

sion) the activation of hidden units z_j is computed as $z_j = \epsilon_j \cdot f(\mathbf{pa}_j)$ with $\epsilon_j \sim p(\epsilon_j) = \text{Bernoulli}(0.5)$, and where the parameters are learned by following the gradient of the log-likelihood lower bound: $\nabla_{\theta} \mathbb{E}_{\epsilon} [\log p_{\theta}(\mathbf{t}^{(i)} | \mathbf{x}^{(i)}, \epsilon)]$; this gradient can sometimes be computed exactly (Maaten et al., 2013) and can otherwise be approximated with a Monte Carlo estimate (Hinton et al., 2012). The two parameterizations explained in section 3.1 offer us a useful new perspective on 'dropout'. A 'dropout' hidden unit (together with its injected noise ϵ) can be seen as the DNCP of latent random variables, whose CP is $z_j | \mathbf{pa}_j \sim p_{\theta}(z_j = \epsilon_j \cdot f(\mathbf{pa}_j) | \mathbf{pa}_j)$. A practical implication is that 'dropout'-type neural networks can therefore be interpreted and treated as hierarchical Bayes nets, which opens the door to alternative approaches to learning the parameters, such as Monte Carlo EM or variational methods.

While 'dropout' is designed as a regularization method, other work on stochastic neural networks exploit the power of stochastic hidden units for generative modeling, e.g. (Frey & Hinton, 1999; Rezende et al., 2014; Tang & Salakhutdinov, 2013) applying (partially) MCMC or (partially) factorized variational approaches to modelling the posterior. As we will see in sections 4 and 6, the choice of parameterization has a large impact on the posterior dependencies and the efficiency of posterior inference. However, current publications lack a good justification for their choice of parameterization. The analysis in section 4 offers some important insight in where the centered or non-centered parameterizations of such networks are more appropriate.

3.4. A differentiable MC likelihood estimator

We showed that many hierarchical continuous latent-variable models can be transformed into a DNCP $p_{\theta}(\mathbf{x}, \epsilon)$, where all latent variables (the introduced auxiliary variables ϵ) are root nodes (see eq. (8)). This has an important implication for learning since (contrary to a CP) the DNCP can be used to form a differentiable Monte Carlo estimator of the marginal likelihood:

$$\log p_{\theta}(\mathbf{x}) \simeq \log \frac{1}{L} \sum_{l=1}^L \prod_j p_{\theta}(\mathbf{x}_j | \mathbf{pa}_j^{(l)})$$

where the parents $\mathbf{pa}_j^{(l)}$ of the observed variables are either root nodes or functions of root nodes whose values are sampled from their marginal: $\epsilon^{(l)} \sim p(\epsilon)$. This MC estimator can be differentiated w.r.t. θ to obtain an MC estimate of the log-likelihood gradient $\nabla_{\theta} \log p_{\theta}(\mathbf{x})$, which can be plugged into stochastic optimization methods such as Adagrad for approximate ML or MAP. When performed one datapoint at a time, we arrive at our on-line Maximum Monte Carlo Likelihood (MMCL) algorithm.

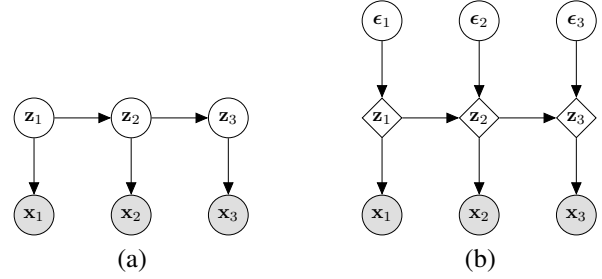


Figure 2. (a) An illustrative hierarchical model in its centered parameterization (CP). (b) The differentiable non-centered parameterization (DNCP), where $z_1 = g_1(\epsilon_1, \theta)$, $z_2 = g_2(z_1, \epsilon_2, \theta)$ and $z_3 = g_3(z_2, \epsilon_3, \theta)$, with auxiliary latent variables $\epsilon_k \sim p_{\theta}(\epsilon_k)$. The DNCP exposes a neural network within the hierarchical model, which we can differentiate efficiently using back-propagation.

Table 1. Limiting behaviour of squared correlations between z and its parent y_i when z is in the centered (CP) and non-centered (DNCP) parameterization.

	$\rho_{y_i, z}^2$ (CP)	$\rho_{y_i, e}^2$ (DNCP)
$\lim_{\sigma \rightarrow 0}$	1	0
$\lim_{\sigma \rightarrow +\infty}$	0	$\frac{\beta w_i^2}{\beta w_i^2 + \alpha}$
$\lim_{\beta \rightarrow 0}$	$\frac{w_i^2}{w_i^2 - \alpha \sigma^2}$	0
$\lim_{\beta \rightarrow -\infty}$	0	1
$\lim_{\alpha \rightarrow 0}$	$\frac{1}{1 - \beta \sigma^2}$	$\frac{\beta \sigma^2}{\beta \sigma^2 - 1}$
$\lim_{\alpha \rightarrow -\infty}$	0	0

4. Effects of parameterizations on posterior dependencies

What is the effect of the proposed reparameterization on the efficiency of inference? If the latent variables have linear-Gaussian conditional distributions, we can use the metric of squared correlation between the latent variable and any of its children in their posterior distribution. If after reparameterization the squared correlation is decreased, then in general this will also result in more efficient inference.

For non-linear Gaussian conditional distributions, the log-PDF can be locally approximated as a linear-Gaussian using a second-order Taylor expansion. Results derived for the linear case can therefore also be applied to the non-linear case; the correlation computed using this approximation is a local dependency between the two variables.

Denote by z a scalar latent variable we are going to reparameterize, and by y its parents, where y_i is one of the parents. The log-PDF of the corresponding conditional dis-