

Parameter Learning in MN

Amr Ahmed

10708-F08 Recitation

Outline

- CRF
- Learning CRF for 2-d image segmentation
- IPF parameter sharing revisited

Log-linear Markov network (most common representation)

- **Feature** is some function $\phi[\mathbf{D}]$ for some subset of variables \mathbf{D}
 - e.g., indicator function
- **Log-linear model** over a Markov network H :
 - a set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables
 - a set of weights w_1, \dots, w_k
 - usually learned from data
 - $$P(X) = \frac{1}{Z} \exp \left[\sum_i w_i \phi_i(D_i) \right]$$

Generative v. Discriminative classifiers

– A review

- **Want to Learn:** $h:\mathbf{X} \mapsto Y$
 - \mathbf{X} – features
 - Y – target classes
- **Bayes optimal classifier** – $P(Y|\mathbf{X})$
- **Generative classifier**, e.g., Naïve Bayes:
 - Assume some **functional form for $P(\mathbf{X}|Y)$, $P(Y)$**
 - Estimate parameters of $P(\mathbf{X}|Y)$, $P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|\mathbf{X}=\mathbf{x})$
 - This is a **'generative' model**
 - **Indirect** computation of $P(Y|\mathbf{X})$ through Bayes rule
 - But, **can generate a sample of the data**, $P(\mathbf{X}) = \sum_y P(y) P(\mathbf{X}|y)$
- **Discriminative classifiers**, e.g., Logistic Regression:
 - Assume some **functional form for $P(Y|\mathbf{X})$**
 - Estimate parameters of $P(Y|\mathbf{X})$ directly from training data
 - This is the **'discriminative' model**
 - Directly learn $P(Y|\mathbf{X})$
 - But **cannot obtain a sample of the data**, because $P(\mathbf{X})$ is not available

Log-linear CRFs

(most common representation)

- **Graph H** : only over hidden vars Y_1, \dots, Y_p
 - **No assumptions about dependency on observed vars X**
 - You must always observe all of X
- **Feature** is some function $\phi[\mathbf{D}]$ for some subset of variables \mathbf{D}
 - e.g., indicator function,
- **Log-linear model** over a CRF H :
 - a set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables
 - a set of weights w_1, \dots, w_k
 - usually learned from data

$$- P(Y | X) = \frac{1}{Z(X)} \exp \left[\sum_i w_i \phi_i(D_i, X) \right]$$

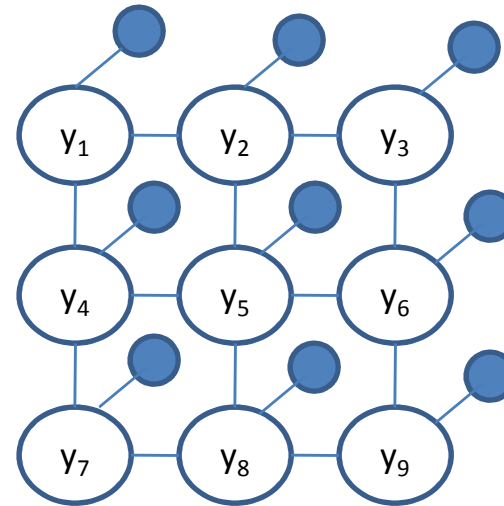
Example: Image Segmentation

- A set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables

$$P(Y | X) = \frac{1}{Z(X)} \exp \left[\sum_i w_i \phi_i(D_i, X) \right]$$

We will define features as follows:

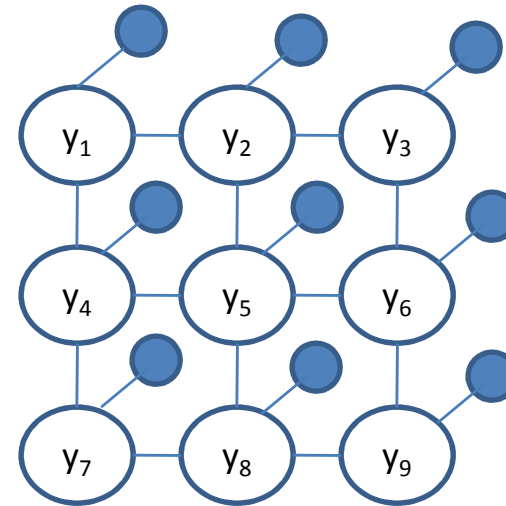
- $\phi(x, y)$: measures compatibility of node color and its segmentation
- A set of indicator features triggered for each edge labeling pair $\{ff, bb, fb, bf\}$
 - This is allowed since we can define many features over the same subset of variables



Example: Image Segmentation

- A set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables

$$P(Y | X) = \frac{1}{Z(X)} \exp \left[\sum_i w_i \phi_i(D_i, X) \right]$$



$$\phi(x, y) = \begin{cases} \log P(x | GMM_b) & y = b \\ \log P(x | GMM_f) & y = f \end{cases}$$

$$\phi_{ff}(y_i, y_j) = \begin{cases} 1 & y_i = f, y_j = f \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_{fb}(y_i, y_j) = \begin{cases} 1 & y_i = f, y_j = b \\ 0 & \text{otherwise} \end{cases}$$

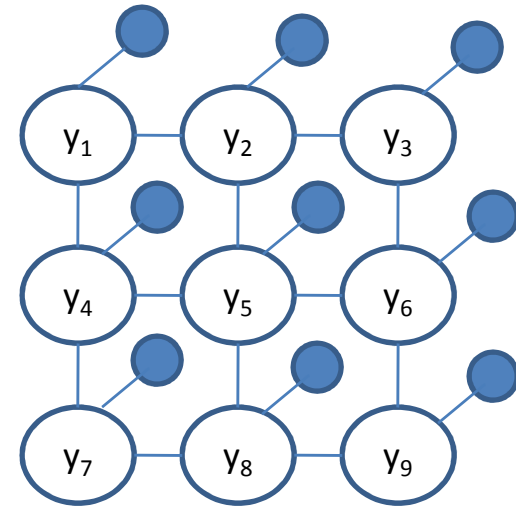
$$\phi_{bf}(y_i, y_j) = \begin{cases} 1 & y_i = b, y_j = f \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_{bb}(y_i, y_j) = \begin{cases} 1 & y_i = b, y_j = b \\ 0 & \text{otherwise} \end{cases}$$

Example: Image Segmentation

- A set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables

$$P(Y | X) = \frac{1}{Z(X)} \exp \left[\sum_i w_i \phi_i(D_i, X) \right]$$



-Now we just need to sum these features

$$P(Y | X) \propto \exp \left[\sum_{i \in V} \phi(x_i, y_i) + \sum_{ij \in E} \sum_{m \in \{ff, fb, bf, bb\}} w_m \phi_m(y_i, y_j) \right]$$

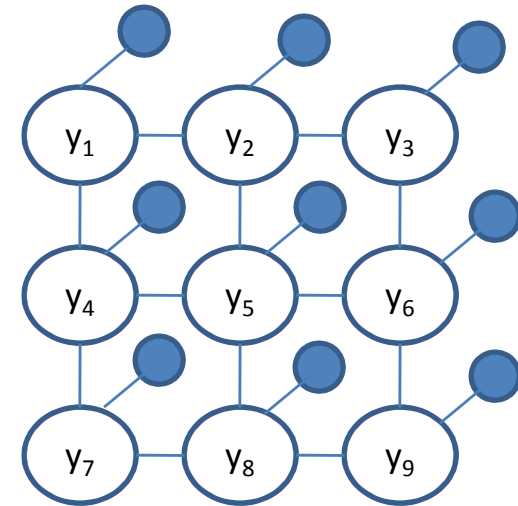
$$P(Y | X) \propto \exp \left[\sum_{i \in V} \phi(x_i, y_i) + \sum_{m \in \{ff, fb, bf, bb\}} w_m C_m \right] \quad C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i, y_j = m)$$

We need to learn parameters w_m

Example: Image Segmentation

$$P(Y | x) \propto \exp \left[\sum_{i \in V} \phi(x_i, y_i) + \sum_{m \in \{ff, fb, bf, bb\}} w_m C_m \right]$$

Given N data points (images and their segmentations)



$$\frac{\partial \ell(\text{Data} : w)}{\partial w_m} = \sum_{n=1}^N C_m[n] - E_w [C_m | X[n]]$$

Count for features m in data n

Requires inference using the current parameter estimates

$$C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i, y_j = m)$$

Example: Inference for Learning

$$P(Y | X) \propto \exp \left[\sum_{i \in V} f(x_i, y_i) + \sum_{m \in \{ff, fb, bf, bb\}} w_m C_m \right]$$

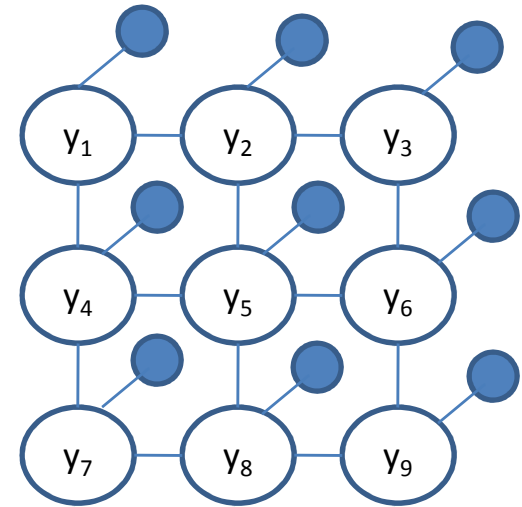
$$\frac{\partial \ell(\text{Data} : w)}{\partial w_m} = \sum_{n=1}^N C_m[n] - E_w [C_m | X[n]]$$

How to compute $E[C_{fb} | X[n]]$

$$E_w [C_{fb} | X[n]] = E \left[\sum_{ij} I(y_i = f, y_j = b) | X[n] \right]$$

$$E_w [C_{fb} | X[n]] = \sum_{ij} E_w [I(y_i = f, y_j = b) | X[n]]$$

$$E_w [C_{fb} | X[n]] = \sum_{ij} P_w (y_i = f, y_j = b | X[n])$$



$$C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i y_j = m)$$

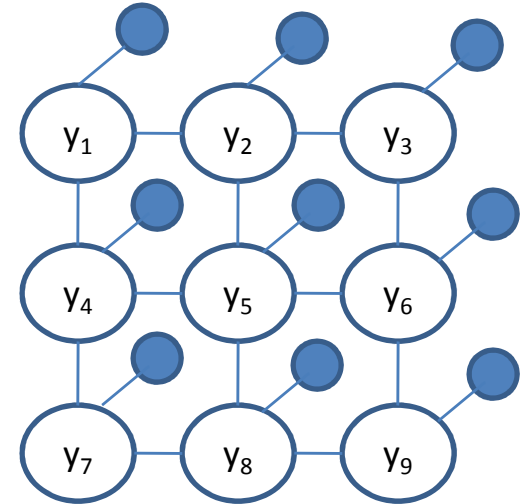
Example: Inference for Learning

$$P(Y | X) \propto \exp \left[\sum_{i \in V} f(x_i, y_i) + \sum_{m \in \{ff, fb, bf, bb\}} w_m C_m \right]$$

$$\frac{\partial \ell(\text{Data} : w)}{\partial w_m} = \sum_{n=1}^N C_m[n] - E_w [C_m | X[n]]$$

How to compute $E[C_{fb} | X[n]]$

$$E_w [C_{fb} | X[n]] = \sum_{ij} P_w (y_i = f, y_j = b | X[n])$$



$$C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i y_j = m)$$

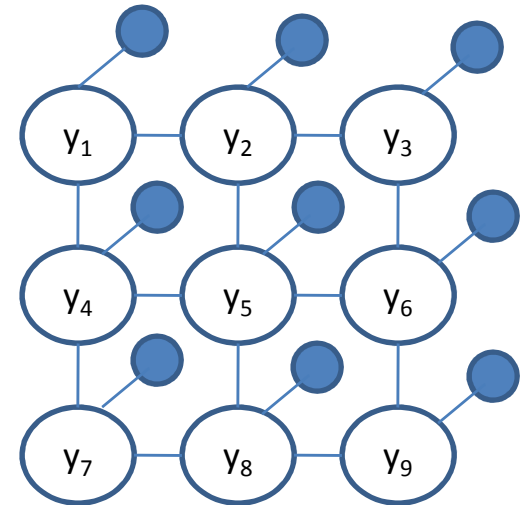
Representation Equivalence

Log linear representation

$$P(Y | X) \propto \exp \left[\sum_{i \in V} \phi(x_i, y_i) + \sum_{m \in \{ff, fb, bf, bb\}} w_m C_m \right]$$

$$C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i, y_j = m)$$

$$\phi(x, y) = \begin{cases} \log P(x | GMM_b) & y = b \\ \log P(x | GMM_f) & y = f \end{cases}$$



Tabular MN representation from HW4

$$P(Y | X) \propto \prod_{i \in V} \Phi(x_i, y_i) \prod_{ij \in E} \psi(y_i, y_j)$$

$$\begin{aligned} \prod_{i \in V} \Phi(x_i, y_i) &= \prod_{i \in V} P(x_i | GMM_{y_i}) \\ &= \prod_{i \in V} \exp[\log P(x_i | GMM_{y_i})] \\ &= \exp \left[\sum_{i \in V} \log P(x_i | GMM_{y_i}) \right] \\ &= \exp \left[\sum_{i \in V} \phi(x_i, y_i) \right] \end{aligned}$$

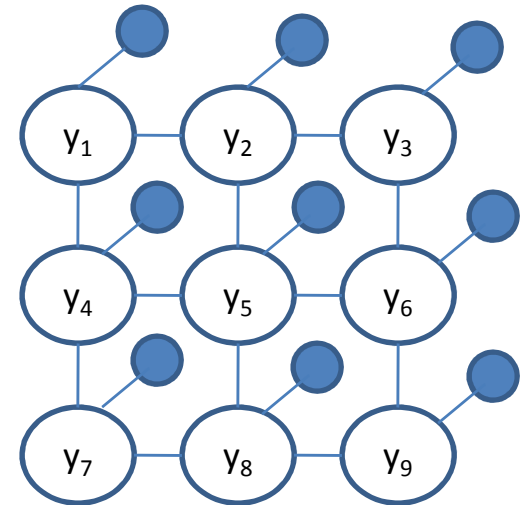
Representation Equivalence

Log linear representation

$$P(Y | X) \propto \exp \left[\sum_{i \in V} \phi(x_i, y_i) + \sum_{m \in \{ff, fb, bf, bb\}} w_m C_m \right]$$

$$C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i y_j = m)$$

$$\phi(x, y) = \begin{cases} \log P(x | GMM_b) & y = b \\ \log P(x | GMM_f) & y = f \end{cases}$$



Tabular MN representation from HW4

Now do it over the edge potential

$$P(Y | X) \propto \prod_{i \in V} \Phi(x_i, y_i) \prod_{ij \in E} \psi(y_i, y_j)$$

$$\psi(y_i, y_j) = \theta_{ff}^{I(y_i y_j = ff)} \theta_{fb}^{I(y_i y_j = fb)} \theta_{bf}^{I(y_i y_j = bf)} \theta_{bb}^{I(y_i y_j = bb)}$$

$$\psi(y_i, y_j) = \prod_{m \in \{ff, bf, fb, bb\}} \theta_m^{I(y_i y_j = m)}$$

This is correct as for every assignment to $y_i y_j$ we select one value from the table

Tabular MN representation from HW4

Now do it over the edge potential

$$P(Y | X) \propto \prod_{i \in V} \Phi(x_i, y_i) \prod_{ij \in E} \psi(y_i, y_j)$$

$$\psi(y_i, y_j) = \prod_{m=\{ff, bf, fb, bb\}} \theta_m^{I(y_i y_j = m)}$$

This is correct as for every assignment to $y_i y_j$ we select one value from the table

$$\psi(y_i, y_j) = \prod_{m=\{ff, bf, fb, bb\}} \theta_m^{I(y_i y_j = m)} = \prod_{m=\{ff, bf, fb, bb\}} \exp[\log \theta_m^{I(y_i y_j = m)}]$$

The cheap exp(log..) trick

$$\psi(y_i, y_j) = \exp \left[\sum_{m=\{ff, bf, fb, bb\}} I(y_i y_j = m) \log \theta_m \right]$$

Just algebra

Now lets combine it over all edge assuming parameter sharing

$$\prod_{ij \in E} \psi(y_i, y_j) = \prod_{ij \in E} \exp \left[\sum_{m=\{ff, bf, fb, bb\}} I(y_i y_j = m) \log \theta_m \right] = \exp \left[\sum_{ij \in E} \sum_{m=\{ff, bf, fb, bb\}} I(y_i y_j = m) \log \theta_m \right]$$

Now use the same C_m trick

$$\exp \left[\sum_{m=\{ff, bf, fb, bb\}} C_m \log \theta_m \right]$$

$$\leftarrow C_m = \sum_{ij \in E} I(y_i y_j = m)$$

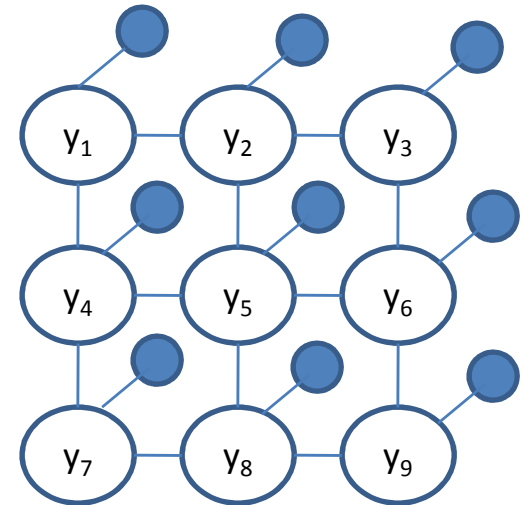
Representation Equivalence

Log linear representation

$$P(Y | X) \propto \exp \left[\underbrace{\sum_{i \in V} \phi(x_i, y_i)}_{\text{red bar}} + \underbrace{\sum_{m \in \{ff, fb, bf, bb\}} w_m C_m}_{\text{blue bar}} \right]$$

$$C_m = \sum_{ij \in E} \phi_m(y_i, y_j) = \sum_{ij \in E} I(y_i y_j = m)$$

$$\phi(x, y) = \begin{cases} \log P(x | GMM_b) & y = b \\ \log P(x | GMM_f) & y = f \end{cases}$$



Tabular MN representation from HW4

Now substitute

$$P(Y | X) \propto \prod_{i \in V} \Phi(x_i, y_i) \prod_{ij \in E} \psi(y_i, y_j)$$

$$P(Y | X) \propto \exp \left[\underbrace{\sum_{i \in V} \phi(x_i, y_i)}_{\text{red bar}} + \underbrace{\sum_{m \in \{ff, fb, bf, bb\}} C_m \log \theta_m}_{\text{blue bar}} \right]$$

Equivalent, with
 $w_m = \log \theta_m$ Where θ is
 the value in the tabular

Outline

- CRF
- Learning CRF for 2-d image segmentation
- IPF parameter sharing revisited

Iterative Proportional Fitting (IPF)

$$\hat{P}(\mathbf{u}) = \frac{\text{Count}(\mathbf{U} = \mathbf{u})}{m}$$

$$\frac{\partial \ell}{\partial \psi_i(\mathbf{c}_i)} = \frac{m \hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - \frac{m P_{\mathcal{F}}^{\psi}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} = 0$$

- Setting derivative to zero:

$$\frac{m \hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - \frac{m P_{\mathcal{F}}^{\psi}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} = 0$$

- Fixed point equation:

$$\psi_i(\mathbf{c}_i) = \psi_i(\mathbf{c}_i) \frac{\hat{P}(\mathbf{c}_i)}{P_{\mathcal{F}}^{\psi}(\mathbf{c}_i)}$$

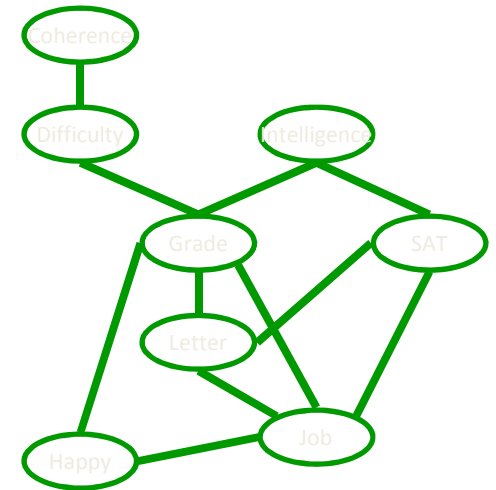
must initialize $\psi_i(\mathbf{c}_i) > 0$

- Iterate and converge to optimal parameters

- Each iteration, must compute:

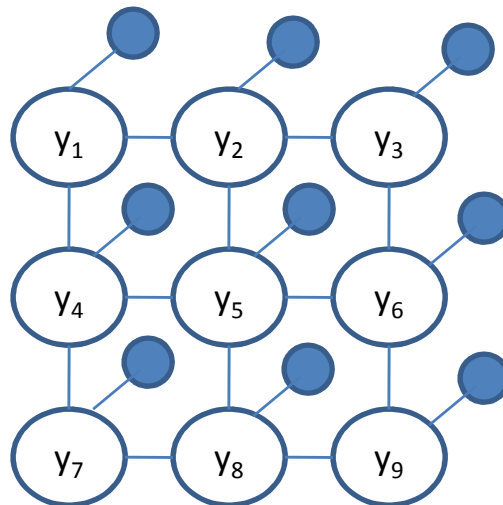
$$\psi_i^{(t+1)}(\mathbf{c}_i) \leftarrow \frac{\hat{P}(\mathbf{c}_i) \psi_i^{(t)}(\mathbf{c}_i)}{P_{\mathcal{F}}^{\psi^{(t)}}(\mathbf{c}_i)}$$

compute $P_{\mathcal{F}}^{\psi^{(t)}}(\mathbf{c}_i)$ at iteration t



Parameter Sharing in your HW

- Note that I am using Y for label
- All edge potentials are shared
- Also we are learning a conditional model



IPF parameter sharing

We only have one data point (image) in this example so we dropped $X[n]$ to only X

$$\psi_i^{t+1}(c_i) = \psi_i^t(c_i) \frac{\hat{p}(c_i | X)}{p(c_i | \psi^t, X)}$$

In total we have 4 parameters as opposed to 4 parameters per edge

How to calculate these quantities using parameter sharing?

$$\hat{p}(c = fb | X) = \frac{\sum_{ij \in E} I(y_i = f, y_j = b | X)}{|E|}$$

We can cancel $|E|$ due to division

$$P(c = fb | \psi^t, X) = \frac{\sum_{ij \in E} P(y_i = f, y_j = b | \psi^t, X)}{|E|}$$

Run lbp, when converged

$$P(y_i, y_j | \psi^t, X) \propto \Phi(x_i, y_i) \Phi(x_j, y_j) \psi(y_i, y_j) \prod_{k \in N(i)-j} \delta_{k \rightarrow i} \prod_{k' \in N(j)-i} \delta_{k' \rightarrow j}$$

