

1 Probability Theory

1.1 Independence

$p(X, Y) = p(X)p(Y) \Leftrightarrow p(X|Y) = p(X) \Leftrightarrow p(Y|X) = p(Y)$

1.2 Conditional Independence

$X \perp\!\!\!\perp Y \mid Z \iff p(X, Y|Z) = p(X|Z)p(Y|Z)$   
 $p(X, Y|Z) = p(X|Y, Z)p(Y|Z) = p(X|Z)p(Y|Z)$

1.3 Sum and Product Rules

$p(X, Y) = p(X)p(Y|X), \quad p(X, Y, Z) = p(X)p(Y|X)p(Z|X, Y)$   
 $p(X) = \sum_Y p(X, Y)$

1.4 Bayes' Theorem

$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}, \quad p(Y|X, Z) = \frac{p(X|Y, Z)p(Y|Z)}{p(X|Z)}$

1.5 Chain Rule

$p(x_a, x_b, x_c) = p(x_c|x_a, x_b)p(x_b|x_a)p(x_a) \quad p(\mathbf{x}) = \prod_{k=1}^K p(x_k|\text{pa}_k)$

1.6 Misc

$p(a, b|c) = \frac{p(a, b, c)}{p(c)} \quad p(a, b, c) = p(a|c)p(b|c)p(c)$

2 Distributions

Binary	Bernoulli	Binomial	Beta
Discrete	Categorical	Multinomial	Dirichlet

2.1 Bernoulli Distribution

$\text{Ber}(x|n) = \mu^x(1 - \mu)^{1-x}, \quad \mathbb{E}[x] = \mu, \quad \text{Var}[x] = \mu - \mu^2,$   
 $P(D, \mu) = \prod_{n=1}^N \mu^{x_n}(1 - \mu)^{1-x_n}, \quad \mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$

2.2 Binomial Distribution

$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m(1 - \mu)^{N-m}, \quad \frac{n!}{k!(n-k)!} = \binom{n}{k},$   
 $\mathbb{E}[m] = N\mu, \quad \text{Var}[m] = N\mu(1 - \mu), \quad \mu_{\text{ML}} = \frac{m}{N}$

2.3 Categorical Distribution

$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_k \mu_k^{x_k}, \quad \boldsymbol{\mu} \in \{0, 1\}^K, \quad \sum_k \mu_k = 1, \quad \boldsymbol{\mu}_{\text{ML}} = \frac{\mathbf{m}}{N},$   
 $m_k = \sum_n x_{nk}, \quad \text{Mult}(m_1, \dots, m_k|N, \boldsymbol{\mu}) = \frac{N!}{m_1! \dots m_k!} \prod_k \mu_k^{m_k}$

2.4 Beta Distribution

$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1}(1 - \mu)^{b-1}, \quad \mathbb{E}[\mu] = \frac{a}{a+b},$   
 $\text{Var}[x] = \frac{ab}{(a+b)^2(a+b+1)}, \quad p(\mu|m, l, a, b) \propto \mu^{m+a-1}(1 - \mu)^{l+b-1}$

2.5 Gamma Distribution

$\text{Gamma}(\tau|a, b) = \frac{b^a}{\Gamma(a)} \tau^{a-1} e^{-b\tau}, \quad \mathbb{E}[\tau] = \frac{a}{b}, \quad \text{Var}[\tau] = \frac{a}{b^2},$   
 $\text{mode}[\tau] = \frac{a-1}{b} \text{ for } a \geq 1, \quad \mathbb{E}[\ln \tau] = \psi(a) - \ln b,$   
 $H(\tau) = \ln \Gamma(a) - (a-1)\psi(a) - \ln b + a$

2.6 Multinomial Distribution

$\mathbf{x} = [0, 0, 0, 0, 1, 0, 0]^\top, \quad \sum_{k=1}^K x_k = 1, \quad p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k},$   
 $\sum_{k=1}^K \mu_k = 1, \quad \mu_k^{\text{ML}} = \frac{m_k}{N}, \quad m_k = \sum_{k=1}^K x_{nk}$

2.7 Dirichlet Distribution

$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} \prod_{k=1}^K \mu_k^{\alpha_k-1}, \quad \alpha_0 = \sum_{k=1}^K \alpha_k$

2.8 Gaussian Distribution

$\mathcal{N}(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x - \mu)^2),$   
 $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$

2.8.1 ML for the Gaussian

$\ln p(X|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}),$   
 $\boldsymbol{\mu}_{\text{ML}} = 1/N \sum_{n=1}^N \mathbf{x}_n, \quad \boldsymbol{\Sigma}_{\text{ML}} = 1/N \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top (\mathbf{x}_n - \boldsymbol{\mu})$

2.8.2 Stochastic gradient descent Gaussian

$\max P(x_1, \dots, x_N|\theta), \quad \theta^N = \theta^{N-1} + \alpha_{N-1} \frac{\partial}{\partial \theta^{N-1}} \ln p(x_N|\theta^{N-1})$

2.8.3 Marginal and Conditional Gaussians

Given  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$  and  $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$ . We get  $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^\top)$  and  $p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}[\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \mathbf{A}\boldsymbol{\mu}], \boldsymbol{\Sigma})$  where  $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$ .

2.9 Student's T distribution

The heavy tail of the student-t distribution makes it more robust against outliers.

$St(x|\mu, \lambda, \nu) = \frac{\Gamma(\nu/2+1/2)}{\Gamma(\nu/2)} (\frac{\lambda^{1/2}}{(\pi\nu)^{D/2}}) [1 + \frac{\lambda(x-\mu)^2}{\nu}]^{-\nu/2-D/2},$

$f_x(x) = \frac{\Gamma((\nu+p)/2)}{\Gamma(\nu/2)\nu^{p/2}\pi^{p/2}|\boldsymbol{\Sigma}|^{1/2}[1+1/\nu(x-\boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(x-\boldsymbol{\mu})]^{(\nu+p)/2}}$   
 $\mathbb{E}(\mathbf{x}) = \frac{\Gamma(D/2+v/2)}{\Gamma(v/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(\pi v)^{D/2}} \int [1 + \frac{(x-\boldsymbol{\mu})^T \boldsymbol{\Lambda} (x-\boldsymbol{\mu})}{v}]^{-D/2-v/2} x d\mathbf{x}$

3 Independent Component Analysis

We have a set of N observations:  
 $D = \{\mathbf{x}^{(n)} \in \mathbb{R}^J\}_{n=1}^N$ , with  $I$  independent sources. Thus:  
 $\mathbf{x} = \mathbf{G}\mathbf{s}$ , where  $G$  is not known. We assume that the latent variables are independently distributed, with marginal distributions  $P(s_i|H) = p_i(s_i)$ . The probability of the observables and the hidden nodes is:  
 $P(\{x^{(n)}, s^{(n)}\}_{n=1}^N | G, H) = \prod_{n=1}^N [P(x^{(n)}|s^{(n)}, G, H)P(s^{(n)}|H)]$   
 $= \prod_{n=1}^N [(\prod_j \delta(x_j^{(n)} - \sum_i G_{ji}s_i^{(n)}))(\prod_i p_i(s_i^{(n)}))]$ .  
The factor  $\mathbf{x}$  is generated without noise. For learning the  $G$  from the data the likelihood is:  $P(D|G, H) = \prod_n P(x^{(n)}|G, H)$ , which is a product of factors if we marginalize over the latent variables. We can express each coefficient of  $\mathbf{G}$ , as a summation over all coefficients multiplied by a delta function such that,  $G_{ji}s_i^{(n)} = \sum_i G_{ji}s_i^{(n)}$ . ( $\mathcal{H}$  denotes the model.)

$p(\mathbf{x}^{(n)}|\mathbf{G}, \mathcal{H}) = \int p(\mathbf{x}^{(n)}|\mathbf{s}^{(n)}, \mathbf{G}, \mathcal{H})p(\mathbf{s}^{(n)}|\mathcal{H})d^I \mathbf{s}^{(n)}$   
 $= \int \prod_j \delta(x_j^{(n)} - G_{ji}s_i^{(n)}) \prod_i p_i(s_i^{(n)})d^I \mathbf{s}^{(n)}$   
 $= \frac{1}{|\det \mathbf{G}|} \prod_i p_i(G_{ji}^{-1}x_j) \Rightarrow$

$\ln p(\mathbf{x}^{(n)}|\mathbf{G}, H) = -\ln |\det \mathbf{G}| + \sum_i \ln p_i(G_{ji}^{-1}x_j)$ , which is the log-likelihood of the data D. Now, to find the gradient of the log-likelihood, we are introducing  $\mathbf{W} = \mathbf{G}^{-1}$ , now the log-likelihood can be written as:

$\ln p(\mathbf{x}^{(n)}|\mathbf{G}, H) = -\ln |\det \mathbf{W}| + \sum_i \ln p_i(W_{ji}x_j)$   
The gradient of the log-likelihood equation will be:  
 $\frac{\partial}{\partial W_{ij}} \ln p(\mathbf{x}^{(n)}|\mathbf{G}, \mathcal{H}) = -\frac{\partial}{\partial W_{ij}} (\ln |\det \mathbf{W}|) + \frac{\partial}{\partial W_{ij}} (\sum_i \ln p_i(W_{ji}x_j))$

4 Generative Models for Discrete Data

We can classify a feature vector  $\mathbf{x}$  using the Bayes rule  
 $p(y = c|\mathbf{x}, \boldsymbol{\theta}) \propto p(\mathbf{x}|y = c, \boldsymbol{\theta})p(y = c|\boldsymbol{\theta})$   
We can use different models for the data when it's discrete, based on

what kind of distribution we expect the data to assume and a respective *conjugate* prior over the model parameters  $\boldsymbol{\theta}$ .

4.1 Beta-Binomial Model

In this model we can observe a series of Bernoulli trials (e.g. coin tosses) or the number of heads (and the number of tails or total number of tosses), which is a Binomial, and it would result in the same **likelihood**:  $p(\mathcal{D}|\boldsymbol{\theta}) = \theta^{N_1}(1 - \theta)^{N_0}$ . A conjugate **prior** for this likelihood is given by  $\text{Beta}(\theta|a, b) \propto \theta^{a-1}(1 - \theta)^{b-1}$ . The **posterior** is then obtained by multiplying the prior with the likelihood,  $p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta) = \text{Bin}(N_1|\theta, N_1 + N_0)\text{Beta}(\theta|a, b) \propto \text{Beta}(\theta|a + N_1, b + N_0)$ . The **evidence** is obtained from  $p(\theta|\mathcal{D}) = \frac{1}{p(\mathcal{D})}p(\mathcal{D}|\theta)p(\theta)$  normalization of posterior is  $1/B(a + N_1, b + N_0)$ <sup>2</sup>, of prior is  $1/B(a, b)$ , hence  $p(\mathcal{D}) = \binom{N}{N_1} B(a + N_1, b + N_0)/B(a, b)$ .

4.2 Dirichlet-Multinomial Model

For example  $N$  dice rolls or *multinomial* events (with  $K$  outcomes). The **likelihood** is  $p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{k=1}^K \theta_k^{N_k}$  where  $N_k$  counts times event  $k$  occurred (it's a suff. statistic). A conjugate **prior** is the Dirichlet distribution. The **posterior** is  $p(\boldsymbol{\theta}|\mathcal{D}) \propto \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha})p(\boldsymbol{\theta}|\mathcal{D}) \propto \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}_1 + N_1, \dots, \boldsymbol{\alpha}_k + N_k)$ . The **evidence** is (obtained as the previous model)  $p(\mathcal{D}) = B(\mathbf{N} + \boldsymbol{\alpha})/B(\boldsymbol{\alpha})$ .

5 Graphical Models

Capture noise, for reasoning, enormous datasets, for causality, for designing models, CIR are encoded in the graph, for inference.

5.1 Discriminative

$p(\mathbf{t}|w) = p(w) \prod_{n=1}^N p(t_n|w)$

5.2 Generative

$T^* = \frac{P(T^*|W, X^*, X_i, T_i, \sigma^2, a)}{P(W|a)P(T^*|X^*, W, \sigma^2)} =, \quad [\prod_{i=1}^N P(X_i|W, \sigma^2)]$

5.3 Directed GMs a.k.a. Bayesian Networks

Nodes are connected by *directed* arrows. The full joint distribution is  $p(\mathbf{x}) = \prod_{k=1}^K p(x_k|\text{pa}(x_k))$   
**Directed Acyclic Graphs** are BNs without *directed* loops.

5.3.1 Blocking Rules

$\bigcirc \rightarrow \bullet \rightarrow \bigcirc \quad \bigcirc \leftarrow \bullet \rightarrow \bigcirc \quad \bigcirc \rightarrow \dots \bigcirc \dots \leftarrow \bigcirc$

5.3.2 D-separation

$A \perp\!\!\!\perp B \mid C$  holds if each path that connects a node in  $A$  with a node in  $B$  is *blocked*, that is  
**a)** the arrows on the path meet either head-to-tail or tail-to-tail at the node and the node is in  $C$ , or  
**b)** the arrows meet head-to-head at the node and either the node nor any of its descendants is in  $C$ .

5.4 Markov Random Fields

Graphical Models with undirected edges (a.k.a. Undirected Graphical Models).

<sup>1</sup>  $\Gamma(x) = \int_0^1 u^{x-1}e^{-u} = 1, \quad \Gamma(x+1) = \Gamma(x)x, \quad \Gamma(x+1) = x!$   
<sup>2</sup>  $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$

### 5.4.1 Conditional Independence

$A \perp\!\!\!\perp B \mid C$  holds if all paths connecting every node in  $A$  to every other node in  $B$  is 'blocked' by a node in  $C$  (blocked means passing through).

### 5.4.2 Cliques and Maximal Cliques

A *clique* is a subset of nodes in the graph such that there exists a link between all pairs of nodes in the subset (the set of nodes in a clique is fully connected). A *maximal clique* is a clique such that it is not possible to include any other nodes from the graph in without it ceasing to be a clique.

### 5.4.3 Factorization

A MRF can be factorized using *potential functions* over its maximal cliques:  $p(\mathbf{x}) = \frac{1}{Z} \prod_C \psi_C(\mathbf{x}_C)$   $Z = \sum_{\mathbf{x}} \prod_C \psi_C(\mathbf{x}_C)$

### 5.4.4 Relation to Directed Graphs

To transform a directed graph into an undirected one we have to perform a *moralization* process of "marrying the parents" of each node (by linking them) and removing all remaining arrows.

### 5.5 Markov Blanket

$p(x_i | x_{\text{MB}_i}, x_{\text{rest}}) = p(x_i | x_{\text{MB}_i})$

**Directed case:** The MB of  $x_i$  consists of: the parents of  $x_i$ , the children of  $x_i$  and the co-parents of the children of  $x_i$ .

**Undirected case:** All neighboring nodes of  $x_i$ .

### 5.6 Naive Bayes

The problem of finding a label  $y^*$  for some previously unobserved vector of features  $\mathbf{x}^*$ , while having previously observed  $\{\mathbf{x}_n, y_n\}_{n=1, \dots, N}$ . We build a *generative* GM, linking the label  $y_n$  with each feature of the feature vector  $\mathbf{x}_n$ . This implies that all features are independent of each other given the label. For a single data case the joint probability is

$$p(y, x_1, \dots, x_D | \boldsymbol{\eta}, \boldsymbol{\theta}_i) = p(y | \boldsymbol{\eta}) \prod_{i=1}^D p(x_i | \boldsymbol{\theta}_i)$$

while the probability of the full dataset is  $\prod_{n=1}^N p(y_n | \boldsymbol{\eta}) \prod_{i=1}^D p(x_i | \boldsymbol{\theta}_i)$  Note that this generative form is chosen as the other option (features pointing towards the label) would imply a very highly parameterized model, as we would considering  $p(y | x_1, \dots, x_D)$ .

Finding a label  $y^*$  implies finding

$$y^* = \arg\max_y \left\{ \ln p(y | \boldsymbol{\eta}) + \sum_{i=1}^D \ln p(x_i^* | y, \boldsymbol{\theta}_i) \right\}$$

### 5.7 Maximum Likelihood Training in BNs

It is fast because the log-likelihood *decomposes* into a sum over all variables  $X_i$ . Learning all parameters reduces into a collection of independent tasks of learning  $p(x_i | \text{pa}_{x_i})$ .

$$p(x_i | \text{pa}_{x_i}) = \frac{N(x_i, \text{pa}_{x_i})}{N(\text{pa}_{x_i})}$$

is the number of times  $x_i$  co-occurred with  $\text{pa}_{x_i}$  divided by the number of times  $\text{pa}_{x_i}$  occurred.

The probability of a set of variables  $\mathbf{x}$  in a Bayes Net is defined as  $p(x) = \prod_i p(x_i | \text{pa}_{x_i})$ . If we define a function  $\theta(x_i, \text{pa}_{x_i})$  and use it instead of  $p(x_i | \text{pa}_{x_i})$ , we have  $p(x) = \prod_i \theta(x_i, \text{pa}_{x_i})$ , if we constrain  $\theta$  such that  $\sum_{x_i} \theta(x_i, \text{pa}_{x_i}) = 1, \forall i$ . We now have that the probability of the full data set:  $p(\{\tilde{x}_{in}\}) = \prod_n \prod_i \prod_{x_i} \prod_{\text{pa}_{x_i}} \theta(x_i, \text{pa}_{x_i})^{\mathbb{I}[x_i = \tilde{x}_{in} \wedge \text{pa}_{x_i} = \tilde{\text{pa}}_{x_i n}]}$

### 5.8 Inference in Graphical Models

In which some of the variables are observed and we wish to compute the posterior distribution of one or more subsets of other variables.

#### 5.8.1 Inference on a chain

$$\begin{aligned} p(\mathbf{x}) &= p(x_1, \dots, x_N) \\ &= \frac{1}{Z} \psi_{1,2}(x_1, x_2) \dots \psi_{N-1,N}(x_{N-1}, x_N) \\ &= \sum_{x_1} \sum_{x_2} \dots \sum_{x_{N-1}} \sum_{x_N} p(\mathbf{x}) \\ &= \sum_{x_1} \dots \sum_{x_{n-1}} \sum_{x_{n+1}} \dots \sum_{x_{N-1}} \sum_{x_N} p(\mathbf{x}) \\ &= \frac{1}{Z} \sum_{x_1} \dots \sum_{x_{n-1}} \psi_{x_1, x_2} \dots \psi_{x_{n-1}, x_n} \mu_{\beta}(x_n) \\ \mu_{\beta}(x_n) &= \frac{1}{Z} \sum_{x_{n+1}} \psi_{x_n, x_{n+1}} \dots \sum_{x_N} \psi_{x_{N-1}, x_N} \\ \mu_{\alpha}(x_n) &= \frac{1}{Z} \sum_{x_{n-1}} \psi_{x_{n-1}, x_n} \dots \sum_{x_1} \psi_{x_2, x_1} \\ p(x_n) &= \frac{1}{Z} \mu_{\alpha}(x_n) \mu_{\beta}(x_n) \quad O(NK^2) \end{aligned}$$

#### 5.8.2 Factor Graphs

A tree is a graph with no loops. Both directed and undirected trees can be converted to a factor graph tree, but a directed tree could result in a non-tree structure when converted to an undirected representation. It is called a poly-tree (and not simply a tree) since its undirected representation (middle graph) includes a loop. The factor graph representation is again a tree. Factor graphs are the most general representation, and since any other tree representation can be easily converted to a factor tree, the sum-product algorithm is defined for factor trees.

#### 5.8.3 Sum-Product Algorithm

It is used to find *individual* marginals.

Probability of the factor graph:  $p(\mathbf{x}) = \frac{1}{Z} \prod_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha})$

**factor**  $\rightarrow$  **variable** message:

$$\mu_{f_{\alpha} \rightarrow x_i}(x_i) = \sum_{\mathbf{x}_{\alpha} \setminus x_i} f_{\alpha}(\mathbf{x}_{\alpha}) \prod_{m \in \text{ne}(f_{\alpha}) \setminus x_i} \mu_{x_m \rightarrow f_{\alpha}}(x_m)$$

**variable**  $\rightarrow$  **factor** message:

$$\mu_{x_m \rightarrow f_{\alpha}}(x_m) = \prod_{\beta \in \text{ne}(x_m) \setminus f_{\alpha}} \mu_{f_{\beta} \rightarrow x_m}(x_m)$$

$$\begin{aligned} \text{leaf nodes} \quad x_l \quad & \text{is variable:} \quad \mu_{x_l \rightarrow f_{\delta}}(x_l) = 1 \\ f_{\varepsilon} \quad & \text{is factor:} \quad \mu_{f_{\varepsilon} \rightarrow x_k}(x_k) = f_{\varepsilon}(x_k) \end{aligned}$$

$$\text{variable marginal: } p(x) = \prod_{\alpha \in \text{ne}(x)} \mu_{f_{\alpha} \rightarrow x}(x)$$

$$\text{factor marginal: } p(\mathbf{x}_{\alpha}) = f_{\alpha}(\mathbf{x}_{\alpha}) \prod_{i \in \text{ne}(f_{\alpha})} \mu_{x_i \rightarrow f_{\alpha}}(x_i)$$

By caching intermediate results and calculating two messages for each edge in a tree-structured factor graph, the Sum-Product algorithm efficiently computes all variable marginals in  $O(2EK^2)$  (where  $E$  is the number of edges and each variable is assumed to have  $K$  possible values).

#### 5.8.4 Max-Sum Algorithm

It is used to find a setting of the variables that has the largest *joint* probability and the value of that probability.

$$\mathbf{x}^{\max} = \arg\max_{\mathbf{x}} \prod_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha}) \quad p(\mathbf{x}^{\max}) = \max_{\mathbf{x}} \frac{1}{Z} \prod_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha})$$

$$\log p(\mathbf{x}^{\max}) = \max_{\mathbf{x}} \left\{ \log Z + \sum_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha}) \right\}$$

**factor**  $\rightarrow$  **variable** message:

$$\mu_{f_{\alpha} \rightarrow x_i}(x_i) = \max_{\mathbf{x}_{\alpha} \setminus x_i} \left[ (\ln f_{\alpha}(\mathbf{x}_{\alpha})) + \sum_{m \in \text{ne}(f_{\alpha}) \setminus x_i} \mu_{x_m \rightarrow f_{\alpha}}(x_m) \right]$$

**variable**  $\rightarrow$  **factor** message:

$$\mu_{x_m \rightarrow f_{\alpha}}(x_m) = \sum_{f_{\beta} \in \text{ne}(x_m) \setminus f_{\alpha}} \mu_{f_{\beta} \rightarrow x_m}(x_m)$$

$$\begin{aligned} \text{leaf nodes} \quad x_l \quad & \text{is variable:} \quad \mu_{x_l \rightarrow f_{\delta}} = 0 \\ f_{\varepsilon} \quad & \text{is factor:} \quad \mu_{f_{\varepsilon} \rightarrow x_k} = \ln f_{\varepsilon}(x_k) \\ \text{at root} \quad & p^{\max} = \max_x \left[ \sum_{f_{\alpha} \in \text{ne}(x)} \mu_{f_{\alpha} \rightarrow x}(x) \right] \\ & x^{\max} = \arg\max_x \left[ \sum_{f_{\alpha} \in \text{ne}(x)} \mu_{f_{\alpha} \rightarrow x}(x) \right] \end{aligned}$$

$$\text{max-marginals} \quad q(x_i) = -\ln Z + \sum_{f_{\alpha} \in \text{ne}(x_i)} \mu_{f_{\alpha} \rightarrow x_i}(x_i)$$

Given max-marginals you have to perform a decoding step (the Viterbi algorithm) in order to find the global optimum. If  $q(x_i)$  has unique maximum, we can use  $x_i^* = \arg\max_{x_i} q(x_i)$ .

#### 5.8.5 Variable Elimination

If we have some clusters  $C_i$  and  $C_j$ , then:

$S_{ij} = C_i \cap C_j$  Define potential functions  $\Psi_i$  for each cluster  $C_i$ :  $\forall C_i : \Psi_i := \prod_{\psi_j \in C_i} \psi_j$   $m_i \rightarrow C_j = \sum_{C_i \setminus S_{ij}} \Psi_i \prod_{k \in \text{ne}(i) \setminus j} m_{k \rightarrow i}(C_i)$

Note that:  $C_i \setminus S_{ij} = C_i \setminus C_j$

Belief:  $\Psi_{\text{root}} \prod_{k \in \text{ne}(\text{root})} m_{k \rightarrow \text{root}}(C_{\text{root}}) \propto p(C_{\text{root}})$

Clusters replace factors, edges replace variables.

**Tree width** The width of a tree is one smaller than the biggest clique in the optimal ordering. Finding the optimal order is NP-hard.

**Running Intersection Property** A clique tree has the running intersection property if  $x \in C_i \wedge x \in C_j \Rightarrow x \in C_k$  on the clique path between  $C_i$  and  $C_j$ .

#### 5.8.6 Loopy Belief Propagation

Is a way to approximate inference in graphs with loops. The idea is simply to apply the sum-product algorithm even though there is no guarantee that it will yield good results. This approach is possible because the message passing rules are purely local. However, because of the loops, information can flow many times around the graph; for some models it will converge, for others it will not. We will say that a node  $a$  has a **pending** message on its link to a node  $b$  if node  $a$  has received any message on any of its other links since the last time it sent a message to  $b$ . Thus, when a node receives a message on one of its links, this creates pending messages on all of its other links. Then we need to send only pending messages otherwise they will be duplicate. When there are no more pending messages, the product of the received messages at every node gives the (exact for tree) marginal. For loopy graphs however the algorithm may not converge and when we stop the marginals are *approximate*.

### 5.9 EM, VM, VB

EM, VEM are learning algorithms, VB is for inference.

If  $Q = P$ , use EM. If  $Q$  is simpler than  $P$ , use VEM (easier computation, might never reach max, biased result)

Bayes: approximate  $P(\theta | X)$  (eg sampling). Maximize both  $\theta$  and  $\alpha$ .

Frequentist:	EM: unbiased, can overfit	VEM: biased, can overfit
Bayesian:	$P(\theta   X)$ : best, but comp.hard, no overfitting bc no fitting variance	VB

Bayes does not overfit because you integrate over parameter space.

## 5.10 Variational EM Algorithm

Variational EM algorithm has two steps:

### E-Step:

Solve  $q_t(z_n|x_n) = p(z_n|x_n, \theta_t) \rightarrow B = l$ , (or increase  $B(\theta_t, q)$  over  $q$ ). This is not easy to solve, but sometimes it can be done (for example in Mixture of Gaussians).

It could be a situation where:

$q(z|t) = \prod_{i=1}^k q(z|x)$  The green dot upon the circle will be the best  $q$  that approximates  $p$  as it is the shortest point from  $p$ .

### M-Step:

Solve  $t+1 = \arg \max_{\theta} \sum_n E q_t[\log p(x_n, z_n|\theta)]$ , (or increase  $B(\theta, q_t)$  over  $\theta$ ). Implicitly,  $\theta$  is in that  $q_t$  of the expectation because it is fixed from the previous iteration. We maximize the lower bounds until we find the maximum. The convergence is rather slow. In terms of speed is not the best algorithm.

Proof of Convergence:

We have to prove that:  $h(\theta_t) \leq \dots \leq h(\theta_{t+1})$

$h(\theta_t) = B(\theta_t, p(z|x, \theta_t))$   $h(\theta_t) \leq B(\theta_{t+1}, p(z|x, \theta_t))$ , (M-Step): We maximise over  $\theta_{t+1}$ . So  $B$  in eq.(13) must be bigger than the one in the (12). In the next step we have:  $h(\theta_t) = l(\theta_{t+1}) - \mathcal{KL}[q_t||p(z|x, \theta_{t+1})]$ , ( $l=B+\mathcal{KL}$ ) The distribution in Kullback-Leibler divergence are not the same so  $\mathcal{KL}$  is not equal to 0 at this occasion, due to the fact that the two terms of  $\mathcal{KL}$  divergence are coming from different iterations. Hence,  $h(\theta_t) \leq h(\theta_{t+1})$ , cause  $\mathcal{KL} \geq 0$ .

**Lower bound  $\mathcal{L}(q)$  and KL-divergence**  $\ln p(X) = \mathcal{L}(q) + \mathcal{KL}(q||p)$

$$\mathcal{L}(q) = \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ$$

$$\mathcal{KL}(q||p) = - \int q(Z) \ln \left\{ \frac{p(Z|X)}{q(Z)} \right\} dZ$$

**Note that**

$$Q(\theta, \theta^{old}) = \sum_z p(Z|X, \theta^{old}) \ln p(X, Z|\theta) = \mathbb{E}_{Z|X, \theta^{old}} [\ln p(X, Z|\theta)]$$

## 6 Sampling Methods

### 6.1 Rejection sampling

We take a  $q(z)$  distribution, choose  $k$  for  $kq(z) \geq p(z)$ ,  $ez \in p(z)$

$$p(accepted) = \int \{p(z)/kq(z)\} q(z) dz = \frac{1}{K} \int p(z) dz$$

### 6.2 Adaptive rejection sampling

When we have a distribution  $p(z)$  where  $\ln p(z)$  has derivatives which are not increasing functions of  $z$ . The envelope distribution is a succession of linear functions.  $q(z) = k_i \lambda_i \exp\{-\lambda_i(z - z_i)\}$ . Then rejection sampling can be applied.

### 6.3 Importance sampling

Provides a framework for approximating distributions but does not provide a mechanism for drawing samples from them. Used when it is not easy to find a function that forms an upper bound. Note that  $q$  does not necessarily have to be a bound on  $p$ , as we weigh by the quotient of the two distributions. Samples are drawn from distribution  $q$  and weighed correspondingly:  $x_i \approx \tilde{q}$ ,  $w_i = \frac{\tilde{p}(x_i)}{\tilde{q}(x_i)}$  Works bad in high dimensions as it is hard to find an  $x$  where both  $q(x)$  and  $p(x)$  are high in the high-dimensional space. If  $p$  and  $q$  overlap just a little, this will give a bad estimate.

### 6.4 Ancestral sampling (Likelihood Weighted sampling)

Given a Bayesian network, we can write down the joint probability as  $p(z_1)p(z_2|z_1)$ . That means there is always a node in the tree that is not dependent on any other variable. We sample this node, and use the sample to draw samples for subsequent nodes. If we observe values ( $x \in Ev$ ,

evidences), we just select the observed value instead of sampling. Works bad in high dimensions as well.

## 6.5 Markov Chain Monte-Carlo (MCMC) sampling

MCMC runs ancestral sampling on a chain. Samples can no longer be drawn independently.  $t \rightarrow 1$ ,  $x_t \sim q_{\infty}(x_{\infty})$  Equilibrium or invariant distribution. We transition from  $q_1$  to  $q_2$  according to transition probability  $T(x_2|x_1)$ .

$$q(x_{t+1}, x_t) = T(x_{t+1}|x_t)q(x_t)$$

**Invariance:** We want to find  $T$  for a given  $p$ , such that  $p = Tp$  holds.

**Detailed balance:**(or reversibility) We want to find  $T$  so that transition from one state to another has probability mass equal to the transition from that next state to the current. Thus,  $p^*(z)T(z, z') = p^*(z')T(z', z)$ . Distribution  $p^*(z)$  is invariant if  $p^*(z) = \sum_{z'} p^*(z')T(z', z)$

$$\text{Proof: } A_k(z', z) = \min \left( 1, \frac{p(z')q(z|z')}{p(z)q(z'|z)} \right)$$

$$p(z)q(z'|z)A_k(z', z) = \min(p(z)q(z'|z), p(z')q(z|z')) \text{ [multiply by denom]}$$

$$= \min(p(z')q(z|z'), p(z)q(z'|z)) \text{ [order irrelevant]}$$

$= p(z')q(z|z')A_k(z, z')$  **Ergodicity:** This algorithm needs ergodicity: A positive probability for every state (so all will be reached).

*For two transition kernels  $T_1, T_2$  that are valid any linear combination of these two is a valid kernel.*

## 6.6 Metropolis-Hastings algorithm

In the  $\tau$  step of the algorithm the current state is  $z^{(\tau)}$ , we draw sample  $z^*$  from  $q_k(z|z^{(\tau)})$ , which it is been accepted with probability:  $A_k(z^*, z^{(\tau)}) = \min(1, \frac{p(z^*)q_k(z^{(\tau)}|z^*)}{p(z^{(\tau)})q_k(z^*|z^{(\tau)})})$ , where  $k$  labels the set of possible transitions. Metropolis-Hastings algorithm work better in high dimensions than previous methods, but is quite slow (as we perform a random walk).

## 6.7 Gibbs sampling

For a distribution  $p(z)$ , for each component of the distribution  $z_i$  draw a sample from the distribution of that variable given the values of the remaining variables:

1. Initialize  $\{z_i : i = 1, \dots, M\}$
2. For  $\tau = 1, \dots, T$ :
  - Sample  $z_1^{\tau+1} \sim p(z_1|z_2^{\tau}, \dots, z_M^{\tau})$
  - Sample  $z_2^{\tau+1} \sim p(z_2|z_1^{\tau+1}, \dots, z_{j-1}^{\tau+1}, z_{j+1}^{\tau}, \dots, z_M^{\tau})$
  - Sample  $z_M^{\tau+1} \sim p(z_M|z_1^{\tau+1}, \dots, z_{M-1}^{\tau+1})$

## 7 HMMs + Kalman filters

### 7.1 Markov models

First order Markov chains hold the Markov property which you whole history is only dependent on the current state. **First order's** number of parameters will be:  $k-1+k(k-1)$ . **Second order:**  $k^2-1+k^2(k-1)$ .

## 7.2 HMMs

Useful in speech recognition, natural language processing, online character recognition... Discrete variables: HMM, Linear-Gaussian: LDS (Linear Dynamical System).

$$p(x, \dots, x_n) = \sum_{z_1} \dots \sum_{z_n} p(z_1) (\prod_{n=2}^N p(z_n|z_{n-1}, A)) \prod_{n=1} p(x_n|z_n)$$

For the first node:

$$p(z_1|\pi) = \prod_{k=1}^K \pi_k^{z_{1,k}}$$

$$p(z_n|z_{n-1}, A) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j} z_{n,k}}$$

$$p(x_n|z_n, \phi) = \prod_{k=1}^K p(x_n|\phi_k)^{z_{n,k}}$$

Where  $A$ ,  $\phi$  and  $\pi$  are the parameters of the model.  $p(x_n|z_n, \phi)$  is model-dependent (e.g.: mixture of Gaussians, multinomial if discrete, etc).

### 7.2.1 Maximum likelihood (EM)

data:  $X$ , latent state:  $Z$ , parameters:  $\theta = (\pi, A, \phi)$ .

$$p(x|\theta) = \sum_z p(x, z|\theta)$$

This represents a sum of KN terms, which quickly becomes intractable, so we can use for instance the EM algorithm.

**E-step:**

$$Q(\theta|\theta^{old}) = \sum_z p(z|x, \theta^{old}) \ln p(x, z|\theta)$$

**M-step:**

$$\theta^{new} = \arg \max_{\theta} Q(\theta|\theta^{old})$$

EM-general:

$$\ln p(X|\theta) = \mathcal{L}(q, \theta) + \mathcal{KL}(q||p) \text{ where } \mathcal{L}(q, \theta) = \sum_Z q(Z) \ln \left( \frac{p(X, Z|\theta)}{q(Z)} \right)$$

$$\text{and } \mathcal{KL}(q||p) = - \sum_Z q(Z) \ln \left( \frac{p(Z|X, \theta)}{q(Z)} \right).$$

E-step: maximize  $\mathcal{L}(q, \theta)$  wrt  $q(Z)$ . M-step: maximize  $\mathcal{L}(q, \theta)$  wrt  $\theta$ .

### 7.2.2 Forward-Backward

$$a(z_n) = p(x_n|z_n) \sum_{z_{n-1}} a(z_{n-1}) p(z_n|z_{n-1})$$

$$\beta(z_n) = \sum_{z_{n+1}} \beta(z_{n+1}) p(x_{n+1}|z_{n+1}) p(z_{n+1}|z_n)$$

**Predictive distribution:** Add  $x_{n+1}$  and  $z_{n+1}$  at the end of the chain and use the calculated parameters to predict the new values.(assuming the chain is homogeneous). By d-separation (if applicable),  $p(x_{N+1}|X) = p(x_{N+1}|x_N)$  (regular Markov Chain, add  $z$ 's to taste).

## 8 Mathematical Tricks

$$-\frac{1}{2}(x - \mu)^{\top} \Sigma^{-1} (x - \mu) = -\frac{1}{2} x^{\top} \Sigma^{-1} x + x^{\top} \Sigma^{-1} \mu + c$$

$$\frac{\partial a^{\top} X b}{\partial X} = ab^{\top}, \quad \frac{\partial a^{\top} X^{\top} b}{\partial X} = ba^{\top}$$

$$\frac{\partial \det X}{\partial X} = \det(X) (X^{-1})^{\top}$$

$$M = M^{\top} \iff u^{\top} M v = v^{\top} M u$$

$$\mathbb{E}[(x - \mu)^{\top} \Sigma^{-1} (x - \mu) \mathcal{N}(x|\mu, \Sigma)] = \text{Tr}(\Sigma^{-1} \Sigma) = D$$

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} x f(x) dx$$

Product of Gaussians:  $\Sigma = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}$ ,  $\mu = \Sigma \Sigma_1^{-1} \mu_1 + \Sigma \Sigma_2^{-1} \mu_2$

Other:  $(P^{-1} + B^{\top} R^{-1} B)^{-1} B^{\top} R^{-1} = P B^{\top} (B P B^{\top} + R)^{-1}$

Woodbury:  $(A + B D^{-1} C)^{-1} = A^{-1} - A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}$

Gradient in exp family wrt for  $\eta$ , for  $p(x|\eta) = h(x)g(\eta)\exp\{\eta^{\top} u(x)\}$ , where  $g(\eta) \int h(x)\exp\{\eta^{\top} u(x)\} dx = 1$ , we have the result:  $-\nabla \ln g(\eta) = \mathbb{E}[u(x)]$