

▶ Supervised Learning: $\langle x_1, y_1 \rangle \dots \langle x_N, y_N \rangle \rightarrow$ feedback immediate

Unsupervised Learning: $\langle x_1, * \rangle \dots \langle x_N, * \rangle \rightarrow$ NO feedback

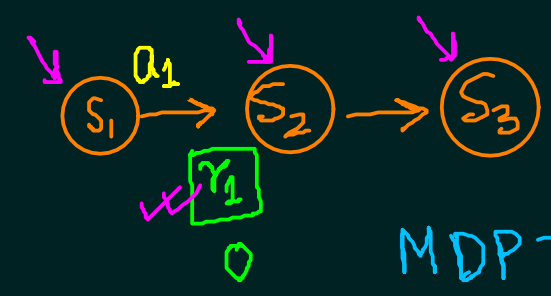
Training Data $\langle x_1, - \rangle \dots \langle W, L, D, \text{feedback delayed} \rangle$ } Ex: Chess
 +1, -1, 0 } Win, Loss, D

Reinforcement Learning

Sequence of steps $\rightarrow \rightarrow \rightarrow \dots \rightarrow$ Goal $\begin{matrix} \nearrow +1 \\ \searrow -1 \end{matrix}$

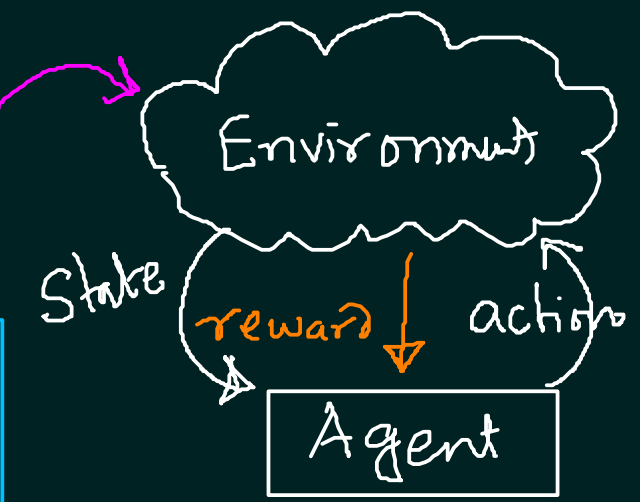
(policy)

$f: X \rightarrow Y$
 $h \approx f$
TE for SL
 $\langle x_1, y_1 \rangle \dots \langle x_N, y_N \rangle$



$$P(S_t | S_{t-1}, a_{t-1})$$

$$P(r_t | S_{t-1}, a_{t-1})$$



$\langle s_1, a_1 \rangle \rightarrow r_1$

$\langle s_1, a_1 \rangle \rightarrow \dots \rightarrow \langle s_N, a_N \rangle \rightarrow r_N$
TE for RL

Learn a policy
 $\pi: S \rightarrow A$

Challenges: ① outcome of action uncertain

Ex: Backgammon 

② Perfect sense of environment

Noise → unknown
Non-det → changing

③ Reward is delayed

④ Reward is stochastic.

⑤ How much do you train?

Task $\Pi: S \rightarrow A$

Cumulative Reward

$$V^\Pi(s_i) = r_i + \gamma V^\Pi(s_{i+1})$$

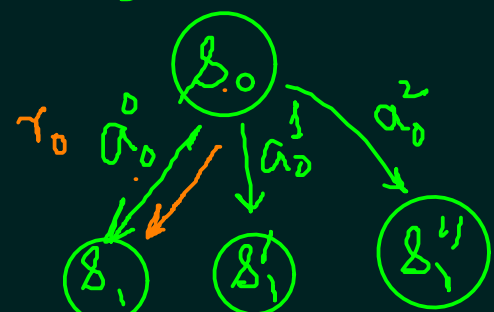
↑
value function

$$\gamma^2 r_{i+2} + \dots$$

$$= \sum_{i=0}^{\infty} \gamma^i r_i$$

$\Pi \rightarrow (a_0 \rightarrow a_1 \rightarrow \dots \rightarrow a_n)$

How to design rewards?



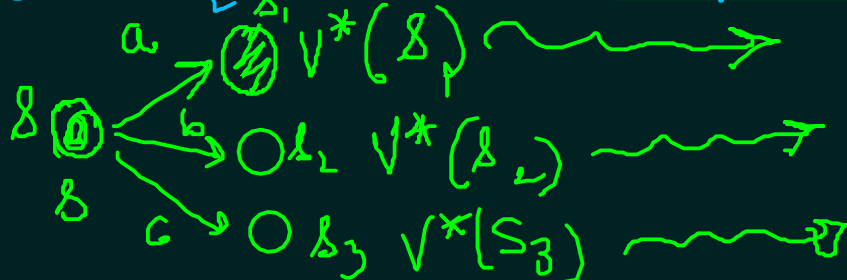
$$\Pi^*(s) =$$

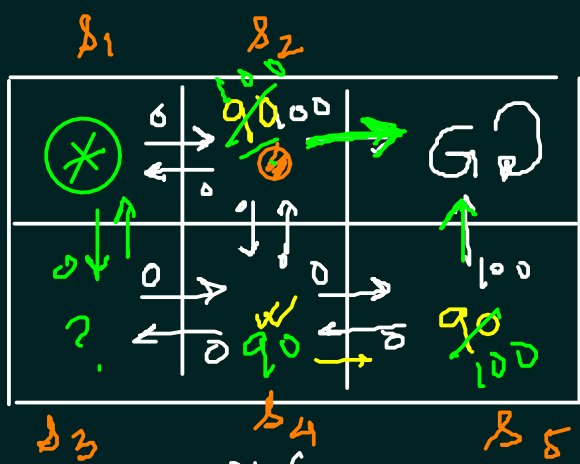
Optimal policy

$$V^*(s)$$

goal

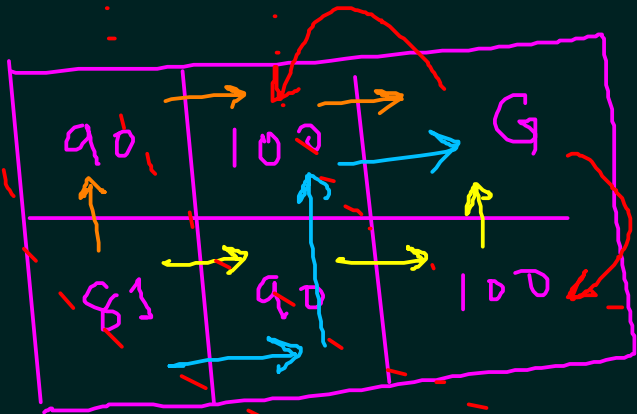
$$\operatorname{argmax}_a \left[r(s,a) + \gamma V^*(s(s,a)) \right]$$





$$V^*(s)$$

$$s \xrightarrow{a} s' \quad r$$



OR

Alt. Best policy

$$V^*(s)$$

$r(s, a)$ } all known
 $\delta(s, a)$ }

$$V^*(s_2) = 100$$

(Value Iteration Policy Iteration)

$$V^*(s_4) = 0 + 0.9 \times \max\{0, 100, 0\} = 90$$

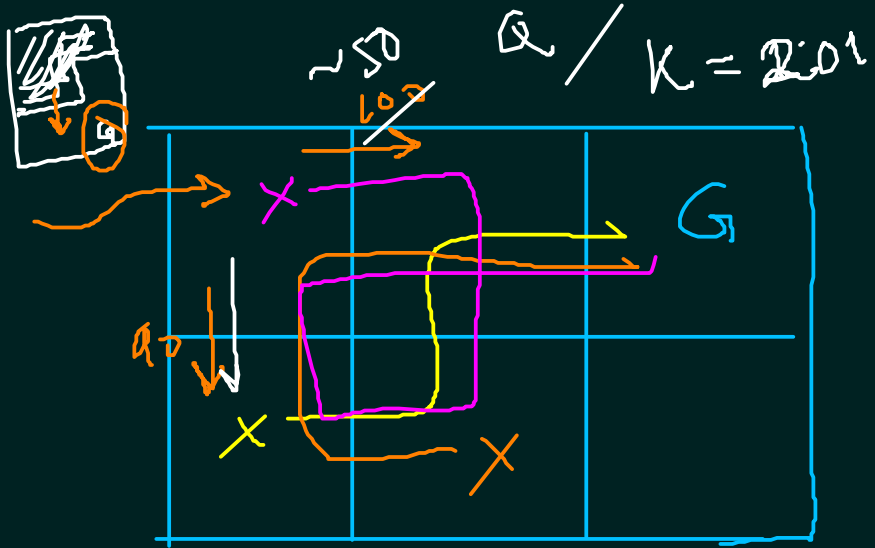
$$r(s_2, a) \quad V^*(s_1) = 0 + 0.9 \times 100 = 90$$

$$V^*(s_3) = 0 + 0.9 \times 90 = 81$$

Learn

$$\langle s_t, r_t \rangle \rightarrow$$





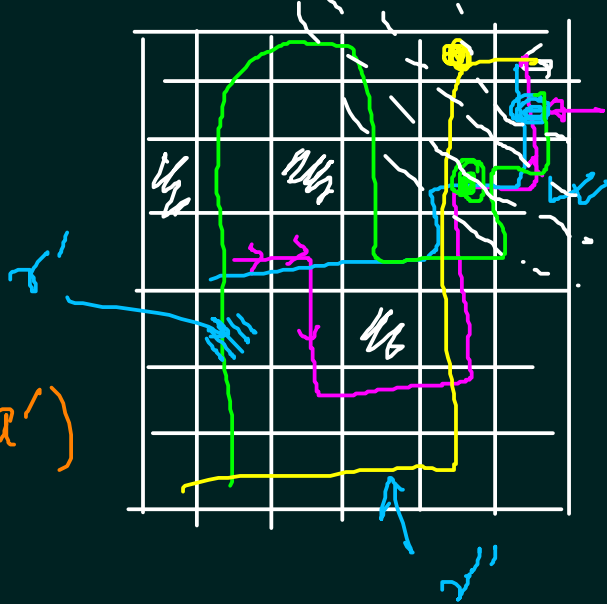
- Convergence \rightarrow
- ① det. Act
 - ② $|r(s,a)| < c$
 - ③ Infinitely often visit
- $\hat{Q}(s,a) \forall s,a$ \rightarrow $Q(s,a)$

locally optimal policy \rightarrow π^*_{NN}

Theroem:

$$\hat{Q}(s,a) = r(s,a) + \gamma V^*(\delta(s,a))$$

$$= r(s,a) + \gamma \max_{a'} \hat{Q}(\delta(s,a), a')$$



Dynamic Prog.

$$\hat{Q}(s,a) = r(s,a) + \gamma \max_{a'} (Q(s', a'))$$

$\delta(s,a) = s'$

