Bellman equation : $Q(s,a) = r_t + \gamma \max_{a'} Q(s',a')$

actions — (above $a$)
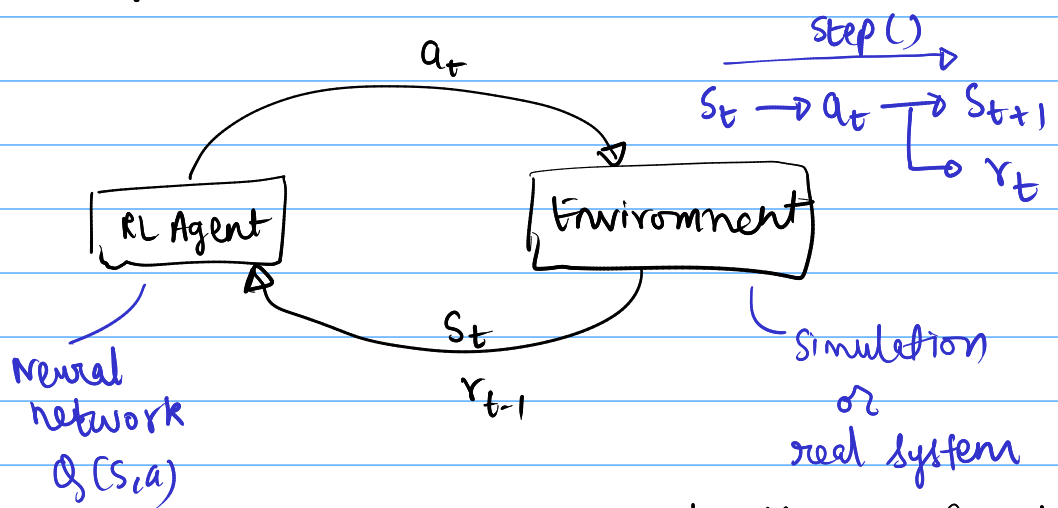
$0 \leq \gamma \leq 1$

current state — (below $s$)

next step reward — (below $r_t$)

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots$$

$\underbrace{\phantom{\gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots}}_{G_{t+1}}$

$$= r_t + \gamma G_{t+1}$$

$$\begin{array}{c c c c c} & a_0 & a_1 & \cdots & a_m \\ s_0 & \boxed{Q(s_0, a_0)} & & & \\ s_1 & & & & \\ s_2 & & & Q(s_k, a_l) & \\ \vdots & & & & \\ s_n & & & & \end{array}$$

① Scale is large $s, a$

② Convergence issues because of lack of time

$a_t$

$\text{step}()$

$s_t \to a_t \to \begin{cases} s_{t+1} \\ r_t \end{cases}$

RL Agent

Environment
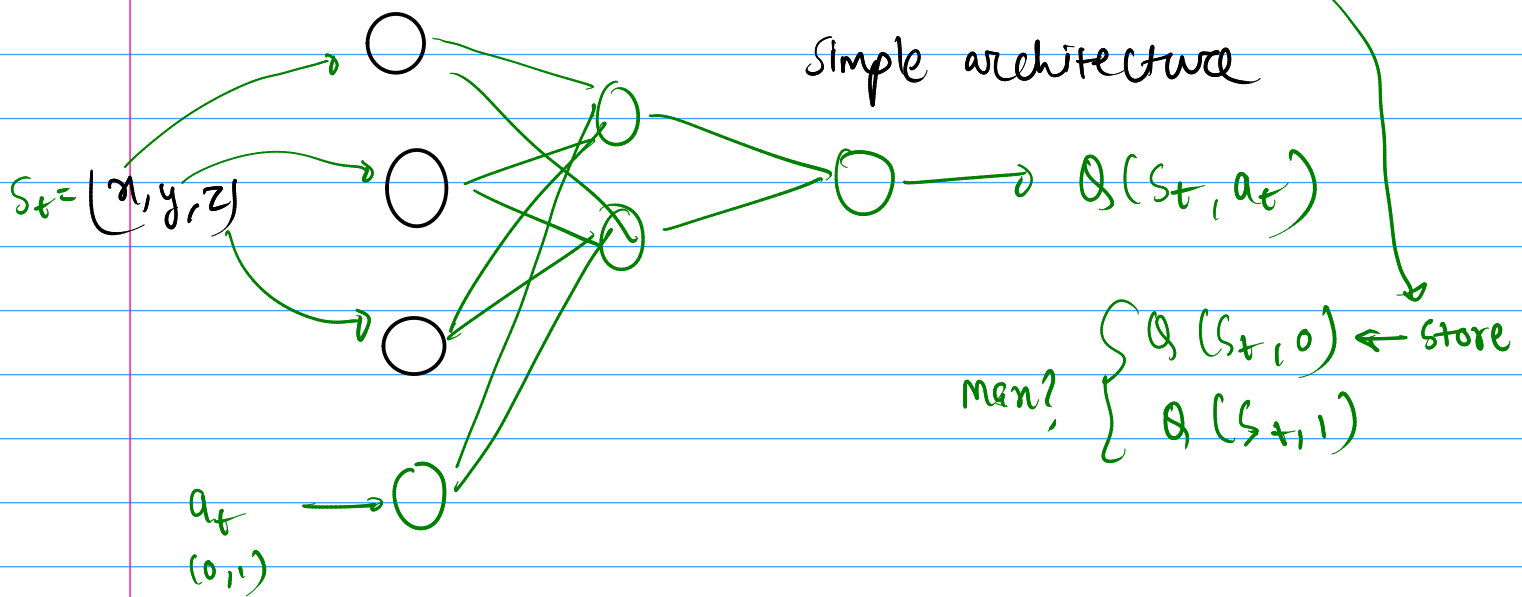
$s_t$

$r_{t-1}$

Neural network $Q(s, a)$

simulation or real system

"sample" : $[s_t, a_t, r_t, s_{t+1}] \; \forall t$

keep them → off-policy

throw away → on-policy
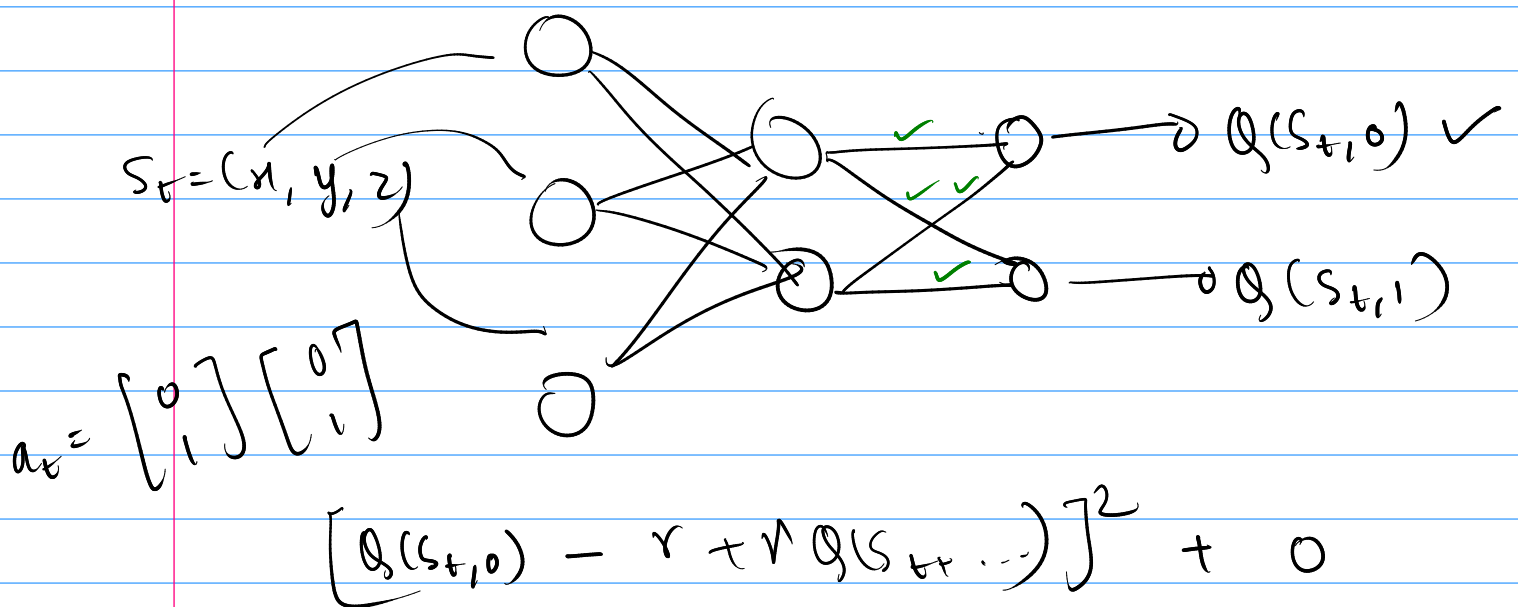
$$Q(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

Predictions or inference

target

$$\text{error} = \left[ \overset{b}{Q}(s_{t}, a_{t}) - r_{t} - \gamma \max Q(s_{t+1}, a_{t+1}) \right]^{2}$$
MSE

Simple architecture

$s_t = (x, y, z)$

$a_t$
$(0, 1)$

$Q(s_t, a_t)$

Man?  $\begin{cases} Q(s_t, 0) \leftarrow \text{store} \\ Q(s_t, 1) \end{cases}$

Practical architecture

$s_t = (x, y, z)$

$a_t = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$
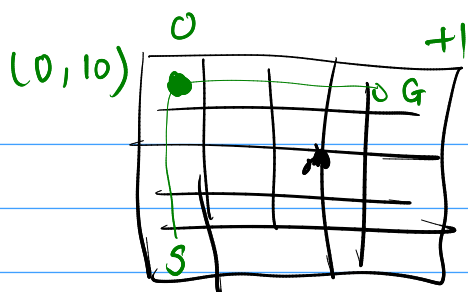
$Q(s_t, 0)$ ✓

$Q(s_t, 1)$

$$\left[ Q(s_t, 0) - r + \gamma Q(s_{t+\cdots}) \right]^{2} + 0$$

"Vanilla DQN"

Issues
① Bootstrapping
② Stabilisation
③ $\Delta Q(s, a)$ implemented only by last layer
④ $\partial b/$ on policy decisions

$(0,10)$    $0$        $+1$

$(x,y)$    $0$    $Q(s_{t+1}, a_{t+1})$

$(s_t, a_t, r_t, s_{t+1})$

$$\begin{bmatrix} s_t, (10,10), a_t, r_t, s_{t+1} \\ \vdots \\ \vdots \end{bmatrix} 0$$

$$\begin{bmatrix} \boxed{s_t, (0,10),} a_t, r_t, s_{t+1} \\ \vdots \\ (0,10) \quad (0,10), a_t, 1 \quad, s_{t+1} \end{bmatrix}$$

within agent

① choose action $\longrightarrow$    State as input
                            compute $Q(s,a)$ by
                                   inference on NN
                          $\epsilon$-greedy action

② Save to memory $\longrightarrow$    $s_t, a_t, r_t, s_{t+1}$

③ train $\longrightarrow$    sample from memory
                       compute loss
                       update NN

$$g(s,a) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots$$
$$= \underbrace{r_t + \gamma r_{t+1}}_{} + \gamma^2 \underbrace{[r_{t+2} \rightarrow \cdots}_{g_{t+2}}$$