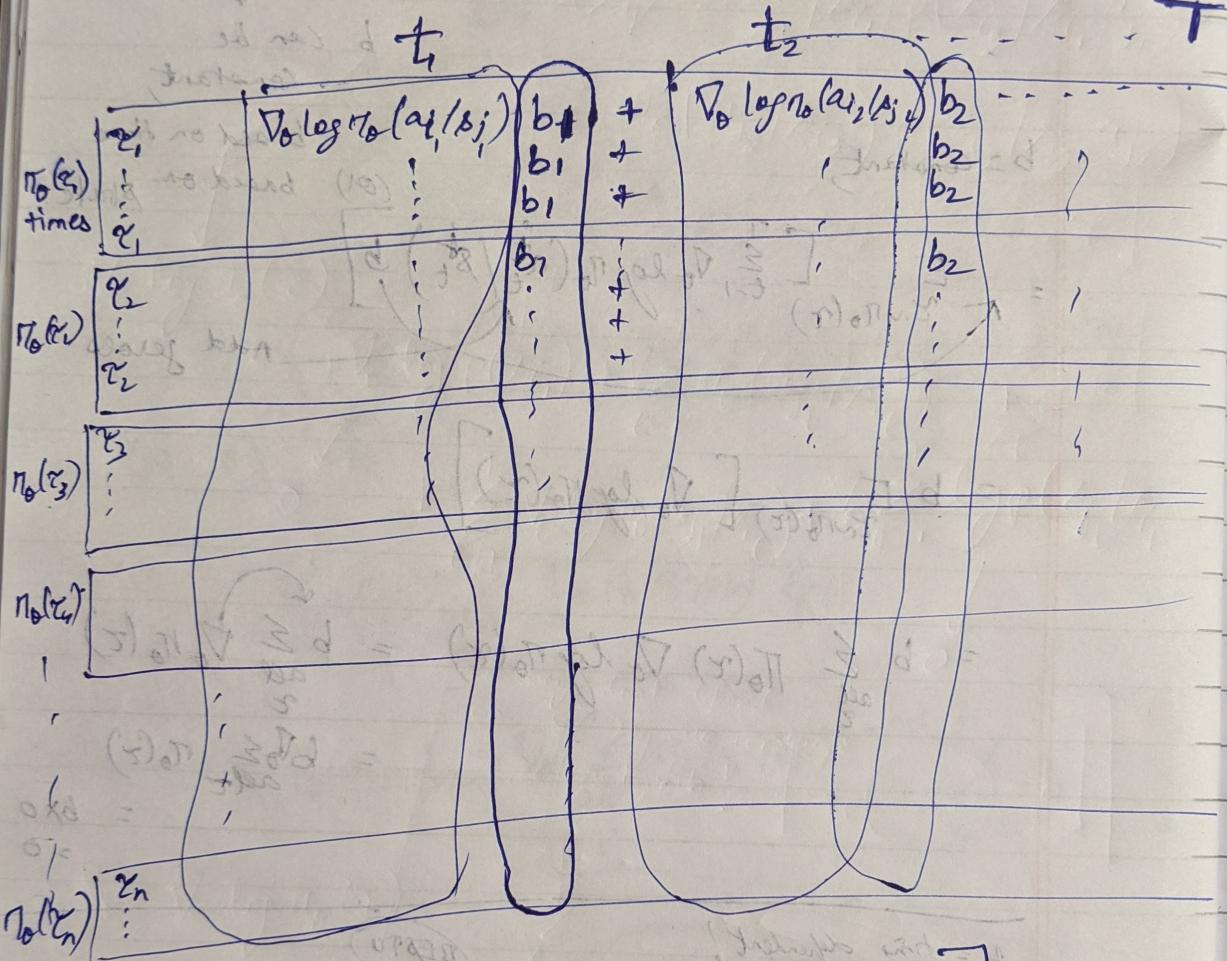


Baseline as function of time

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}(z)} \left[\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) [\tilde{Q}(a_t^*, s_t^*) - b_t] \right]$$

↓
Visualize the double sum for only
the second term



$$\sum_{t=1}^T b_t \left[\sum_{\text{all } a_i, a_j} p_{\theta}(a_i, a_j) \nabla_{\theta} \log \pi_{\theta}(a_j | s_i) \right]$$

$$\sum_{t=1}^T b_t \left[\sum_{\text{all } a_i, a_j} p_{\theta}(a_i) \frac{\nabla_{\theta} (a_j | s_i) \nabla_{\theta} \log \pi_{\theta}(a_j | s_i)}{\nabla_{\theta} \pi_{\theta}(a_j | s_i)} \right]$$

$$\sum_{t=1}^T b_t \sum_{\text{all } a_i, a_j} p_{\theta}(a_i) \nabla_{\theta} \pi_{\theta}(a_j | s_i)$$

$$\sum_{t=1}^T b_t \sum_{\text{all } a_i} p_{\theta}(a_i) \left[\sum_{\text{all } a_j} \nabla_{\theta} \pi_{\theta}(a_j | s_i) \right]$$

Taking p_{θ} common!

$$\nabla_{\theta} J(\theta) = 0$$

softmax & NN

= 0 !

Baseline as function of [both] state + time

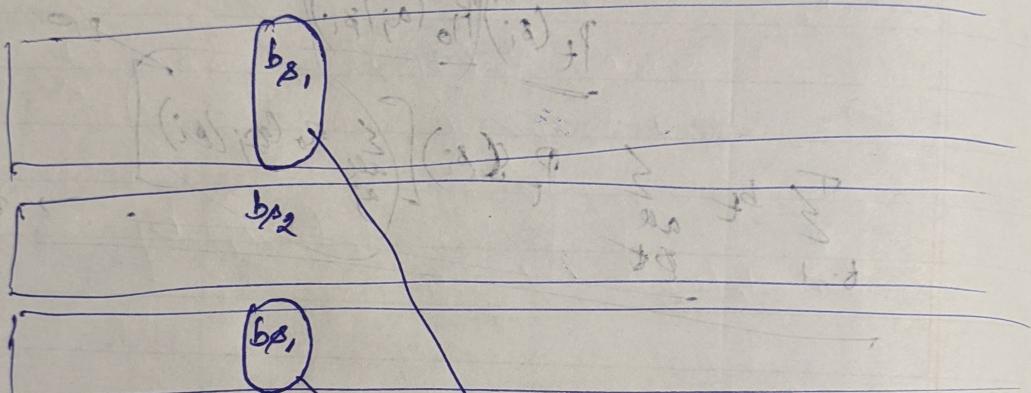
$$E_{\text{softmax}}(z) \left[\sum_{t=1}^T \nabla_b \log \pi_0(a_i | s_j) [Q(a_i, s_j) - b_{s,t}] \right]$$

visualizing double sum

for only the 2nd term

t_1

s_1



grouped together!
under each timestep!

~~$\sum_{t=1}^T \sum_{all \beta_j} b_{s,t} P_0(s_j | a_i)$~~

$$\sum_{t=1}^T \sum_{all \beta_j} b_{s,j,t} \sum_{all a_i} \uparrow \nabla_b \log \pi_0(a_i | s_j) \times$$

$$= \sum_{t=1}^T \sum_{all \beta_j} b_{s,j,t} \sum_{all a_i} \underbrace{P_0(s_j)}_{\pi_0(a_i | s_j)} \left[\nabla_b \log \pi_0(a_i | s_j) \nabla_b \log \pi_0(a_i | s_j) \right]$$

$$= \sum_{t=1}^T \sum_{all \beta_j} b_{s,j,t} P_0(s_j) \underbrace{\sum_{all a_i} \nabla_b \log \pi_0(a_i | s_j)}_{\substack{\text{sum of softmax of} \\ \text{a NN} = 1}} \rightarrow$$

$$= \boxed{0} !$$

$$\nabla_b(1) = 0$$