

Chow et Liu dans le cas gaussien

November 7, 2017

Abstract

Approximating Gaussian distribution with dependence Trees

1 Notations and definitions

Let \mathbf{x} be a vector of n random continuous variables, $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{R}^n$. Let \mathcal{S} be a collection of samples of \mathbf{x} : $\mathcal{S} = \{\mathbf{x}^1, \dots, \mathbf{x}^s\}$.

We consider a tree structure, each node of the tree being a variable. All nodes i except for one have an ancestor node, indexed a_i . The structure of a tree is recorded in the vector $\boldsymbol{\alpha} = (a_1, \dots, a_n)$.

We assume Gaussian densities :

$$f(x_i | \mu_i, \sigma_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right)$$

$$f(x_i, x_{a_i} | \mu_i, \mu_{a_i}, \Sigma_{ia_i}) = \frac{1}{2\pi\sqrt{\det(\Sigma_{ia_i})}} \exp\left(-\frac{1}{2}(x_i - \mu_i, x_{a_i} - \mu_{a_i})^T \Sigma_{ia_i}^{-1} (x_i - \mu_i, x_{a_i} - \mu_{a_i})\right)$$

$$\text{Where } \Sigma_{ia_i} = \begin{pmatrix} \sigma_i^2 & \sigma_{ia_i} \\ \sigma_{a_i i} & \sigma_{a_i}^2 \end{pmatrix}.$$

The gaussian dependence-tree is parametrised by $\boldsymbol{\alpha}$:

$$\begin{aligned} f(\mathbf{x} | \boldsymbol{\alpha}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= f(x_1 | \mu_1, \sigma_1) \prod_{i=2}^n f(x_i | x_{a_i}, \mu_i, \mu_{a_i}, \Sigma_{ia_i}) \\ &= f(x_1 | \mu_1, \sigma_1) \prod_{i=2}^n \frac{f(x_i, x_{a_i} | \mu_i, \mu_{a_i}, \Sigma_{ia_i})}{f(x_{a_i} | \mu_{a_i}, \sigma_{a_i})} \end{aligned}$$

2 Likelihood maximisation

2.1 Likelihood

We want to maximise the likelihood function, which writes

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \log(f(\mathbf{x} | \boldsymbol{\alpha}, \boldsymbol{\mu}, \boldsymbol{\Sigma})) \\ &= \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \left[\log(f(x_1 | \mu_1, \sigma_1)) + \sum_{i=2}^n \log\left(\frac{f(x_i, x_{a_i} | \mu_i, \mu_{a_i}, \Sigma_{ia_i})}{f(x_{a_i} | \mu_{a_i}, \sigma_{a_i})}\right) \right] \\ &= \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \left[\sum_{i=1}^n \log(f(x_i | \mu_i, \sigma_i)) + \sum_{i=2}^n \log\left(\frac{f(x_i, x_{a_i} | \mu_i, \mu_{a_i}, \Sigma_{ia_i})}{f(x_{a_i} | \mu_{a_i}, \sigma_{a_i}) \times f(x_i | \mu_i, \sigma_i)}\right) \right] \end{aligned}$$

2.2 Estimation

The maximum-likelihood estimators for the parameters μ_i , σ_i and σ_{ia_i} are known :

$$\begin{aligned}\hat{\mu}_i &= \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} x_i, \\ \hat{\sigma}_i^2 &= \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} (x_i - \hat{\mu}_i)^2, \\ \hat{\sigma}_{ia_i}^2 &= \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} (x_i - \hat{\mu}_i)(x_{a_i} - \hat{\mu}_{a_i})\end{aligned}$$

Using these estimators, we finally get

$$\mathcal{L}(\alpha, \hat{\mu}, \hat{\Sigma}) = \sum_{i=1}^n \frac{1}{2} (1 + \log(2\pi\hat{\sigma}_i^2)) + \sum_{i=2}^n \underbrace{-\frac{1}{2} \log \left(1 - \frac{\hat{\sigma}_{ia_i}}{\hat{\sigma}_i^2 \hat{\sigma}_{ia_i}^2} \right)}_{\mathcal{I}(f(\cdot|\hat{\mu}_i, \hat{\sigma}_i), f(\cdot|\hat{\mu}_{a_i}, \hat{\sigma}_{a_i}))}$$

The first term of this quantity is independant of α , meaning it is independant of the structure of the dependance-tree. We then only need to maximise the second term to obtain an MLE estimator for the tree structure.

The last term is known as the Shannon information between two variables i and a_i . As the Shannon information is increasing with the correlation between these variables, it is a weight on each branch of the tree. The structure is then optimized by the maximum-weight spanning tree :

$$\hat{\alpha} = \arg \max_{\alpha} \left\{ \sum_{i=2}^n \mathcal{I}(f(\cdot|\hat{\mu}_i, \hat{\sigma}_i), f(\cdot|\hat{\mu}_{a_i}, \hat{\sigma}_{a_i})) \right\}$$