# Finite Markov Decision Processes

- STATE, REWARD, ACTION   • AGENT/ENVIRONMENT   • POLICY $\pi_T$, $\pi_T(u|s)$ PROBABILITY, MAPS STATES TO PROBABILITY OF SELECTING ACTION

- EXPRESS THE GOAL THROUGH **REWARD SIGNAL**. WHAT TO ACHIEVE NOT HOW TO ACHIEVE

- **CUMULATIVE REWARD:** A SPECIFIC FUNCTION OF THE REWARD SEQUENCE, IE SUM

- **EPISODIC TASKS:** THERE IS A TERMINATING STATE   • **CONTINUING TASKS**. NO TERMINAL STATE $\rightarrow$ **DISCOUNTED RETURNS:** $G_T = R_{T+1} + \gamma R_{T+2} + \gamma^2 R_{T+3} + \sum \gamma^u R_{T+u+1}$
  - $\hookrightarrow$ CAN BE SEEN AS CONTINUING WHERE TERMINAL IS A ABSORBING STATE WITH R=0

  $0 \leq \gamma \leq 1$ DISCOUNT RATE

- **GENERAL RETURN** $G_T = \sum \gamma^u R_{T+u+1}$   • **STATE SIGNAL:** USUALLY SET TO HAVE MARKOV PROPERTY
  $\begin{cases} 0 & \text{ONLY IMMEDIATE REWARDS} \\ 1 & \text{FARSIGHTEDNESS} \end{cases}$

- RL TASK + MARKOV PROPERTY $\longrightarrow$ FINITE MARKOV DECISION PROCESS
  - $\begin{cases} \textbf{EXPECTED REWARDS } R(s;A) = \sum\limits_{s'\in S} R \sum P(s'_i R|s_i A) \\ \textbf{STATE TRANSITIONS: } P(s'|s,A) = \sum P(s',R|s,A) \\ \textbf{EXPECTED REWARDS } s\text{-}A\text{-}NS: R(s,A,s') = \dfrac{\sum\limits_R R \cdot P(s'|R|s,A)}{P(s'|s,A)} \end{cases}$

- DIAGRAM WITH STATE + ACTION NODES. $\sum P_{out}(a) = 1$

- **VALUE FCN:** $V_\pi(s) = E_\pi[G_T|S_{T=S}] = E_\pi\left[\sum \gamma^u R_{T+u+1}|S_T=s\right] \rightarrow$ VALUE OF STATE UNDER POLICY $\pi$. EXPECTED RETURN OF STARTING IN S AND FOLLOWING $\pi$
  **STATE-VALUE FCN FOR POLICY $\pi$**

  - $Q_\pi(s,A) = E_\pi[G_T|S_T=s, A_T=a] = E_\pi\left[\sum \gamma^u R_{T+u+1}\Big|S_T=s, A_T=a\right] \rightarrow$ VALUE OF TAKING A, UNDER STATE S, AND FOLLOWING $\pi$
    **ACTION-VALUE FCN FOR POLICY $\pi$**

**BELLMAN EQUATION:** $V_\pi(s) = \ldots = \sum\limits_a \pi(u|s)\sum\limits_{s',R} P(s',R|s,a)[R + \gamma V_\pi(s')]$   RELATIONSHIP BETWEEN VALUE OF STATES AND ITS SUCCESSORS  • FUNDAMENTAL IDENTITY
  $\approx T^\pi V^\pi = V^\pi$, $T^\pi$ BELLMAN OPERATOR
  $\hookrightarrow$ CONTRACTION FOR LEAST POINTS OF
  $\rightarrow$ CONVERGENCE UNIQUE SOLUTION

**OPTIMAL VALUE FCN:** FINITE MDP HAVE CLOSED-FORM OPTIMAL POLICY. VALUE FCN INDUCE PARTIAL ORDERING OVER POLICIES.

  $\pi \geq \pi'$ IFF $V_\pi(s) \geq V_{\pi'}(s) \forall s \in S$. MULTIPLE OPTIMAL POLICIES $\longrightarrow$ EQUIVALENT, SAME $V_*(s)$ AND $Q_*(s)$

  $V_*(s) = \max\limits_\pi V_\pi(s)$     $Q_*(s,A) = \max\limits_\pi Q_\pi(s,A) = E[R_{T+1} + \gamma V_*(S_{T+1})|S_T=s, A_T=a]$

**BELLMAN OPTIMALITY EQUATION:**
- FOR FINITE MDP, UNIQUE SOLUTION INDEPENDENT OF POLICY
  $\begin{cases} V_*(s) = \max \sum P(s',R|s,A)[R + \gamma V_*(s')] \\ Q_*(s,A) = \sum\limits_{s,R} P(s',R|s,A)[R + \gamma \max Q_*(s',A')] \end{cases}$
  • VALUE OF STATE UNDER OPTIMAL POLICY MUST EQUAL EXPECTED RETURN FOR BEST ACTION FROM THAT STATE

- IS $|S|$ SYSTEM OF EQUATION IN $|S|$ UNKNOWNS, CAN SOLVE IF DYNAMICS ARE KNOWN  • **FROM $V_*$** BEST ACTIONS AFTER 1-STEP SEARCH, OPTIMAL IN LONG RUN

- **FROM $Q_*$** EVEN EASIER, FOR ANY S $\rightarrow$ IS A ARGMAX $Q_*(s,A)$. NOT EVEN 1 STEP SEARCH.   $\rightarrow$ GREEDY BUT OPTIMAL
  IS LIKE ALREADY CACHING THE RESULTS. WE DON'T HAVE TO KNOW ANYTHING ABOUT FUTURE STATES AND THEIR VALUE, NO DYNAMICS NEEDED!!!

- **IN PRACTICE:** APPROXIMATIONS BECAUSE OPTIMAL COMPUTATIONS REQUIRE EXHAUSTIVE SEARCH. BOOM. RL IS ONLINE $\rightarrow$ WE CAN 'PRUNE' LOW OCCURRING STATES
  $\hookrightarrow$ CAN COMPUTE INDIVIDUALS