \* Number of variational parameters:  $2m$
    - \* Essentially, need to use the  $n$  repeated measures from each patient  $i$  to estimate the variational parameters  $m_i$  and  $s_i$ , so  $n$  cannot be too small.
- In our notations, their model is equivalent to

$$\begin{aligned} \mathbf{Y}_i \mid \mathbf{Z}_i &\sim \text{Poisson}(e^{\mathbf{Z}_i}) \quad , \quad Y_{ij} \mid Z_{ij} \text{ indep.} \\ \mathbf{Z}_i &= \mathbf{X}_i^\top \mathbf{B} + \boldsymbol{\epsilon}_i \quad , \quad \boldsymbol{\epsilon}_i = \epsilon_i \mathbf{1}_p \quad , \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad iid \end{aligned}$$

v.s. our model

$$\begin{aligned} \mathbf{Y}_i \mid \mathbf{Z}_i &\sim \text{Poisson}(e^{\mathbf{Z}_i}) \quad , \quad Y_{ij} \mid Z_{ij} \text{ indep.} \\ \mathbf{Z}_i &= \mathbf{X}_i^\top \mathbf{B} + \boldsymbol{\epsilon}_i \quad , \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}_p, \boldsymbol{\Sigma}_{p \times p}) \quad iid \end{aligned}$$

- Number of model parameters:  $dp + p(p+1)/2$
- Number of variational parameters:  $2np \rightarrow$  too large, won't be able to estimate properly

- Their model is similar (though not exactly the same) to our model if we assume  $\Sigma$  has rank  $q = 1$ . In that case number of variational parameters can be reduce to  $2nq = 2n$ . Run a simulation to check: fix  $rank = 1$ , generate data with varying  $p$  and  $n$ , check bias.  $\rightarrow$  verified that with large enough  $n$ , bias goes to 0 as  $p$  grows.

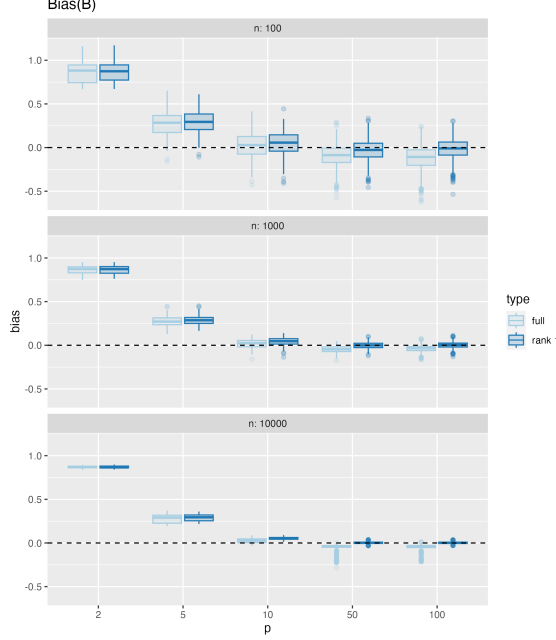


Figure 2: *full*: VGA with  $2np$  variational parameters; *rank 1*: VGA with  $2n$  variational parameters. Here I'm only showing bias of  $B$  but the bias of  $\Sigma$  has a similar pattern.

## Potential remedies?

### 1. With Known $\Sigma$ , modify VGA

Our model:

$$\mathbf{Y}_i | \mathbf{Z}_i \sim \text{Poisson}(e^{\mathbf{Z}_i}) \quad , \quad Y_{ij} | Z_{ij} \text{ indep.}$$

$$\mathbf{Z}_i = \mathbf{X}_i^\top \mathbf{B} + \mathbf{V} \mathbf{\Lambda}^{1/2} \mathbf{W}_i \quad ,$$

$$\mathbf{W}_i \sim \mathcal{N}(0_q, \mathbf{I}_{q \times q}) \quad iid$$

- Re-parameterize:  $\mathbf{W}_i = \mathbf{\Lambda}^{-1/2} \mathbf{V}^\top (\mathbf{Z}_i - \mathbf{X}_i^\top \mathbf{B}) \sim \mathcal{N}(0_q, \mathbf{I}_{q \times q}) \quad iid$
- Question: can we leverage the iid nature of  $\mathbf{W}_i$  to use less variational parameters to estimate  $f(\mathbf{W}_i | \mathbf{Y}_i; \mathbf{\Lambda}, \mathbf{V})$ ? For now, consider  $q = p$ ,  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$  known,  $\mu_1, \dots, \mu_p$  unknown,
  - For each  $j = 1, \dots, p$ , let  $p_\theta$  be pdf of  $\text{poisson}(\theta)$  and  $\phi$  be pdf of  $\text{Normal}(0, 1)$ ,

$$\begin{aligned} f(W_{ij} | \mathbf{Y}_i; \mathbf{\Lambda}, \mathbf{V}) &= \frac{\left\{ f(Y_{ij} | W_{ij}) \prod_{j' \neq j} f(Y_{ij'}) \right\} f(W_{ij})}{f(\mathbf{Y}_i)} \\ &= \frac{f(Y_{ij} | W_{ij}) f(W_{ij})}{f(Y_{ij})} \\ &= \frac{p_{\exp(\mu_j + \sigma_j^2 W_{ij})}(Y_{ij}) \phi(W_{ij})}{\int_{\mathbb{R}} p_{\exp(\mu_j + \sigma_j^2 w)}(Y_{ij}) \phi(w) dw} \end{aligned}$$

- In this case, the integral can be computed numerically, e.g, through quadrature or sampling, so no VGA is needed. If so, our estimators should be unbiased, and knowing  $\sigma_j$  is expected to help reduce the variance in our estimation.
- However in general, VGA is still needed for approximation of  $p(\mathbf{W}_i | \mathbf{Y}_i; \mathbf{\Lambda}, \mathbf{V})$ . In that case:
  - \* For  $j \neq k$ ,  $f(W_{ij} | \mathbf{Y}_i; \mathbf{\Lambda}, \mathbf{V}) \neq f(W_{ik} | \mathbf{Y}_i; \mathbf{\Lambda}, \mathbf{V})$  if  $\mu_j \neq \mu_k$  or  $\sigma_j \neq \sigma_k$ , so we still need  $2p$  parameters ( $p$  for mean and  $p$  for variance) for each  $i = 1, \dots, n$ .
  - \* Knowing  $\sigma_j$  will not help reduce the bias of the estimation since the number of variational parameters to be estimated stays the same [pending: ]. It might still help reduce the variance.

## 2. VGA with growing $p$ in low rank

Based on the intuition from ?, VGA for our model might be consistent if we allow  $p$  to grow along with  $n$  while controlling  $q$  to be small relative to  $p$ .

- Run simulation to check intuition: For fixed  $n = 1000$ ,  $p = 50$  or  $100$ , increase  $q$  from 1 to  $p$ .
- *full* is using VGA with  $2np$  variational parameters, and *true rank* is using VGA with  $2nq$  variational parameters.

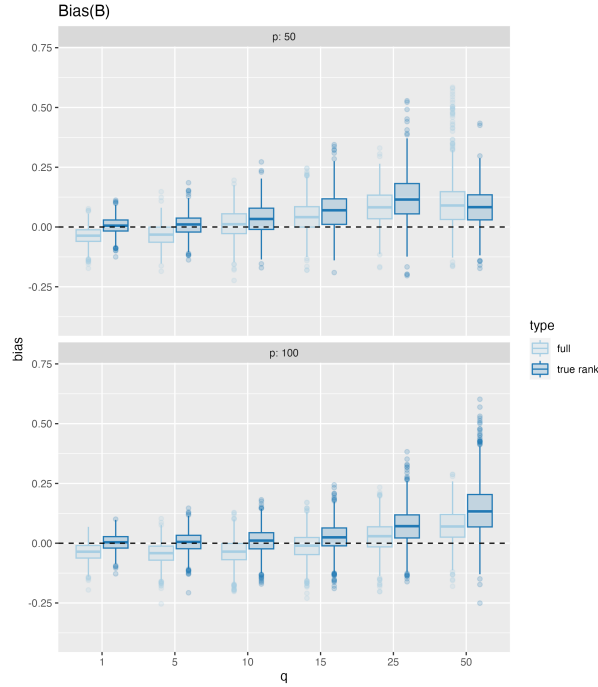


Figure 3: *full*: VGA with  $2np$  variational parameters; *true rank*: VGA with  $2nq$  variational parameters. Here I'm only showing bias of  $B$  but the bias of  $\Sigma$  has a similar pattern.

- As expected, bias for *true rank* with  $2nq$  variational parameters is small with small  $q$ , but increases as  $q$  grows. We do not know how to control the bias for *full* with  $2np$  parameters.
- No theoretical results on consistency rate currently exists under this setting.

### 3. Other methods

**BBVI:** *Black Box Variational Inference*, ?

- Algorithm:

---

**Algorithm 2** Black Box Variational Inference (II)

---

**Input:** data  $x$ , joint distribution  $p$ , mean field variational family  $q$ .  
**Initialize**  $\lambda_{1:n}$  randomly,  $t = 1$ .  
**repeat**  
    // Draw  $S$  samples from the variational approximation  
    **for**  $s = 1$  **to**  $S$  **do**  
         $z[s] \sim q$   
    **end for**  
    **for**  $i = 1$  **to**  $n$  **do**  
        **for**  $s = 1$  **to**  $S$  **do**  
             $f_i[s] = \nabla_{\lambda_i} \log q_i(z[s]|\lambda_i)(\log p_i(x, z[s]) - \log q_i(z[s]|\lambda_i))$   
             $h_i[s] = \nabla_{\lambda_i} \log q_i(z[s]|\lambda_i)$   
        **end for**  
         $\hat{a}_i^* = \frac{\sum_{s=1}^S \text{Cov}(f_i^s, h_i^s)}{\sum_{s=1}^S \text{Var}(h_i^s)}$   
         $\hat{\nabla}_{\lambda_i} \mathcal{L} \triangleq \frac{1}{S} \sum_{s=1}^S f_i[s] - \hat{a}_i^* h_i[s]$   
    **end for**  
     $\rho = t$ th value of a Robbins Monro sequence  
     $\lambda = \lambda + \rho \hat{\nabla}_{\lambda} \mathcal{L}$   
     $t = t + 1$   
**until** change of  $\lambda$  is less than 0.01.

---

- Advantage:
  - Flexible variational family, not limited to Gaussian distribution, as long as
    - \* we know the score function of the variational distribution, and
    - \* we can sample from that distribution
  - It should be computationally fast.
- It is currently unclear whether it can help resolve the consistency issue:
  - No matter what distribution we choose, we still need to estimate  $2np$  variational parameters.
  - The sampling scheme is used for estimating the gradients of the ELBO. Increasing the number of samples  $S$  might make this estimation better, which will likely make convergence faster. However, it does not necessarily help estimate the variational distribution better as essentially we are not gaining any information from this step. Therefore it might not be able to fix our consistency issue.
  - Currently no theoretical guarantee found.
- Still looking into this and trying to understand whether/how it can help.

**MCMC-EM:** A multivariate Poisson-log normal mixture model for clustering transcriptome sequencing data ?

- Algorithm: sample from  $f(\mathbf{W} \mid \mathbf{Y}; \mathbf{A}, \mathbf{V})$  using MCMC, then calculate expected values based on samples
  - Using **RStan** package



- Advantage:
  - If  $p$  not too large when the MCMC can work well, parameter estimates should have consistency rate of mle.
- Concerns:
  - Will not work well for large  $p$ . In the paper they only tested the alg on  $n = 200, p = 6$ . Highly doubt that it will work well in our setting where  $p$  can be in the hundreds.
  - In the paper they didn't show any result on the performance of parameter estimation.
  - Probably too slow if directly applied to our datasets with much larger  $n$  and  $p$ .
- Also looking into other methods based on sampling or approximation that might work...

## 4. Two-step smoothing

**Idea:** approximate  $p(W | X)$  by  $q_x(W) = \mathcal{N}(g(x), h(x))$ , where  $g(x)$  and  $h(x)$  are modeled by NN.

### 4.1 Simulation

With  $n = 10000, 50000, p = 50, 100, 200, q = 5, 10, 20$ ,

#### Data Generation

1. Generate  $\Sigma = AR_1(0.5)$ . Find its eigen-decomposition  $\Sigma = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ .
2. Let  $\tilde{\Sigma} = \mathbf{V}_{[:,1:q]}\mathbf{\Lambda}_{[1:q,1:q]}\mathbf{V}_{[:,1:q]}^\top$ .
3. Generate  $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \tilde{\Sigma})$  and  $\mathbf{Y} \sim \text{Poisson}(e^{\mathbf{Z}})$ .

#### Parameter estimation

1. Find optimal  $\mathbf{B}, \Sigma, \mathbf{M}, \mathbf{S}$  that maximize the ELBO (??) using blockwise gradient descent (?).
2. Fit  $\mathbf{M} \sim \mathbf{X}$  and  $\mathbf{S} \sim \mathbf{X}$  using two-layer FCNN to obtain  $\hat{\mathbf{M}}, \hat{\mathbf{S}}$ 
  - Hidden layer: 30 nodes each, ReLU activation
  - Output layer:  $q$  nodes, linear activation for  $\mathbf{M}$  and softplus activation for  $\mathbf{S}$  (to ensure positiveness)
3. Plug in  $\hat{\mathbf{M}}$  and  $\hat{\mathbf{S}}$  and use blockwise gradient descent to find  $\hat{\mathbf{B}}$  and  $\hat{\Sigma}$

## Results

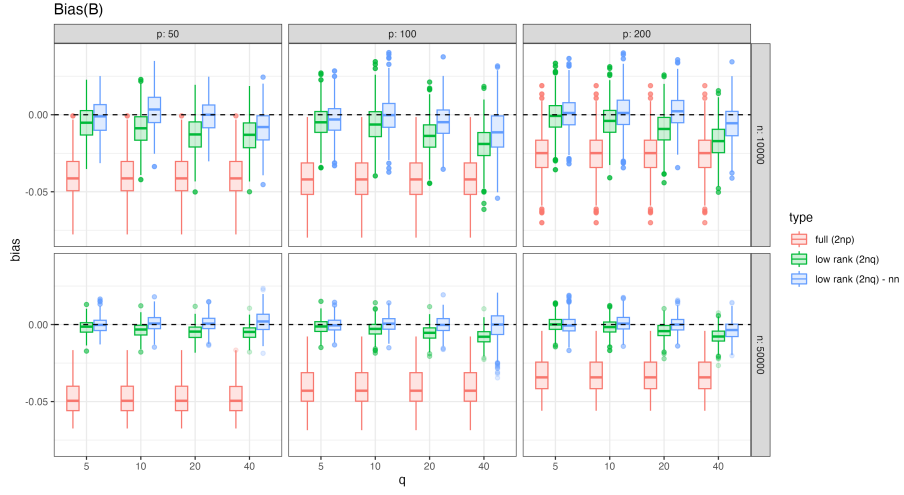


Figure 4: *full*: VGA with  $2np$  variational parameters; *true rank*: VGA with  $2nq$  variational parameters; *true rank*: with additional NN smoothing step. Here I'm only showing bias of  $B$  but the bias of  $\Sigma$  has a similar pattern.

### 3 Notes 4/23/2023

**Model:**

$$\mathbf{X}_i \mid \mathbf{Z}_i \sim \text{Poisson}(e^{\mathbf{Z}_i}) \quad , \quad Y_{ij} \mid Z_{ij} \text{ indep.}$$

$$\mathbf{Z}_i = \mathbf{D}_i^\top \mathbf{B} + \mathbf{V} \mathbf{\Lambda}^{1/2} \mathbf{W}_i \quad ,$$

$$\mathbf{W}_i \sim \mathcal{N}(0_q, \mathbf{I}_{q \times q}) \quad iid$$

**Idea:** approximate  $p(W \mid X)$  by  $q_x(W) = \mathcal{N}(g(x), h(x))$ , where  $g(x)$  and  $h(x)$  are modeled by NN.

#### 3.1 Simulation

With  $n = 10000, 50000$ ,  $p = 50, 100, 200$ ,  $q = 5, 10, 20$ ,

##### 3.1.1 Data Generation

1. Generate  $\mathbf{\Sigma} = AR_1(0.5)$ . Find its eigen-decomposition  $\mathbf{\Sigma} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$ .
2. Let  $\tilde{\mathbf{\Sigma}} = \mathbf{V}_{[:,1:q]} \mathbf{\Lambda}_{[1:q,1:q]} \mathbf{V}_{[:,1:q]}^\top$ .
3. Generate  $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{\Sigma}})$  and  $\mathbf{Y} \sim \text{Poisson}(e^{\mathbf{Z}})$ .

##### 3.1.2 Parameter estimation

1. Find optimal  $\mathbf{B}, \mathbf{\Sigma}, \mathbf{M}, \mathbf{S}$  that maximize the ELBO (??) using blockwise gradient descent (?).
2. Fit  $\mathbf{M} \sim \mathbf{X}$  and  $\mathbf{S} \sim \mathbf{X}$  using two-layer FCNN to obtain  $\widehat{\mathbf{M}}, \widehat{\mathbf{S}}$ 
  - Hidden layer: 30 nodes each, ReLU activation
  - Output layer:  $q$  nodes, linear activation for  $\mathbf{M}$  and softplus activation for  $\mathbf{S}$  (to ensure positiveness)
3. Plug in  $\widehat{\mathbf{M}}$  and  $\widehat{\mathbf{S}}$  and use blockwise gradient descent to find  $\widehat{\mathbf{B}}$  and  $\widehat{\mathbf{\Sigma}}$

##### 3.1.3 Results

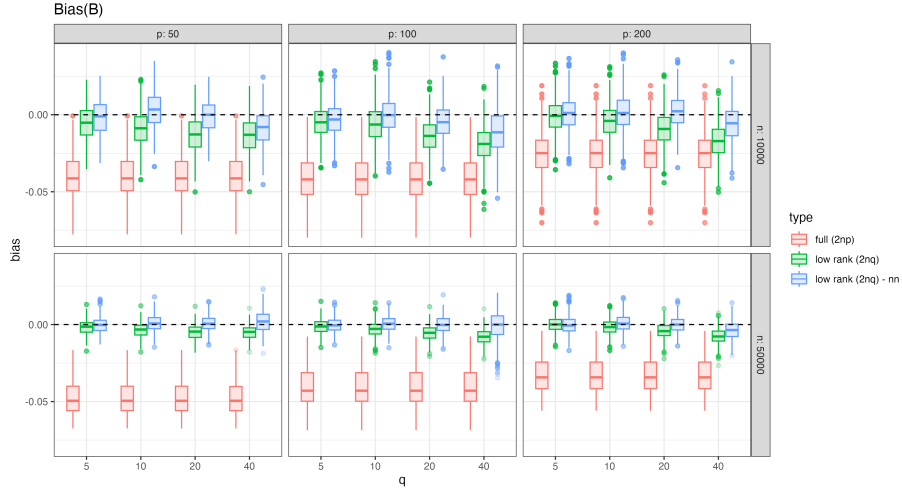
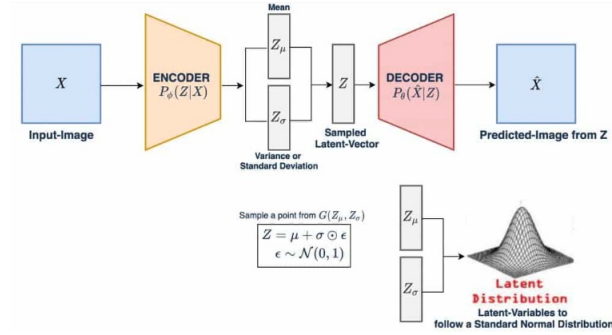


Figure 5: *full*: VGA with  $2np$  variational parameters; *true rank*: VGA with  $2nq$  variational parameters; *true rank*: with additional NN smoothing step. Here I'm only showing bias of  $B$  but the bias of  $\Sigma$  has a similar pattern.

## 4 Notes on Variational Autoencoder (VAE)

### 4.1 VAE Framework



$$p(z) = \mathcal{N}(0, I)$$

$$p(x | z) = \mathcal{N}(f(z), cI) \quad , \quad f \in F$$

Approximate  $p(z | x)$  by a Gaussian distribution

$$q_x(z) = \mathcal{N}(g(x), h(x)) \quad , \quad g \in G \quad h \in H.$$

We are looking for the optimal  $g^*, h^*, f^*$  such that

$$\begin{aligned} (g^*, h^*) &= \arg \min_{(g, h) \in G \times H} KL(q_x(z), p(z | x)) \\ &= \arg \max_{(g, h) \in G \times H} \left\{ \mathbb{E}_{z \sim q_x} \left( -\frac{\|x - f(x)\|^2}{2c} \right) - KL(q_x(z), p(z)) \right\} \\ f^* &= \arg \max_{f \in F} \mathbb{E}_{z \sim q_x^*} (\log p(x | z)) \\ &= \arg \max_{f \in F} \mathbb{E}_{z \sim q_x^*} \left( -\frac{\|x - f(x)\|^2}{2c} \right) \end{aligned}$$

i.e.

$$(f^*, g^*, h^*) = \arg \max_{(f, g, h) \in F \times G \times H} \left\{ \underbrace{\mathbb{E}_{z \sim q_x} \left( -\frac{\|x - f(x)\|^2}{2c} \right)}_{\text{"reconstruction term"}} - \underbrace{KL(q_x(z), p(z))}_{\text{"approximation/regularization term"}} \right\}$$

- the constant  $c$  rules the balance between the two terms: with higher  $c$  we assume high variance around  $f(z)$  for the decoder in the model so we favor the second term over the first term.
- the second term is often referred to as the "regularization term" because it forces the distributions returned by the encoder to be close to a standard normal distribution

– We expect the latent space to be "regular"

- \* Continuity: two close points in the latent space should not give two completely different contents once decoded

- \* Completeness: for a chosen distribution, a point sampled from the latent space should give meaningful content once decoded
- To achieve this we want to regularise both the covariance matrix and the mean of the distributions returned by the encoder
- By forcing it to be standard normal, we require covariance to be close to identity, preventing distributions with tiny variances; and mean close to 0, preventing encoded distributions to be too far apart from each others.

## 4.2 Our Idea

Model:

$$\begin{aligned}
 p(w) &= \mathcal{N}(0, I) \\
 z &= x^\top B + V\Lambda^{1/2}w \\
 p(x \mid z) &= \text{Poisson}(e^z)
 \end{aligned}$$

Let  $\theta = (B, \Sigma)$  denotes the model parameters. If we were to put the model under the VAE framework:  
 Approximate  $p(w \mid x)$  by a Gaussian distribution

$$q_x(w) = \mathcal{N}(g(x), h(x)) \quad , \quad g \in G \quad h \in H.$$

We can find optimal  $(f^*, g^*, \theta^*)$  by maximizing the ELBO, similar as VAE but with different forms since the model assumptions are different. The encoder is the same but for decoder we no longer need NN - simply update by gradient descent.