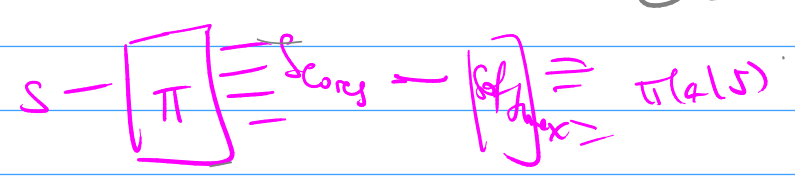


Policy gradient
 $\pi(a|s, \theta)$

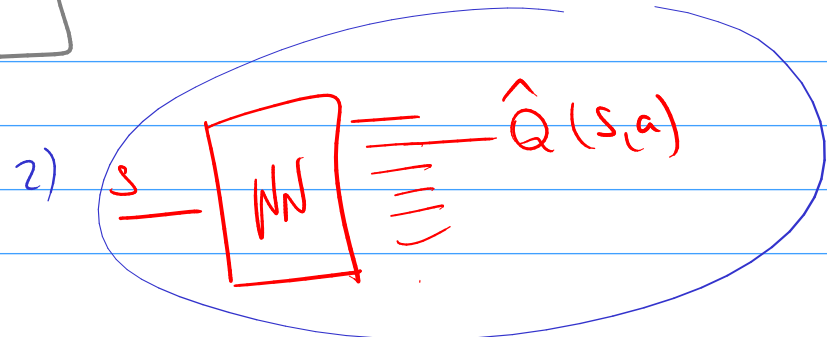
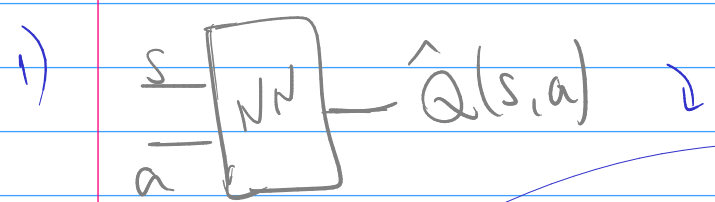


$$J(\theta) = V_{\pi_\theta}(s_0) \xrightarrow{\theta} \max$$

$$\nabla_\theta J(\theta) = \nabla_\theta [E_{\pi_\theta} [G_0]] = \dots = \text{Policy gradient theorem}$$

$$\nabla_\theta J(\theta) \propto E_\pi [\cdot G_t \cdot \nabla_\theta [\ln \pi_\theta(A_t | S_t)]]$$

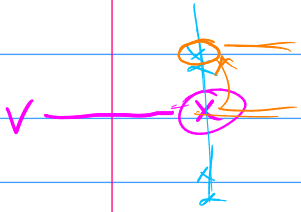
π: — — — — — REINFORCE



$$\bar{\theta} = \bar{\theta} + \alpha \left(z + \gamma \cdot \max_{a'} \underbrace{Q(s', a')}_{\theta} - \underbrace{Q(s, a)}_{\theta} \right) \cdot \nabla_{\theta} Q$$

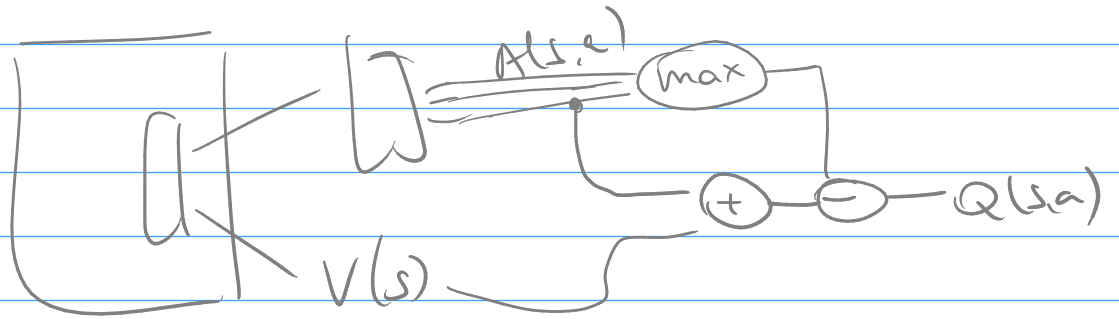
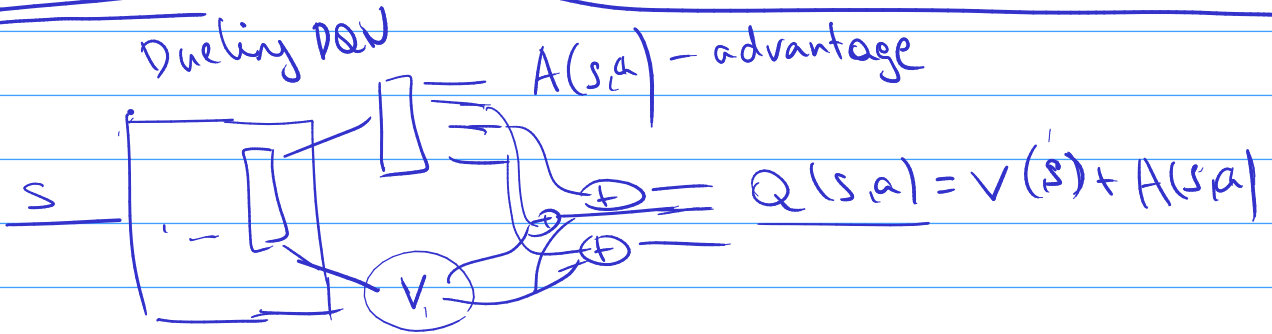
Winner's curse

$$\hat{Q}(s', a') \approx Q(s', a')$$



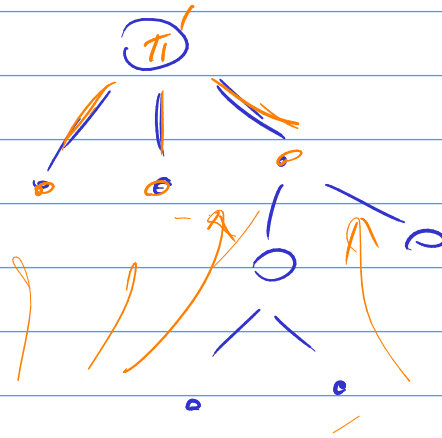
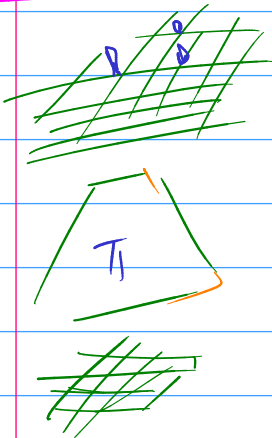
$$v_i \sim \mathcal{N}(v_i | v, \sigma_i^2)$$

3)



$$Q(s, a) = V(s) + (A(s, a) - \max_{a'} A(s, a'))$$

$$= \frac{1}{|A|} \sum A(s, a')$$



$$\pi \rightarrow \pi' > \pi$$

$$\bar{p}_t \approx \pi(a_t | \text{history})$$

$$v_t \approx V(s | \text{history})$$

$$z_t \approx R_{t+1}$$

