

### 3. Graphical models on directed acyclic graphs

- $S = \{S_i \mid i \in V\}$  - collection of  $K$ -valued random variables
- $(V, E)$  - directed acyclic graph (DAG)

Definition 1  $S$  is a Bayesian network (aka belief network) w.r.t. the DAG  $(V, E)$  if its joint distribution is

$$p(S) = \prod_{i \in V} p(s_i \mid S_{pa(i)}),$$

where  $pa(i) \subset V$  denotes the parents of node  $i$ . ■

Let  $de(i) \subset V$  denote all descendants of a node  $i \in V$ . Then  $S_i$  conditioned on  $S_{pa(i)}$  is independent of all  $S_j$ ,  $j \in V \setminus de(i)$ , i.e.

$$S_i \perp\!\!\!\perp S_{V \setminus de(i)} \mid S_{pa(i)}.$$

Remark 1 Notice that BNs do not imply causality. For example, a Markov model on a chain is both, a Markov model on an undirected graph (chain) and a BN on the directed chain. ■

### A. Stochastic neural networks with binary units

- $X = \{X_i \mid i \in V\}$  a collection of  $\pm 1$  valued r.v.
- $(V, E)$  -  $M$ -partite DAG with subsets  $V_0, V_1, \dots, V_M$

Denote  $X^m = \{X_i \mid i \in V_m\}$  and consider the conditional Bayesian network

$$p(X^M, X^{M-1}, \dots, X^1 \mid X^0) = p(X^M \mid X^{M-1}) \circ \dots \circ p(X^1 \mid X^0),$$

where

$$p(x^m | x^{m-1}) = \prod_{i \in V_m} p(x_i^m | x^{m-1})$$

and

$$p(x_i^m | x^{m-1}) = \frac{e^{x_i^m \langle w_i^m, x^{m-1} \rangle}}{2 \operatorname{ch} \langle w_i^m, x^{m-1} \rangle}$$

This is a sigmoid belief network

Remark 2 A sigmoid belief network on a  $M$ -partite DAG is not a MRF on the corresponding undirected graph. ■

Remark 3 Computing the probabilities of the output nodes  $x^M$  given the input  $x^0$ , requires to marginalise over all hidden (latent) layers  $x_1^1, \dots, x^{M-1}$  and is hard. ■

Remark 4 Sampling a realisation  $x_1^1, \dots, x^M$  given  $x^0$  is easy (linear in model size). ■

## B. Stochastic neural networks with Gaussian units

- $Z = \{z_i | i \in V\}$  a collection of real valued r.v.
- $(V, E)$  a  $M$ -partite DAG with subsets  $V_0, V_1, \dots, V_M$

Denote  $Z^m = \{z_i | i \in V_m\}$  and consider the conditional BN

$$p(z^M, z^{M-1}, \dots, z^1 | z^0) = p(z^M | z^{M-1}) \cdot \dots \cdot p(z^1 | z^0),$$

where

$$p(z^m | z^{m-1}) \sim \mathcal{N}(\mu | \theta, z^{m-1}), C(\theta, z^{m-1})$$

with a diagonal covariance matrix  $C$

Remark 5 To compute  $p(z^M | z^0)$ , we need to solve the integrals

$$p(z^M | z^0) = \int dz^{M-1} \dots \int dz^1 p(z^M | z^{M-1}) \dots p(z^1 | z^0).$$

This is hard. □

Sampling a realisation  $z^1, \dots, z^M$  given  $z^0$  is easy (linear in model size). Derivatives w.r.t. to model parameters can be computed by sampling and using the identities

$$\nabla_{\mu_i} \mathbb{E}_{\mathcal{N}(\mu, C)} f(z) = \mathbb{E}_{\mathcal{N}(\mu, C)} \frac{\partial f(z)}{\partial z_i}$$

$$\nabla_{C_{ij}} \mathbb{E}_{\mathcal{N}(\mu, C)} f(z) = \frac{1}{2} \mathbb{E}_{\mathcal{N}(\mu, C)} \frac{\partial^2 f(z)}{\partial z_i \partial z_j}$$