



Bayesian Hierarchical Dynamic Factor Models

Anthony M. Thomas, Jr.
anthony.thomas3@mavs.uta.edu

Department of Mathematics
The University of Texas at Arlington

September 8, 2020

Contents

1 Background

Bayesian Inference

Factor Analysis

Dynamic Factor Analysis

2 Hierarchical Dynamic Factor Analysis

State Space Representation

Data and Model Structure

Gibbs Sampler

3 Variational Bayesian Inference

Coordinate Ascent Variational Inference

Variational Bayesian EM

4 Future Work

5 References

Section 1

1 Background

Bayesian Inference

Factor Analysis

Dynamic Factor Analysis

2 Hierarchical Dynamic Factor Analysis

State Space Representation

Data and Model Structure

Gibbs Sampler

3 Variational Bayesian Inference

Coordinate Ascent Variational Inference

Variational Bayesian EM

4 Future Work

5 References

Bayesian Inference

Bayesian Inference can be described by two parts:

- ① Build a model based on data \mathbf{X} and parameters $\Theta = \{\Theta_1, \Theta_2\}$
 - Likelihood: $p(\mathbf{X}|\Theta)$
 - Prior: $p(\Theta)$
- ② Compute the posterior
 - Posterior:

$$p(\Theta|\mathbf{X}) = \frac{p(\mathbf{X}|\Theta)p(\Theta)}{p(\mathbf{X})}$$

- Report summaries, e.g. posterior expectations

$$\mathbb{E}[h(\Theta)|\mathbf{X}]$$

or marginal posterior expectations

$$\mathbb{E}[h(\Theta_i)|\mathbf{X}]$$

Factor Analysis

- ▶ **Factor Analysis** (FA) is a method that assumes that the covariance structure of a set of cross-sectional observations can be described in terms of a linear combination of latent variables called factors
- ▶ A sample of P observations are related to a set of factors through the equation

$$\mathbf{X}_i = \boldsymbol{\Lambda} \mathbf{F}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, P \quad (1)$$

where

- $\mathbf{X}_i = (X_{1i}, \dots, X_{Ni})^\top$ denotes a vector of observations for variable i
- $\mathbf{F}_i = (F_{1i}, \dots, F_{Ki})^\top$ denotes a vector of factors for variable i
- $\boldsymbol{\epsilon}_i = (\epsilon_{1i}, \dots, \epsilon_{Ni})^\top$ denotes a vector of measurement errors and idiosyncratic (unique) factors for variable i
- $\boldsymbol{\Lambda} = [\lambda_{nk}]_{N \times K}$ denotes a matrix of factor loadings

- ▶ The following assumptions are typically made in FA:

- ① $\text{rank}(\boldsymbol{\Lambda}) = K$
- ② $\mathbb{E}[\mathbf{X}_i] = \mathbb{E}[\mathbf{e}_i] = \mathbf{0}_N$ and $\mathbb{E}[\mathbf{F}_i] = \mathbf{0}_K \quad \forall i$
- ③ $\text{Var}(\mathbf{F}_i) = \mathbf{I}_K$ and $\text{Var}(\mathbf{e}_i) = \boldsymbol{\Sigma}$ where $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_N^2) \quad \forall i$
- ④ $\text{Cov}(\mathbf{F}_i, \mathbf{e}_i) = \mathbf{0}_{K \times N} \quad \forall i$

- ▶ Under these assumptions it follows that

$$\text{Var}(\mathbf{X}_i) = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Sigma} \quad \forall i$$

- ▶ Typical uses of FA:

- ① Dimension reduction: explain the covariation between N variables using $K < N$ factors
- ② Data interpretation: find factors that explain the covariation
- ③ Theory testing: test whether a hypothesized factor structure fits observed data

Time Series Factor Analysis

- ▶ In [1] FA is extended to time series data as **time series factor analysis** (TSFA)
- ▶ A sample of T time series observations are related to the factors through the equation

$$\mathbf{X}_t = \boldsymbol{\Lambda} \mathbf{F}_t + \boldsymbol{\epsilon}_t \quad t = 1, \dots, T \quad (2)$$

where

- $\mathbf{X}_t = (X_{1t}, \dots, X_{Nt})^\top$ denotes a vector of observations at time t
- $\mathbf{F}_t = (F_{1t}, \dots, F_{K_F t})^\top$ denotes a vector of factors at time t
- $\boldsymbol{\epsilon}_t$ is a vector of measurement errors and idiosyncratic factors at time t
- $\boldsymbol{\Lambda} = [\lambda_{nk}]_{N \times K}$ denotes a matrix of factor loadings

Dynamic Factor Analysis

- ▶ In **dynamic factor analysis** (DFA) the factors are assumed to not only affect the observations contemporaneously, but affect them through their lags as well:

$$X_{nt} = \lambda^n(L)\mathbf{F}_t + \mathbf{e}_t \quad n = 1, \dots, N \quad (3)$$

where

$$\lambda^n(L) = \lambda_0^n + \lambda_1^n L + \dots + \lambda_q^n L^q$$

$$L^s \mathbf{F}_t = \mathbf{F}_{t-s} \quad \forall s \geq 0$$

is a distributed lag polynomial of factor loadings in the lag operator L for the n th series

Dynamic Factor Analysis

- ▶ In DFA the factors are modeled as a time series process
- ▶ The time series process is commonly taken to be a vector autoregressive process, i.e.

$$\Psi(L)\mathbf{F}_t = \boldsymbol{\varepsilon}_t \quad (4)$$

where

$$\Psi(L) = \mathbf{I}_K - \Psi_1 L - \dots - \Psi_p L^p$$

is a matrix polynomial of autocorrelation coefficients in the lag operator L

Section 2

1 Background

Bayesian Inference

Factor Analysis

Dynamic Factor Analysis

2 Hierarchical Dynamic Factor Analysis

State Space Representation

Data and Model Structure

Gibbs Sampler

3 Variational Bayesian Inference

Coordinate Ascent Variational Inference

Variational Bayesian EM

4 Future Work

5 References

Hierarchical Dynamic Factor Analysis

- ▶ Suppose you organize a large panel of data into B blocks (e.g. Production, Employment, Demand, etc.) and each block b has N_b series
- ▶ $N = \sum_{b=1}^B N_b$
- ▶ A block may be divided into subblocks (e.g. Demand: Retail Sales, Auto Sales)

Hierarchical Dynamic Factor Model

- ▶ In [2], the authors generalize the dynamic factor model by positing that for each t , the n th series in a given block b , denoted by X_{bnt} , has three sources of variation:
 - ① idiosyncratic
 - ② block-specific
 - ③ common

Hierarchical Dynamic Factor Model

- ▶ A three-level representation of the data for $b = 1, \dots, B$ and $n = 1, \dots, N_b$ is given as

$$X_{bnt} = \boldsymbol{\lambda}_{G.b}^n(L) \mathbf{G}_{bt} + e_{Xbnt} \quad (5)$$

$$\mathbf{G}_{bt} = \boldsymbol{\Lambda}_{F.b}(L) \mathbf{F}_t + e_{Gbt} \quad (6)$$

$$\boldsymbol{\Psi}_F(L) \mathbf{F}_t = \boldsymbol{\epsilon}_{Ft}, \quad (7)$$

where

- $\boldsymbol{\lambda}_{G.b}^n(L)$ denotes a distributed lag polynomial of block-level factor loadings
- $\boldsymbol{\Lambda}_{F.b}(L)$ denotes a distributed lag matrix polynomial of common factor loadings
- $\mathbf{G}_{bt} = (G_{b1t}, \dots, G_{bK_{\mathcal{G}_b}t})^\top$ denotes the block-level factors
- $\mathbf{F}_t = (F_{1t}, \dots, F_{K_F t})^\top$ denotes the common factors

Hierarchical Dynamic Factor Model

- ▶ For some blocks, it may be appropriate to break up the data into subblocks, which adds another source of variation
- ▶ Let Z_{bsnt} be the n th series in subblock s of block b
- ▶ A four-level representation of the subblock data is given as

$$Z_{bsnt} = \boldsymbol{\lambda}_{H.bs}^n(L) \mathbf{H}_{bst} + e_{Zbsnt} \quad (8)$$

$$\mathbf{H}_{bst} = \boldsymbol{\Lambda}_{G.bs}(L) \mathbf{G}_{bt} + \mathbf{e}_{Hbst} \quad (9)$$

$$\mathbf{G}_{bt} = \boldsymbol{\Lambda}_{F.b}(L) \mathbf{F}_t + \mathbf{e}_{Gbt} \quad (10)$$

$$\boldsymbol{\Psi}_F(L) \mathbf{F}_t = \boldsymbol{\epsilon}_{Ft} \quad (11)$$

where

- $\boldsymbol{\lambda}_{H.bs}^n(L)$ denotes a distributed lag polynomial of subblock-level factor loadings
- $\boldsymbol{\Lambda}_{G.bs}(L)$ denotes a distributed lag matrix polynomial of block-level factor loadings
- $\mathbf{H}_{bst} = (H_{bs1t}, \dots, H_{bsK_{Hbst}})^\top$ denotes the subblock-level factors

Hierarchical Dynamic Factor Model

- The idiosyncratic components, the subblock-specific, block-specific, and common factors are assumed to be stationary, Gaussian autoregressive processes of orders q_{Zbsn} , q_{Xbn} , q_{Hbsi} , q_{Gb_j} , and q_{Fk} , respectively, i.e.

$$\begin{aligned}\psi_{Z.bsn}(L)e_{Zbsnt} &= \epsilon_{Zbsnt}, \quad \epsilon_{Zbsnt} \sim \mathcal{N}(0, \sigma_{Zbsn}^2) \quad n = 1, \dots, N_{bs} \\ \psi_{X.bn}(L)e_{Xbnt} &= \epsilon_{Xbnt}, \quad \epsilon_{Xbnt} \sim \mathcal{N}(0, \sigma_{Xbn}^2) \quad n = 1, \dots, N_b \\ \psi_{H.bsi}(L)e_{Hbsit} &= \epsilon_{Hbsit}, \quad \epsilon_{Hbsi} \sim \mathcal{N}(0, \sigma_{Hbsi}^2) \quad i = 1, \dots, K_{Hbs} \\ \psi_{G.bj}(L)e_{Gbjt} &= \epsilon_{Gbjt}, \quad \epsilon_{Gbjt} \sim \mathcal{N}(0, \sigma_{Gb_j}^2) \quad j = 1, \dots, K_{Gb} \\ \psi_{F.k}(L)F_{kt} &= \epsilon_{Fkt}, \quad \epsilon_{Fkt} \sim \mathcal{N}(0, \sigma_{Fk}^2) \quad k = 1, \dots, K_F\end{aligned}$$

Hierarchical Dynamic Factor Model

- ▶ Not all series need to belong to blocks and subblocks
- ▶ In general, the data used in a four-level model are a mixture of Z_{bsnt} , X_{bnt} , and X_{nt}

State Space Representation

- Observed data (vector/matrix form):

$$\begin{aligned}\mathbf{Z}_{bst} &= \boldsymbol{\Lambda}_{H.bst}(L) \mathbf{H}_{bst} + \mathbf{e}_{Zbst} \\ \boldsymbol{\Psi}_{Z.bst}(L) \mathbf{e}_{Zbst} &= \boldsymbol{\epsilon}_{Zbst}\end{aligned}$$

implies that the measurement equation is

$$\begin{aligned}\boldsymbol{\Psi}_{Z.bst}(L) \mathbf{Z}_{bst} &= \boldsymbol{\Psi}_{Z.bst}(L) \boldsymbol{\Lambda}_{H.bst}(L) \mathbf{H}_{bst} + \boldsymbol{\epsilon}_{Zbst} \\ \tilde{\mathbf{Z}}_{bst} &= \tilde{\boldsymbol{\Lambda}}_{H.bst} \vec{\mathbf{H}}_{bst} + \boldsymbol{\epsilon}_{Zbst}, \quad \boldsymbol{\epsilon}_{Zbst} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{Zbst})\end{aligned}\tag{12}$$

with transition equation

$$\vec{\mathbf{H}}_{bst} = \vec{\boldsymbol{\alpha}}_{G.bst} + \tilde{\boldsymbol{\Psi}}_{Hbst} \vec{\mathbf{H}}_{bst-1} + \vec{\boldsymbol{\epsilon}}_{Hbst}, \quad \vec{\boldsymbol{\epsilon}}_{Hbst} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}\mathbf{1}^\top \otimes \boldsymbol{\Sigma}_{Hbst})\tag{13}$$

where $\boldsymbol{\alpha}_{G.bst} = \boldsymbol{\Psi}_{H.bst}(L) \boldsymbol{\Lambda}_{G.bst}(L) \mathbf{G}_{bt}$

State Space Representation

- Subblock-level factors (vector/matrix form):

$$\begin{aligned}\mathbf{H}_{bt} &= \boldsymbol{\Lambda}_{G.b}(L)\mathbf{G}_{bt} + \boldsymbol{\epsilon}_{Hbt} \\ \boldsymbol{\Psi}_{H.b}(L)\mathbf{e}_{Hbt} &= \boldsymbol{\epsilon}_{Hbt}\end{aligned}$$

implies that the (pseudo) measurement equation is

$$\begin{aligned}\boldsymbol{\Psi}_{H.b}(L)\mathbf{H}_{bt} &= \boldsymbol{\Psi}_{H.b}(L)\boldsymbol{\Lambda}_{G.b}(L)\mathbf{G}_{bt} + \boldsymbol{\epsilon}_{Hbt} \\ \tilde{\mathbf{H}}_{bt} &= \tilde{\boldsymbol{\Lambda}}_{G.b}\vec{\mathbf{G}}_{bt} + \boldsymbol{\epsilon}_{Hbt}, \quad \boldsymbol{\epsilon}_{Hbt} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{Hb})\end{aligned}\tag{14}$$

with transition equation

$$\vec{\mathbf{G}}_{bt} = \boldsymbol{\alpha}_{F.bt} + \tilde{\boldsymbol{\Psi}}_{G.b}\vec{\mathbf{G}}_{bt-1} + \vec{\boldsymbol{\epsilon}}_{Gbt}, \quad \vec{\boldsymbol{\epsilon}}_{Gbt} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}\mathbf{1}^\top \otimes \boldsymbol{\Sigma}_{Gb})\tag{15}$$

where $\boldsymbol{\alpha}_{F.bt} = \boldsymbol{\Psi}_{G.b}(L)\boldsymbol{\Lambda}_{F.b}(L)\mathbf{F}_t$

State Space Representation

- Block-level factors (vector/matrix form):

$$\begin{aligned}\mathbf{G}_t &= \boldsymbol{\Lambda}_F(L)\mathbf{F}_t + \boldsymbol{\epsilon}_{Gt} \\ \boldsymbol{\Psi}_G(L)\mathbf{e}_{Gt} &= \boldsymbol{\epsilon}_{Gt}\end{aligned}$$

implies that the (pseudo) measurement equation is

$$\begin{aligned}\boldsymbol{\Psi}_G(L)\mathbf{G}_t &= \boldsymbol{\Psi}_G(L)\boldsymbol{\Lambda}_F(L)\mathbf{F}_t + \boldsymbol{\epsilon}_{Gt} \\ \tilde{\mathbf{G}}_t &= \tilde{\boldsymbol{\Lambda}}_F\vec{\mathbf{F}}_t + \boldsymbol{\epsilon}_{Gt}, \quad \boldsymbol{\epsilon}_{Gt} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_F)\end{aligned}\tag{16}$$

with transition equation

$$\vec{\mathbf{F}}_t = \tilde{\boldsymbol{\Psi}}_F\vec{\mathbf{F}}_{t-1} + \vec{\boldsymbol{\epsilon}}_{Ft}, \quad \vec{\boldsymbol{\epsilon}}_{Gbt} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}\mathbf{1}^\top \otimes \boldsymbol{\Sigma}_{Gb})\tag{17}$$

Data and Model Structure

- ▶ The authors organized a dataset consisting of $N = 445$ series giving $T = 227$ observations into 5 blocks

Table 1: Block Structure

Block	Subblock	Source	N	K_{Hbs}
Production	CU	Fed	25	1
	IP	Fed	38	1
	DG	Census	60	2
Employment	ES	BLS	82	2
	HS	BLS	92	1
Consumption	WT	Census	54	1
	RS	Census	30	1
Housing		Census	29	
Manufacturing surveys		ISM, Fed	35	

Data and Model Structure

- The distributed lag (matrix) polynomials of factor loadings are assumed to be constant (order 0), i.e.

$$\lambda_{H.bs}^n(L) = \lambda_{H.bs0}^n$$

$$\Lambda_{G.bs}(L) = \Lambda_{G.bs0}$$

$$\Lambda_{F.b}(L) = \Lambda_{F.b0}$$

- The (matrix) polynomials of autocorrelation coefficients are assumed to be of order 1, i.e.

$$\psi_{Z.bsn}(L) = 1 - \psi_{Z.bsn1} L$$

$$\psi_{X.bn}(L) = 1 - \psi_{X.bn1} L$$

$$\psi_{H.bsi}(L) = 1 - \psi_{H.bsi1} L$$

$$\psi_{G.bj}(L) = 1 - \psi_{G.bj1} L$$

$$\psi_{F.k}(L) = 1 - \psi_{F.k1} L$$

Data and Model Structure

- ▶ They estimate one common factor, i.e. $K_F = 1$, one common factor per block, i.e. $K_{Gb} = 1 \quad \forall b$, and one or two factors per subblock, i.e. $K_{Hbs} = 1, 2$
- ▶ For the cases where $K_{Hbs} = 2$, the factor loading matrices are assumed to be lower triangular with 1's along the diagonal, i.e.

$$\boldsymbol{\Lambda}_{H.bs} = \begin{bmatrix} 1 & 0 \\ \lambda_{H.bs_{2,1}} & 1 \\ \lambda_{H.bs_{3,1}} & \lambda_{H.bs_{3,2}} \\ \vdots & \vdots \\ \lambda_{H.bs_{N_{bs},1}} & \lambda_{H.bs_{N_{bs},2}} \end{bmatrix}$$

Parameter Priors

- ▶ All free factor loadings and autocorrelation coefficients are assigned independent standard Gaussian priors, i.e.

$$\begin{aligned}\lambda_{(\cdot)} &\sim \mathcal{N}(0, 1) \\ \psi_{(\cdot)} &\sim \mathcal{N}(0, 1)\end{aligned}$$

- ▶ All variance parameters are assigned independent scaled-inverse chi squared distributions with 4 degrees of freedom and a scale of 0.01, i.e.

$$\sigma_{(\cdot)}^2 \sim \text{scale-inv-}\chi^2(4, 0.01^2)$$

Gibbs Sampler

- ▶ The authors implement a Gibbs sampling algorithm based on a degenerate linear Gaussian state-space representation of the model to obtain samples from the posterior and compute posterior means
- ▶ 50,000 draws are used as a burn-in, i.e. discarded
- ▶ 50,000 more draws are obtained while storing every 50th draw to obtain a posterior sample size of 1,000
- ▶ The main idea of the algorithm is presented next

Gibbs Sampler

- ▶ Let $\boldsymbol{\Lambda} = (\boldsymbol{\Lambda}_H, \boldsymbol{\Lambda}_G, \boldsymbol{\Lambda}_F)$, $\boldsymbol{\Psi} = (\boldsymbol{\Psi}_F, \boldsymbol{\Psi}_G, \boldsymbol{\Psi}_H, \boldsymbol{\Psi}_Z)$, and $\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_F, \boldsymbol{\Sigma}_G, \boldsymbol{\Sigma}_H, \boldsymbol{\Sigma}_Z)$
- ① Organize the data into blocks and subblocks to get Z_{bst} for $b = 1, \dots, B, s = 1, \dots, B_S$. Get initial values for $\{H_{bst}\}$, $\{G_{bt}\}$, and $\{F_t\}$ using principal components. Use these to produce initial values for $\boldsymbol{\Lambda}, \boldsymbol{\Psi}, \boldsymbol{\Sigma}$
- ② Conditional on $\boldsymbol{\Lambda}, \boldsymbol{\Psi}, \boldsymbol{\Sigma}, \{G_{bt}\}, \{F_t\}$ and the data, draw $\{H_{bst}\} \quad \forall b, s$
- ③ Conditional on $\boldsymbol{\Lambda}, \boldsymbol{\Psi}, \boldsymbol{\Sigma}, \{H_{bst}\}, \{F_t\}$ and the data, draw $\{G_{bt}\} \quad \forall b$
- ④ Conditional on $\boldsymbol{\Lambda}, \boldsymbol{\Psi}, \boldsymbol{\Sigma}, \{H_{bst}\}, \{G_{bt}\}$ and the data, draw $\{F_t\}$
- ⑤ Conditional on $\{H_{bst}\}, \{G_{bt}\}, \{F_t\}$ and the data, draw $\boldsymbol{\Lambda}, \boldsymbol{\Psi}, \boldsymbol{\Sigma}$
- ⑥ Return to 2

Section 3

1 Background

Bayesian Inference

Factor Analysis

Dynamic Factor Analysis

2 Hierarchical Dynamic Factor Analysis

State Space Representation

Data and Model Structure

Gibbs Sampler

3 Variational Bayesian Inference

Coordinate Ascent Variational Inference

Variational Bayesian EM

4 Future Work

5 References

Why Variational Bayesian Inference?

- ▶ Consider a model with data \mathbf{X} and latent variables \mathbf{Z} (model parameters can be included)
- ▶ The goal is to compute the joint posterior of the latent variables given the data

$$p(\mathbf{Z}|\mathbf{X}) = \frac{p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})}{p(\mathbf{X})} \quad (18)$$

- Likelihood: $p(\mathbf{X}|\mathbf{Z})$
- Prior: $p(\mathbf{Z})$
- Evidence: $p(\mathbf{X})$

$$p(\mathbf{X}) = \int p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} \quad (19)$$

Why Variational Bayesian Inference?

- ▶ For complex models (18) and (19) either have no closed-form or require high-dimensional integration which causes the “inference problem”
- ▶ As a result, the joint posterior has to be approximated
- ▶ Markov chain Monte Carlo (MCMC) methods have been the gold standard to solve this problem
 - ① Construct an ergodic Markov chain on \mathbf{Z} whose stationary distribution is the joint posterior $p(\mathbf{Z}|\mathbf{X})$
 - ② Sample from the chain to collect samples from the stationary distribution
 - ③ Approximate the posterior with an empirical estimate constructed from a subset of the collected samples
 - ④ Use the subset of collected samples to estimate expectations of interest

Why Variational Bayesian Inference?

- ▶ MCMC uses sampling to solve the inference problem
- ▶ MCMC methods eventually produce accurate results, but are usually computationally intensive for interesting models
- ▶ Variational Bayes (VB) uses optimization to solve the inference problem
- ▶ VB typically obtains results much faster

Variational Bayesian Inference

- ▶ The main idea of the VB framework is to
 - ① Posit a family of “nice” approximate densities \mathcal{Q}
 - ② Find a member of that family that is “closest” to the exact posterior

$$q^*(\mathbf{Z}) = \arg \min_{q(\mathbf{Z}) \in \mathcal{Q}} \text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z}|\mathbf{X})) \quad (20)$$

where

$$\text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z}|\mathbf{X})) = \mathbb{E}_q[\log q(\mathbf{Z})] - \mathbb{E}_q[\log p(\mathbf{Z}|\mathbf{X})] \quad (21)$$

Variational Bayesian Inference

- ▶ It follows that

$$\text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z}|\mathbf{X})) = \mathbb{E}_q[\log q(\mathbf{Z})] - \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z})] + \log p(\mathbf{X})$$

which reveals that the objective in (20) depends on the evidence, thus it cannot be computed directly

Evidence Lower Bound

- ▶ The reason why $\text{KL}(q||p)$ is a desirable measure of “closeness” is because it leads to a lower bound on $\log p(\mathbf{X})$ called the **evidence lower bound** (ELBO)

$$\begin{aligned}\log p(\mathbf{X}) &= \log \int q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= \log \left(\mathbb{E}_q \left[\frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right] \right) \\ &\geq \mathbb{E}_q \left[\log \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right], \quad \text{by Jensen's Inequality} \\ &= \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z})] - \mathbb{E}_q[\log q(\mathbf{Z})] \\ &= \text{ELBO}(q)\end{aligned}$$

Evidence Lower Bound

- ▶ It follows that

$$\text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z}|\mathbf{X})) = \log p(\mathbf{X}) - \text{ELBO}(q)$$

- ▶ Since $\log p(\mathbf{X})$ is constant wrt $q(\mathbf{Z})$ we can redefine the objective function as

$$\arg \min_{q(\mathbf{Z}) \in \mathcal{Q}} \text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z}|\mathbf{X})) = \arg \max_{q(\mathbf{Z}) \in \mathcal{Q}} \text{ELBO}(q)$$

Variational Bayesian Inference

- ▶ The VB framework is now
 - ➊ Posit a family of “nice” approximate densities \mathcal{Q}
 - ➋ Find a member of that family that is “closest” to the exact posterior, i.e.

$$q^*(\mathbf{Z}) = \arg \max_{q(\mathbf{Z}) \in \mathcal{Q}} \text{ELBO}(q) \quad (22)$$

- ▶ Solving this optimization problem is still difficult in general
- ▶ Using the mean-field assumption can make it easier

Mean-Field Assumption

- ▶ The mean-field assumption says to:
 - ➊ Partition the latent variables into M groups, say $\mathbf{Z}_1, \dots, \mathbf{Z}_M$
 - ➋ Assume that the distributions in \mathcal{Q} factorize across the groups, i.e.

$$\mathcal{Q} = \left\{ q : q(\mathbf{Z}) = \prod_{m=1}^M q_m(\mathbf{Z}_m) \right\}$$

- ▶ Learning the optimal q now reduces to learning the optimal q_1, \dots, q_M
- ▶ Straightforward to optimize via coordinate ascent
- ▶ This is **NOT** a modeling assumption

Mean-Field Assumption

- ▶ Interestingly, under the mean-field assumption, the optimization problem for a single q_m has the solution:

$$q_m(\mathbf{Z}_m) = \frac{\exp\{\mathbb{E}_{-m}[\log p(\mathbf{X}, \mathbf{Z})]\}}{\int \exp\{\mathbb{E}_{-m}[\log p(\mathbf{X}, \mathbf{Z})]\} d\mathbf{Z}_m} \quad \forall m \quad (23)$$

- ▶ This establishes what is called the **coordinate ascent variational inference** (CAVI) algorithm

Coordinate Ascent Variational Inference

Algorithm 1: Coordinate ascent variational inference

Input: Model $p(\mathbf{X}, \mathbf{Z})$, Data \mathbf{X} ,

Output: Variational density $q(\mathbf{Z}) = \prod_{m=1}^M q_m(\mathbf{Z}_m)$

Initialize: Variational densities $q_m(\mathbf{Z}_m)$

while the ELBO has not converged **do**

for $m \in \{1, 2, \dots, M\}$ **do**

 | Set $q_m(\mathbf{Z}_m) \propto \exp\{\mathbb{E}_{-\mathbf{Z}_m}[\log p(\mathbf{X}, \mathbf{Z})]\}$

end

 Compute ELBO(q)

end

return $q(\mathbf{Z})$

Conjugate-Exponential Models

- ▶ Conjugate-exponential models satisfy two conditions:
 - ➊ The complete-data likelihood is in the exponential family, i.e.

$$p(\mathbf{X}_i, \mathbf{Z}_i | \boldsymbol{\theta}) = g(\boldsymbol{\theta}) f(\mathbf{X}_i, \mathbf{Z}_i) \exp\left\{ \phi(\boldsymbol{\theta})^\top u(\mathbf{X}_i, \mathbf{Z}_i) \right\} \quad (24)$$

where $\phi(\boldsymbol{\theta})$ is the vector of natural parameters and u is a vector of sufficient statistics

- ➋ The parameter prior is conjugate to the complete-data likelihood, i.e.

$$p(\boldsymbol{\theta} | \eta, \boldsymbol{\nu}) = h(\eta, \boldsymbol{\nu}) g(\boldsymbol{\theta})^\eta \exp\left\{ \phi(\boldsymbol{\theta})^\top \boldsymbol{\nu} \right\} \quad (25)$$

where η and $\boldsymbol{\nu}$ are hyperparameters of the prior

- ▶ The author in [3] generalizes the expectation-maximization (EM) to a VB-EM algorithm for conjugate-exponential models (HMMs, MFA, and SSMs)

Variational Bayesian EM

- Given an i.i.d. dataset $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ the ELBO(q) can be maximized iteratively

- The VBE step:

$$q(\mathbf{Z}_i) \propto f(\mathbf{X}_i, \mathbf{Z}_i) \exp\left\{\bar{\phi}^\top u(\mathbf{X}_i, \mathbf{Z}_i)\right\} \quad \forall i \quad (26)$$

with

$$\bar{\phi} = \mathbb{E}_{\boldsymbol{\theta}}[\phi(\boldsymbol{\theta})]$$

- The VBM step:

$$q(\boldsymbol{\theta}) = h(\tilde{\eta}, \tilde{\boldsymbol{\nu}}) g(\boldsymbol{\theta})^{\tilde{\eta}} \exp\left[\phi(\boldsymbol{\theta})^\top \tilde{\boldsymbol{\nu}}\right] \quad (27)$$

with

$$\tilde{\eta} = \eta + N$$

$$\tilde{\boldsymbol{\nu}} = \boldsymbol{\nu} + \sum_{i=1}^N \bar{u}(\mathbf{X}_i)$$

$$\bar{u}(\mathbf{X}_i) = \mathbb{E}_{\mathbf{Z}_i}[u(\mathbf{X}_i, \mathbf{Z}_i)]$$

Section 4

① Background

Bayesian Inference

Factor Analysis

Dynamic Factor Analysis

② Hierarchical Dynamic Factor Analysis

State Space Representation

Data and Model Structure

Gibbs Sampler

③ Variational Bayesian Inference

Coordinate Ascent Variational Inference

Variational Bayesian EM

④ Future Work

⑤ References

Future Work

- ▶ Ultimate goal is to develop a VB framework to handle the three- and four-level hierarchical dynamic factor model in [2]
- ▶ Developing a VB framework according to CAVI is straightforward, but very tedious
 - Still an open problem
- ▶ Developing a VB framework according to the VBEM using the degenerate state space representation is not as straightforward
 - The state transition covariance matrix is singular
 - The state transition equations have time-varying intercepts
 - Still an open problem

Section 5

1 Background

Bayesian Inference

Factor Analysis

Dynamic Factor Analysis

2 Hierarchical Dynamic Factor Analysis

State Space Representation

Data and Model Structure

Gibbs Sampler

3 Variational Bayesian Inference

Coordinate Ascent Variational Inference

Variational Bayesian EM

4 Future Work

5 References

Selected References I

-  Paul D Gilbert, Erik Meijer, et al. “Time series factor analysis with an application to measuring money”. In: *University of Groningen, Research School SOM Research Report 05F10* (2005).
-  Emanuel Moench, Serena Ng, and Simon Potter. “Dynamic Hierarchical Factor Models”. en. In: *Review of Economics and Statistics* 95.5 (Dec. 2013), pp. 1811–1817. ISSN: 0034-6535, 1530-9142. DOI: [10.1162/REST_a_00359](https://doi.org/10.1162/REST_a_00359).
-  Matthew James Beal. “Variational Algorithms for Approximate Bayesian Inference”. en. In: (), p. 281.