

How to deal with missing data in supervised deep learning?

Niels Bruun Ipsen, Pierre-Alexandre Mattei, Jes Frellsen

July 15, 2020

How to approach $p(\mathbf{y}|\mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}})$, assuming MAR?

0-imputation:
$$\iota_0(\mathbf{x}^{\text{obs}}) = \mathbf{x} \odot \mathbf{s} + \mathbf{0} \odot (1 - \mathbf{s}) \quad (1)$$

learnable imputation:
$$\iota_\lambda(\mathbf{x}^{\text{obs}}) = \mathbf{x} \odot \mathbf{s} + \boldsymbol{\lambda} \odot (1 - \mathbf{s}) \quad (2)$$

concatenation in separate channels:
$$\iota_0(\mathbf{x}^{\text{obs}}), \boldsymbol{\lambda}, \text{ and } \mathbf{s} \quad (3)$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathcal{X}^n$ contain n i.i.d. copies of the random variable $\mathbf{x} \in \mathcal{X}$ and the positions of observed entries in the data matrix \mathbf{X} are contained in a mask matrix $\mathbf{S} = (\mathbf{s}_1, \dots, \mathbf{s}_n)^T \in \{0, 1\}^{n \times p}$ and $\mathbf{x} = (\mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}})$.

Deep Latent Variable Model, DLVM: $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$ (4)

Joint DLVM and discriminative model: $p(\mathbf{z})p(\mathbf{x}^{\text{obs}}|\mathbf{z})p(\mathbf{x}^{\text{miss}}|\mathbf{z})p(\mathbf{y}|\mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}})$ (5)

Lower bound for training

$$\mathcal{L}_K = \mathbb{E} \left[\log \left(\frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{z}_k)p(\mathbf{x}^{\text{obs}}|\mathbf{z}_k)p(\mathbf{y}|\mathbf{x}^{\text{obs}}, \mathbf{x}_k^{\text{miss}})}{q(\mathbf{z}_k|\mathbf{x}^{\text{obs}}, \mathbf{s})} \right) \right] \leq \log p(\mathbf{y}, \mathbf{x}^{\text{obs}}) \quad (6)$$

where $q(\mathbf{z}_k|\mathbf{x}^{\text{obs}}, \mathbf{s})$ is the *variational distribution* (learnable proposal) and $(\mathbf{z}_k, \mathbf{x}_k^{\text{miss}})_{k \in \{1, \dots, K\}}$ are i.i.d. samples from $p(\mathbf{x}^{\text{miss}}|\mathbf{z})q(\mathbf{z}|\mathbf{x}^{\text{obs}}, \mathbf{s})$.

Prediction: self-normalized importance sampling

$$p(\mathbf{y}|\mathbf{x}^{\text{obs}}) \approx \sum_{i=1}^K w_k p(\mathbf{y}|\mathbf{x}^{\text{obs}}, \mathbf{x}_k^{\text{miss}}), \quad (7)$$

where

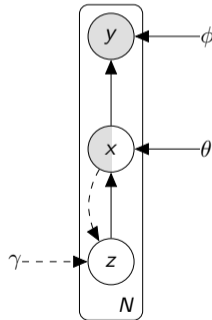
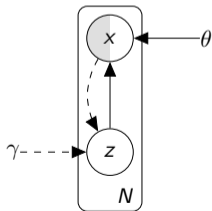
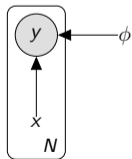
$$w_k = \frac{r_k}{r_1 + \dots + r_K}, \quad \text{and } r_k = \frac{p(\mathbf{z}_k)p(\mathbf{x}^{\text{obs}}|\mathbf{z}_k)}{q(\mathbf{z}_k|\mathbf{x}^{\text{obs}}, \mathbf{s})}, \quad (8)$$

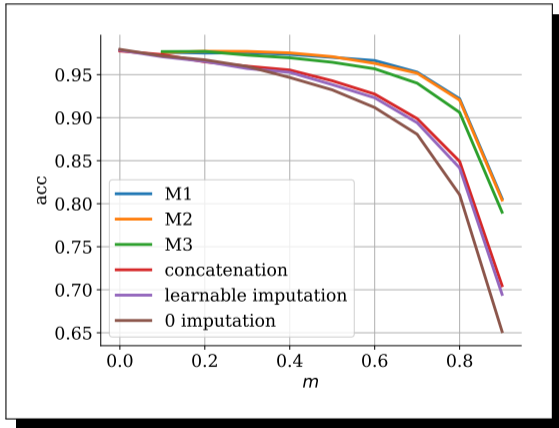
and $(\mathbf{z}_k, \mathbf{x}_k^{\text{miss}})_{k \in \{1, \dots, K\}}$ are i.i.d. samples from $p(\mathbf{x}^{\text{miss}}|\mathbf{z})q(\mathbf{z}|\mathbf{x}^{\text{obs}}, \mathbf{s})$.

Discriminative model: $p_{\phi}(\mathbf{y}|\mathbf{x})$ (9)

Deep Latent Variable Model, DLVM: $p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$ (10)

Joint DLVM and discriminative model: $p(\mathbf{z})p_{\theta}(\mathbf{x}^{\text{obs}}|\mathbf{z})p_{\theta}(\mathbf{x}^{\text{miss}}|\mathbf{z})p_{\phi}(\mathbf{y}|\mathbf{x}^{\text{obs}}, \mathbf{x}^{\text{miss}})$ (11)





M1, joint model, trained jointly
M2, joint model, DLVM and discriminative model trained separately
M3, DLVM used for imputing. Discriminative model trained on imputed dataset.

Thank You