

Probabilistic Neural Networks

- $X = \{X_i | i \in V\}$ - collection of K -valued random variables
- (V, E) - directed acyclic graph (DAG)

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

$$P(x) = \prod_{i \in V} P(x_i | x_{pa(i)}),$$

where $pa(i) \subset V$ denotes the parents of node $i \in V$ □

Let $de(i)$ denote all descendants of a node $i \in V$. Then X_i conditioned on $X_{pa(i)}$ is independent of all X_j with $j \in V \setminus de(i)$, i.e.

$$\& X_i \perp\!\!\!\perp X_{V \setminus de(i)} | X_{pa(i)}$$

Remark 1 Notice that BNs do not imply causality.

For example, a Markov chain model is both, a Markov model on an undirected graph (chain) and a BN on a directed chain.

A. Probabilistic networks

In probabilistic networks all units are considered random. The units are organised in layers. Each layer $k=1, \dots, \ell$ is described by a random vector X^k with dimension n_k (number of layer units). The network \mathcal{N} is a conditional Bayesian network with pdf

$$p(x^{1, \dots, \ell} | x^0) = \prod_{k=1}^{\ell} p(x^k | x^{k-1})$$

where it is usually assumed that the conditional distributions factorise as

$$p(x^k | x^{k-1}) = \prod_{i=1}^{n_k} p(x_i^k | x^{k-1})$$

1C. Approximate inference for PNN

Computing $p(x^k | x^0)$ for probabilistic NN is not tractable.

Approximate $p(x^k | x^0) \approx \prod_{i=1}^{n_k} q(x_i^k)$ with some simple distributions q . This factorised approximation can be computed layer by layer

$$\begin{aligned} q(x_i^k) &= \mathbb{E}_{q(x^{k-1})} [p(x^k | x^{k-1})] \\ &= \int p(x^k | x^{k-1}) \prod_{i=1}^{n_{k-1}} q(x_i^{k-1}) dx^{k-1} \end{aligned}$$

Linear layer $A = \langle w, x \rangle$, where $\mathbb{E}[x] = \mu$, $V[x] = \sigma^2 \mathbb{I}$

Then $\mu' = \mathbb{E}[A] = \langle w, \mu \rangle$

$$\sigma'^2 \approx \sum_i w_i^2 \sigma_i^2$$

Non-linear (coordinate-wise) mappings A is a scalar r.v. with statistic (μ, σ^2) . $y = f(A - z)$ with independent noise. Assume $A' = A - z$ is normally (logistically) distributed with (μ', σ'^2) . We can approximate the expectation and variance of $y = f(A')$

$$\tilde{\mu} = \mathbb{E}_{q(A')} [f(A')], \quad \tilde{\sigma}^2 = \mathbb{E}_{q(A')} [f^2(A')] - \tilde{\mu}^2$$

Example 2 (Heaviside nonlinearity with noise)

We consider $y = h(A - z)$, where z is standard logistic noise. Statistic of $A' = A - z$ is given by

$$\mu' = \mu, \quad \sigma'^2 = \sigma^2 + \sigma_s^2, \quad \text{where } \sigma_s^2 = \pi^2/3$$

The posterior distribution of each layer $k > 0$ given the observations x^0 expresses recurrently ~~as~~ as

$$p(x^k | x_0) = \int p(x^k | x^{k-1}) p(x^{k-1} | x_0) dx^{k-1}$$

Remark 2 Computing these (conditional) marginal distributions is usually NP-hard. However, sampling a realisation x^1, \dots, x^k given x^0 is often easy (linear in model size).

B. Standard NN \rightarrow Probabilistic NN

Standard (deterministic) artificial neuron

$$y = f\left(\underbrace{\sum_i w_i x_i + b}_a\right) = f(a)$$

Inject noise Z with known distribution (e.g. normal distr.)

$$Y = f(A - Z)$$

We have then $F_Y(u) = \mathbb{P}(f(A - Z) < u)$

Example 1 (Stochastic binary unit)

Let Y be a binary r.v. $Y = h(A - Z)$, where h is the Heaviside step function and Z is noise with cdf. F_Z . Then

$$\mathbb{P}(Y=1 | A) = \mathbb{P}(A - Z > 0 | A) = \mathbb{P}(Z < A) = F_Z(A).$$

If for instance Z has standard logistic distribution, then $\mathbb{P}(Y=1 | A) = S(A)$, where S is the sigmoid function $S(a) = (1 + e^{-a})^{-1}$

Assuming A' to have logistic distribution, we get the statistic of Y by

$$\begin{aligned}\hat{\mu} &= \mathbb{E}[h(A')] = \mathbb{P}(A' > 0) = \mathbb{P}\left(\frac{A' - \mu'}{\sigma' \sigma_s} \geq -\frac{\mu'}{\sigma' \sigma_s}\right) = \\ &= \mathcal{S}\left(\frac{\mu'}{\sigma' \sigma_s}\right) = \mathcal{S}\left(\frac{\mu}{\sqrt{\sigma^2 \sigma_s^2 + 1}}\right)\end{aligned}$$

□

Using such approximations for ReLU, max-pooling and for categorical layers, we can forward propagate means and variances through the layers of probabilistic NN and then apply back-propagation.