

# Mean Field Asymptotics in Statistical Learning.

Apr 7

## Lecture 20. Derivation of AMP. I.

[Montanari, 2011]

Graphical model concepts in CS.

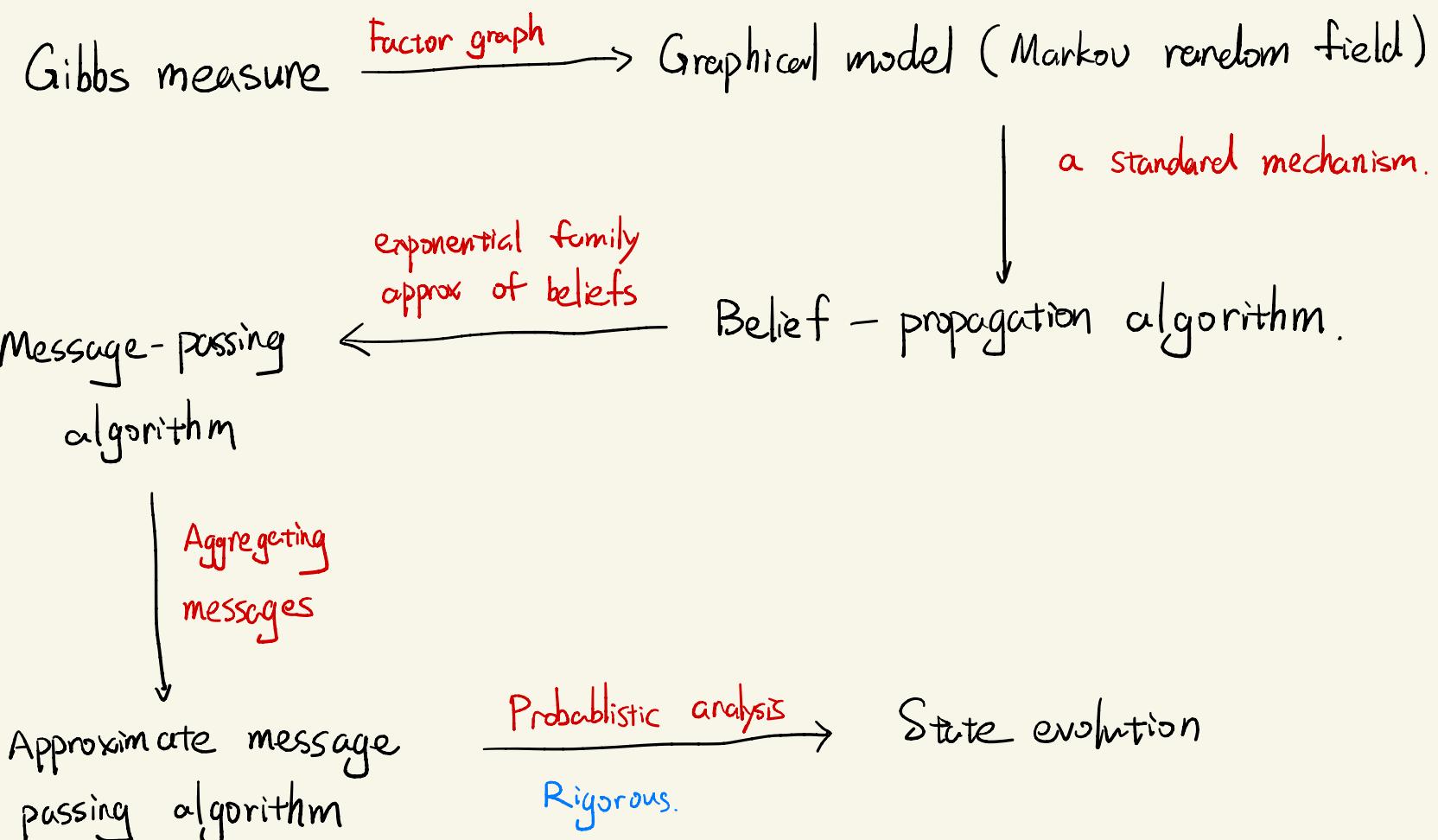
[Rangan, Schniter, Fletcher, 2016]

Vector AMP

[Donoho, Matalki, Montanari, 2009]

MP algorithm for CS : I, motivation and construction.

### ① The overall methodology.



### ② Markov random field.

$$V = [N] \quad F = [M].$$

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

$$\mu(x) = \frac{1}{Z} \prod_{a \in F} \psi_a(x_{\partial a}) \prod_{i \in V} \psi_i(x_i).$$

$\partial a \subseteq V, \forall a \in F$   
 $x_{\partial a} = (x_1, \dots, x_{|\partial a|}) \in \mathbb{R}^{|\partial a|}$

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

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

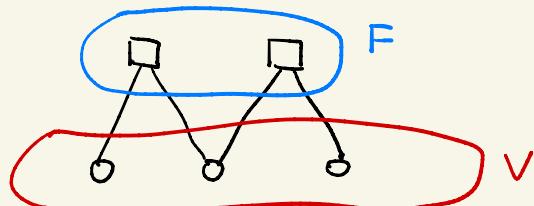
$$a \in F, \quad \partial a$$

$$i \in V, \quad \partial i$$

#### 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.

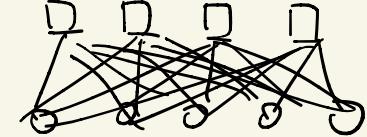
$$\Omega = \{\pm 1\}^3$$

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

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

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

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



Example 2: Bayes linear model.

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

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

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

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

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

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

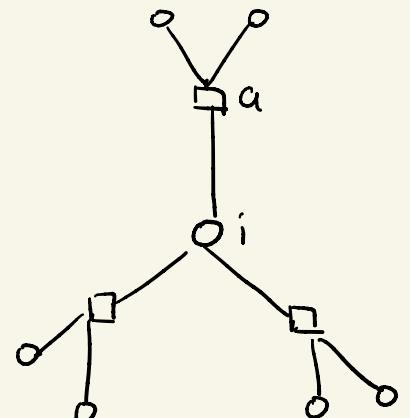
$$\psi(x) = \mu(x)$$

$$\mu(x) \propto \psi(x)$$

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"

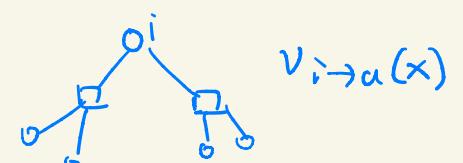
$F_{i \rightarrow a}$  - - - - factors starting at "i" by blocking "a".

Define

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



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



$$x_{\sim i} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$$

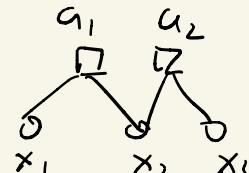
$$\begin{aligned}\hat{\mu}_{a \rightarrow i}(x_i) &\propto \sum_{x_{\sim i}} \hat{v}_{a \rightarrow i}(x) \in \mathcal{P}(X) \\ \mu_{i \rightarrow a}(x_i) &\propto \sum_{x_{\sim i}} v_{i \rightarrow a}(x) \in \mathcal{P}(X).\end{aligned}$$

Claim:

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

$$\mu_i(x_i) = \int \mu(x) dx_{\sim i} \propto \mu_{i \rightarrow a}(x) \hat{\mu}_{a \rightarrow i}(x) \propto \psi_i(x_i) \prod_{b \in \partial i} \hat{\mu}_{b \rightarrow i}(x)$$

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



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

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

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

$$\begin{aligned}\hat{\mu}_{a_2 \rightarrow 3}(x_3) &\propto \sum_{x_2} \psi_{a_2}(x_2, x_3) \mu_{2 \rightarrow a_2}(x_2) \\ &= \sum_{x_2} \psi_{a_2}(x_2, x_3) \sum_{x_1} \psi_{a_1}(x_1, x_2).\end{aligned}$$

$$\mu_3(x_3) \propto \psi_3(x_3) \hat{\mu}_{a_2 \rightarrow 3}(x_3)$$

Def (Belief propagation algorithm, on trees).

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

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

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

$$\text{Extract marginal: } \mu_{i \rightarrow a}^{k+1}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial i} \hat{\mu}_{b \rightarrow i}^k(x)$$

Thm: BP algorithm on trees will be exact after  $2K$  iteration  
 $K$  is the diameter of the tree.

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

#### ④ Loopy BP algorithm on general graphs.

Def (Loopy BP).

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

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

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

$$\text{Extract marginal: } \mu_i^{k+1}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial i} \hat{\mu}_{b \rightarrow i}^k(x_i)$$

Example: LASSO. w/ temperature  $\beta$ .

$$\mu_\beta(x) = \prod_{a=1}^n \exp\left\{-\left(\frac{\beta}{2}\right) (y_a - \langle A_a, x \rangle)^2\right\} \cdot \prod_{i=1}^d \exp\{-\beta \lambda |x_i|\}.$$

$$\begin{cases} \hat{\mu}_{a \rightarrow i}^{k+1}(x_i) \propto \int_{\mathbb{R}^{d-1}} \exp\left\{-\left(\frac{\beta}{2}\right)(y_a - \langle A_a, x \rangle)^2\right\} \prod_{j \in \partial a \setminus i} \hat{\mu}_{j \rightarrow a}^k(x_j) \prod_{j \neq i} dx_j \\ \mu_{i \rightarrow a}^{k+1}(x_i) \propto \exp\{-\beta \lambda |x_i|\} \times \prod_{b \in \partial i \setminus a} \hat{\mu}_{b \rightarrow i}^k(x_i) \\ \mu_i^{k+1}(x_i) \propto \prod_{b \in \partial i} \hat{\mu}_{b \rightarrow i}^k(x_i) \times \exp\{-\beta \lambda |x_i|\}. \end{cases}$$

Remark:  $\otimes$  We hope  $\mu_i^k(x_i) \rightarrow \tilde{\mu}_i(x_i)$  as  $k \rightarrow \infty$ .

$\otimes$  This will happen in practice, in many case.

$\otimes$  Practical when  $\psi_a, \psi_i$ , are simple.

$\otimes$  Not very practical for LASSO.

⑤ From BP to message passing. ( $X = \{\pm 1\}, X = \mathbb{R}^k$ )

Consider  $X = \mathbb{R}$

Gaussian approx.

Wrong intuition: the beliefs are approximately Gaussian.

Intuition: In the update rule, only mean and variance of incoming beliefs are important.

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

Real beliefs are still non-Gaussian.

Def (Message passing).

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

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

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

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

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

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

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

Extract marginals:

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