

1) Calculating  $\nabla_w \log \pi_w(A_t | s_t)$

Gradient w.r.t weights  $\Rightarrow s = \begin{bmatrix} a=4 \\ -p(a_j | s_1) \\ -p(a_j | s_2) \\ -p(a_j | s_3) \\ -p(a_j | s_4) \end{bmatrix}$

$$\Rightarrow \nabla_w \left[ \log \left( \frac{e^{w_{st,a_t}/\gamma}}{\sum_{a'} e^{w_{st,a'}/\gamma}} \right) \right]$$

$$\Rightarrow \nabla_w \left[ \underbrace{\log(e^{w_{st,a_t}/\gamma})}_i - \underbrace{\log\left(\sum_{a'} e^{w_{st,a'}/\gamma}\right)}_{ii} \right]$$

While  $w$  is a  $4 \times 4$  matrix, we only have gradients through the specific row (representing state  $s_t$  of our choosing). Therefore:

For some  $i$ :

$$i) \nabla_{w_{st,i}} \log(e^{w_{st,a_t}/\gamma}) = \nabla_{w_{st,i}} \frac{w_{st,a_t}}{\gamma} = \begin{cases} 1/\gamma & i=a_t \\ 0 & i \neq a_t \end{cases}$$

$$ii) \nabla_{w_{st,i}} \log\left(\sum_{a'} e^{w_{st,a'}/\gamma}\right) = \frac{1}{\sum_{a'} e^{w_{st,a'}/\gamma}} \times \frac{\partial \sum_{a'} e^{w_{st,a'}/\gamma}}{\partial w_{st,i}}$$

$$\Rightarrow \frac{e^{w_{st,i}}}{\sum_{a'} e^{w_{st,a'}/\gamma}}$$

$$\therefore \nabla_{w_{s',a'}} = \begin{cases} 1 - \frac{e^{w_{st,a_t}}}{\sum_{a'} e^{w_{st,a'}/\gamma}} & s'=s_t, a'=a_t \\ -\frac{e^{w_{st,a_t}}}{\sum_{a''} e^{w_{st,a''}/\gamma}} & s'=s_t, a' \neq a_t \\ 0 & \text{o.t.w} \end{cases}$$

$$2) \pi_k(a|s) = \frac{1}{\sigma\sqrt{2\pi}} \times e^{-\frac{(a-ks)^2}{2\sigma^2}}$$

$$\log \pi_k(a|s) = \log\left(\frac{1}{\sigma\sqrt{2\pi}} \times e^{-\frac{(a-ks)^2}{2\sigma^2}}\right)$$

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

$$\therefore \nabla_k \log \pi_k(a|s) = \frac{1}{2\sigma^2} \nabla_k (a-ks)^2$$

$$\Rightarrow \frac{a-ks}{\sigma^2} \times s$$

## 2.1 Linear Policy

$$N(s) = k s \Rightarrow s = \begin{bmatrix} x\text{-pos} \\ \text{velocity} \end{bmatrix} \quad k = \begin{bmatrix} 1 \\ k_1 \\ k_2 \end{bmatrix}$$

Mean is a linear combination of features  $k_1 x_{\text{pos}} + k_2 \text{velocity}$

$$i) \nabla$$