

# Notes on Probabilistic Graphical Models

Shoichiro Yamanishi

May 1, 2020

## 1 Overview

This document briefly goes through some important points regarding the probabilistic graphical models. A graphical model expresses the factoring of a joint distribution over the random variables. Each node represents a random variable, except for a factor graph in which node is either a factor or a random variable. A model can be programmatically drawn using Python package Daft which works with matplotlib. See [2].

- Directed Acyclic Graph (Bayesian Network) : This expresses conditional dependencies among the random variables by directed edges.
- Undirected Graph (Markov Random Field) : This expresses factoring by maximal cliques.
- Factor Graph : This expresses the exact factoring. Any path on the graph visits factors and random variables alternately.

Chap. 8 of PRML [1] is a good source of information.

## 2 Directed Acyclic Graph (Bayesian Network)

### 2.1 Chain

Two interesting cases: One is for a chain of 1-of-K discrete variables and the other is for gaussian distribution. The random variables on a chain are usually latent variables.

#### 2.1.1 Discrete Variables Chain

Please see Figure 1 for a chain of Dirichlet distribution for latent variables  $\mathbf{z}_i \in \mathcal{R}^K$ .

$$\begin{aligned} p(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3, \dots, \mathbf{z}_{N-1}, \mathbf{z}_N | \boldsymbol{\alpha}) &= p(\mathbf{z}_N | \mathbf{z}_{N-1}) \cdots p(\mathbf{z}_3 | \mathbf{z}_2) p(\mathbf{z}_2 | \mathbf{z}_1) p(\mathbf{z}_1 | \boldsymbol{\alpha}) \\ &= \text{Dir}(\mathbf{z}_N | \mathbf{z}_{N-1}) \cdots \text{Dir}(\mathbf{z}_3 | \mathbf{z}_2) \text{Dir}(\mathbf{z}_2 | \mathbf{z}_1) \text{Dir}(\mathbf{z}_1 | \boldsymbol{\alpha}) \end{aligned} \tag{1}$$

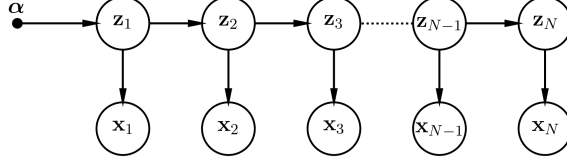


Figure 1: Chain of Dirichlet Distribution

where  $\text{Dir}(\mathbf{z}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_{j=1}^K \alpha_j)}{\prod_{j=1}^K \Gamma(\alpha_j)} \prod_{j=1}^K z_j^{\alpha_j-1}$ . The corresponding random variables can be one-of-K variable  $\mathbf{x}_i \in \{0, 1\}^K$  such that  $\mathbf{x}_i \sim \text{Mult}(\mathbf{x}_i|\mathbf{z}_i) = \prod_{j=1}^K z_{i_j}^{x_{i_j}}$ .

### 2.1.2 Linear Gaussian Model

Please see Figure 2 for a chain of Gaussian distribution for latent variables  $\mathbf{z}_i$ . This is a conceptual model and not really useful. A useful example that includes inference to the covariance matrices is Kalman filter described in a separate document.

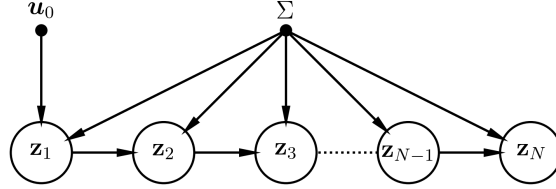


Figure 2: Chain of Gaussian Distribution

$$\begin{aligned}
 p(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{N-1}, \mathbf{z}_N | \boldsymbol{\alpha}) &= p(\mathbf{z}_N | \mathbf{z}_{N-1}) \cdots p(\mathbf{z}_2 | \mathbf{z}_1) p(\mathbf{z}_1) \\
 &= \mathcal{N}(\mathbf{z}_N; \mathbf{u}_{N-1}(\mathbf{z}_{N-1}), \Sigma) \cdots \mathcal{N}(\mathbf{z}_2; \mathbf{u}_1(\mathbf{z}_1), \Sigma) \mathcal{N}(\mathbf{z}_1; \mathbf{u}_0, \Sigma)
 \end{aligned}
 \tag{2}$$

## 2.2 Parameter Reduction

Consider a probability distribution  $p(y|x_1, x_2, \dots, x_N)$  where  $\mathbf{x}_i \in \{0, 1\}^K$ ,  $\sum_{j=1}^K x_{i_j} = 1$ , i.e.,  $\mathbf{x}_i$  is a 1-of-K variable. The possible combinations will add up to  $K^N$ , which could be unmanageable. To reduce the complexity we can use a linear combination of  $\mathbf{x}_i$  as follows.  $p(y|x_1, x_2, \dots, x_N, \mathbf{w}_0, W) = p(y|\mathbf{w}_0 + \sum_{i=1}^N W \mathbf{x}_i)$  where  $\mathbf{w}_0 \in \mathcal{R}^K$ ,  $W$  is a  $K \times N$  matrix. See figure 3.

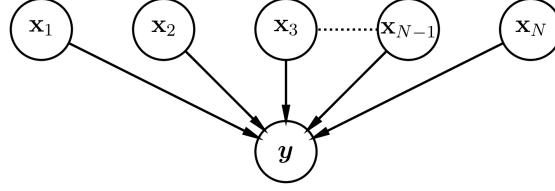


Figure 3: Parameter Reduction

### 2.3 I.I.D. and Naive Bayes

Consider the graphical models depicted in figure 4. This can manifest as *independent and identically distributed* assumption for  $N$  samples, or *Naive Bayes* model of a random variable  $\mathbf{x} \in \mathcal{R}^N$  with  $N$  features that depends on the class  $C$ , where each feature is assumed to be independent of each other conditioned on  $C$ .

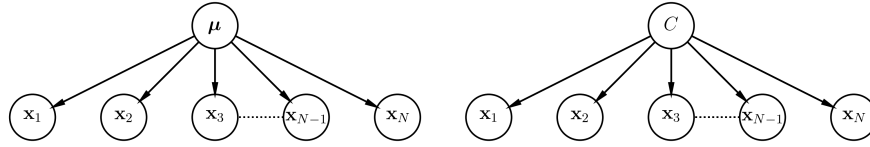


Figure 4: I.I.D. and Naive Bayes

### 2.4 Conditional Independence, D-Separaton, and Markov Blanket

#### 2.4.1 Conditional Independence

BTW, the math symbols  $\perp\!\!\!\perp$  and  $\not\perp\!\!\!\perp$  for Text are `\upmodels` and `\nupmodels` in the package *MnSymbol*.

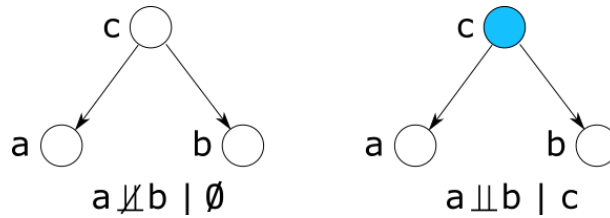


Figure 5: Conditional Independence Case 1

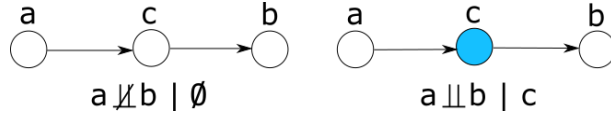


Figure 6: Conditional Independence Case 2

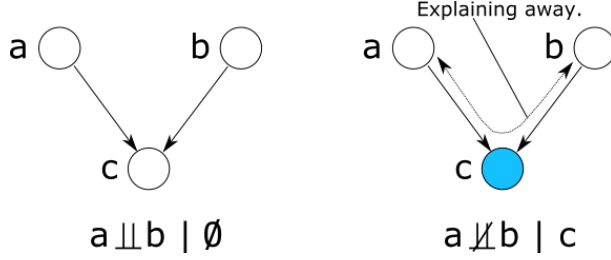


Figure 7: Conditional Independence Case 3

#### 2.4.2 D-Separation

Let  $V_s$  be a set of random variables which you want to make conditionally independent from the rest. Given a set of observed random variables, how to determine if  $V_s$  are conditionally independent from the rest? D-Separation gives an answer. Let  $G_s$  be a subgraph induced by  $V_s$ . Also, Let  $N^+(v)$  and  $N^-(v)$  denote the in-neighbor and out-neighbor of  $v$  respectively. Then Let  $N^+(G_s) = (\bigcup_{v \in G_s} N^+(v)) \setminus V(G_s)$  and  $N^-(G_s) = (\bigcup_{v \in G_s} N^-(v)) \setminus V(G_s)$  where  $V(G_s)$  denotes the vertices in the graph  $G_s$ . In other words,  $N^+(G_s)$  is the immediate in-neighbor of  $G_s$  and  $N^-(G_s)$  is the immediate out-neighbor of  $G_s$ . Let  $D(G_s)$  be the subgraph induced by the nodes reachable from  $G_s$ . In other words,  $D(G_s)$  is the descendent subgraph of  $G_s$ . Then in order to make  $N(G_s)$  conditionally independent, the following must hold.

- $N^+(G_s)$  must be observed.
- if any  $v \in D(G_s)$  is observed, then  $N^+(v) \setminus V(G_s)$  must be observed.

This is illustrated in figure 8 If the set of observed random variables meet the criteria above, then  $V_s$  are conditionally independent from the rest.

#### 2.4.3 Markov Blanket

Again, let  $V_s$  be a set of random variables which you want to make conditionally independent from the rest. What other random variables must be observed to guarantee the conditional independence of  $V_s$ ? Markov Blanket gives an answer. Please see figure 9. It's basically  $N^+(G_s)$ ,  $N^-(G_s)$ , and its outneighbors.

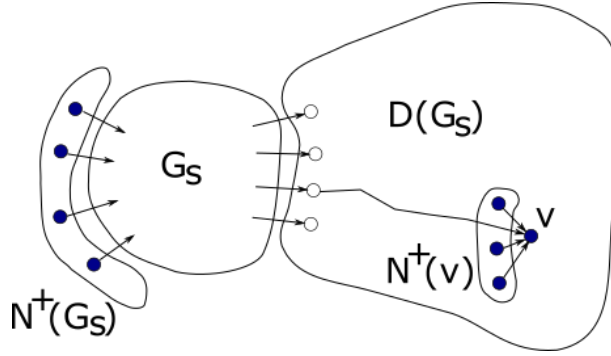


Figure 8: D-Separation

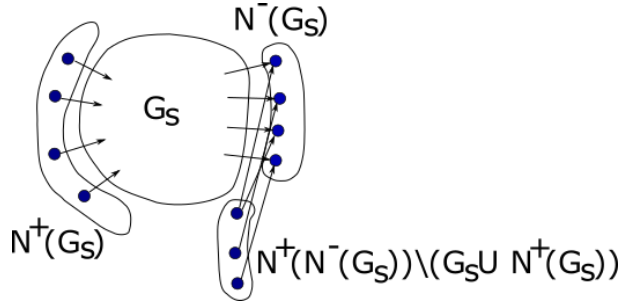


Figure 9: D-Separation

### 3 Undirected Graph (Markov Random Field) and Belief Propagation

We consider an exact inference called belief propagation in a tree. There are two types to consider. One is to find a marginal distribution of a particular random variable, and the other is find the values for all the random variables that together maximizes the joint probability as in finding a mode in a MAP estimation.

We first consider the simple case of a chain and expand the discussion to a tree.

#### 3.1 Inference in a Chain : Simple Case

Consider the joint probability of discrete random variables  $\mathbf{x}_i \in \{0, 1\}^K$ ,  $\sum_{j=1}^K x_{i_j} = 1$ , which is factorized into a chain as in figure 10.

The joint probability is factorized into the following.

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}, \mathbf{x}_N) = \frac{1}{Z} \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \dots \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \quad (3)$$

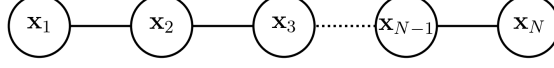


Figure 10: Undirected Chain

### 3.1.1 Marginal distribution $p(\mathbf{x}_N)$

Consider first the marginal distribution  $p(\mathbf{x}_N)$ .

$$\begin{aligned}
 p(\mathbf{x}_N) &= \sum_{\mathbf{x}_1} \sum_{\mathbf{x}_2} \cdots \sum_{\mathbf{x}_{N-1}} \frac{1}{Z} \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \cdots \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \\
 &= \frac{1}{Z} \left[ \sum_{\mathbf{x}_{N-1}} \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \left[ \sum_{\mathbf{x}_{N-2}} \phi_{N-2,N-1}(\mathbf{x}_{N-2}, \mathbf{x}_{N-1}) \left[ \cdots \sum_{\mathbf{x}_2} \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \left[ \sum_{\mathbf{x}_1} \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \right] \cdots \right] \right] \right] \right] \quad (4)
 \end{aligned}$$

On the first line, there are  $K^N$  terms. On the second, I abused the notation such that each square brackets pair indicates a function of one variable. We define the following recursive definition.

$$\phi'_i(\mathbf{x}_i) = \sum_{\mathbf{x}_{i-1}} \phi_{i-1,i}(\mathbf{x}_{i-1}, \mathbf{x}_i) \phi'_{i-1}(\mathbf{x}_{i-1}) \quad (5)$$

Evaluation of each of  $K$  elements of  $\phi'_i(\mathbf{x}_i)$  takes  $K$  evaluations of  $\phi(\mathbf{x}_{i-1}, \mathbf{x}_i)$ ,  $K$  multiplications, and  $K - 1$  additions. Hence the construction of one  $\phi'_i(\mathbf{x}_i)$  is  $O(K^2)$ .

Then  $p(\mathbf{x}_N)$  is constructed sequentially by evaluating from  $\phi'_1(\mathbf{x}_2)$  up to  $\phi'_1(\mathbf{x}_N)$  as follows.

$$\begin{aligned}
 p(\mathbf{x}_N) &= \frac{1}{Z} \left[ \sum_{\mathbf{x}_{N-1}} \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \left[ \sum_{\mathbf{x}_{N-2}} \phi_{N-2,N-1}(\mathbf{x}_{N-2}, \mathbf{x}_{N-1}) \left[ \cdots \sum_{\mathbf{x}_2} \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \left[ \sum_{\mathbf{x}_1} \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \right] \cdots \right] \right] \right] \right] \\
 &= \frac{1}{Z} \left[ \sum_{\mathbf{x}_{N-1}} \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \left[ \sum_{\mathbf{x}_{N-2}} \phi_{N-2,N-1}(\mathbf{x}_{N-2}, \mathbf{x}_{N-1}) \left[ \cdots \sum_{\mathbf{x}_2} \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \phi'_2(\mathbf{x}_2) \right] \right] \right] \\
 &= \frac{1}{Z} \left[ \sum_{\mathbf{x}_{N-1}} \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \left[ \sum_{\mathbf{x}_{N-2}} \phi_{N-2,N-1}(\mathbf{x}_{N-2}, \mathbf{x}_{N-1}) [\cdots \phi'_3(\mathbf{x}_3)] \right] \right] \\
 &= \frac{1}{Z} \left[ \sum_{\mathbf{x}_{N-1}} \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \left[ \sum_{\mathbf{x}_{N-2}} \phi_{N-2,N-1}(\mathbf{x}_{N-2}, \mathbf{x}_{N-1}) \phi'_{N-2}(\mathbf{x}_{N-2}) \right] \right] \\
 &= \frac{1}{Z} \left[ \sum_{\mathbf{x}_{N-1}} \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \phi'_{N-1}(\mathbf{x}_{N-1}) \right] \\
 &= \frac{1}{Z} \phi'_N(\mathbf{x}_N) \quad (6)
 \end{aligned}$$

This derivation can be viewed as a message passing from  $\mathbf{x}_2$  up to  $\mathbf{x}_N$  with the message defined by  $\phi'_i(\mathbf{x}_i)$ . This is the simplest form of belief propagation.

Please note that  $Z = \sum_{\mathbf{x}_N} \phi'_N(\mathbf{x}_N)$ . Overall the complexity of finding  $p(\mathbf{x}_N)$  is  $O(NK^2)$ .

### 3.1.2 Finding a Mode : $\text{argmax}\{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}, \mathbf{x}_N)\}$ for MAP etc

Now we want to find the values for each random variables that maximizes the joint probability. This occurs in a MAP estimate, where you want to get the highest mode of the probability distribution. Here we define an operation called *amax* to be a combined operation of max and argmax. I.e., it stores both the maximum value as well as the element that holds the maximum value.

$$\begin{aligned} (p^*, \mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_N^*) &= \text{amax}_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N} \{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}, \mathbf{x}_N)\} \\ &= \text{amax}_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N} \left\{ \frac{1}{Z} \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \dots \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \right\} \end{aligned} \quad (7)$$

We can exploit the chain structure to distribute the *amax* operation as follows.

$$\begin{aligned} &\text{amax}_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N} \left\{ \frac{1}{Z} \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \dots \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \right\} \\ &= \text{amax}_{\mathbf{x}_1} \left\{ \text{amax}_{\mathbf{x}_2} \left\{ \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \text{amax}_{\mathbf{x}_3} \left\{ \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \dots \text{amax}_{\mathbf{x}_N} \{ \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \} \dots \right\} \right\} \right\} \end{aligned} \quad (8)$$

where I abused the notation of  $\text{amax}\{\}$  be a function of  $\mathbf{x}_i$ .

We define the following recursive definition.

$$\phi'_{i-1}(\mathbf{x}_{i-1}) = \text{amax}_{\mathbf{x}_i} \{ \phi_{i-1,i}(\mathbf{x}_{i-1}, \mathbf{x}_i) \phi'_{i-1}(\mathbf{x}_{i-1}) \} \quad (9)$$

Evaluation of each of  $K$  elements of  $\phi'_{i-1}(\mathbf{x}_{i-1})$  takes takes  $K$  evaluations of  $\phi(\mathbf{x}_{i-1}, \mathbf{x}_i)$  and  $K$  comparisons. Hence the construction of one  $\phi'_{i-1}(\mathbf{x}_{i-1})$  is  $O(K^2)$ . Then,

$$\begin{aligned} &\text{amax}_{\mathbf{x}_1} \left\{ \text{amax}_{\mathbf{x}_2} \left\{ \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \text{amax}_{\mathbf{x}_3} \left\{ \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \dots \text{amax}_{\mathbf{x}_N} \{ \phi_{N-1,N}(\mathbf{x}_{N-1}, \mathbf{x}_N) \} \dots \right\} \right\} \right\} \\ &= \text{amax}_{\mathbf{x}_1} \left\{ \text{amax}_{\mathbf{x}_2} \left\{ \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \text{amax}_{\mathbf{x}_3} \{ \phi_{2,3}(\mathbf{x}_2, \mathbf{x}_3) \dots \phi'_{N-1}(\mathbf{x}_{N-1}) \dots \} \right\} \right\} \\ &= \text{amax}_{\mathbf{x}_1} \left\{ \text{amax}_{\mathbf{x}_2} \{ \phi_{1,2}(\mathbf{x}_1, \mathbf{x}_2) \phi'_2(\mathbf{x}_2) \} \right\} \\ &= \text{amax}_{\mathbf{x}_1} \{ \phi'_1(\mathbf{x}_1) \} \\ &= (p^*, \mathbf{x}_1^*) \end{aligned} \quad (10)$$

This derivation can be viewed as a message passing from  $\mathbf{x}_{N-1}$  up to  $\mathbf{x}_1$  with the message defined by  $\phi'_i(\mathbf{x}_i)$ . To obtain  $(\mathbf{x}_2^*, \dots, \mathbf{x}_N^*)$ , we can back track the operation from  $\mathbf{x}_1$  to  $\mathbf{x}_{N-1}$ . Overall the complexity of finding  $p(\mathbf{x}_N)$  is  $O(NK^2)$ .

### 3.2 Inference in a Tree

We extend the idea described in the chain to trees. There are two types of extensions in the factor graph.

- One factor can take more than two random variables.
- One variable can be associated to more than two factors.

In order to accommodate those extensions, we define the recursive relation in two phases. Please see figure 11.

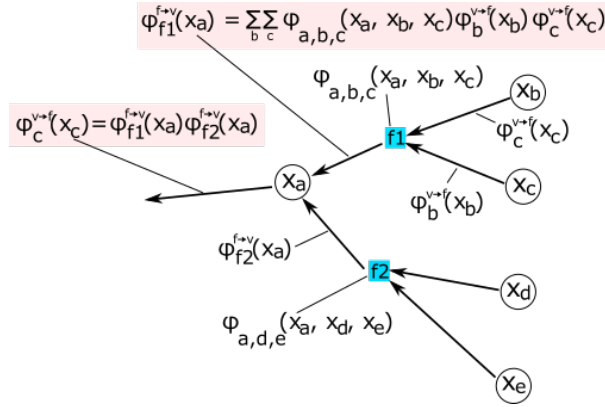


Figure 11: Propagation in a Tree

We have to nominate one random variable node  $\mathbf{x}_{root}$  in the factor graph as the root, and form a rooted oriented tree by orienting the edges from the root toward leaves. Such a tree is unique up to the edge orientation. In case of marginalization,  $\mathbf{x}_{root}$  is the random variable for which, the rest of the random variables are marginalized to obtain  $p(\mathbf{x}_{root})$ .

$\phi_{\mathbf{f}}^{f \rightarrow v}(\mathbf{x})$  is from a factor to a variable, and  $\phi_{\mathbf{x}}^{v \rightarrow f}(\mathbf{x})$  is from a variable to a factor.

$\phi_{\mathbf{f}_a}^{f \rightarrow v}(\mathbf{x}_i)$  represent the subtree under  $\mathbf{f}_a$  as a function of  $\mathbf{x}_i$ .  $\phi_{\mathbf{x}_i}^{v \rightarrow f}(\mathbf{x}_i)$  represent the subtree under  $\mathbf{x}_i$  as a function of  $\mathbf{x}_i$ , i.e.,

$$\phi_{\mathbf{x}_i}^{v \rightarrow f}(\mathbf{x}_i) = \prod_{\mathbf{f}_j \in N^-(\mathbf{x}_i)} \phi_{\mathbf{f}_j}^{f \rightarrow v}(\mathbf{x}_i)$$

Now,  $\phi_{\mathbf{f}}^{f \rightarrow v}(\mathbf{x}_a)$  is defined as a marginalization or maximization of  $\phi_{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}}(\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL})$  where  $\mathbf{x}_a \in \{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}\}$ . Then

$$\phi_{\mathbf{x}_a}^{f \rightarrow v}(\mathbf{x}_a) = \sum_{\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}\} \setminus \mathbf{x}_a} \left( \phi_{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}}(\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}) \prod_{\mathbf{x}_b}^{\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}\} \setminus \mathbf{x}_a} (\phi_{\mathbf{x}_b}^{v \rightarrow f}(\mathbf{x}_b)) \right) \quad (11)$$



or,

$$\phi_{\mathbf{x}_a}^{f \rightarrow v}(\mathbf{x}_a) = \operatorname{amax}_{\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}\} \setminus \mathbf{x}_a} \left\{ \phi_{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}}(\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}) \prod_{\mathbf{x}_b}^{\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iL}\} \setminus \mathbf{x}_a} (\phi_{\mathbf{x}_b}^{v \rightarrow f}(\mathbf{x}_b)) \right\} \quad (12)$$

### 3.3 Inference in a Grid

We can't use a belief propagation if the graph contains a cycle. If we force a belief propagation to such a graph and assume the propagation to converge, then it is called *loopy belief propagation*.

Assume the latent variables are discrete, i.e.  $\mathbf{x}_i = \{0, 1\}^K$ ,  $\sum_{j=1}^K x_{ij} = 1$  and the joint probability is represented as follows.

$$\begin{aligned} p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) &= \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{z}_i) p(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) \\ &= \prod_{i=1}^N p(\mathbf{x}_i | \mathbf{z}_i) \prod_{\{u,v\} \in \mathcal{C}} p(\mathbf{z}_u, \mathbf{z}_v) \\ &= \prod_{i=1}^N \operatorname{Mult}(\mathbf{x}_i | \mathbf{z}_i) \prod_{\{u,v\} \in \mathcal{C}} \frac{1}{Z_{\{u,v\}}} \exp(-\phi(\mathbf{z}_u, \mathbf{z}_v)) \\ &= \prod_{i=1}^N \prod_{j=1}^K (z_{ij}^{x_{ij}}) \frac{1}{Z_w} \exp\left(-\sum_{\{u,v\} \in \mathcal{C}} \phi(\mathbf{z}_u, \mathbf{z}_v)\right) \\ &= \frac{1}{Z} \exp\left(\sum_{i=1}^N \sum_{j=1}^K x_{ij} \ln(z_{ij}) - \sum_{\{u,v\} \in \mathcal{C}} \phi(\mathbf{z}_u, \mathbf{z}_v)\right) \\ &= \frac{1}{Z} \exp\left(\sum_{i=1}^N U_i(\mathbf{z}_i) - \sum_{\{u,v\} \in \mathcal{C}} P_{u,v}(\mathbf{z}_u, \mathbf{z}_v)\right) \end{aligned} \quad (13)$$

where  $\mathcal{C}$  is the set of edges in the grid. Please see figure 12

After observing  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ , the posterior will be:

$$\begin{aligned} p(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) &= \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N)}{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)} \\ &\propto p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) \\ &= \frac{1}{Z} \exp\left(\sum_{i=1}^N U_i(\mathbf{z}_i) - \sum_{\{u,v\} \in \mathcal{C}} P_{u,v}(\mathbf{z}_u, \mathbf{z}_v)\right) \end{aligned} \quad (14)$$

Then the MAP estimate will be:

$$\operatorname{argmax}_{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N} \{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N)\} = \operatorname{argmax}_{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N} \left\{ \sum_{i=1}^N U_i(\mathbf{z}_i) - \sum_{\{u,v\} \in \mathcal{C}} P_{u,v}(\mathbf{z}_u, \mathbf{z}_v) \right\} \quad (15)$$

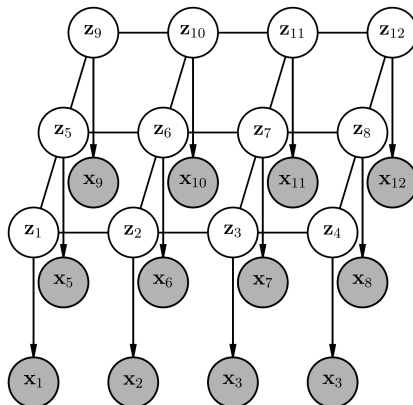


Figure 12: Generative Grid

If you want a MAP estimate, then we can use a max-flow min-cut paradigm, but it will be impractical if  $K$  is large, as you have to form a graph whose number of nodes is  $K + 1$  times the number of latent variables, and between two latent variables there are  $2K + K^2$  edges as a complete bipartite graph. Please see section 12.2 and 12.3 of [3].

The same technique can be applied for a model in conditional random field as in figure 13. The probability distribution is formulated as follows.

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N) = \frac{1}{Z} \exp \left( - \sum_{\{u,v\} \in \mathcal{C}} \phi(\mathbf{z}_u, \mathbf{z}_v) - \sum_{i=1}^N \eta(\mathbf{x}_i, \mathbf{z}_i) \right) \quad (16)$$

## References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007.
- [2] David S. Fulford, Dan Foreman-Mackey, and David W. Hogg. Daft : Beautifully rendered probabilistic graphical models. <https://docs.daft-pgm.org/en/latest>. Accessed: 2020-03-30.
- [3] Simon J. D. Prince. *Computer Vision: Models, Learning, and Inference*. Cambridge University Press, 2012.

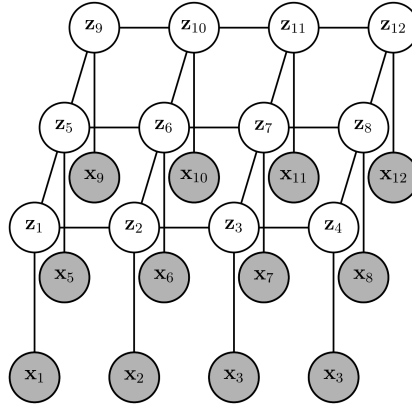


Figure 13: Conditional Random Field Grid