Softmax Policy Gradient

(1)  $\pi_w(a|s) = \dfrac{e^{w_{s,a}/\tau}}{\sum_{a'} e^{w_{s,a'}/\tau}}$

$$\nabla \log(\pi_w(a|s)) = \nabla_w\left[ w_{s,a}/\tau - \log\left(\sum_{a'} e^{w_{s,a'}/\tau}\right)\right]$$

$C_1:\ \nabla_{w_{s,b}}\left[ w_{s,a}/\tau - \log\left(e^{w_{s,a}/\tau}\right)\right] \qquad b = a$

$= \dfrac{1}{\tau} - \nabla_{w_{s,b}} \log\left(\sum_{a'} e^{w_{s,a'}/\tau}\right)$

$= 1/\tau - \dfrac{1}{\sum_{a'} e^{w_{s,a'}/\tau}} \nabla_{w_{s,b}} \left(\sum_{a'} e^{w_{s,a'}/\tau}\right)$

$= 1/\tau - \dfrac{\nabla_{w_{s,b}}\left(e^{w_{s,a'}/\tau}\right)}{\sum_{a'} e^{w_{s,a}/\tau}}$

$= \dfrac{1}{\tau} - \dfrac{e^{w_{s,a}/\tau}\left(\frac{1}{\tau}\right)}{\sum_{a'} e^{w_{s,a}/\tau}}$

for $\tau = 1$,

$= \boxed{1 - \dfrac{e^{w_{s,a}}}{\sum_{a'} e^{w_{s,a'}/\tau}}}$

$C_2:\ \nabla_{w_{s,b}}\left[ w_{s,a}/\tau - \log\left(\sum_{a'} e^{w_{s,a'}/\tau}\right)\right],\quad b \neq a.$

$= 0 - \nabla_{w_{s,b}} \log\left(\sum_{a'} e^{w_{s,a'}/\tau}\right)$

$= -\dfrac{1}{\sum_{a'} e^{w_{s,a'}/\tau}} \nabla_{w_{s,b}}\left(\sum_{a'} e^{w_{s,a'}/\tau}\right)$

$= \boxed{\dfrac{e^{w_{s,b}/\tau}\left(\frac{1}{\tau}\right)}{\sum_{a'} e^{w_{s,a'}/\tau}}}$

$C_3:\ \nabla_{w_{s,b}}\left[ w_{s,a}/\tau - \log\left(\sum_{a'} e^{w_{s,a'}/\tau}\right)\right],\quad s^{\dagger} \neq s,\ b \neq a.$
$= 0.$

Linear Policy Gradient

(2.1) $\pi_\mu(a|s) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(a-\mu_s)^2}{2\sigma^2}}$

$\log \pi_\mu(a|s) = \log\left(\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(a-\mu_s)^2}{2\sigma^2}}\right)$

$= \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \log\left(e^{-\frac{(a-\mu_s)^2}{2\sigma^2}}\right)$

$\nabla_\mu \log \pi_\mu(a|s) = \nabla_\mu\left[\log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \log\left(e^{-\frac{(a-\mu_s)^2}{2\sigma^2}}\right)\right]$

$= \overset{0}{\nabla_\mu \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right)} + \nabla_\mu\left(-\frac{(a-\mu_s)^2}{2\sigma^2}\right)$

$= -\frac{1}{2\sigma^2} \nabla_\mu (a-\mu_s)^2$

$= -\frac{1}{2\sigma^2} (s)(a-\mu_s(s))$

$= \frac{(a-\mu_s)}{\sigma^2} \cdot s \qquad$, where $\mu, s$ are vectors.

Neural Network Gradient Derivation.

(2-3) $x$ = input                  alpha/k/ - denotes learning rate

$z = wx_1 + b_1$                    * (Derivation does not include
$h = ReLu(z)$                        sampling from the current distribution
$\theta = Uh + b_2$                  - Here $J(\theta)$ would denote the unseen
$\hat{y} = Softmax(\theta)$                                $u(S$

$J(\theta) = (CE(y,\hat{y}))(U+1)(k)$ →

Here $CE$ denotes the cross entropy loss function

∴ we maximize the likelihood the estimator is close
to the actual output of $\hat{y}$

$CE = \sum_i y_i \cdot log(\hat{y}_i)$

We assume that $b_1, b_2$ (the biases in our network $= 0$)
$\theta = Uh$ and $z = Wx_1$

We compute the following gradients:

$$\frac{\partial J}{\partial u}, \quad \frac{\partial J}{\partial w}, \quad \frac{\partial J}{\partial x}$$

$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta} \cdot \frac{\partial \theta}{\partial u}$$

$$\frac{\partial J}{\partial \theta} = (\hat{y}-y)(U+1)(k) = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta}$$

$$\frac{\partial \theta}{\partial u} = h$$

$ReLu(x) = max(x,0)$

$$\frac{\partial J}{\partial u} = (\hat{y}-y)h,$$

$$\frac{\partial ReLu(x)}{\partial x} \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \frac{\partial ReLu(w)}{\partial b} = sgn(ReLu(b))$$

$$\frac{\partial \hat{y}}{\partial t} = \frac{\partial ReLu(z)}{\partial z} \cdot sgn(ReLu(z))$$

$z$ is a vector.

$$\frac{\partial J}{\partial w} = \left(\frac{\partial J}{\partial \theta}\right)\left(\frac{\partial \theta}{\partial u}\right)\left(\frac{\partial u}{\partial t}\right)\frac{\partial t}{\partial w}$$

$$= (\hat{y}-y)(U+1)(k)(sgn(h))(k)$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial \theta} \cdot \frac{\partial \theta}{\partial u} \cdot \frac{\partial u}{\partial t} \cdot \frac{\partial t}{\partial x} = (\hat{y}-y)(U+1)(k)(sgn(h))(w)$$