

which means that

$$E_{p(\mathbf{y})} \left[ (J(\mathbf{y}) - b(\boldsymbol{\omega})) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] = E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] - b(\boldsymbol{\omega}) E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right] \quad (20.66)$$

$$= E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \right]. \quad (20.67)$$

Furthermore, we can obtain the optimal  $b(\boldsymbol{\omega})$  by computing the variance of  $(J(\mathbf{y}) - b(\boldsymbol{\omega})) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}}$  under  $p(\mathbf{y})$  and minimizing with respect to  $b(\boldsymbol{\omega})$ . What we find is that this optimal baseline  $b^*(\boldsymbol{\omega})_i$  is different for each element  $\omega_i$  of the vector  $\boldsymbol{\omega}$ :

$$b^*(\boldsymbol{\omega})_i = \frac{E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \right]}{E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \right]^2}. \quad (20.68)$$

The gradient estimator with respect to  $\omega_i$  then becomes

$$(J(\mathbf{y}) - b(\boldsymbol{\omega})_i) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \quad (20.69)$$

where  $b(\boldsymbol{\omega})_i$  estimates the above  $b^*(\boldsymbol{\omega})_i$ . The estimate  $b$  is usually obtained by adding extra outputs to the neural network and training the new outputs to estimate  $E_{p(\mathbf{y})} [J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i}]$  and  $E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})}{\partial \omega_i} \right]^2$  for each element of  $\boldsymbol{\omega}$ . These extra outputs can be trained with the mean squared error objective, using respectively  $J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i}$  and  $\frac{\partial \log p(\mathbf{y})}{\partial \omega_i}^2$  as targets when  $\mathbf{y}$  is sampled from  $p(\mathbf{y})$ , for a given  $\boldsymbol{\omega}$ . The estimate  $b$  may then be recovered by substituting these estimates into equation 20.68. Mnih and Gregor (2014) preferred to use a single shared output (across all elements  $i$  of  $\boldsymbol{\omega}$ ) trained with the target  $J(\mathbf{y})$ , using as baseline  $b(\boldsymbol{\omega}) \approx E_{p(\mathbf{y})} [J(\mathbf{y})]$ .

Variance reduction methods have been introduced in the reinforcement learning context (Sutton *et al.*, 2000; Weaver and Tao, 2001), generalizing previous work on the case of binary reward by Dayan (1990). See Bengio *et al.* (2013b), Mnih and Gregor (2014), Ba *et al.* (2014), Mnih *et al.* (2014), or Xu *et al.* (2015) for examples of modern uses of the REINFORCE algorithm with reduced variance in the context of deep learning. In addition to the use of an input-dependent baseline  $b(\boldsymbol{\omega})$ , Mnih and Gregor (2014) found that the scale of  $(J(\mathbf{y}) - b(\boldsymbol{\omega}))$  could be adjusted during training by dividing it by its standard deviation estimated by a moving average during training, as a kind of adaptive learning rate, to counter the effect of important variations that occur during the course of training in the