# Chapter 3: Markov decision processes (MDP)
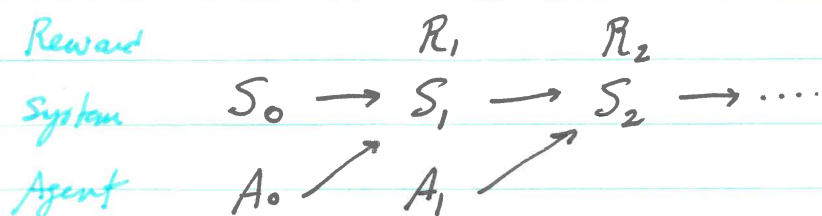
## 3.1. Model

- Environment state : $S_t \in S$ $\qquad t = 0, 1, \dots$
    - System to control
    - Info for making decision

- Agent state : $A_t \in A(s)$ $\qquad t = 0, 1, \dots$
    - Controller
    - Action state / decision $\qquad$ available actions from given state
    - State space can depend on current state
    - Simplification: $A(s) = A \quad \forall s$

- Reward : $R_t \in R$
    - Signal from system / environment
    - Defines good / bad actions
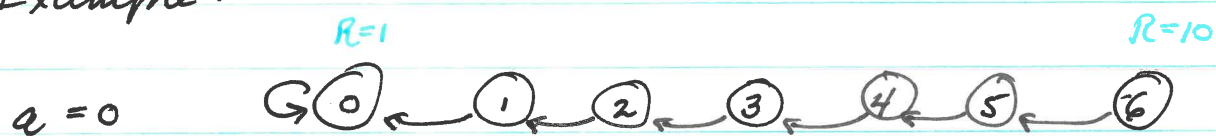      $\qquad\quad +\qquad -$

    - Zero / baseline not important
    - RV correlated with states / actions
    - Sometimes expressed as $r(s, a, s')$

- Transition probability (dynamics):

$$p(s', r \mid s, a) = P(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$$

| Reward | | $R_1$ | | $R_2$ | | |
|---|---|---|---|---|---|---|
| System | $S_0$ | $\rightarrow$ | $S_1$ | $\rightarrow$ | $S_2$ | $\rightarrow \dots$ |
| Agent | $A_0$ | $\nearrow$ | $A_1$ | $\nearrow$ | | |

- Assumed homogeneous (time-independent, stationary)

· Example:



- 2 deterministic actions: left, right
- Deterministic reward given state · and doesn't depend on action (same for a=0,1)

· State transition probability:

$$p(s'|s,a) = \sum_r p(s',r|s,a)$$
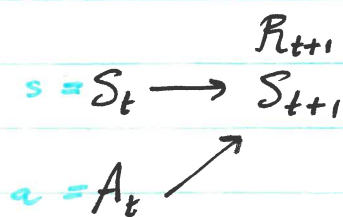
Transition matrix for each a
$P_a(j|i)$   $P_a$

· Reward probability:

$$p(r|s,a) = \sum_{s'} p(s',r|s,a)$$

compare with MRP

$\uparrow$ action

· Expected reward:   $\rho(s,a) = E[R_{t+1}|S_t=s, A_t=a]$

$$s = S_t \xrightarrow{R_{t+1}} S_{t+1}$$
$$a = A_t \nearrow$$
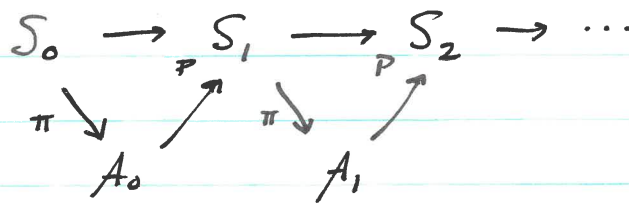
$$\rho(s,a) = \sum_r r\, p(r|s,a)$$

$$= \sum_{r,s'} r\, p(s',r|s,a)$$

- One-step reward from state s and action a
- Doesn't depend on time (stationary)

## 3.2. Policy

· How to choose action/control/decision at given state?

· Policy transition probability:

$$\pi(a|s) = P(A_t = a | S_t = s)$$

$\uparrow$ action $\qquad$ $\uparrow$ state



$S_0 \xrightarrow{} S_1 \xrightarrow{} S_2 \rightarrow \cdots$

$\pi \downarrow \quad A_0 \quad \pi \downarrow \quad A_1$

} noise in environment
} exploration

· Mapping state → action
· Probabilistic policy: actions ~~can be~~ are random
· Assumed homogeneous (time-independent, stationary)

· Deterministic policy: $\pi(s) = a$

one state $\qquad$ one action

· Joint probability: $p(s', r, a | s) = p(s', r | s, a) \, \pi(a|s)$

· System transition probability:

$$P_\pi(s'|s) = \sum_a p(s'|s,a) \pi(a|s) \qquad \text{Average over actions}$$

$$= \sum_{a,r} p(s', r | s, a) \pi(a|s)$$

· Effective dynamics under policy $\pi$
· Transition matrix: $P_\pi \qquad (P_\pi)_{ij} = P_\pi(j|i)$

$$S_0 \xrightarrow{P_\pi} S_1 \xrightarrow{P_\pi} S_2 \xrightarrow{P_\pi} \cdots$$

· Comparison :

- No control: $S_0 \xrightarrow{P} S_1 \xrightarrow{P} S_2 \to \dots$  $p(s'|s)$

- Open-loop control:

$S_0 \xrightarrow{P} S_1 \xrightarrow{P} S_2 \to \dots$  $p(s'|s,a)$

$A_0 \nearrow \quad A_1 \nearrow$

- Closed-loop control: feedback

$S_0 \xrightarrow{P} S_1 \xrightarrow{P} S_2 \to \dots$  $p(s'|s,a)$
$\quad \pi \downarrow \nearrow \quad \pi \downarrow \nearrow \quad$  $\pi(a|s)$
$\quad A_0 \quad A_1$

- Reduced dynamics: $S_0 \xrightarrow{P_\pi} S_1 \xrightarrow{P_\pi} S_2 \xrightarrow{P_\pi} \dots$  $p_\pi(s'|s)$


· Expected reward under policy $\pi$ :

$$\rho_\pi(s) = E_\pi\left[R_{t+1} \mid S_t = s\right]$$

$$= E_\pi \, E\left[R_{t+1} \mid S_t = s, A_t\right]$$

$$= E_\pi\left[\rho(s, A_t)\right]$$

$$\rho_\pi(s) = \sum_a \rho(s, a)\, \pi(a|s) \qquad \textcolor{teal}{average over actions}$$

$$= \sum_{a,r} r\, p(r|s,a)\, \pi(a|s)$$

$$= \sum_{a,r,s'} r\, p(s', r|s, a)\, \pi(a|s)$$

- One-time/step reward from state $s$ under policy $\pi$
- Doesn't depend on time (stationary)
- $\rho(s,a)$ doesn't depend on $\pi$
- Deterministic: $\rho_\pi(s) = \rho(s, \pi(s))$
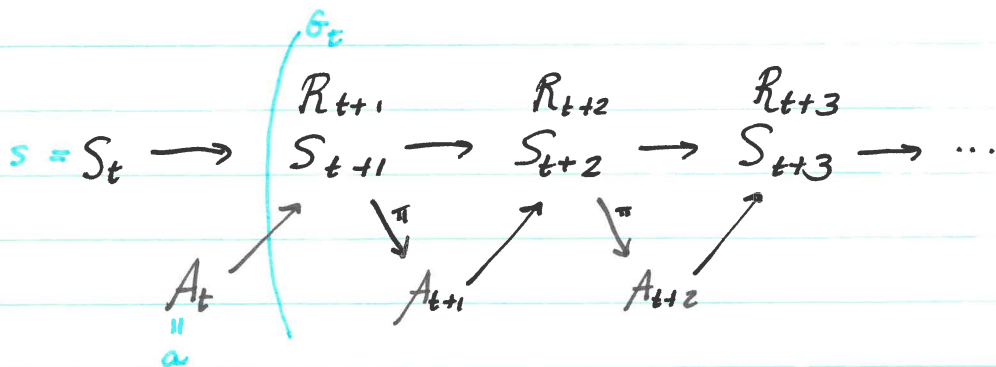
## 3.3. Value functions

- State value function:

$$v_\pi(s) = E_\pi\left[G_t \mid S_t = s\right]$$

return



- Future return / cost to go from state $s$ + policy $\pi$
- Long term value of $s$ under $\pi$
- Function of $\pi(\cdot \mid s)$
- Doesn't depend on time $t$ ( stationary dynamics / infinite return )

- State-action value function:

$$q_\pi(s,a) = E_\pi\left[G_t \mid S_t = s, A_t = a\right]$$



- Return / cost to go if action $a$ taken from state $s$
- Value of action $a$ when in $s$ if policy $\pi$ followed after
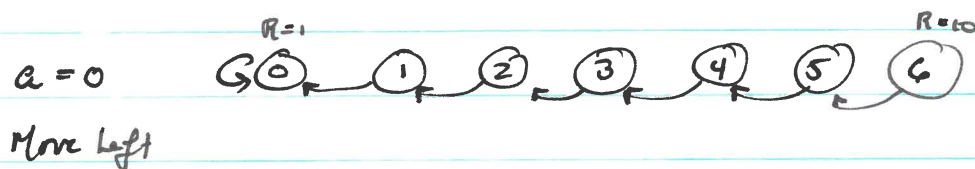- Doesn't depend on time $t$

· Connection:

$$V_\pi(s) = E_\pi \left[ q_\pi(s, A_t) \mid S_t = s \right] \quad \text{average over actions}$$
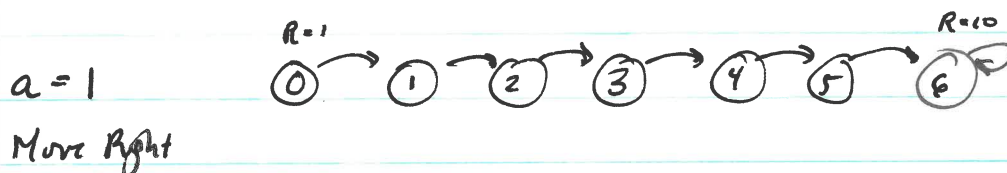
$$= \sum_a q_\pi(s, a)\, \pi(a/s)$$

$$V_\pi(s) = q_\pi(s, \pi(s)) \quad \text{for deterministic policy}$$
$$a = \pi(s)$$

· Example:

$a = 0$
More Left

$R=1$ ... $R=10$ $\quad P(s'|s, a=0)$

$a = 1$
More Right

$R=1$ ... $R=10$ $\quad P(s'|s, a=1)$

· Policy 1:  $\pi_1(s) = 0 \quad \forall s$  move left from all state   deterministic

· Policy 2:  $\pi_2(s) = 1 \quad \forall s$  moves right   "

· Policy 3:  $\pi_3 = \frac{1}{2}\pi_1 + \frac{1}{2}\pi_2$   random policy/control

CW: Calculate $V_{\pi_1}$, $V_{\pi_2}$, $V_{\pi_3}$

· Note:  $\rho_\pi(s) = 1$ or $10 \quad s=0, s=6 \quad \forall \pi$   reward only depends on $s$

$\quad\quad\quad \rho(s,a) = 1$ or $10 \quad s=0, s=6 \quad \forall a$

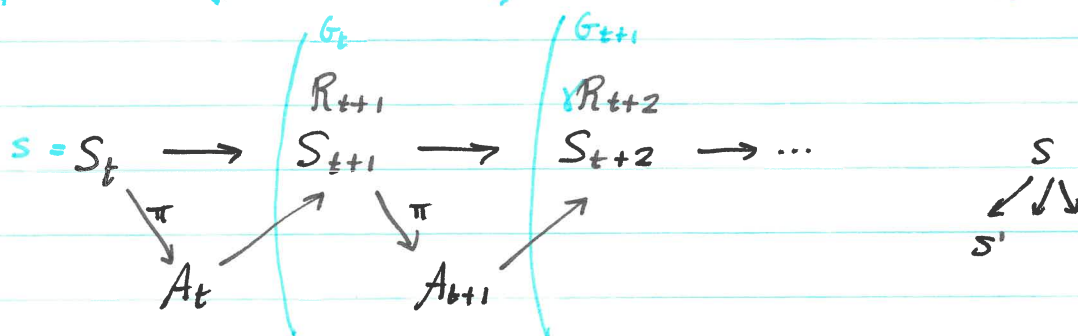· Note:  2 actions per state $\Rightarrow 2^7$ possible policies

## 3.4. Bellman equations (BE)

· Return : $\quad G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = R_{t+1} + \gamma G_{t+1}$

· Bellman equation for Value function :

$$V_\pi(s) = E_\pi \left[ R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s \right] \qquad \text{Compare with} \\ \text{MRP}$$

$$s = S_t \longrightarrow S_{t+1} \longrightarrow S_{t+2} \longrightarrow \cdots$$

$$V_\pi(s) = \sum_{s',a} p(s,a)\, \pi(a\mid s) + \gamma \sum_{s',a} V_\pi(s')\, p(s'\mid s,a)\, \pi(a\mid s)$$

$$= \sum_{s',r,a} r\, p(s',r\mid s,a)\, \pi(a\mid s)$$

$$+ \gamma \sum_{s',r,a} V_\pi(s')\, p(s',r\mid s,a)\, \pi(a\mid s)$$

· Value = immediate reward + discounted Value of
  $\quad\; s \qquad\qquad\qquad s \qquad\qquad\qquad$ successor state
  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad s'$

· Linear equation: $\quad V_\pi = \rho_\pi + \gamma P_\pi V_\pi$

· Solution : $\quad V_\pi = (I - \gamma P_\pi)^{-1} \rho_\pi$

· BE for action-value function:

$$q_\pi(s,a) = E_\pi\left[ R_{t+1} + \gamma\, q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s,\ A_t = a \right]$$

$$s = S_t \longrightarrow \overset{\overset{\displaystyle G_t}{\overset{\big|}{R_{t+1}}}}{S_{t+1}} \longrightarrow \overset{\overset{\displaystyle G_{t+1}}{\overset{\big|}{\gamma R_{t+2}}}}{S_{t+2}} \longrightarrow \cdots$$

$$\overset{A_t}{\underset{\overset{\|}{a}}{\nearrow}} \qquad \overset{\pi\searrow}{A_{t+1}}\nearrow$$

$$s,a$$
$$\downarrow$$
$$S_{t+1}, A_{t+1}$$
$$s' \underset{\pi}{\curvearrowright} a'$$

$$q_\pi(s,a) \qquad\qquad q_\pi(S_{t+1}, A_{t+1})$$

$$q_\pi(s,a) = \rho(s,a) + \gamma \sum_{a',s'} q_\pi(s',a')\, p(s'|s,a)\, \pi(a'|s')$$

$$= \sum_{s',r} r\, p(s',r|s,a)$$

$$\qquad\qquad + \gamma \sum_{a',s',r} q_\pi(s',a')\, p(s',r|s,a)\, \pi(a'|s')$$

· Value = immediate reward + discounted value of
  $(s,a)$        $(s,a)$      successor state, action
                                            $(s', a')$

· Linear equation for "matrix" $q_\pi$

· Two ways of defining MDPs

① ②

$p(s', r \mid s, a)$

$\qquad \hookrightarrow p(r \mid s, a)$

$\qquad\qquad \hookrightarrow p(s, a) = E[R_{t+1} \mid S_t = s, A_t = a]$   $\rho$

$\qquad \hookrightarrow p(s' \mid s, a)$

$\qquad\qquad \hookrightarrow P_a$   $P_a$

$\pi(a \mid s)$

$\qquad \hookrightarrow \rho_\pi(s) = E_\pi[R_{t+1} \mid S_t = s]$   $\rho_\pi$
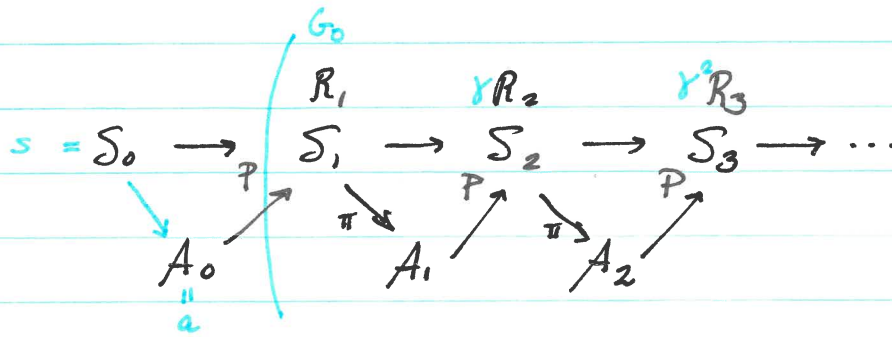
$\hookrightarrow p_\pi(s' \mid s)$

$\qquad \hookrightarrow P_\pi$   $P_\pi$

$$V_\pi = \rho_\pi + \gamma P_\pi V_\pi$$

$$q_\pi = \rho + \gamma (P_a \pi \, q_\pi)$$

## 3.5 Optimal policies

$$s = S_0 \xrightarrow[P]{} S_1 \xrightarrow[P]{R_1} S_2 \xrightarrow[P]{\gamma R_2} S_3 \xrightarrow{\gamma^2 R_3} \cdots$$

$G_0$

$A_0 \quad A_1 \quad A_2$

$a$

| Dynamics | Policy /control | Reward |
|---|---|---|
| $p(s'\mid s,a)$ | $\pi(a\mid s)$ | $p(s',r\mid s,a)$ |

$$V_\pi(s) = E_\pi[G_t \mid S_t = s] = E_\pi[G_0 \mid S_0 = s]$$

$$q_\pi(s,a) = E_\pi[G_t \mid S_t = s, A_t = a] = E_\pi[G_0 \mid S_0 = s, A_0 = a]$$

*infinite horizon return*

- Goal : Find $\pi$ with max expected return

- Optimal value function :

$$V_*(s) = \max_\pi V_\pi(s)$$

  - Best actions / policies from $s$
  - At least one solution

- Optimal action-value function :

$$q_*(s,a) = \max_\pi q_\pi(s,a)$$

  - Optimal policy for successor states from $s,a$
  - Same solution $\pi_*$ as $V_*$

- Optimal policy :  $\pi_* = \arg\max_\pi V_\pi(s)$

- **General results** (Silver: L2): Fn any finite MDPs
  - There exists at least one optimal policy $\pi_*$
  - All optimal $\pi_*$ achieve $V_*$
  - "       "       "       "    $q_*$
  - Optimal policy is deterministic (one action per state,
  - "       "       "   stationary

$$\Rightarrow \quad \pi_*(s) = a$$

- Optimal policy:

$$\pi_*(s) = \arg\max_a q_*(s,a) = \text{best action from } s$$

- Optimal value function:

$\downarrow$ best action

$$V_*(s) = V_{\pi_*}(s) = q_*(s, \pi_*(s))$$

$$V_*(s) = \max_a q_*(s,a)$$

$\rightarrow$ best action

$$\pi_* \iff q_* \iff V_*$$

- Note: $V_\pi(s) = \sum_a q_\pi(s,a)\, \pi(a|s)$

See p.3-6

$$= q_\pi(s, \pi(s)) \qquad \text{deterministic action}$$

· Example:

R=1                                        R=10

⑩    ? ①  →  ②    ③    ④    ⑤    ⑥
    ?↙

↑
start

$\gamma \approx 0$ : Better to move left to get R=1 reward

$\gamma \approx 1$ :    "    "    "   right  "   "   R=10    "

Optimal policy : for each s, either left or right

$p(s'|s, a=0)$    $p(s'|s, a=1)$
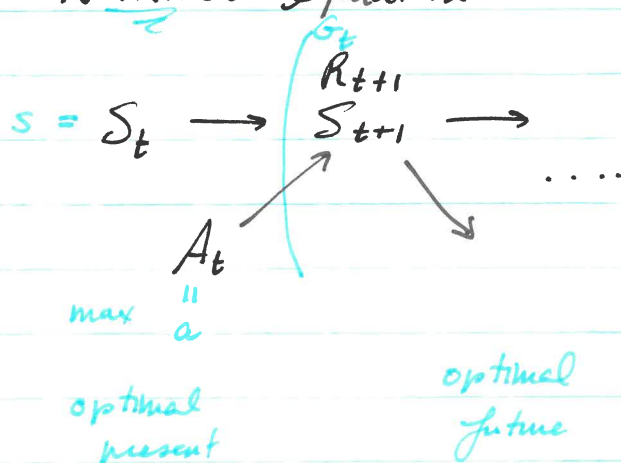
$\pi_*$ among $2^7$ possible policies
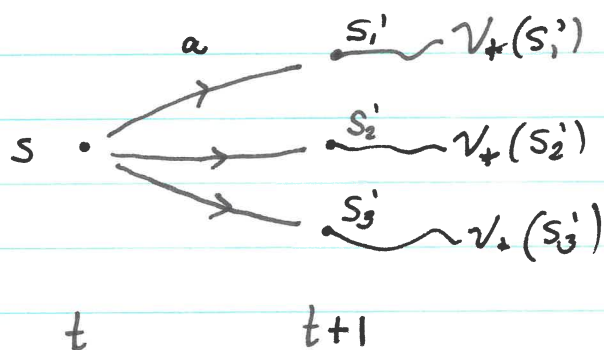
## 3.6 Bellman optimality equations

See 3-12b

· Value function :

$$V_*(s) = \max_a E\left[R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a\right]$$

- · Optimal for s = optimal action from s and then optimal policy after

- · Expectation on one step — no $E_\pi$
- · max outside expectation — influences reward
- · arg max defines $\pi_*(s)$

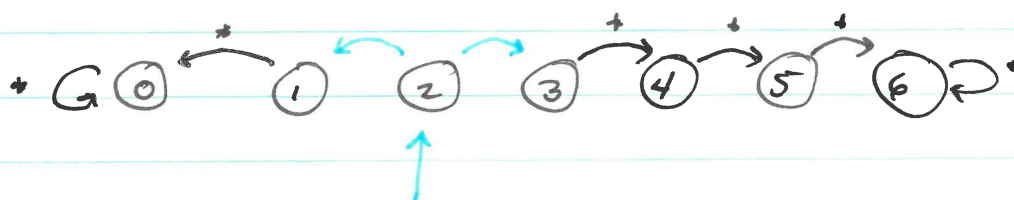- · Nonlinear equation because of max

$s = S_t \longrightarrow \begin{matrix} S_t \\ R_{t+1} \\ S_{t+1} \end{matrix} \longrightarrow \cdots$

$A_t$

max $a$

optimal present

optimal future

$V_*(s) = \max_a \rho(s,a) + \gamma \sum_{s'} V_*(s') p(s'|s,a)$

$$t \qquad\qquad t+1$$

- Whatever action chosen at time $t$, rest of dynamics is optimal

- To be optimal at $S_t = s$, choose optimal action and then follow optimal policy from $S_{t+1}$

$$\Rightarrow \quad v_*(s) = \max_a E\left[ R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a \right]$$

- Example:

· Action Value Function:

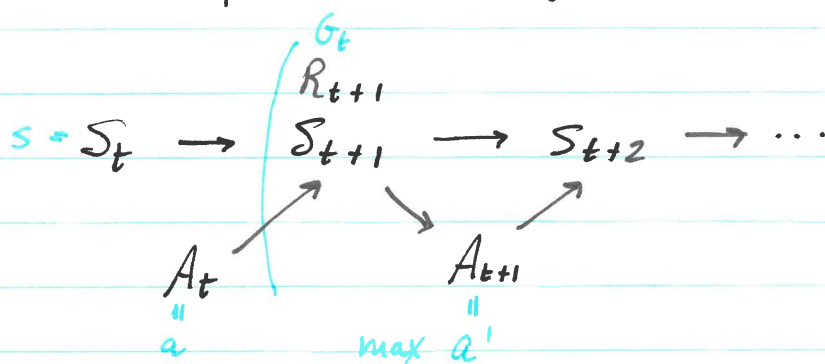$$q_*(s,a) = E\left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a\right]$$

action after

$$= E\left[R_{t+1} + \gamma \, v_*(S_{t+1}) \mid S_t = s, A_t = a\right]$$

i.e.

$$q_*(s,a) = \rho(s,a) + \gamma \max_{a'} \sum_{s'} q(s',a') \, p(s' \mid s,a)$$

$$= \rho(s,a) + \gamma \sum_{s'} v_*(s') \, p(s' \mid s,a)$$



$G_t$
$R_{t+1}$
$s = S_t \longrightarrow S_{t+1} \longrightarrow S_{t+2} \longrightarrow \cdots$
$A_t$  $A_{t+1}$
$\overset{\shortparallel}{a}$  $\overset{\shortparallel}{\max a'}$

· Optimal from $s,a$ = optimal after reaching $S_{t+1}$ by taking optimal action after

· Max inside — doesn't influence reward
· arg max defines $\pi_*(s')$

· Nonlinear equation because of max

· Bellman's optimality principle:
An optimal policy has the property that whatever the initial state/decision are, the remaining decisions constitute an optimal policy with regard to the state resulting from the first decision.

$s$  $s'$  optimal
max $a$  max $a'$
$t$  $t+1$  $t+2$ $\cdots$