
Neural Variational Inference and Learning in Belief Networks: Supplementary Material

Andriy Mnih
Karol Gregor
DeepMind Technologies Ltd.

ANDRIY@DEEPMIND.COM
KAROL@DEEPMIND.COM

A. Algorithm for computing NVIL gradients

Algorithm 1 provides an outline of our implementation of NVIL gradient computation for a minibatch of n randomly chosen training cases. The exponential smoothing factor α used for updating the estimates of the mean c and variance v of the inference network learning signal was set to 0.8 in our experiments.

Algorithm 1 Compute gradient estimates for the model and the inference network

```

 $\Delta\theta \leftarrow 0, \Delta\phi \leftarrow 0, \Delta\psi \leftarrow 0$ 
 $\mathcal{L} \leftarrow 0$ 
{Compute the learning signal and the bound}
for  $i \leftarrow 1$  to  $n$  do
     $x_i \leftarrow$  random training case
    {Sample from the inference model}
     $h_i \sim Q_\phi(h_i|x_i)$ 
    {Compute the unnormalized learning signal}
     $l_i \leftarrow \log P_\theta(x_i, h_i) - \log Q_\phi(h_i|x_i)$ 
    {Add the case contribution to the bound}
     $\mathcal{L} \leftarrow \mathcal{L} + l_i$ 
    {Subtract the input-dependent baseline}
     $l_i \leftarrow l_i - C_\psi(x_i)$ 
end for
{Update the learning signal statistics}
 $c_b \leftarrow \text{mean}(l_1, \dots, l_n)$ 
 $v_b \leftarrow \text{variance}(l_1, \dots, l_n)$ 
 $c \leftarrow \alpha c + (1 - \alpha)c_b$ 
 $v \leftarrow \alpha v + (1 - \alpha)v_b$ 
for  $i \leftarrow 1$  to  $n$  do
     $l_i \leftarrow \frac{l_i - c}{\max(1, \sqrt{v})}$ 
    {Accumulate the model parameter gradient}
     $\Delta\theta \leftarrow \Delta\theta + \nabla_\theta \log P_\theta(x_i, h_i)$ 
    {Accumulate the inference net gradient}
     $\Delta\phi \leftarrow \Delta\phi + l_i \nabla_\phi \log Q_\phi(h_i|x_i)$ 
    {Accumulate the input-dependent baseline gradient}
     $\Delta\psi \leftarrow \Delta\psi + l_i \nabla_\psi C_\psi(x_i)$ 
end for

```

B. Derivation of the inference network gradient

Differentiating the variational lower bound w.r.t. to the inference network parameters gives

$$\begin{aligned}
\nabla_\phi \mathcal{L}(x) &= \nabla_\phi E_Q[\log P_\theta(x, h) - \log Q_\phi(h|x)] \\
&= \nabla_\phi \sum_h Q_\phi(h|x) \log P_\theta(x, h) - \\
&\quad \nabla_\phi \sum_h Q_\phi(h|x) \log Q_\phi(h|x) \\
&= \sum_h \log P_\theta(x, h) \nabla_\phi Q_\phi(h|x) - \\
&\quad \sum_h (\log Q_\phi(h|x) + 1) \nabla_\phi Q_\phi(h|x) \\
&= \sum_h (\log P_\theta(x, h) - \log Q_\phi(h|x)) \nabla_\phi Q_\phi(h|x),
\end{aligned}$$

where we used the fact that $\sum_h \nabla_\phi Q_\phi(h|x) = \nabla_\phi \sum_h Q_\phi(h|x) = \nabla_\phi 1 = 0$. Using the identity $\nabla_\phi Q_\phi(h|x) = Q_\phi(h|x) \nabla_\phi \log Q_\phi(h|x)$, then gives

$$\begin{aligned}
\nabla_\phi \mathcal{L}(x) &= \sum_h (\log P_\theta(x, h) - \log Q_\phi(h|x)) \\
&\quad \times Q_\phi(h|x) \nabla_\phi \log Q_\phi(h|x) \\
&= E_Q [(\log P_\theta(x, h) - \log Q_\phi(h|x)) \nabla_\phi \log Q_\phi(h|x)].
\end{aligned}$$