

Mean Field Asymptotics in Statistical Learning.

Apr 7

Lecture 20. Derivation of AMP. I.

[Montanari, 2011]

Graphical model concepts in CS.

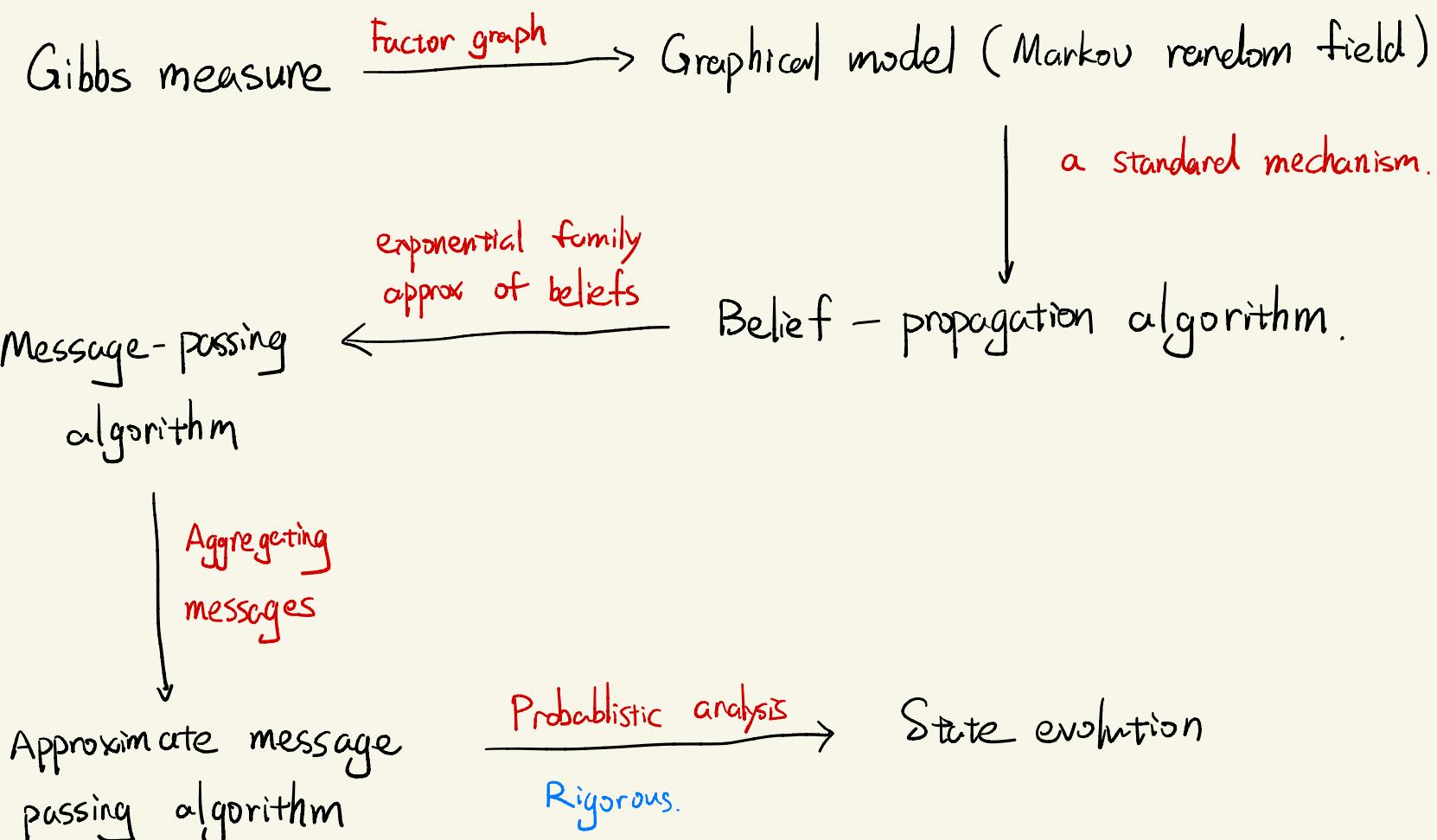
[Rangan, Schniter, Fletcher, 2016]

Vector AMP

[Donoho, Matalki, Montanari, 2009]

MP algorithm for CS : I, motivation and construction.

① The overall methodology.



② Markov random field.

Probability measure on $\Omega = \bigotimes_{i \in [N]} X_i = X^{\otimes N}$ $\mu(x) \propto \exp\{-\beta H(x)\}$

$$\mu(x) = \frac{1}{Z} \prod_{a \in F} \psi_a(x|_a) \prod_{i \in V} \psi_i(x_i), \quad \forall a \subseteq V, \forall a \in F$$

$$x|_a = (x_1, \dots, x_{|a|}) \in \mathbb{R}^{|a|}$$

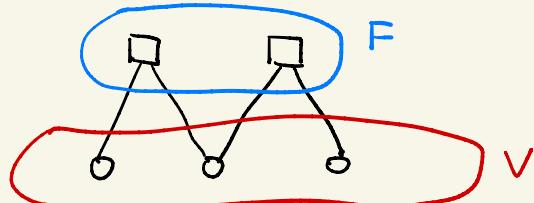
$$\psi_a: X^{\otimes |a|} \rightarrow \mathbb{R}_{\geq 0}, \quad a \in F$$

$$\psi_i: X \rightarrow \mathbb{R}_{\geq 0}, \quad i \in V.$$

Factor graph visualization

$$G = (V, F, E), \quad |V| = N, \quad |F| = M,$$

$$E \subseteq \{(i, a) : i \in V, a \in F\}.$$



Example 1: 1-D Ising model with 3 spins.

$$\Omega = \{\pm 1\}^3$$

$$\mu(x_1, x_2, x_3) \propto e^{-\beta(x_1x_2 + x_2x_3)} = \psi_{a_1}(x_1, x_2) \times \psi_{a_2}(x_2, x_3).$$

$$\partial a_1 = \{1, 2\}, \quad \partial a_2 = \{2, 3\}.$$

$$\psi_{a_1}(x_1, x_2) = e^{-\beta x_1 x_2}, \quad \psi_{a_2}(x_2, x_3) = e^{-\beta x_2 x_3}$$

$$\psi_1(x_1) = \psi_2(x_2) = \psi_3(x_3) = 1.$$

Example 2: Bayes linear model.

$$y = Ax_0 + w \in \mathbb{R}^n$$

$$x_0 \in \mathbb{R}^d, \quad A \in \mathbb{R}^{n \times d}, \quad w \in \mathbb{R}^n,$$

$$x_{0,i} \sim \text{iid } P_0(x). \quad w_i \sim \text{iid } N(0, \sigma^2).$$

$$\begin{aligned} \mu(x) = P(x | A, y) &\propto \exp \left\{ - \frac{\|y - Ax\|_2^2}{2\sigma^2} \right\} \prod_{i=1}^d P_0(x_i) \\ &= \prod_{a=1}^n \exp \left\{ - \frac{(y_a - \langle A_a, x \rangle)^2}{2\sigma^2} \right\} \prod_{i=1}^d P_0(x_i) \end{aligned}$$

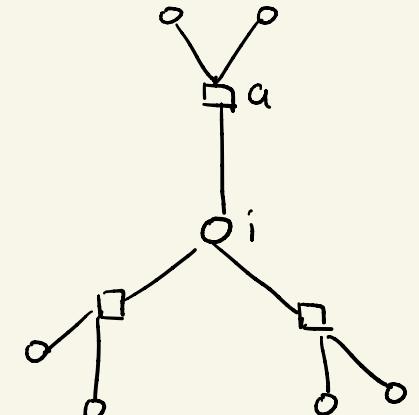
$$\psi_a(x) = \exp \left\{ - \frac{(y_a - \langle A_a, x \rangle)^2}{2\sigma^2} \right\}$$

$$\psi_i(x_i) = P_0(x_i).$$

Remark: For $\mu \in \mathcal{P}(\Omega)$, the MRF representation is not unique.

Task: Compute the marginal distributions

$$\forall i \in [d], \quad \mu_i(x_i) \equiv \int \mu(x) \prod_{j \neq i} dx_j$$



(3) Belief propagation algorithm on tree graphs

Let G be a tree factor graph.

Denote $V_{a \rightarrow i}$ be all reachable vertices starting at "a" by blocking "i",

$V_{i \rightarrow a}$ - - - - - - - - - "i" by blocking "a".

$F_{a \rightarrow i}$ - - - - factors starting at "a" by blocking "i"

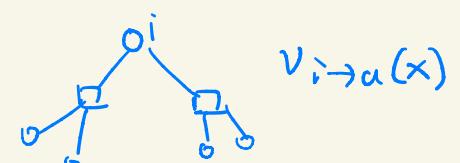
$F_{i \rightarrow a}$ - - - - factors starting at "i" by blocking "a".

Define

$$\hat{\nu}_{a \rightarrow i}(x) \propto \prod_{b \in F_{a \rightarrow i}} \psi_b(x_{ab}) \prod_{j \in F_{a \rightarrow i}} \psi_j(x_j)$$



$$\nu_{i \rightarrow a}(x) \propto \prod_{b \in F_{i \rightarrow a}} \psi_b(x_{ab}) \prod_{j \in F_{i \rightarrow a}} \psi_j(x_j)$$



$$x_{\setminus i} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_d).$$

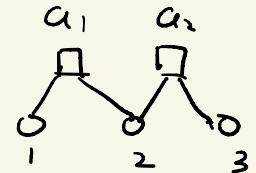
Define $\hat{\mu}_{a \rightarrow i}(x_i) \propto \sum_{x_{\setminus i}} \hat{v}_{a \rightarrow i}(x) \propto \sum_{x_{\setminus i}} \prod_{b \in F_{a \rightarrow i}} \psi_b(x_{a b}) \prod_{j \in F_{a \rightarrow i}} \psi_j(x_j) \in \mathcal{P}(X)$

$$\mu_{i \rightarrow a}(x_i) \propto \sum_{x_{\setminus i}} v_{i \rightarrow a}(x) \propto \sum_{x_{\setminus i}} \prod_{b \in F_{i \rightarrow a}} \psi_b(x_{a b}) \prod_{j \in F_{i \rightarrow a}} \psi_j(x_j) \in \mathcal{P}(X)$$

Claim : $\left\{ \begin{array}{l} \mu_{i \rightarrow a}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial_i \setminus a} \hat{\mu}_{b \rightarrow i}(x_i) \\ \hat{\mu}_{a \rightarrow i}(x_i) \propto \sum_{x_{\setminus a \setminus i}} \psi_a(x_{a a}) \prod_{j \in \partial a \setminus i} \mu_{j \rightarrow a}(x_j) \end{array} \right.$

} *

$$\mu_i(x_i) \propto \mu_{i \rightarrow a}(x_i) \hat{\mu}_{a \rightarrow i}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial_i} \hat{\mu}_{b \rightarrow i}(x_i)$$



Example : $\mu(x_1, x_2, x_3) \propto \psi_{a_1}(x_1, x_2) \cdot \psi_{a_2}(x_2, x_3)$

$$\mu_{1 \rightarrow a_1}(x_1) \propto 1.$$

$$\hat{\mu}_{a_1 \rightarrow 2}(x_2) \propto \sum_{x_1} \psi_{a_1}(x_1, x_2) \quad \mu_{1 \rightarrow a_1}(x_1) = \sum_{x_1} \psi_{a_1}(x_1, x_2)$$

$$\mu_{2 \rightarrow a_2}(x_2) \propto \hat{\mu}_{a_1 \rightarrow 2}(x_2) \psi_{a_2}(x_2) = \hat{\mu}_{a_1 \rightarrow 2}(x_2) = \sum_{x_1} \psi_{a_1}(x_1, x_2)$$

$$\hat{\mu}_{a_2 \rightarrow 3}(x_3) \propto \sum_{x_2} \psi_{a_2}(x_2, x_3) \quad \mu_{2 \rightarrow a_2}(x_2) = \sum_{x_2} \psi_{a_2}(x_2, x_3) \sum_{x_1} \psi_{a_1}(x_1, x_2)$$

$$\begin{aligned} \mu_3(x_3) &\propto \psi_3(x_3) \hat{\mu}_{a_2 \rightarrow 3}(x_3) = \sum_{x_2} \psi_{a_2}(x_2, x_3) \sum_{x_1} \psi_{a_1}(x_1, x_2) \\ &= \sum_{x_1, x_2} \psi_{a_1}(x_1, x_2) \psi_{a_2}(x_2, x_3). \end{aligned}$$

Def (BP on trees). (Make * iterative updates)

$$\left\{ \hat{\mu}_{a \rightarrow i}^k, \mu_{i \rightarrow a}^k \right\}_{a \in F, i \in V, k \geq 0} \subseteq \mathcal{P}(X) \text{ beliefs.}$$

Update rule: Calculate $\left\{ \hat{\mu}_{a \rightarrow i}^{k+1}, \mu_{i \rightarrow a}^{k+1} \right\}_{a,i}$ using $\left\{ \hat{\mu}_{a \rightarrow i}^k, \mu_{i \rightarrow a}^k \right\}_{a,i}$.

$$\left\{ \begin{array}{l} \mu_{i \rightarrow a}^{k+1}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial_i \setminus a} \hat{\mu}_{b \rightarrow i}^k(x_i) \\ \hat{\mu}_{a \rightarrow i}^{k+1}(x_i) \propto \sum_{x_{\setminus a \setminus i}} \psi_a(x_{a a}) \prod_{j \in \partial a \setminus i} \mu_{j \rightarrow a}^k(x_j) \end{array} \right.$$

Extract marginals: $\mu_i^k(x_i) \propto \mu_{i \rightarrow a}^k(x_i) \hat{\mu}_{a \rightarrow i}^k(x_i) \propto \psi_i(x_i) \prod_{b \in \partial_i} \hat{\mu}_{b \rightarrow i}^k(x_i)$.

Initialization: $\hat{\mu}_{a \rightarrow i}^0(x_i), \mu_{i \rightarrow a}^0(x_i)$: arbitrary, often uniform.

Thm: BP algorithm on trees will give exact results after $2K$ iterations.

where K is diameter of the tree (length of longest path)

$$\mu_i^{2K}(x_i) = \mu_i(x_i) \quad \mu_{i \rightarrow a}^{2K}(x_i) = \mu_{i \rightarrow a}(x_i), \quad \hat{\mu}_{a \rightarrow i}^{2K}(x_i) = \hat{\mu}_{a \rightarrow i}(x_i).$$

④ Loopy belief propagation on general graphs.

Def (Loopy BP).

$$\{\hat{m}_{a \rightarrow i}^k, m_{i \rightarrow a}^k\}_{a \in F, i \in V, k \geq 0} \subseteq \mathcal{P}(X) \text{ beliefs.}$$

Update rule:

$$\begin{cases} m_{i \rightarrow a}^{k+1}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial i \setminus a} \hat{m}_{b \rightarrow i}^k(x_i) \\ \hat{m}_{a \rightarrow i}^{k+1}(x_i) \propto \sum_{x_{\partial a \setminus i}} \psi_a(x_{\partial a}) \prod_{j \in \partial a \setminus i} m_{j \rightarrow a}^k(x_j) \end{cases}$$

Extract marginals: $m_i^k(x_i) \propto m_{i \rightarrow a}^k(x_i) \hat{m}_{a \rightarrow i}^k(x_i) \propto \psi_i(x_i) \prod_{b \in \partial i} \hat{m}_{b \rightarrow i}^k(x_i)$.

Example: LASSO w/ temperature β

$$\begin{cases} m_{i \rightarrow a}^{k+1}(x_i) \propto \prod_{b \in \partial i \setminus a} \hat{m}_{b \rightarrow i}^k(x_i) \times \exp\{-\beta \lambda |x_i|\} \\ \hat{m}_{a \rightarrow i}^{k+1}(x_i) \propto \int dx_{\partial i} \exp\left\{-\frac{\beta}{2} (y_a - \langle A_a, x \rangle)^2\right\} \prod_{j \in \partial a \setminus i} m_{j \rightarrow a}^k(x_j) \\ m_i^k(x_i) \propto \prod_{b \in \partial i} \hat{m}_{b \rightarrow i}^k(x_i) \times \exp\{-\beta \lambda |x_i|\} \end{cases}$$

Remark: \otimes We hope $m_i^k(x_i) \rightarrow m_i(x_i)$ as $k \rightarrow \infty$.

\otimes For many models, especially mean field models, this happens in practice.

\otimes Practical when there is closed form update.

\otimes When ψ_a, ψ_i 's are not simple, no closed form update, Not practical.

⑤ From BP to message passing algorithm.

Consider the case when $X = \mathbb{R}$.

Gaussian approximation of (exponential family approx).

Wrong intuition: the beliefs are approximately Gaussian.

Intuition: In the update rule, only the mean and var of input beliefs are important.

Input beliefs can be approximated by Gaussian dist. in the update rule.

Real belief is still non-Gaussian.

Def (Message passing algorithm)

mean and variance of beliefs.

$$\{m_{i \rightarrow a}^k, v_{i \rightarrow a}^k, \hat{m}_{a \rightarrow i}^k, \hat{v}_{a \rightarrow i}^k\}_{a \in F, i \in V, k \geq 0} \subseteq \mathbb{R}. \text{ Messages.}$$

Update rule: Calculating $\{m_{i \rightarrow a}^{k+1}, v_{i \rightarrow a}^{k+1}, \hat{m}_{a \rightarrow i}^{k+1}, \hat{v}_{a \rightarrow i}^{k+1}\}$ using $\{m_{i \rightarrow a}^k, v_{i \rightarrow a}^k, \hat{m}_{a \rightarrow i}^k, \hat{v}_{a \rightarrow i}^k\}$.

$$p_{i \rightarrow a}^k(x_i) = \frac{1}{\sqrt{2\pi v_{i \rightarrow a}^k}} \exp \left\{ -\frac{(x_i - m_{i \rightarrow a}^k)^2}{2v_{i \rightarrow a}^k} \right\}. \quad \text{density of } N(m_{i \rightarrow a}^k, v_{i \rightarrow a}^k)$$

$$\hat{p}_{a \rightarrow i}^k(x_i) = \frac{1}{\sqrt{2\pi \hat{v}_{a \rightarrow i}^k}} \exp \left\{ -\frac{(x_i - \hat{m}_{a \rightarrow i}^k)^2}{2\hat{v}_{a \rightarrow i}^k} \right\}. \quad \text{density of } N(\hat{m}_{a \rightarrow i}^k, \hat{v}_{a \rightarrow i}^k)$$

$$r_{i \rightarrow a}^{k+1}(x_i) \propto 4_i(x_i) \prod_{b \in \partial_i \setminus a} \hat{p}_{b \rightarrow i}^k(x_i)$$

$$\hat{r}_{a \rightarrow i}^{k+1}(x_i) \propto \sum_{x_{\partial a \setminus i}} 4_a(x_{\partial a}) \prod_{j \in \partial a \setminus i} p_j^k(x_j)$$

$$(m_{i \rightarrow a}^{k+1}, v_{i \rightarrow a}^{k+1}) = (\text{mean, variance}) \text{ under } r_{i \rightarrow a}^{k+1}(x_i)$$

$$(\hat{m}_{a \rightarrow i}^{k+1}, \hat{v}_{a \rightarrow i}^{k+1}) = (\text{mean, variance}) \text{ under } \hat{r}_{a \rightarrow i}^{k+1}(x_i)$$

Extract marginals:

$$r_i^k(x_i) \propto \psi_i(x_i) \prod_{b \in \partial_i} \hat{p}_{b \rightarrow i}^{k-1}(x_i)$$

$\{p_{i \rightarrow a}^k, \hat{p}_{a \rightarrow i}^k\}$ input beliefs, assumed to be Gaussian

$\{r_{i \rightarrow a}^k, \hat{r}_{a \rightarrow i}^k\}$ output beliefs, possibly non-Gaussian

Example: LASSO w/ temperature β .

$$\begin{aligned} \otimes \quad r_{i \rightarrow a}^{k+1}(x_i) &\propto \exp\{-\beta \lambda |x_i|\} \times \exp\left\{-\sum_{b \neq a} \frac{(x_i - \hat{m}_{b \rightarrow i}^k)^2}{2\hat{v}_{b \rightarrow i}^k}\right\} \\ &\propto \exp\left\{-\beta \left(\frac{(x_i - \theta_{i \rightarrow a}^k)^2}{2\zeta_{i \rightarrow a}^k} + \lambda |x_i| \right)\right\}. \end{aligned}$$

$$\left\{ \begin{array}{l} \beta(\zeta_{i \rightarrow a}^k)^{-1} = \sum_{b \neq a} (\hat{v}_{b \rightarrow i}^k)^{-1} \\ \beta(\zeta_{i \rightarrow a}^k)^{-1} \theta_{i \rightarrow a}^k = \sum_{b \neq a} (\hat{v}_{b \rightarrow i}^k)^{-1} \hat{m}_{b \rightarrow i}^k \\ m_{i \rightarrow a}^{k+1} = \mathbb{E}_{X_i \sim \pi(\beta, \lambda, \theta_{i \rightarrow a}^k, \zeta_{i \rightarrow a}^k)} [X_i] \\ v_{i \rightarrow a}^{k+1} = \text{Var}_{X_i \sim \pi(\beta, \lambda, \theta_{i \rightarrow a}^k, \zeta_{i \rightarrow a}^k)} [X_i] \end{array} \right.$$

$$\text{where } \pi(\beta, \lambda, \theta, \zeta) \sim \exp\left\{-\beta \left(\frac{(x - \theta)^2}{2\zeta} + \lambda |x| \right)\right\}.$$

$$\textcircled{X} \quad \hat{Y}_{a \rightarrow i}^{k+1}(x_i) \propto \int_{\mathbb{R}^{d-1}} \exp \left\{ -(\beta/2) (y_a - \langle A_a, x \rangle^2) \right\} \\ \times \exp \left\{ - \sum_{j \neq i} \frac{(x_j - m_{j \rightarrow a}^k)^2}{2 \hat{\nu}_{j \rightarrow a}^{k+1}} \right\} dx_{-i}.$$

This is a Gaussian integration, which gives a Gaussian density.

$$\propto \exp \left\{ - \frac{(x_i - \hat{m}_{a \rightarrow i}^{k+1})^2}{2 \hat{\nu}_{a \rightarrow i}^{k+1}} \right\}.$$

$$\begin{cases} \hat{m}_{a \rightarrow i}^{k+1} = \frac{y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k}{A_{ai}} \\ \hat{\nu}_{a \rightarrow i}^{k+1} = \left(\sum_{j \neq i} A_{aj}^2 \nu_{j \rightarrow a}^k + 1/\beta \right) / A_{ai}^2. \end{cases}$$

\textcircled{X} Simplification as $\beta \rightarrow \infty$.

$$\lim_{\beta \rightarrow \infty} \mathbb{E}_{x \sim \pi(\beta, \lambda, \theta, \zeta)} [x] = \arg \min_x \left\{ \frac{(x - \theta)^2}{2s} + \lambda |x| \right\} \\ = \eta(\theta; \lambda \zeta).$$

$$\lim_{\beta \rightarrow \infty} \beta \cdot \text{Var}_{x \sim \pi(\beta, \lambda, \theta, \zeta)} [x] = \begin{cases} \zeta, & \text{if } |\theta| > \lambda \zeta \\ 0, & \text{if } |\theta| < \lambda \zeta \end{cases} \\ = \zeta \cdot \eta'(\theta; \lambda \zeta).$$

We are short of alphabets. Abuse notation:

$$\hat{m}_{a \rightarrow i} \leftarrow \lim_{\beta \rightarrow \infty} \hat{m}_{a \rightarrow i}, \quad m_{i \rightarrow a} \leftarrow \lim_{\beta \rightarrow \infty} m_{i \rightarrow a} \\ \hat{\nu}_{a \rightarrow i} \leftarrow \lim_{\beta \rightarrow \infty} \beta \hat{\nu}_{a \rightarrow i}, \quad \nu_{i \rightarrow a} \leftarrow \lim_{\beta \rightarrow \infty} \nu_{i \rightarrow a}.$$

$$\begin{cases} (\zeta_{i \rightarrow a}^k)^{-1} = \sum_{b \neq a} (\hat{\nu}_{b \rightarrow i}^k)^{-1} \\ (\zeta_{i \rightarrow a}^k)^{-1} \theta_{i \rightarrow a}^k = \sum_{b \neq a} (\hat{\nu}_{b \rightarrow i}^k)^{-1} \hat{m}_{b \rightarrow i}^k \\ m_{i \rightarrow a}^{k+1} = \eta(\theta_{i \rightarrow a}^k; \lambda \zeta_{i \rightarrow a}^k) \\ \nu_{i \rightarrow a}^{k+1} = \zeta_{i \rightarrow a}^k \cdot \eta'(\theta_{i \rightarrow a}^k, \lambda \zeta_{i \rightarrow a}^k) \\ \hat{m}_{a \rightarrow i}^{k+1} = \frac{y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k}{A_{ai}} \\ \hat{\nu}_{a \rightarrow i}^{k+1} = \left(\sum_{j \neq i} A_{aj}^2 \nu_{j \rightarrow a}^k + 1 \right) / A_{ai}^2. \end{cases}$$

$$\text{Define } T_{a \rightarrow i}^{k+1} = A_{ai}^2 \hat{v}_{a \rightarrow i}^{k+1} = \sum_{j \neq i} A_{aj}^2 v_{j \rightarrow a}^k + 1$$

$$z_{a \rightarrow i}^{k+1} = A_{ai} \hat{m}_{a \rightarrow i}^{k+1} = y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k .$$

$$\Rightarrow \theta_{i \rightarrow a}^k = \frac{\sum_{b \neq a} A_{bi} z_{b \rightarrow i}^k / T_{b \rightarrow i}^k}{\sum_{b \neq a} A_{bi}^2 / T_{b \rightarrow i}^k}$$

$$\varsigma_{i \rightarrow a}^k = \frac{1}{\sum_{b \neq a} A_{bi}^2 / T_{b \rightarrow i}^k}$$

