

REINFORCEMENT LEARNING

7/31/2023

THE GOAL IS TO CONTROL AN AGENT (NOT JUST MAKING PREDS)
IN AN ENVIRONMENT. EACH EXAMPLE INFLUENCES EACH OTHER, SEQUENTIAL DECISION MAKING

MARKOV - DECISION PROCESS

TUPLE $(S, A, \{P_{sa}\}, \gamma, R)$

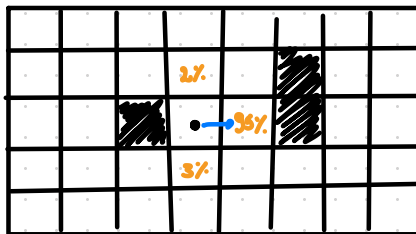
S - SET OF STATES EX $(\underbrace{x, y, z}_{\text{POSITION}}, \underbrace{v_x, v_y, v_z}_{\text{VELOCITY}}, \underbrace{\theta, \phi, \alpha}_{\text{DIRECTIONS}}, \dots) \rightarrow \text{HELICOPTER}$

WE WILL Δ^{ST} CONSIDER IT DISCRETE

A - SET OF ACTIONS EX ('MOVE LEFT', 'MOVE RIGHT')

P_{sa} - STATE TRANSITION PROBABILITIES

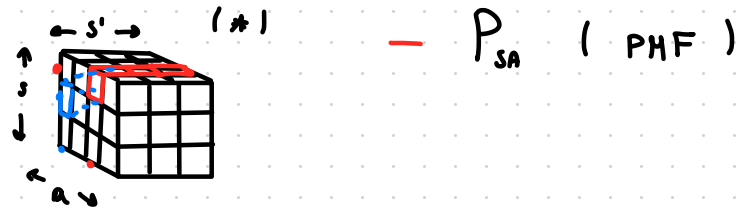
EX



RANDOMNESS IN THE ENVIRONMENT

DIFFERENT DISTRIBUTION FOR EACH
 (s, a)

WE CAN THINK OF $\{P_{sa}\}$ AS A 3D STRUCTURE (TENSOR)



γ - DISCOUNT FACTOR, $\gamma \in [0, 1)$ HOW MUCH DO WE DESIRE TO GET THE R VS LATER

R - REWARD, OF BEING IN A STATE, $R(s) \in \mathbb{R}$

Ex



WE WANT TO MAXIMIZE THE EXPECTATION OF (*)

$$\max \mathbb{E} \left[\underbrace{R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots + \gamma^d R(s_{d+1})}_{(*)} \right]$$

↓
WILL ALWAYS CONVERGE!

WE NEED TO DEFINE

- POLICY π : FUNCTION THAT MAPS CURRENT STATE TO ACTION (AN AGENT EXECUTES POLICY)

$$\pi: S \mapsto A$$

- VALUE $V_\pi(s)$: FUNCTION THAT GIVES EXPECTED VALUE FOR BEING IN s AND FOLLOWING π

$$V_\pi(s) = \mathbb{E} \left[R(s_0) + \gamma R(s_1) + \dots + \gamma^{K-1} R(s_K) \right]$$

$$V_\pi(s) = \mathbb{E} \left[R(s_0) + \dots + \gamma^{K-1} R(s_K) \mid \pi \right]$$

ONCE REACHED K , π IS GIVEN

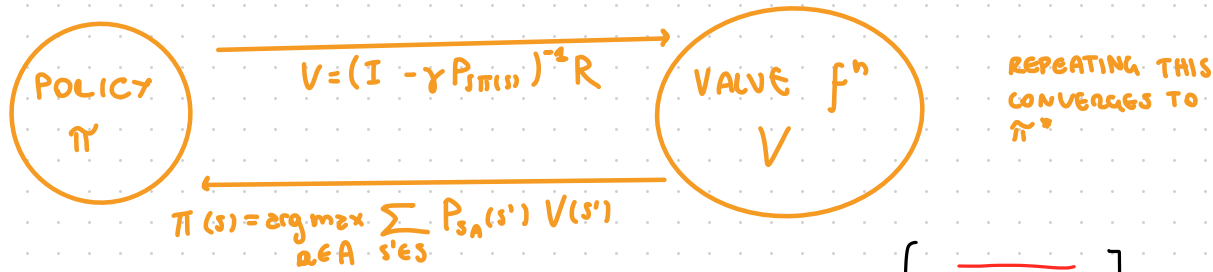
$$V_\pi(s) = R(s_0) + \mathbb{E} \left[\gamma R(s_1) + \dots + \gamma^{K-1} R(s_K) \mid \pi \right]$$

AND, SINCE THIS HAS A RECURSIVE STRUCTURE

$$V_\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V_\pi(s')$$

BELLMAN
EQUATION

$$V_\pi(s) = R(s) + \mathbb{E}_{s \sim P_{s\pi(s)}} [V_\pi(s')]$$



$$V(s) = [V(s_0), V(s_1), V(s_2), \dots]$$

$$R(s) = [R(s_0), \dots]$$

$$P_{S\pi(s)} = \begin{bmatrix} \text{---} \\ \text{---} \\ (*) \end{bmatrix}$$

$$\gamma \in \mathbb{R}$$

SYSTEM OF K
LINEAR EQUATIONS
(K VARS)

$$\begin{bmatrix} V(s_0) \\ V(s_1) \\ \vdots \end{bmatrix} = \begin{bmatrix} R(s_0) \\ R(s_1) \\ \vdots \end{bmatrix} + \gamma \begin{bmatrix} P_{S\pi(s)} \end{bmatrix} \begin{bmatrix} V(s'_0) \\ V(s'_1) \end{bmatrix}$$

$$V = R + \gamma P_{S\pi(s)} V \mid (I - \gamma P_{S\pi(s)}) V = R \mid V = (I - \gamma P_{S\pi(s)})^{-1} R$$

THEN WE EXPRESS THE RATIONAL POLICY

$$\pi(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s')$$

AND THEN FIND V W/ THE GREEDY POLICY AND BELLMAN EQUATION

OPTIMAL VALUE FUNCTION

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

$$V^*(s) = R(s) + \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

WE ARE NO MORE FOLLOWING A POLICY, BUT DIRECTLY MAXIMIZING THROUGH ACTIONS
SO WE CAN DEFINE AN OPTIMAL POLICY π^*

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

ALL POSSIBLE POLICIES' DOMAIN IS $|A|^{|S|}$, so A LOT

$$V^*(s) = V^{\pi^*}(s) \geq V^{\pi}(s)$$

ALGORITHMS

1. VALUE ITERATION
2. POLICY ITERATION

VALUE ITERATION

I. SET $V(s) = 0 \quad \forall s$

II. REPEAT UNTIL CONVERGENCE {

$$\forall s \quad V(s) := R(s) + \max_{a \in A} \gamma \sum_{s'} P_{sA}(s') V(s')$$

}

ex on 1st iter.  THIS IS $R(s)$

ex $V^{(0)}(s) = 0 \rightarrow V^{(1)}(s) = R(s) \rightarrow V^{(2)}(s) = R(s) + \gamma R(s) + \dots \rightarrow V^{(T)}(s) = V^*(s)$

$$V^{(t)}(s) \rightarrow V^{(t+1)}(s)$$

BELLMAN BACKUP
OPERATOR

WILL ALWAYS CONVERGE TO V^*, π^*

POLICY ITERATION

I. INIT π RANDOMLY

II. REPEAT UNTIL CONVERGENCE {

(a) LET $V = V^\pi$

(b) FOR EVERY s , LET

$$\pi(s) = \operatorname{argmax}_{a \in A} \sum_{s'} P_{sA}(s') V(s')$$

}

ALWAYS CONVERGE TO V^*, π^*

LEARNING THE MDP

UNTIL NOW WE KNEW R AND P_{SA}

IF WE WERE GIVEN ONLY A SUBSET OF THE MDP

$MDP = (S, A, \{P_{SA}\}, \gamma, R)$ \rightarrow WE DON'T KNOW THE RULES OF THE GAME, SO WE NEED TO LEARN THE GAME WHILE GETTING GOOD AT IT

EX

EPISODES $t_1: S_0^{(1)} \xrightarrow{a_0^{(1)}} S_1^{(1)} \xrightarrow{a_1^{(1)}} S_2^{(1)} \longrightarrow \dots$

$t_2: S_0^{(2)} \xrightarrow{a_0^{(2)}} S_1^{(2)} \longrightarrow \dots$

RUN TRIALS OVER AND OVER (EXPLORING STATE), THEN WE CAN DEFINE

$$P_{SA}(s') = \frac{\text{\#TIMES WE TOOK } a \text{ IN } s \text{ AND GOT IN } s'}{\text{\#TIME WE TOOK } a \text{ IN } s}$$

AND W/ GOOD COVERAGE WE CAN GET GOOD ESTIMATES

AT THE BEGINNING WE ASSIGN $P_{SA}(s) = 0, 1$ (FOR EXAMPLE), BECAUSE IT WOULD BE $\frac{0}{0}$

AND REWARDS CAN BE DEFINED AS THE MEAN FROM THE TRIALS

$$R(s) = \frac{1}{n} \sum_{i=1}^n R^{(i)}$$

ALGORITHM

I. INIT π RANDOMLY

II. REPEAT {

(a) EXECUTE π IN MDP FOR SOME TRIALS

(b) USE ACCUMULATED EXPERIENCE AND UPDATE P_{SA} AND R

LEARN THE GAME

(c) APPLY VALUE-ITERATION W/ EST. P_{SA}, R

(d) UPDATE π TO BE GREEDY W.R.T. V_{π}

BECOME GOOD AT IT

OPTIMIZE

WHEN INIT V IN $V-I$ WE INIT V AS THE OBTAINED R (?)

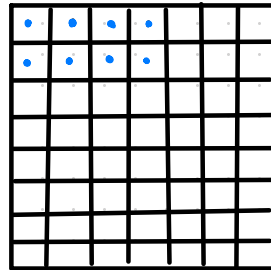
AS LONG AS R AND P_{SA} CONVERGE TO TRUE VALUES, THE ALGORITHM WILL CONVERGE

RL IN CONTINUOUS SPACES

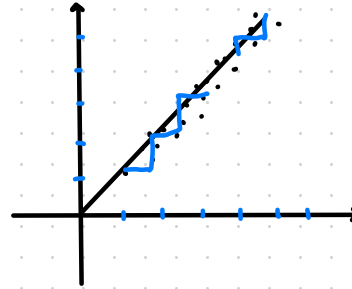
ex

$s = (x, y, \theta, \dot{x}, \dot{y}, \dot{\theta})$ IS A CONTINUOUS STATE SPACE

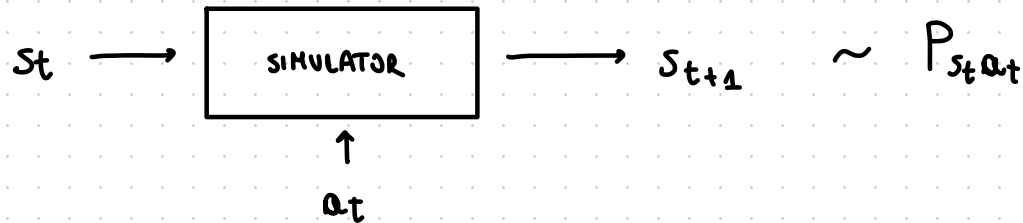
WE CAN DISCRETIZE!



STATE SPACE



TYPICALLY WE HAVE A SIMULATOR



WE RUN AN n NUMBER OF TRIALS

$$\begin{aligned} s_0^{(1)} &\xrightarrow{a_0^{(1)}} s_1^{(1)} \xrightarrow{a_1^{(1)}} s_2^{(1)} \xrightarrow{a_2^{(1)}} \dots \xrightarrow{a_{T-1}^{(1)}} s_T^{(1)} \\ s_0^{(2)} &\xrightarrow{a_0^{(2)}} s_1^{(2)} \xrightarrow{a_1^{(2)}} s_2^{(2)} \xrightarrow{a_2^{(2)}} \dots \xrightarrow{a_{T-1}^{(2)}} s_T^{(2)} \\ &\dots \\ s_0^{(n)} &\xrightarrow{a_0^{(n)}} s_1^{(n)} \xrightarrow{a_1^{(n)}} s_2^{(n)} \xrightarrow{a_2^{(n)}} \dots \xrightarrow{a_{T-1}^{(n)}} s_T^{(n)} \end{aligned}$$

WE APPLY A LEARNING ALGORITHM TO PREDICT s_{t+1} GIVEN s_t, a_t
TYPICALLY A LINEAR ONE S.T.

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

$$s_{t+1} = A s_t + B a_t$$

HERE THE PARAMETERS OF THE MODEL ARE A, B AND ESTIMATE THEM
W/ DATA FROM n TRIALS

$$\arg \min_{A, B} \sum_{i=1}^n \sum_{t=0}^{T-1} \| s_{t+1}^{(i)} - (A s_t^{(i)} + B a_t^{(i)}) \|_2^2$$

HERE WE CAN CHOOSE TO MAKE THE MODEL

- DETERMINISTIC
- STOCHASTIC

$$s_{t+1} = A s_t + B a_t + \epsilon_t \quad (**)$$

$$\epsilon \sim N(0, \Sigma)$$

$$s_{t+1} = A \phi_s(s_t) + B \phi_a(a_t)$$

$$s_{t+1} \sim N(A \phi_s(s_t), B \phi_a(a_t), \Sigma)$$

ϕ_s, ϕ_a NON LINEAR FEATURE MAPPINGS OF s AND a

WE CAN NOW DESCRIBE FITTED VALUE ITERATION, THAT APPROX.
 V OF A CONTINUOUS STATE MDP

$$V(s) := R(s) + \gamma \max_a \int_{s'} P_{sa}(s') V(s') ds'$$

$$= R(s) + \gamma \max_a \mathbb{E}_{s' \sim P_{sa}} [V(s')]$$

$$= R(s) + \gamma \max_a \frac{1}{K} \sum_{i=1}^K V(s_i) \quad s_i \sim P_{SA}$$

MONTÉ-CARLO EST. FROM \mathbb{E}

$$V(s) = \Theta^T \phi(s) \quad (*)$$

∞ NUMBER OF STATES, WE CANNOT STORE THEM. SO WE FIT A MODEL THAT APPROX THE RIGHT VALUE TO THIS $(*)$ MODELING (PARAMETRIC) ASSUMPTION GIVEN THE ASSUMPTION, WE CAN DEFINE THIS NEW ALGORITHM

ALG

I. SAMPLE n STATES $s^{(1)}, s^{(2)}, \dots, s^{(n)}$

II. INIT $\theta = 0$ ($V(s) = \theta^T \phi(s)$)

III. REPEAT FOR $i = 1, \dots, n$

FOR EACH $a \in A$

SAMPLE $s'_1, \dots, s'_K \sim P_{SA}$

SET $q(a) = \frac{1}{K} \sum_{i=1}^K R(s^{(i)}) + \gamma V(s^{(i)})$

LEARNED IN PHASE 1 $(**)$

$q(a)$ EMPIRICAL AVERAGE

}

$$\text{SET } y^{(i)} = \max_a q(a)$$

}

$y^{(i)}$ IS THE O: OF BELLMAN
OPERATOR (EST. OF $R(s^{(i)}) + \gamma \max_a \mathbb{E}_{s' \sim P_{s,a}} [V(s')]$)

$$(s^{(i)}, y^{(i)})_{i=1}^n \quad V^{(t+1)}(s^{(i)}) = y^{(i)}_{(*)}$$

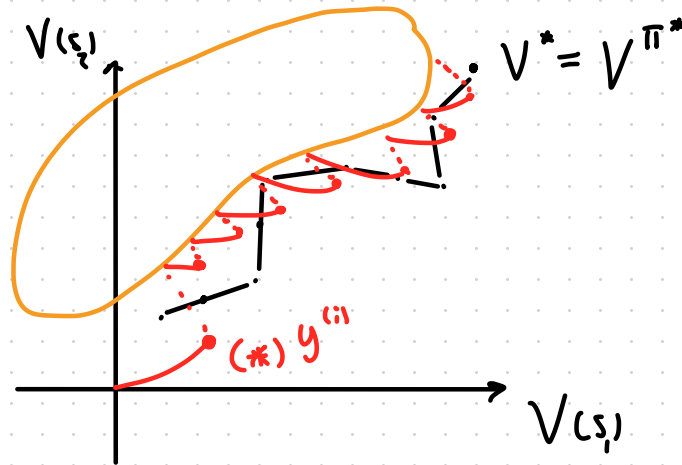
$$\text{SET } \Theta = \arg \min_{\Theta} \sum_{i=1}^n (\Theta^T \phi(s^{(i)}) - y^{(i)})^2$$

BASICALLY, BY APPLYING THE BELLMAN BACKUP OPERATOR, WE WILL CONVERGE TO V^*

$$V(s) = \Theta^T (\phi(s))$$

WE KNOW HOW TO
REPRESENT THIS
SUBSPACE (AND ITS F)

THIS SPACE DEPENDS
ON THE FEATURE MAP



∞ - DIMENSIONAL SPACE
(CONTINUOUS)

WE PROJECT y''' ONTO THE SUBSPACE WE KNOW HOW TO REPRESENT
AT SOME POINT WE OBTAIN THE APPROX VALUE $V^{(T)}$