

## Models

- (0,1)
- (2,1)-(2,4)
- (3,1)-(3,3)
- (3,11)-(3,17)

11/21/2025

## Derivation of EBMF Models

- Factor analysis Bayesian model (add TTE survival component)

(1.1)

- Matrix factorization represent  $n \times p$  data matrix  $Y$  as:

$$Y = L^T F + E \quad (1.1)$$

where  $L$  is  $K \times n$  matrix,  $F$  is  $K \times p$  matrix,  $E$  is  $n \times p$  matrix of residuals (assume normally distributed)

- $L$ : "loadings"
- $F$ : "factors"

- simplest approaches to estimating  $L$  and  $F$  are based on maximum likelihood or least squares. Include sparsity assumptions on  $L$  and  $F$ .
- some methods have sparsity in loadings ( $L$ ) only, others in both loadings + factors.

- Empirical Bayes approach exploits variational approximation methods to obtain simple algorithms that jointly estimate the prior distributions for both loadings and factors, as well as the loadings and factors themselves.  
 $\Rightarrow$  Empirical Bayes approach to matrix factorization

$\hookrightarrow$  fitting EB MF with any prior family can be reduced to repeatedly solving the "empirical bayes normal means" (EBNM) problem w/some prior

- $K$ -factor Empirical Bayes Matrix Factorization (EBMF) model

2.1 - linear structure

• Definition:  $Y = \sum_{k=1}^K l_k f_k^T + E$  (instead of  $Y = L^T F + E$ )

$\Downarrow \sum_{k=1}^K l_k f_k^T$

$\Downarrow \sum_{k=1}^K l_k f_k^T$

$\Downarrow \sum_{k=1}^K l_k f_k^T$

- condition represents low-rank assumption. Assumes that high-dimensional observed data  $Y_{n \times p}$  is composed of a relatively small number of underlying processes/components ( $K$ ).

- Matrix vs. summation notation: Standard factor analysis is often written as  $Y = L F^T$ . This is explicitly decomposed into a sum of  $K$  rank-1 matrices. If  $l_k$  is an  $n \times 1$  vector (column) and  $f_k^T$  is a  $1 \times p$  vector (row), then their product  $l_k f_k^T$  is an  $n \times p$  matrix.

$$Y_{ij} \propto (l_1 f_1^T)_{ij} + (l_2 f_2^T)_{ij} + \dots + (l_K f_K^T)_{ij} \quad \begin{array}{c} \text{effect of} \\ \text{factor 1} \end{array} \quad \begin{array}{c} \text{effect of} \\ \text{factor 2} \end{array} \quad \dots \quad \begin{array}{c} \text{effect of} \\ \text{factor } K \end{array} \quad \left[ \begin{array}{c} \text{summation of} \\ K \text{ components} \end{array} \right]$$

↳ cont: This summation notation helps w/ clarity because the fitting algorithm (algs 1 and 2) updates the model one factor ( $k$ ) at a time. The algorithm treats the contributions of other factors ( $k' \neq k$ ) as fixed offsets.

## 2.2 & 2.3 - Hierarchical priors

(2.2) loadings:  $l_{11}, \dots, l_{nn} \stackrel{iid}{\sim} g_{lk}, g_{lk} \in G_l$  (weights)

(2.3) factors:  $f_{11}, \dots, f_{nn} \stackrel{iid}{\sim} g_{fk}, g_{fk} \in G_f$  (latent variables)  
 ↳ selected "factors"

- define Empirical Bayes nature of the model

- The notation  $iid$  implies exchangeability. It assumes that, a priori, we have no reason to distinguish sample  $i$  and sample  $i+1$  regarding their loading (weight) on factor  $k$ , other than what the distribution  $g_{lk}$  dictates.

- In classical Bayesian factor analysis, we might fix the prior  $l_{ik} \sim N(0, 1)$ . In EBMF, the prior  $g_{lk}$  is unknown. The "derivation" here relies on adaptive shrinkage.

- If factor  $k$  is sparse, the estimated  $g_{lk}$  might look like a "spike" at zero.
- If factor  $k$  is dense,  $g_{lk}$  might be a broad Normal distribution.
- different factors ( $k$ ) have different priors ( $g_{lk}$  vs.  $g_{l'k}$ ). This allows the model to handle a mix of sparse and dense signals simultaneously.

- prior family ( $G$ ): user chooses a family  $G$  (all unimodal distributions at 0 or "spike-and-slab"). The algorithm then searches within this family to find the best  $g$  (prior).

## 2.4 - The Noise Model

$$E_{ij} \sim N(0, \frac{1}{\pi_{ij}}) \text{ with } \boldsymbol{\gamma} := (\gamma_{ij}) \in \Gamma$$

- defines the likelihood function



- Gaussian Assumption: Assuming the errors  $\epsilon_{ij}$  are Gaussian leads to the squared error loss function used in the optimization.

$$P(Y_{ij} | L, F, \gamma) = \sqrt{\frac{\pi}{2\gamma}} \exp \left[ -\frac{\gamma_{ij}}{2} (Y_{ij} - \sum_k f_{ik})^2 \right] \quad (\text{derivation to follow})$$

- Precision structure: paper derives flexibility by constraining the precision matrix  $\gamma$ . The term  $\gamma$  represents the set of allowed precision structures:
  - constant variance:  $\gamma_{ij} = \gamma$ , simplest case
  - column-specific variance:  $\gamma_{ij} = \gamma_j$ , standard in "fMRI", assumes some features are noisier than others
  - row-specific variance:  $\gamma_{ij} = \gamma_i$

- Derivation of how (2+) Gaussian assumption transforms into "Weighted Squared Error" loss function used in optimization.

1) model assumes errors are Normally distributed

$$\epsilon_{ij} = Y_{ij} - (\sum_k f_{ik}) \sim N(0, \frac{1}{\gamma_{ij}})$$

$\hookrightarrow$  observed data  $Y_{ij}$ , given the latent factors, follows a Gaussian distribution with mean  $\mu_{ij} = \sum_k f_{ik}$  and variance  $\sigma_{ij}^2 = \frac{1}{\gamma_{ij}}$

2) probability density function (PDF)

"probability density for a single point  $Y_{ij}$  is"

$$P(Y_{ij} | L, F, \gamma) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \exp \left( -\frac{(Y_{ij} - \mu_{ij})^2}{2\sigma_{ij}^2} \right)$$

• Substitute precision  $\gamma_{ij} = 1/\sigma_{ij}^2$

$$P(Y_{ij} | L, F, \gamma) = \sqrt{\frac{\gamma_{ij}}{2\pi}} \exp \left( -\frac{\gamma_{ij}}{2} (Y_{ij} - \sum_k f_{ik})^2 \right)$$

3) Joint likelihood

- assuming iid errors, probability of entire data matrix  $Y$  is product of individual probabilities

$$L(\gamma) = \prod_{i=1}^I \prod_{j=1}^J \left[ \sqrt{\frac{\gamma_{ij}}{2\pi}} \exp \left( -\frac{\gamma_{ij}}{2} (Y_{ij} - \sum_k f_{ik})^2 \right) \right]$$

4) log-likelihood

- take a log of both sides (product  $\rightarrow$  sum)

$$\log[L(\gamma)] = \sum_{i,j} \left[ \frac{1}{2} \log(\gamma_{ij}) - \frac{1}{2} \log(2\pi) - \frac{\gamma_{ij}}{2} (Y_{ij} - \sum_k f_{ik})^2 \right]$$

5) re-write as "loss"

- maximizing likelihood is equivalent to minimizing negative log-likelihood,
- consider NLL, drop terms that are constant wrt L and F

$$\text{Loss}(L, F) \propto -\log(L(\theta)) \approx \sum_{ij} \frac{\tau_{ij}}{2} (Y_{ij} - \sum_k f_{ik} f_{kj})^2 - \sum_{ij} \frac{1}{2} \log(\tau_{ij})$$

• holding  $\tau$  (precision) constant, minimizing loss is equivalent to minimizing:

$$SSE = \sum_{ij} \tau_{ij} (Y_{ij} - \hat{Y}_{ij})^2$$

- weighted sum of squared errors,
- $\tau_{ij}=1$  (constant noise) → standard least squares
- $\tau_{ij}=0$  (infinite variance) → ignore data point (missing data imputation)

↳ connect to (A.6) where a function  $F(\gamma)$  is defined  
to update the precision parameters:

$$F(\gamma) = \frac{1}{2} \sum_{ij} \left[ \log(\gamma_{ij}) + \gamma_{ij} R_{ij}^2 \right] \rightarrow \begin{matrix} \text{expected squared} \\ \text{residuals} \end{matrix}$$

$R$  + or - ? tips?

### Notes

- $Y$ : observed data matrix w/ dimension  $n \times p$
- $l_k$ :  $k$ th set of "badges" ( $n$ -vector)
- $f_k$ :  $k$ th "factor" ( $p$ -vector)
- $g_e, g_f$ : pre-specified families of distributions
- $g_{lk}, g_{fk}$ : unknown "prior" distributions to be estimated
- $E$ :  $n \times p$  matrix of independent error terms
- $\tau$ : unknown  $n \times p$  matrix of precisions ( $\tau_{ij}$ ) (constant/column-specific/mu-specific)

• "Empirical Bayes" in EBME means fitting (2.1)-(2.4) by obtaining point estimates for the priors  $g_{lk}, g_{fk}$ , ( $k=1, \dots, t$ ) and approximate the posterior distributions for the parameters  $l_k, f_k$  given those point estimates.

### 3 - Fitting the EBME Model

- need to estimate all of  $g_e, g_f, l, f, \tau$
- 2 step approach:
  - estimate  $g_e, g_f, \tau$  b) maximizing likelihood (product of all  $g_e$  and  $g_f$ ) over

$$L(g_e, g_f, \tau) := \prod_{i=1}^n p(Y_i | l, f, \tau) g_e(dl_1) \cdots g_e(dl_n) g_f(df_1) \cdots g_f(df_p) \underbrace{g_e(g_f, \tau)}_{\text{over}}$$

• estimate  $l$  and  $f$  using posterior distribution:  $p(l, f | Y, \tilde{g}_l, \tilde{g}_f, \tilde{\tau})$

\* Derivation ↴

• Marginal likelihood:

$$L(g_l, g_f, \tau) := \iint p(Y | l, f, \tau) g_l(dl) \cdots g_f(df)$$

- integrate out unknown latent variables  $l$  and  $f$

- find how "likely" observed data  $Y$  is given hyperparameters  $(g_l, g_f, \tau)$

1) probability of data given hyperparameters

- find estimated priors  $\tilde{g}_l$ ,  $\tilde{g}_f$ , and precision  $\tilde{\tau}$  that make observed data most likely. Maximize

$$p(Y | g_l, g_f, \tau)$$

2) Law of Total Probability

- data  $Y$  depends on latent factors  $l$  (loadings) and  $f$  (factors) (not observed), can't directly compute  $p(Y)$ . Need to integrate over all possible values that  $l$  and  $f$  take.

$$P(Y | g_l, g_f, \tau) = \iint p(Y | l, f | g_l, g_f, \tau) dl df$$

3) factorize joint distribution

- factorize joint probability inside integral

$$p(Y, l, f) = \underbrace{p(Y | l, f, \tau)}_{\text{likelihood}} \cdot \underbrace{p(l | g_l)}_{\text{prior on } l} \cdot \underbrace{p(f | g_f)}_{\text{prior on } f}$$

4) expand priors

- consider rank 1 case ( $k=1$ ).  $l$  and  $f$  vectors consist of  $n$  and  $p$  independent draws:

$$\cdot p(l | g_l) = \prod_{i=1}^n g_l(l_i)$$

$$\cdot p(f | g_f) = \prod_{j=1}^p g_f(f_j)$$

Note:  $g(\cdot)$  allows for distributions that are partly continuous (slap) and partly discrete (spike at zero)

↳ computational difficulty motivates use of variational approximation to avoid solving MLE (3.1) directly!

Next: use variational approximation to approximate estimates for  $g_e$ ,  $g_f$ , and  $\gamma$  by maximizing likelihood (3.1).

goal: transform (3.1) "impossible integration"  $\rightarrow$  "solvable optimization"!

### Variational Approach

• 3.2: write log of likelihood (3.1) as:

$$l(g_e, g_f, \gamma) := \log L(g_e, g_f, \gamma)$$

- in factor analysis (maximum likelihood), find loadings ( $L$ ) and factors ( $F$ ) that minimize the error:  $\min \|Y - LF\|^2$

- In 3.2,  $L$  and  $F$  are missing from likelihood parameters ( $l(g_e, g_f, \gamma)$ ). EBML is a "Type II Maximum Likelihood" method where we don't care about the specific values of  $L$  and  $F$ . We are interested in finding the distributions (priors  $g_e, g_f$ ) that generated them ( $L$  and  $F$ ).

- Do this by integrating out the latent variables  $L$  and  $F$ .

$$L(g_e, g_f, \gamma) = \int p(Y, L, F | g_e, g_f, \gamma) dL dF$$

$L > (3.2)$  says: "find the shape of the priors ( $g$ ) that makes the observed

data ( $Y$ ) most probable, averaging over all possible factor configurations"

• Take log to turn product of probabilities into a sum of expectations (separate  $L$  and  $F$  terms)

• 3.3

• Target:  $l(g_e, g_f, \gamma) = F(g_e, g_f, \gamma, \gamma) + D_{KL}(q || p)$

#### 1) marginal likelihood

Let  $\Theta = \{g_e, g_f, \gamma\}$  by hyperparameters, Let  $Z = \{L, F\}$  be the latent variables, the log of the marginal likelihood is:  $\log [P(Y|\Theta)]$

#### 2) Expectation

$\log P(Y|\Theta)$  is constant wrt  $Z$ . Thus, if we take expectation w.r.t any distribution  $q(z)$ , the value remains unchanged.

$$\log P(Y|\Theta) = \int q(z) \log P(Y|z) dz$$

$$\text{Note: } \int q(z) \cdot c dz = c \underbrace{\int q(z) dz}_{=1} = c$$

### 3) Bayes Rule Expansion

- consider conditional probability:

$$P(z|Y, \mathcal{D}) = \frac{P(Y, z|\mathcal{D})}{P(Y|\mathcal{D})}$$

- rearrange:

$$P(Y|\mathcal{D}) = \frac{P(Y, z|\mathcal{D})}{P(z|Y, \mathcal{D})}$$

- sub fraction into integral from (2)  Plug in

$$\log P(Y|\mathcal{D}) = \int q(z) \log \left( \frac{P(Y, z|\mathcal{D})}{P(z|Y, \mathcal{D})} \right) dz$$

### 4) variational distribution q

- multiply and divide by approximation  $q(z)$

$$\log P(Y|\mathcal{D}) = \int q(z) \log \left[ \frac{P(Y, z|\mathcal{D})}{P(z|Y, \mathcal{D})} \cdot \frac{q(z)}{q(z)} \right] dz$$

### 5) rearrange terms

- group numerator + denominator to match definitions of ELBO and KL divergence.

$$\bullet \text{ELBO: } \frac{P(Y, z|\mathcal{D})}{q(z)} \quad \bullet \text{KL: } \frac{q(z)}{P(z|Y, \mathcal{D})}$$

$$\hookrightarrow \log P(Y|\mathcal{D}) = \int q(z) \log \left[ \underbrace{\frac{P(Y, z|\mathcal{D})}{q(z)}}_{\text{ELBO}} \cdot \underbrace{\frac{q(z)}{P(z|Y, \mathcal{D})}}_{\text{KL}} \right] dz$$

### 6) Split up logarithm

$$\log(A \cdot B) = \log(A) + \log(B)$$

$$\hookrightarrow \log P(Y|\mathcal{D}) = \int q(z) \left[ \overbrace{\log \left( \frac{P(Y, z|\mathcal{D})}{q(z)} \right)}^{\text{ELBO}} + \overbrace{\log \left( \frac{q(z)}{P(z|Y, \mathcal{D})} \right)}^{\text{KL}} \right] dz$$

### 7) Separate integral (ELBO + KL components)

- definition of  $F(q, p)$  (ELBO) from (3,4)

$$\int q(z) \log \frac{P(Y, z|\mathcal{D})}{q(z)} dz$$

- definition of  $D_{KL}(q||p)$  (KL divergence) from (3,5)

$$\int q(z) \log \frac{q(z)}{P(z|Y, \mathcal{D})} dz$$

Note:  $D_{KL}(q||p) = \int q \log \left( \frac{q}{p} \right) = - \int q \log \frac{p}{q}$ .

Based on the separation of the integral we get:

$$\underbrace{l(g_e, g_f, \gamma)}_{\text{likelihood}} = \underbrace{F(\ell, g_e, g_f, \gamma)}_{\text{ELBO}} + \underbrace{D_{KL}(q || p)}_{\text{KL divergence}} \rightarrow \text{this is (3.3)}$$

### 3.4 - Evidence Lower Bound

- derived for 3.2  $\rightarrow$  3.3

$$F(\ell, g_e, g_f, \gamma) = \int q(\ell, f) \log \frac{p(Y, h, f | g_e, g_f, \gamma)}{q(\ell, f)} d\ell df$$

- defined directly from first term of log-likelihood splitting derived for (3.3)

1) expanded marginal log-likelihood:

$$l(\gamma) = \int q(z) \log \left( \frac{p(z|\gamma)}{p(z)} \right) dz + \int q(z) \log \left( \frac{q(z)}{p(z|\gamma)} \right) dz$$

- 2)  $F$  is defined to be exactly the first integral

$$F(\ell, \gamma) := \int_a^b q(z) \log p(Y, z | \gamma) - \log(q(z)) dz$$

- a) places probability mass where likelihood is high (good fit)
- b) Shannon entropy of  $q$ , encourages  $q$  to be spread out, preventing the model from collapsing to a single point estimate.

### 3.5 - KL divergence

- derived for 3.2  $\rightarrow$  3.3

$$D_{KL}(q || p) = - \int q(\ell, f) \log \frac{p(\ell, f | Y, g_e, g_f, \gamma)}{q(\ell, f)} d\ell df$$

- 2nd term of split integral in (3.3)

$\int q(z) \log \left( \frac{q(z)}{p(z|Y, \gamma)} \right) dz \rightarrow$  definition of KL divergence from distribution  $q$  to distribution  $p(\cdot | Y)$

- negative sign on integral is identical to  $\int q \log \left( \frac{q}{p} \right) \rightarrow$  maximizing term inside  $\log \Rightarrow$  minimizing divergence

- KL divergence is non-negative so  $F(q, g_e, g_f, \gamma)$  is lower bound for log-likelihood (3.6)  $\rightarrow$  proof

### 3.6 - Variational Inequality

- $F(a, g_e, g_f, \gamma)$  is a lower bound for the log-likelihood!

$$l(g_e, g_f, \gamma) \geq F(a, g_e, g_f, \gamma)$$

with equality when  $q(h|f) = p(h|f | Y, g_e, g_f, \gamma)$

#### • Proof (via Jensen's Inequality)

- to prove that  $F$  is a lower bound, prove that the term  $D_{KL}$  is non-negative

- definition of K-L Divergence,

$$- D_{KL}(q||p) = \int q(z) \log\left(\frac{p(z|Y)}{q(z)}\right) dz$$

• let  $f(x) = \log(x)$  (concave)

• apply Jensen's inequality:  $E[\log(x)] \leq \log(E[x])$

$$\hookrightarrow \int q(z) \log\left[\frac{p(z|Y)}{q(z)}\right] dz \leq \log\left(\int q(z) \frac{p(z|Y)}{q(z)} dz\right)$$

- simplify integral inside log

$$\int q(z) \frac{p(z|Y)}{q(z)} dz = \underbrace{\int p(z|Y) dz}_\text{PdS integrates to 1} = 1 \quad (p(z|Y) \text{ integrates to 1!})$$

• plug into log()

$$\log\left(\int q(z) \frac{p(z|Y)}{q(z)} dz\right) = \log\left(\int p(z|Y) dz\right) = \log(1) = 0$$

• therefore (consider inequality)

$$- D_{KL}(q||p) = \int q(z) \log\left(\frac{p(z|Y)}{q(z)}\right) dz \leq 0$$

$$\text{So, } - D_{KL}(q||p) \leq 0 \Rightarrow D_{KL}(q||p) \geq 0$$

• Now consider (3.3) again,

$$l(\mathbb{Q}) = F(a, \mathbb{Q}) + D_{KL} \quad (\text{and } D_{KL}(q||p) \geq 0)$$

such that it must be true that

$$l(\mathbb{Q}) \geq F(a, \mathbb{Q}) \quad (3.6)$$

### 3.7 - maximize objective

Based on (3.6),  $\ell(g_e, g_f, \gamma) = \max_q F(q, g_e, g_f, \gamma)$  with maximization over all possible distributions  $q(l, f)$ .

- Can't minimize  $D_{KL}(q||p)$  directly because calculating  $p$  (posterior) requires the integral  $p(y)$ .
- $\ell(\theta)$  is fixed wrt  $q \Rightarrow \ell(\theta) = F(q) + D_{KL}(q||p) \Rightarrow D_{KL}(q||p) = \ell(\theta) - F(q)$
- so maximizing  $F(q)$  is identical to minimizing  $\ell(\theta) - D_{KL}(q||p)$ .  $\ell(\theta)$  is constant so this is equivalent to minimizing  $D_{KL}(q||p)$
- So finding  $q^*$  that maximizes  $F$  means that  $q^*$  is the distribution closest to the true posterior.

$$\Rightarrow F(q) = \ell(\theta) - D_{KL}(q||p)$$

make as small  
as possible.

Note: variational approach simplifies problem by maximizing  $F$  but restricting the family of distributions for  $q$ .

### 3.8 - mean field assumption

- restrict  $q$  to family  $Q$  of distributions that "fully-factorize":

$$Q = \{q: q(l, f) = \prod_{i=1}^n q_{l_i}(l_i) \prod_{j=1}^m q_{f_j}(f_j)\}$$

- constraint to make maximization for (3.7) feasible.

1) factorization indices:

•  $\prod_{i=1}^n q_{l_i}(l_i)$ : assumes loading for Sample 1 is independent of the loading for Sample 2

•  $\prod_{j=1}^m q_{f_j}(f_j)$ : assumes factor value for gene 1 is independent of gene 2.

$\hookrightarrow$  implication:  $q(l, f) = q(l)q(f)$  (can factor)

2) independence assumption

- expectation of the likelihood term in  $F$  simplifies. Note the cross-term in the Gaussian likelihood:  $f_i^T f_j$ .

$$E[f_i^T f_j] = \sum_l \sum_f q(l, f) (f_i^T f_j) dldf \Rightarrow$$

assume loadings + factors  
are independent

Assuming,  $q(l,f) = q(l)q(f)$ , the integral separates to be:

$$\left( \int q(l) l_i dl \right)^T \left( \int q(f) f_i df \right) = E_q[l_i]^T E_q[f_i]$$

Separation means that when  $q(l)$  is updated, the factors ( $f$ ) look like constants (expectations). This converts the complex optimization into an iterative EBNF problem

$\hookrightarrow$  To optimize  $F(a_l, b_f, g_e, g_f, \gamma)$  by alternating between optimizing over variables related to  $l$  [ $a_l, g_e$ ], over variables related to  $f$  [ $a_f, g_f$ ], and over  $\gamma$ .  
 $\hookrightarrow$  each step is guaranteed to increase  $F$ .

### EBNF Algo 1

initial values:  $q_l^{(0)}, q_f^{(0)}, g_e^{(0)}, g_f^{(0)}$

1)  $t \leftarrow 0$

2) repeats

3)  $t \leftarrow t + 1$

4)  $\gamma^{(t)} \leftarrow \arg \max_{\gamma} F(a_l^{(t-1)}, q_f^{(t-1)}, g_e^{(t-1)}, g_f^{(t-1)}, \gamma)$

5)  $q_l^{(t)}, g_e^{(t)} \leftarrow \arg \max_{q_l, g_e} F(a_l, b_f^{(t-1)}, g_e, g_f^{(t-1)}, \gamma^{(t)})$

6)  $q_f^{(t)}, g_f^{(t)} \leftarrow \arg \max_{q_f, g_f} F(a_l^{(t)}, b_f, g_e^{(t)}, g_f, \gamma^{(t)})$

7) until converged

8) return  $q_l^{(t)}, q_f^{(t)}, g_e^{(t)}, g_f^{(t)}, \gamma^{(t)}$

maximization  
steps

Step 4, update of  $\gamma$ , involves computing expected squared residuals!

$$\bar{R}_{ij}^2 := E_{q_l, q_f} [(Y_{ij} - l_i f_j)^2] = [Y_{ij} - E_q(l_i) E_q(f_j)]^2 - E_q(l_i)^2 E_q(f_j)^2 + E_q(l_i^2) E_q(f_j^2)$$

### 3.9 - definition

ELBO contains a term for the expected log-likelihood of the data. Noise (error) is Gaussian, so maximizing likelihood is equivalent to minimizing the Expected Squared Error.

$$\bar{R}_{ij}^2 := E_{q_l, q_f} [(Y_{ij} - l_i f_j)^2]$$



### 3.10 -

• We can't compute  $(Y_{ij} - lf_j)^2$  directly because  $l$  and  $f$  are distributions.

1) expand quadratic: Let  $x = lf_j$ :

$$(Y_{ij} - x)^2 = Y_{ij}^2 - 2Y_{ij}x + x^2$$

2) expectation

$$E[(Y_{ij} - lf_j)^2] = Y_{ij}^2 - 2Y_{ij}E[lf_j] + E[(lf_j)^2]$$

3) apply independence assumption ( $a(l, f) = g(l)g(f)$ )

$$\cdot E[lf_j] = \sum l_i E[f_j]$$

$$\cdot E[(lf_j)^2] = \sum l_i^2 E[f_j^2] = \sum l_i^2 \sum f_j^2$$

Sub in:

$$\bar{R}_{ij}^2 = Y_{ij}^2 - 2Y_{ij} \sum l_i \sum f_j + \sum l_i^2 \sum f_j^2$$

4) rearrange to align w/ B10)

let  $m_l = \sum l_i$ ,  $m_f = \sum f_j$ , consider the squared error using just the means:

$$(Y_{ij} - m_l m_f)^2 = Y_{ij}^2 - 2Y_{ij} m_l m_f + m_l^2 m_f^2$$

true expected error (step 3)

$$\text{Error}_{\text{true}} = Y_{ij}^2 - 2Y_{ij} m_l m_f + \sum l_i^2 \sum f_j^2$$

now write:

$$\bar{R}_{ij}^2 = (Y_{ij} - m_l m_f)^2 - m_l^2 m_f^2 + \sum l_i^2 \sum f_j^2$$

↳ accounts for the variance of the factors, if the model is very

uncertain about a loading  $f_j$  the term  $\sum f_j^2$  grows, increasing expected residual  $\bar{R}^2$ .

→ model increases estimated noise variance ( $1/\epsilon$ ) or shrink the priors to reduce uncertainty.

### 3.3 - The EBMM Problem

• observations  $x = (x_1, \dots, x_n)$

• underlying quantities  $\theta = (\theta_1, \dots, \theta_n)$

• independent Gaussian errors w/ known sd  $s = (s_1, \dots, s_n)$

• elements of  $\theta$  are iid from some distribution,  $g \in \mathcal{G}$

### 3.11 - observation model

$$x | \theta \sim N_n(\theta, \text{diag}(s_1^2, \dots, s_n^2))$$

• Assume error around the true parameter is Gaussian and independent for each data point.

### 3.12 - prior model

$\theta_1, \dots, \theta_n$  iid  $g$ ,  $g \in \mathcal{G}$

- assume all  $n$  parameters ( $\theta_j$ ) are drawn from the same underlying distribution
- $g$  belongs to family  $\mathcal{G}$ , must estimate the specific shape of  $g$  based on data  $x$ .

Now solve the FBML problem by fitting (3.11)-(3.12) by the following:

### 3.13 - solve for prior

- estimate prior  $g$  from data (Empirical Bayes)

$$\hat{g} = \arg \max_{g \in \mathcal{G}} \prod_j \int p(x_j | \theta_j, s_j) g(d\theta_j)$$

#### derivation

1) Likelihood: maximize probability of observing the data  $x$ . (depends on unobserved  $\theta$ )

2) Marginalize  $\theta$ : integrate out unknown  $\theta$ .

$$p(x_j | g) = \int p(x_j | \theta_j, s_j) p(\theta_j | g) d\theta_j$$

write as measure  $g(d\theta)$  for point masses at zero

3) Independence: since  $\theta_j$  are iid, probability of full vector  $x$  is product of

individual probabilities

$$L(g) = \prod_{j=1}^n p(x_j | g)$$

✓ how?

4) optimization: find  $\hat{g}$  ( $\in \mathcal{G}$ ) that maximizes  $L(g)$

• two approaches for maximizing  $L(g)$ :

- 1) parametric optimization
- 2) convex optimization  $\rightarrow$  numerical optimization algs

### 3.14 - posterior distribution

- estimate prior  $\hat{g} \rightarrow$  treat it as true prior and proceed w/ standard Bayesian inference.

$$\bullet P(\theta | x, s, \hat{g}) \propto \prod_j \hat{g}(\theta_j) p(x_j | \theta_j, s_j)$$

- comes from Bayes' Theorem!, posterior  $\propto$  prior  $\times$  likelihood

• prior:  $\hat{g}(\theta_j)$

• likelihood:  $p(x_j | \theta_j, s_j) \rightarrow$  normal dist. from 3.11

• posterior:  $P(\theta | x, s, \hat{g})$

• posterior: belief about

$\theta_j$  after viewing data  $x_j$ .

### 3.15 + 3.16 - Moments

- EBNF doesn't need full posterior distribution object (only 1st + 2nd moments)

(3.15) posterior mean (1st moment)

$$\bar{S}_j := \mathbb{E}(Q_j | x, s, \hat{g}) = \underbrace{\int}_{\mathbb{E}[\cdot]} \underbrace{Q_j p(Q_j | x, s, \hat{g}) dQ_j}_{\text{posterior}}$$

(3.16) 2nd moment

$$\bar{S}^2 := \mathbb{E}(Q_j^2 | x, s, \hat{g}) = \underbrace{\int}_{\mathbb{E}[\cdot^2]} \underbrace{Q_j^2 p(Q_j | x, s, \hat{g}) dQ_j}_{\text{posterior}}$$

↳ (3.15) had  $\mathbb{E}[(Y - f)^2]$ , need  $\mathbb{E}[f^2]$ ,  $\mathbb{E}[f^2]$  to compute residual error.

$$\Rightarrow \mathbb{E}[S^2] = \mathbb{E}[Q^2] + \text{Var}(Q)$$

### 3.17 - mapping

- EBNM( $x, s$ ) =  $(\hat{g}, p)$

• inputs

- $x$ : noisy data from matrix residuals
- $s$ : noise standard deviation (from precision  $\gamma$ )

• outputs

- $\hat{g}$ : fitted prior distribution
  - $p$ : posterior (moments  $\bar{S}$  and  $\bar{S}^2$ )
- Note: calculation of  $x, s$  in 3.4

### 3.4 - Connecting EBNF and EBNM Problems

- prove task (complex) of updating a single factor ( $f$ ) in Matrix Factorization is identical to solving the simple EBNM problem (Section 3.3)
- To update new model, need to calculate two vectors - pseudo data ( $\tilde{l}$ ) and pseudo standard errors ( $\tilde{s}_e$ ) → plug into solver

3.18 - main result :  $\underset{\tilde{l}, \tilde{s}_e}{\operatorname{argmax}} f(\tilde{l}, \tilde{s}_e) = \text{EBNM}(\tilde{l}, \tilde{s}_e)$

- update posterior ( $q_e$ ) and prior ( $g_e$ ) for the loadings vector  $\tilde{l}$ , holding factors ( $f$ ) and noise ( $\gamma$ ) fixed. Maximizing ELBO ( $f$ ) to do this is the same as doing EBNM on pseudo observations

3.19 + 3.20 - derive pseudo parameters

- match terms in ELBO w/ Gaussian log-likelihood.

1) isolate terms in the ELBO w/  $\hat{l}_{ik}$

$$\text{likelihood} = \text{E}_{\theta, f|f} [\log P(Y|l, f, \epsilon)]$$

'sub in likelihood'

$$= \text{E}_{\theta, f|f} \left[ \frac{1}{2} \sum_{ij} \gamma_{ij} (Y_{ij} - \sum_k \hat{l}_{ik} f_{jk})^2 \right]$$

2) consider single factor ( $k$ )

- update one factor at a time (update  $\hat{l}_{ik}$  holding on  $i$ th person ( $f_{ik}$ ))

$$R_{ij}^{-k} = Y_{ij} - \sum_{k \neq k} \hat{l}_{ik} f_{jk} \Rightarrow Y_{ij} \propto R_{ij}^{-k} + \hat{l}_{ik} f_{jk}$$

expand square w/  $\hat{l}_{ik}$  term:

$$\propto -\frac{1}{2} \sum_j \gamma_{ij} \text{E}_{\theta} [(R_{ij}^{-k} - \hat{l}_{ik} f_{jk})^2]$$

$$\propto -\frac{1}{2} \sum_j \gamma_{ij} \text{E}_{\theta} [(R_{ij}^{-k})^2 - 2\hat{l}_{ik} f_{jk} R_{ij}^{-k} + \hat{l}_{ik}^2 f_{jk}^2]$$

3) expectation and group by  $\hat{l}_{ik}$

- independence  $\Rightarrow \text{E}[l_i f] = \text{E}[l_i] \text{E}[f]$

$$\propto \sum_j \gamma_{ij} \left[ \hat{l}_{ik} R_{ij}^{-k} \text{E}[f_{jk}] - \frac{1}{2} \hat{l}_{ik}^2 \sum_j \text{E}[f_{jk}^2] \right]$$

arrange as quadratic equation in  $\hat{l}_{ik}$ :

$$\text{objective}(\hat{l}_{ik}) \propto -\frac{1}{2} \hat{l}_{ik}^2 \left( \sum_j \gamma_{ij} \text{E}[f_{jk}^2] \right) + \hat{l}_{ik} \left( \sum_j \gamma_{ij} R_{ij}^{-k} \text{E}[f_{jk}] \right)$$

4) align w/ NM problem

- log-likelihood of single data point  $x$  from  $N(\hat{l}_{ik}, s^2)$

$$\log N(x | \hat{l}_{ik}, s^2) \propto -\frac{1}{2s^2} (x - \hat{l}_{ik})^2 \Rightarrow -\frac{1}{2s^2} \hat{l}_{ik}^2 + \frac{x}{s^2} \hat{l}_{ik} - \frac{x^2}{2s^2}$$

5) compare coefficients ( $\hat{l}_{ik}^2$  and  $\hat{l}_{ik}$ )

$$\bullet \hat{l}_{ik}^2 : \frac{1}{2} = \sum_j \gamma_{ij} \text{E}[f_{jk}^2] \Rightarrow S_{\hat{l}}(\cdot) = \left( \sum_j \gamma_{ij} \text{E}[f_{jk}^2] \right)^{-\frac{1}{2}} \quad (3.20)$$

$$\bullet \hat{l}_{ik} : \frac{x}{s^2} = \sum_j \gamma_{ij} R_{ij}^{-k} \text{E}[f_{jk}] \xrightarrow{\text{solve for } \hat{l}_{ik}} \hat{l}_{ik} = S_{\hat{l}} x_B = \frac{\sum_j \gamma_{ij} R_{ij}^{-k} \text{E}[f_{jk}]}{\sum_j \gamma_{ij} \text{E}[f_{jk}^2]} \Rightarrow \hat{l}(N) = \frac{\sum_j \gamma_{ij} Y_{ij} \text{E}[f_{jk}]}{\sum_j \gamma_{ij} \text{E}[f_{jk}^2]} \quad (3.19)$$

### 3.21 + 3.22 - Moments

- need 1st + 2nd moments of factors ( $f_j$ ), defined as:

$$\bullet \bar{f} := (\text{E}_{\text{af}}[f_j])$$

$$\bullet \bar{f}^2 := (\text{E}_{\text{af}}[f_j^2])$$

• Note 1 (3.19) uses  $\bar{f}^2$ , not  $(\bar{f})^2$

$$\bullet \bar{f}^2 = (\bar{f})^2 + \text{Var}(f)$$

• If  $f$  is uncertain  $\rightarrow \bar{f}^2$  is large

• large denominator  $\rightarrow$  loading ( $\hat{l}$ ) shrinks  $\rightarrow 0$

### 3.23 - 3.26 - known factor intuition

- what if we knew  $f$  and  $\gamma$ ?

↳ model becomes  $n$  independent regressions of the rows of  $Y$  on  $f$ ,  
and the maximum likelihood estimate for  $\lambda$  has  $n$  elements!

(3.23) Regression estimate is  $f$  is fixed constants, estimating  $\lambda_i$  is a

weighted least squares (WLS) regression of the data row  $y_i$  on the factor  $f$ .

WLS solution is  $\hat{\beta}_i = \frac{\sum w_{ij} y_i}{\sum w_{ij} f_j}$  where  $w = \gamma$ ,  $x = f$ ,  $y = y_i$

$$\hat{\lambda}_i = \frac{\sum j \gamma_{ij} y_i f_j}{\sum j \gamma_{ij} f_j^2}$$

### (3.24) Regression Standard Error

- SE of WLS estimate is  $(x^T w x)^{-1/2}$

$$\Rightarrow s_i = (\sum j \gamma_{ij} f_j^2)^{-1/2}$$

### (3.25) + (3.26)

- goal: prove  $\hat{\lambda}_i \sim N(\lambda_i, s_i^2)$

• estimator from WLS solution is  $\hat{\lambda}_i = \sum j \frac{\gamma_{ij} y_i f_j}{\sum j \gamma_{ij} f_j^2}$

• distribution shape:  $y_{ij}$  is drawn from a Normal distribution, and  $\hat{\lambda}_i$  is a linear combination of the  $y_{ij}$  values so  $\hat{\lambda}_i \sim \text{Normal}$  (find Mean and Variance)

• prove mean is unbiased ( $E(\hat{\lambda}_i) = \lambda_i$ ) true model:  $y_{ij} = \lambda_i f_j + \epsilon_{ij}$

$$\hat{\lambda}_i = \sum j \frac{\gamma_{ij} (\lambda_i f_j + \epsilon_{ij})}{\sum j \gamma_{ij} f_j^2} \Rightarrow \hat{\lambda}_i = \sum j \frac{\gamma_{ij} \lambda_i f_j^2}{\sum j \gamma_{ij} f_j^2} + \sum j \frac{\gamma_{ij} \epsilon_{ij} f_j}{\sum j \gamma_{ij} f_j^2}$$

• factor out constant  $\lambda_i \Rightarrow \hat{\lambda}_i = \lambda_i \left[ \frac{\sum j \gamma_{ij} f_j^2}{\sum j \gamma_{ij} f_j^2} \right] + " " \Rightarrow \hat{\lambda}_i = \lambda_i + "$

→ consider expectation:

$$E[\hat{\lambda}_i] = E[\hat{\lambda}_i] + E\left[\frac{\sum_j \gamma_{ij} f_j}{\sum_j \gamma_{ij} f_j^2}\right] \quad E[\sum_j \gamma_{ij}] = 0$$

$$\text{So, } E[\hat{\lambda}_i] = \hat{\lambda}_i + 0$$

- show  $\text{Var}(\hat{\lambda}_i) = s_i^2$ : Note  $\text{Var}(cx) = c^2 \text{Var}(x)$  w/  $y_{ij}$  independent errors,
  - consider denominator  $D = \sum_j \gamma_{ij} f_j^2$

$$\text{Var}(\hat{\lambda}_i) = \text{Var}\left(\frac{\sum_j \gamma_{ij} f_j y_{ij}}{D}\right) = \frac{1}{D^2} \sum_j (\gamma_{ij} f_j)^2 \text{Var}(y_{ij})$$

$$\text{Sub in } \text{Var}(y_{ij}) = \frac{1}{\gamma_{ij} f_j^2},$$

$$= \frac{1}{D^2} \sum_j \gamma_{ij}^2 f_j^2 (1/\gamma_{ij} f_j^2)$$

$$= \frac{1}{D^2} \sum_j \gamma_{ij}^2 f_j^2 \text{ where } D = \sum_j \gamma_{ij} f_j^2$$

$$\text{So, } \text{Var}(\hat{\lambda}_i) = \frac{1}{D^2} \cdot D = \frac{1}{D} \Rightarrow \text{Var}(\hat{\lambda}_i) = \frac{1}{\sum_j \gamma_{ij} f_j^2}$$

$$\text{such that } \text{SD}(\hat{\lambda}_i) = (\sum_j \gamma_{ij} f_j^2)^{-1/2} \quad \checkmark$$

Therefore, we showed mean + variance so we have  $\hat{\lambda}_i \sim N(\hat{\lambda}_i, s_i^2)$  (3.25)

(3.26) - prior

$$\cdot \lambda_1, \dots, \lambda_n \text{ iid } g_\lambda, g \in \mathcal{G} \quad (3.26)$$

So if we combine (3.26) and (3.25) (prior and likelihood), we get the structure for each loading:

1) observation:  $\hat{\lambda}_i | \lambda_i \sim N(\lambda_i, s_i^2)$

2) prior:  $\lambda_i \sim g_\lambda$

which is the definition of the EBMM problem