

Fundamentals of Reinforcement Learning

LESSON 01 / THE K-ARMED BANDIT PROBLEM

- * understanding the temporal nature of the bandit problem.
- * define K-armed bandit.
- * define action-values.
- * define reward.

LESSON 02 / WHAT TO LEARN? ESTIMATING ACTION VALUES

- * define action-value estimation methods.
- * define exploration and exploitation.
- * select actions greedily using an action-value function.
- * define online learning.
- * understand a simple online action-value estm. method.
- * define the general online equation.
- * understand const. step size in case of non stationarity.

LESSON 03 / EXPLORATION VS EXPLOITATION TRADEOFF

- * cmp short term benefits of exploit vs long term of explore.
- * understand optimistic initial values.
- * describe the benefits of opt. init. val. for early exploration.
- * explain criticism of optimistic initial values.
- * desc. the upper confidence bound action selection method.
- * define optimism in the face of uncertainty.

WEEK 01 - PAGES 25-36 OF RL-BOOK → LECTURES

PAGES 42-43 OF RL-BOOK → BEFORE ASSIGNMENTS

RL - explore alternatives, find which one probably maxes rewards.

SL - from historical data, try to generalize with least error & no bias.

q^* is not known to the agent, just like the doctor doesn't know the effectiveness of each treatment. So, estimate it.

$$\begin{array}{|c|c|c|} \hline \text{100% = 6 to 12} & \text{78% = 10} & \text{92% = 10} \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline \text{22% = 12} & & \text{8% = 144} \\ \hline \end{array}$$

REPEATEDLY TRY ACTIONS
FIND MEAN OF REWARDS

An action is per unit of time
so prior actions is $t-1$.

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i}{t-1}$$

q^* is the mean of the distributions for each action

$$q^*(\tau_1) = .5 \times .11 + .5 \times 9 = 1$$

$$q^*(\tau_2) = 1$$

$$q^*(\tau_3) = 3$$

HEALTH IMPROVEMENT?

HARD TO MEASURE

WHAT ABOUT IMPROVEMENTS IN BLOOD PRES?

CAN FOLLOW ≠ DISTRIBUTIONS

bernoulli

binomial

uniform

$q^*(\tau_1)$

$q^*(\tau_2)$

$q^*(\tau_3)$

$\tau_1 \dots \tau_3$

<p

Incremental Action Value estimation /

Say you've got millions of requests. It takes a lot of space to store past rewards and be computing again the action value estimates. You really want to only store one single value per treatment, and update it on the go every time an explore is done.

THE INCREMENTAL UPDATE RULE IS AN INSTANCE OF A MORE GENERAL RULE.

THE GENERAL UPD. RULE CAN BE USED TO SOLVE NON-STATIONARY BANDIT PROBLEM.

Non-Stationary Bandit Problem /

What if one of the treatments suddenly became more effective under certain conditions, like seasonality / winter, which could change the reward distribution.

* USE A FIXED α_n : e.g. "0.1" the most recent reward will weigh more than all previous (which have $0.1 * 0.1 \dots$) weight fades exponentially with time.

Decaying Past Rewards / the most recent contributes most the current estimate.

$$Q_{n+1} = Q_n + \alpha_n (R_n - Q_n) \text{ GENERAL RULE/FORM.}$$

$$= \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2} + \dots + (1-\alpha)^{n-1} \alpha R_1 + (1-\alpha)^n \alpha R_0 Q_0$$

$$= (1-\alpha)^n Q_0 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

INITIAL PREDICTION

REWARDS FURTHER BACK IN TIME CONTRIBUTE EXPONENTIALLY LESS.

Optimistic Initial Values /

This encourages early exploration, instead of assuming the original estimated action values Q to be 0, use a large overestimation in all of them to make them appealing to greedy selection.

$$Q_{n+1} = Q_n + \alpha (R_n - Q_n) \quad \alpha = 0.5 \text{ (EXPONENTIAL DECAY)}$$

$$\textcircled{1} \quad 0.2 \cdot 1 + 1 = Q = 1.5$$

ENCOURAGES GREEDY- CHOOSING ALL THREE TREATMENTS EARLY ON TIME IT IS A VERY USEFUL HEURISTIC!

$$\textcircled{2} \quad 0.2 \cdot 1 + 0 = Q = 1.0$$

$$+ 1 = Q = 1.5 + 0 = Q = 0.625$$

$$\textcircled{3} \quad 0.2 \cdot 1 + 1 = Q = 1.5 + 1 = Q = 0.625$$

IN FIRST ROUND, EITHER EXPLOIT OR EXPLORE WILL CHOOSE AT RANDOM

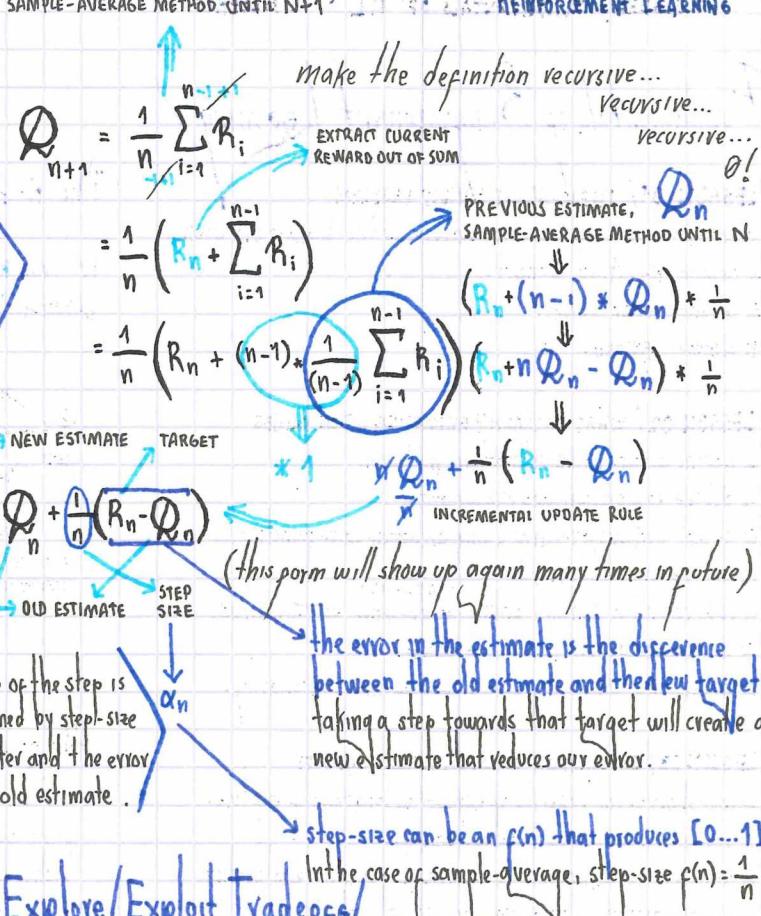
THE REWARD WAS $+1$, BECAUSE THE INITIAL ESTIMATE WAS LARGER THAN POSSIBLE REWARDS. THIS WILL BRING ESTIMATE DOWN AND MAKE THE OTHER TREATMENTS APPEALING TO GREEDY.



THE 10-ARM TEST BED

from textbook, 10 armed bandit, 10 reward distributions sampled from Normal Dist, an initial estimated value of 5 is highly optimistic, all vals in plot are less than 3. whenever agent selects an action, the observed reward will likely be smaller than original estimation, thus bringing it down, other actions will become more appealing in comparison.

- * ONLY DRIVES EARLY EXPLORATION, STOPS LATER
- * HAS ISSUES WITH NON-STATIONARY PROBLEMS
- * WE MAY NOT KNOW HOW HIGH TO SET THIS VALUE



Explore/Exploit Tradeoff /

EXPLORATION IMPROVE KNOWLEDGE FOR LONG-TERM BENEFIT,

INCREASE THE ACCURACY OF THE ESTIMATED ACTION VALUES

EXPLOITATION ACT UPON CURRENT KNOWLEDGE OF BEST FOR SHORT-TERM

BENEFIT, CHOOSE THE GREEDY ACTION WHICH VERY EARLY ON IS NOT SO LIKELY TO HAVE THE BEST REW. DISTRIBUTION.

Epsilon Greedy Action Selection /

to exploit knowledge but not get in a trap of not learning, we can have an exploration ratio, or epsilon greedy, to randomly decide



$$A_t \leftarrow \begin{cases} \operatorname{argmax} Q_t(a) & \text{WITH PROBABILITY } 1 - \epsilon \\ a \sim \operatorname{Uniform}\{\{a, \dots, a_k\}\} & \text{WITH PROBABILITY } \epsilon \end{cases}$$

ACTION SELECTED AT TIME STEP t

OPTIMAL ACTION

Percentage of time agent chooses optim. action, avg'd over runs

100%

OPTIMISTIC, greedy

$$Q=5, \epsilon=0, \alpha=0.1$$

REALISTIC, E-greedy

$$Q=0, \epsilon=0.1, \alpha=0.1$$

In early learning, optimistic agent performs worse, because it explores more, its exploration decreases w/ time, because the optimism in its estim. washes out w/ more samples, but in long run it selects most of the time the most optimal action compared with realistic agent because it explored more early on time.

steps E-greedy just keeps exploring

Practice Quiz/

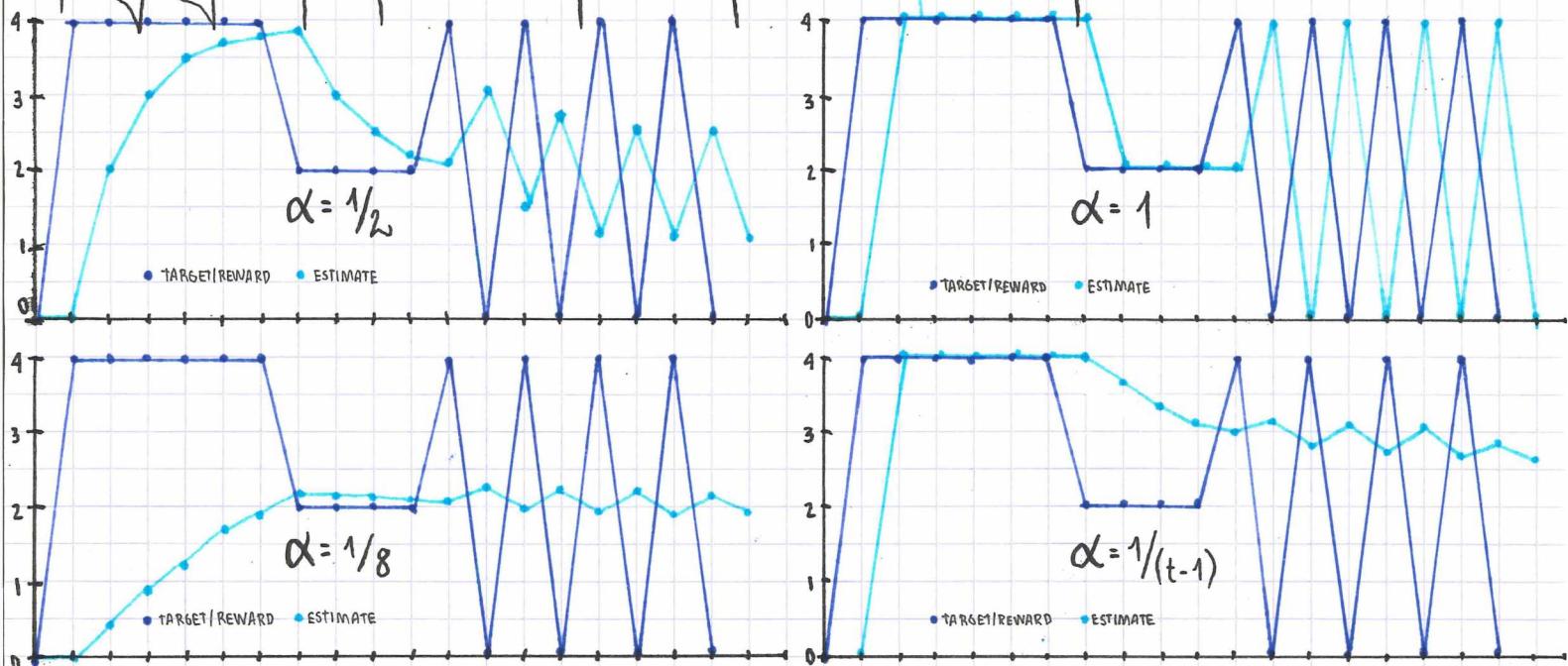
What is the incremental rule (sample average) for action values? What is the exploration/exploitation tradeoff?

$$R/ \quad Q_{n+1} = Q_n + \frac{1}{n} (\text{TARGET} - Q_n)$$

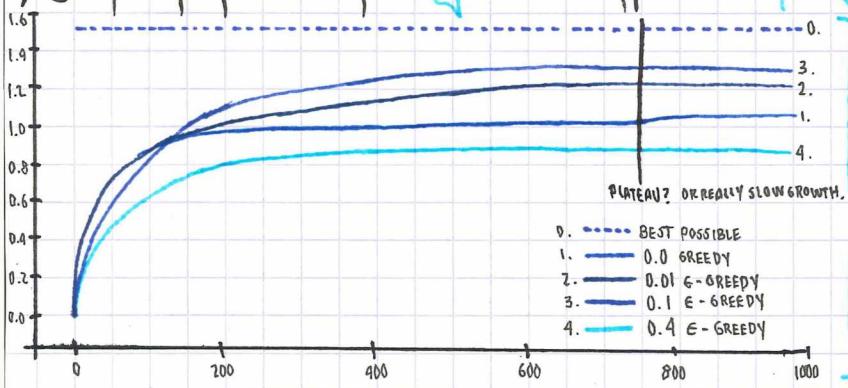
At each time step the agent moves its prediction in the direction of the error by the step size α in the case of the sample-average by $\frac{1}{n}$.

R/ the agent wants to explore to get more accurate estimates of its action values, the agent also wants to exploit its knowledge to get more reward, but it can't choose to do both things simultaneously.

By using the general form of the estimate-update rule, find the α (step size) that updates estimates based on rewards.



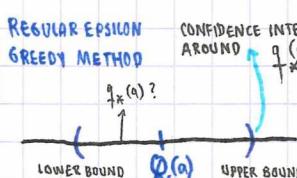
Compare performance of an agent with different ϵ (epsilon greedy values).



- 0 is the best possible
- 1 doesn't explore, at first performs better but long run, has not much growth. 2 starts similarly but better later.
- 3 is not so close to 0,0, starts slow, later really good.
- 4 explores so much it never quite takes off, need to exploit more the information it has acquired.

Upper-Confidence Bound (UCB) Action Selection/

$$A_t \leftarrow \begin{cases} \underset{a}{\operatorname{argmax}} Q_t(a) & \text{prob } 1-\epsilon \\ a \sim \text{Uniform}\{\{a_1, \dots, a_K\}\} & \text{prob } \epsilon \end{cases}$$



HIGHEST ESTIMATED VALUE

UCB ACTION SELECTION

$$A_t = \underset{a}{\operatorname{argmax}} [Q_t(a) + C \sqrt{\frac{\ln t}{N_t(a)}}]$$

C: USER SPECIFIED PARAMETER, CONTROLS EXPLORATION.

ln t: ln of amount of timesteps t

N_t(a): number of times action a has been taken until time t

UPPER CONFIDENCE BOUND EXPLORATION TERM. SO MIXES SOMEHOW THE EXPLORATION AND EXPLOIT.

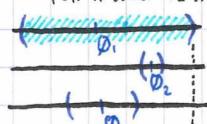
$$C \sqrt{\frac{\ln 10000}{5000}} \Rightarrow 0.0430$$

$$C \sqrt{\frac{\ln 100}{100}} \Rightarrow 0.303 \approx 10x \text{ value!}$$

The exploratory actions are selected uniformly. Could it be done better? If we had a notion of uncertainty in our value estimates, we could potentially select actions in a more intelligent way.

UCB → OPTIMISM IN FACE OF UNCERTAINTY

IF WE ARE UNCERTAIN ABOUT SOMETHING, WE SHOULD OPTIMISTICALLY ASSUME IT IS GOOD. EXPLORE THE ACTION WITH HIGHEST UPPER BOUND. IF VALUE IS BIG, GREAT, IF NOT, WE LEARNED ABOUT IT & NARROWED.



UCB REDUCES EXPLORATION OVER TIME AS INTERVALS ARE NARROWED, BETTER PERFORMANCE LONG TERM, LIKES TO EXPLORE LESS TAKEN ACTIONS.

Introduction to Markov Decision Processes/

LESSON 01 / INTRODUCTION TO MARKOV DECISION PROCESSES

- * understand MARKOV DECISION PROCESSES, or MDPs.
- * describe how the dynamics of an MDP are defined.
- * understand the graphical representation of an MDP.
- * explain how diverse processes can be written in terms of the MARKOV DECISION PROCESS framework.

LESSON 02 / GOAL OF REINFORCEMENT LEARNING

- * describe how rewards relate to the goal of an agent.
- * understand episodes and identify episodic tasks.

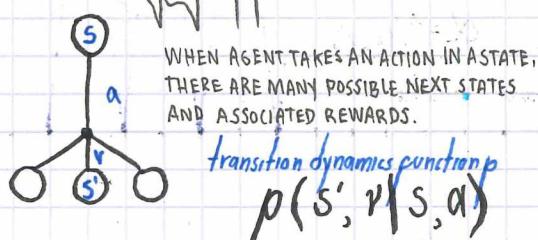
LESSON 03 / CONTINUING TASKS

- * formulate returns for cont. tasks using discounting.
- * describe how returns at successive time steps are related.
- * understand when to formulate a task as episodic or cont.

WEEK 02 - PAGES 47 - 56 OF RL-BOOK → LECTURES

The Dynamics of an MDP/

As in bandits, the outcomes are stochastic, so we have to use the language of probabilities.



Given a state s and an action a , the function p tells us the joint probability of next state s' and reward r above. In this course we will typically assume that the set of states, actions, rewards are finite, but will learn algorithms that handle inf or uncountable sets.

The Markov Property/

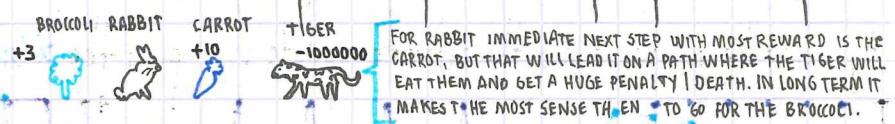
Since p is a probability distribution, it must be non-negative and its sum over all possible next states and rewards must equal 1. Future state and reward only depends on the current state and action. the present state is sufficient and remembering earlier states would not improve predictions about the future.

$$p: S \times R \times S \times A \rightarrow [0, 1]$$

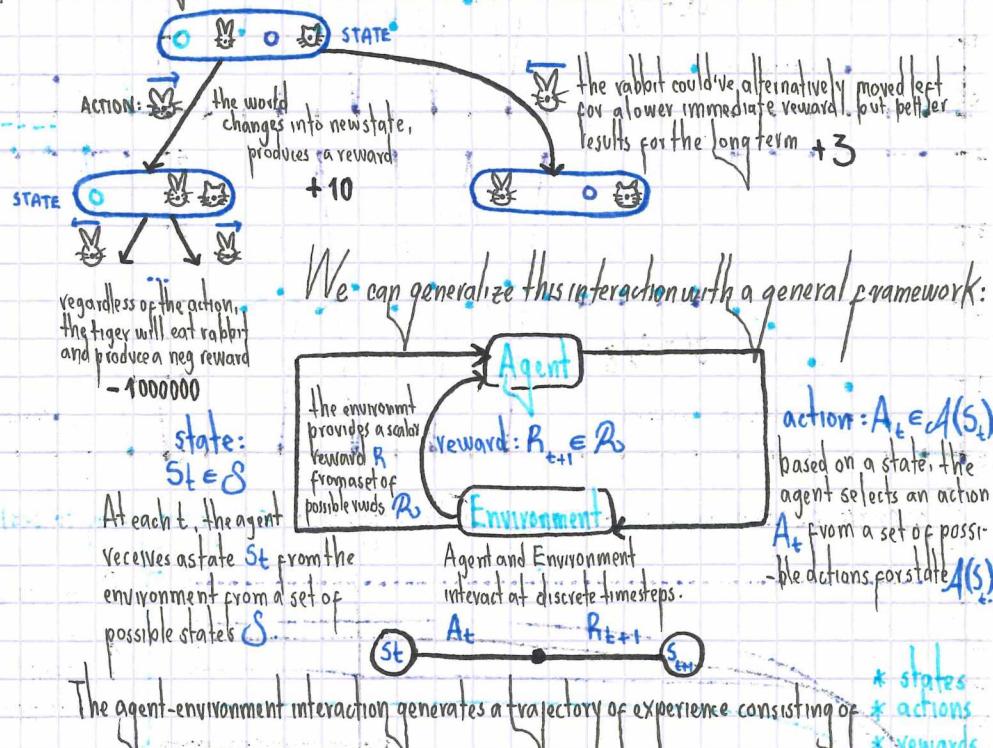
$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \quad \forall s \in S, a \in A(s)$$

Markov Decision Processes/

The **K-armed bandit** problem doesn't include many aspects of real-world problems. The **agent** is presented with same situation, and each time the **action** is optimal. In many real-life problems, different situations call for different responses. The **action** we choose now, affect the amount of **reward** we can get in the **future**. The **Markov Decision Process** formalism captures this two aspects of real-world problems.



A **bandit rabbit** will only care about **immediate reward**, but we can make a better decision by considering the **long-term impact** of our decisions. The **situation** changes as the rabbit takes **actions**, we'll call the situations "**states**".



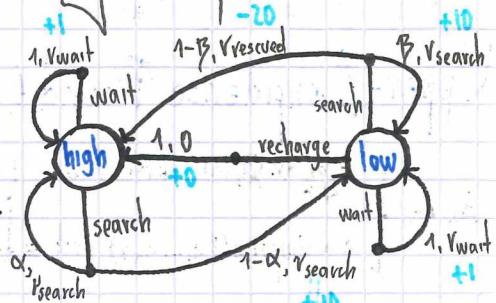
Recycling Robot MDP/

Robot collects empty cans in a office environment. It can detect cans, grab them, and drop them in a recycling bin. The robot runs on rechargeable battery. It seeks to collect as much cans.

$$A = A(\text{low}) = \{\text{search, wait, recharge}\}$$

$$A(\text{high}) = \{\text{search, wait, recharge}\}$$

ONLY ALLOW RECHARGING FROM THE LOW ENERGY STATE.



$$\alpha: \text{probability of remaining high after search.}$$

$$1-\alpha: \text{prob of going to low from high after search.}$$

$$\beta: \text{probability of remaining low after search.}$$

$$1-\beta: \text{prob of needing to be rescued after search in low, takes you back to high.}$$

Search: done in any state, can bring a reward of 10 but can trigger a state change.

Wait: done in any state, can bring a reward of 1, doesn't trigger a state change.

Recharge: no reward, done in low to go to high, no alternatives

Rescued: if run out of battery will be recharged but involves a penalty.

The Goal of Reinforcement Learning

In RL, the agent's objective is to maximize future rewards. Bandits maximize immediate reward, but this won't work on an MDP. Let's formally define what maximizing total future reward is:

RANDOM VARIABLE BECAUSE DYNAMICS OF MDP CAN BE STOCHASTIC.

$$\text{return } G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$$

the return at time step t is the sum of rewards obtained onwards

$$E[G_t] = E[R_{t+1} + R_{t+2} + R_{t+3} + \dots] + R_T$$

due to randomness in individual rewards and state transitions, let's maximize the expected return, many diff traj. from same state are possible. FOR THIS TO BE WELL DEFINED, SUM OF REWARDS MUST BE FINITE what happens when the interaction ends?

EPISODIC TASKS.



In Chess, each game starts at the same state with all pieces reset, it ends with check, draw, or resignation. Each game constitutes an episode.

EPISODIC TASKS

- * INTERACTION BREAKS NATURALLY INTO EPISODES.
- * EACH EPISODE ENDS IN A TERMINAL STATE.
- * EPISODES ARE INDEPENDENT

$$G_t = \sum_{i=0}^n R_{t+i+1} \mid R_{t+n+1} = R_T \quad \text{TERMINAL STATE}$$

CONTINUING TASKS

- * INTERACTION GOES ON CONTINUALLY.
- * THERE ARE NO TERMINAL STATES.
- * LIKE A THERMOSTAT ALWAYS INTERACTING WITH ENVIRONMENT, STATE IS TEMP, ACTION IS ON/OFF, NEG REWARD IF MANUAL CHANGE.

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid 0 \leq \gamma < 1$$

Practice Quiz

* the learner and decision maker is the agent.

* at each time step the agent takes an action.

* MDP stands for MARKOV DECISION PROCESS.

IMAGINE THE AGENT IS LEARNING IN AN EPISODIC PROBLEM:

the number of steps is stochastic, each episode can have a different number of steps occurring.

WHAT IS THE DIFFERENCE BETWEEN SMALL & LARGE GAMMA:

with a larger discount factor the agent is more far-sighted and considers rewards further into the future.

WHAT IS THE REWARD HYPOTHESIS:

goals and purposes can be thought of as the maximization of the expected value of the cumulative sum of rewards received.

FORMULATE AS: when agent-environment interaction naturally breaks into sequences, each sequence begins independently of how the episode ended.

FORMULATE AS: when agent-environment interaction does not naturally break into sequences, each continuing: new episode begins independently of how the previous episode ended.

Continuing Tasks/

In many problems, the agent-environment interaction continues without end, such problems can be formulated as **CONTINUING TASKS**. How do we formulate the return as we did for episodic tasks, without letting it become infinite?

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

$$\text{return } G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{k-1} R_{t+k} + \dots$$

the return is the sum of discounted rewards obtained onwards. Immediate rewards contribute more to the sum, far into future less. \$1 today is worth more than \$1 in a year.

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad \text{FINITE AS LONG AS } 0 \leq \gamma < 1$$

assume R_{\max} is the max reward an agent can receive, so return is bounded like:

$$\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \sum_{k=0}^{\infty} \gamma^k R_{\max} \quad \text{CONSTANT} \rightarrow R_{\max} \sum_{k=0}^{\infty} \gamma^k \rightarrow R_{\max} \frac{1}{1-\gamma}$$

$$R_{\max} * \frac{1}{1-\gamma} \quad * \text{IS FINITE} \\ * \text{IS AN UPPER BOUND ON } G_t \quad \therefore G_t \text{ IS FINITE.}$$

Effect of γ on Agent Behavior/

$\gamma=0$ $G_t = R_{t+1}$, agent only cares about immediate reward, short-sighted.

$\gamma=1$ immediate and future R are weighted nearly equally, more far-sighted.

Recursive Nature of The Returns/

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ = R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots)$$

$$G_t = R_{t+1} + \gamma G_{t+1} \quad \text{IN FUTURE WE'LL EXPLOIT THIS EQUATION TO DESIGN LEARNING ALGORITHMS.}$$

SUPPOSE $\gamma=0.8$ AND REWARDS $R_1=-3, R_2=5, R_3=2, R_4=7, R_5=1$. $T=5$, WHAT IS G_0 ?

$$-3 + 5 * 0.8 + 2 * 0.8^2 + 7 * 0.8^3 + 1 * 0.8^4 = 6,273.6$$

SUPPOSE $\gamma=0.5$ AND REWARD SEQUENCE $R_1=5$ FOLLOWED BY INF 10s. WHAT IS G_0 ?

$$\text{think: } G_2 = \frac{10}{1-0.5} = 50. \quad G_1 = 10 + 0.5 * (50) = 50 \quad G_0 = 5 + 0.5 * (50) = 45$$

IMAGINE AN AGENT IN A MAZE-LIKE GRID WORLD, WANT THEM TO FIND GOAL ASAP. GIVE A REWARD ON GOAL OF +1. DISCOUNT IS 1 (EPISODIC). WHEN RUN, AGENT FINDS GOAL, BUT DOES NOT CARE HOW LONG IT TOOK. HOW TO FIX THIS?

* give -1 reward at each step * use a discount $\neq 1$, something like 0.9. This will incentivize agent to reach goal faster to maximize the return.

IMAGINE YOU ARE A VISION SYSTEM, WHEN 1ST TURNED ON AN IMAGE FEEDS CAMERA. CAN SEE LOTS OF THINGS BUT NOT ALL (OCCLUDED, BEHIND, ETC). AFTER FIRST SCENE, DO YOU HAVE ACCESS TO MARKOV STATE OF ENVIRONMENT? SUPPOSE CAMERA WAS BROKEN THAT DAY, NO IMG AT ALL DAY, DO YOU HAVE ACCESS TO THE MARKOV STATE THEN?

* you have access to MARKOV STATE before and after damage. There's no history before first image, it has markov property. The property doesn't tell state tells all that's useful to know, only that hasn't forgotten anything useful. The case with camera broken, future is impoverished but still has markov property. All futures are same (black img), nothing needed to remember to predict them.

$$G_t = 1 - \frac{1}{1-\gamma} \quad \text{CONVERGING GEOMETRIC SERIES WHEN } \gamma < 1.$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k * 1 \quad \text{CONSTANT}$$

Policies and Value Functions/

LESSON 01 / POLICIES AND VALUE FUNCTIONS

- * recognize that a policy is a distribution over actions for each possible state.
- * describe the similarities and differences between stochastic & deterministic policies.
- * identify the characteristics of a well defined policy.
- * generate examples of valid policies for a given MDP.
- * describe role of state-value and action-value functions in V .
- * describe the relationship between value functions and policies.
- * create examples of valid value functions for a given MDP.

LESSON 02 / BELLMAN EQUATIONS

- * derive the Bellman Equation for state-value functions.
- * derive the Bellman Equation for action-value functions.
- * understand how Bellman Equations relate current and future values.
- * use the Bellman Equations to compute value functions.

LESSON 03 / OPTIMALITY (OPTIMAL POLICIES & VALUE FUNCTIONS)

- * define an optimal policy.
- * understand how a policy can be at least as good as every other policy in every state.
- * identify an optimal policy for given MDPs.
- * derive the Bellman optimality equation for state-value functions.
- * derive the Bellman optimality equation for action-value functions.
- * understand how Bellman optimality leads relate to previous Bellman eqns.
- * understand the connection between the optimal value fn & optimal policies.
- * verify the optimal value function for given MDPs.

WEEK 02 - PAGES 58-67 OF RL-BOOK → LECTURES

PAGES 68-69 OF RL-BOOK → SUMMARY BEFORE QUIZZES

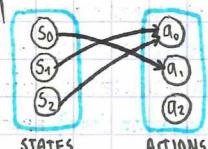
Specifying Policies/

"Reinforcement learning is a problem formulation for sequential decision making under uncertainty. The agent's role in this interaction is to choose an action on each time step. The choice of action has an immediate impact in both immediate reward and the next step. Policies are how an agent selects these actions."

Deterministic Policies/

In simplest case, a policy maps each state to a single action.

$$\pi(s) = a$$



THE AGENT CAN SELECT THE SAME ACTION IN MULTIPLE STATES, AND SOME ACTION MIGHT NOT BE SELECTED IN ANY STATE.

Stochastic Policies/

$$\pi(a|s)$$

$$\sum_{a \in \mathcal{A}(s)} \pi(a|s) = 1$$

THE SUM OVER ALL ACTION PROBABILITIES MUST BE ONE IN EACH STATE NOMATTER THE STATE.

π of a given s , the probability of selecting action a in state s .

multiple actions could be selected with non-zero probability.

IN GENERAL, SET OF ACTIONS \mathcal{A} IS DIFF IN EACH STATE.

$$\pi(a|s) \geq 0$$

EACH ACTION PROB. MUST NON-NEG.

DISTRIBUTION OVER ACTIONS, SPECIFIES A SEPARATE DISTRIBUTION FOR EACH DIFFERENT STATE.

It's important that policies depend only on current state, not in time or prev. states. The state defines all the information used to select the current action, think of this as a requirement on the state, not a limitation on the agent.

* a policy maps the current state onto a set of probabilities for taking each action.

* policies can only depend on the current state, not time or previous states.

Bellman Equations/

"In every day we learn a lot without getting explicit positive or negative feedback. If you're riding a bike and hit a rock, you fall almost stand recover, you may learn to avoid rocks in the future & react more quickly if you do hit one. How do we know hitting a rock is bad even if nothing bad happened this time? We recognize losing balance is bad even without falling/hurting, perhaps we had similar exp. in past. In RL A SIMILAR IDEA ALLOWING TO RELATE THE VALUE OF CURRENT STATE TO VALUE OF FUTURE STATES WITHOUT WAITING TO OBSERVE ALL REWARDS. FORMALIZED IN BELLMAN EQNs."

State-Value Bellman Equation/

"relationship between value of a state and the value of possible successor states"

$$V_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] \rightarrow \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

EXPAND EXPECTED RETURN AS A SUM OF POSSIBLE ACTION CHOICES BY AGENT.

$$\sum_a \pi(a|s) \sum_{s'} p(s'|s,a) [r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s']]$$

EXPAND OVER POSSIBLE REWARDS AND NEXT STATES CONDITIONED ON STATE s .

$$\pi(s') : t+1, \text{ BUT DOESN'T MATTER, AND ACTION } a.$$

BECAUSE NEITHER POLICY OR p DEPEND ON TIME.

Action-Value Bellman Equation/ "the value of a state-action pair in terms of its possible successor state-action pairs"

$$q_\pi(s,a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

$$\sum_{s'} \sum_a p(s'|s,a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']] q_\pi(s',a)$$

CHANGE EXPECTATION TO BE CONDITIONED ON NEXT STATE AND NEXT ACTION AND THEN SUM OVER ALL POSSIBLE ACTIONS.

$$\sum_{s'} \sum_a p(s'|s,a) [r + \gamma \sum_a \pi(a|s') \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s', A_{t+1} = a]]$$

$$V(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

STATE-VALUE FN EXPECTED RETURN FROM STATE.

Recall that...

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

RETURN: DISCOUNTED SUM OF FUTURE REWARDS.

Action-Value Function/ "Expected return if an agent selects action a and then follows policy π "

ACTION-VALUE FN

$$q_\pi(s,a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

...at a given state s ...

VALUE FUNCTIONS ARE CRUCIAL, ALLOW AN AGENT TO QUERY QUALITY OF ITS CURRENT SITUATION, INSTEAD OF WAITING TO OBSERVE LONG-TERM OUTCOME.

THE VALUE FUNCTION SUMMARIZES ALL POSSIBLE FUTURES BY AVG OVER RETURNS.

THEY ALSO ENABLE US TO JUDGE THE QUALITY OF DIFFERENT POLICIES.

Why Bellman Equations / the 'magic' of value functions is we can use them as stand-in for the average of an infinite num of possible futures.

STATE-VALUE BELLMAN EQUATION

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} \sum_v p(s', v | s, a) [r + \gamma V_{\pi}(s')]$$

SUM OVER POSSIBLE ACTION CHOICES MADE BY THE AGENT.

SUM OVER POSSIBLE REWARDS AND NEXT STATES CONDITIONED ON STATE S AND ACTION a.

ACTION-VALUE BELLMAN EQUATION

$$q_{\pi}(s, a) = \sum_{s'} \sum_v p(s', v | s, a) [r + \gamma \sum_a \pi(a|s) q_{\pi}(s', a)]$$

DOES NOT BEGIN WITH POLICY SELECTING AN ACTION, IT IS ALREADY FIXED IN THE S, A PAIR.

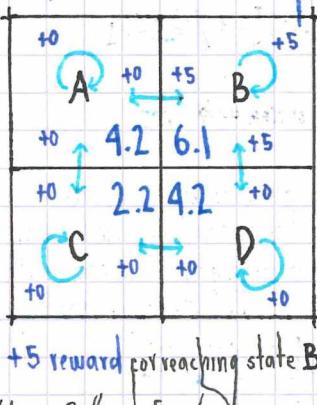
SELECT THE IMMEDIATE REWARD R AND THE NEXT STATE S'.

EXPECTED RETURN ON NEXT STEP IS A SUM OF THE AGENT POSSIBLE ACTION CHOICES, EXPECTATION CONDITIONED ON NEXT STATE, NEXT ACTION, THEN SUM OVER ALL POSSIBLE ACTIONS.

relationship between value of a state and value of possible successor states.

the value of a state-action pair in terms of possible successor state-action pairs.

Gridworld Example /



- * only four states, A, B, C, D on a grid
- * the action space consists of moving up, down, left, right; actions which would move agent off-grid just bounce.
- * uniform random policy, moves in every direction 25% of the time. π .

Using Bellman Equation, we can write down an expression for the value of state A in terms of the sum of the four possible actions and the resulting possible successor states. In this case can be simplified because for each action there's only one possible associated state and reward.

$$V_{\pi}(A) = \mathbb{E}_{\pi}[G_t | S_t = A] = \sum_a \pi(a|s) \sum_{s'} \sum_v p(s', v | s, a) [r + \gamma V_{\pi}(s')]$$

$$V_{\pi}(B) = \mathbb{E}_{\pi}[G_t | S_t = B]$$

$$V_{\pi}(C) = \mathbb{E}_{\pi}[G_t | S_t = C]$$

$$V_{\pi}(D) = \mathbb{E}_{\pi}[G_t | S_t = D]$$

value function is defined as expected return under policy π . It's an avg over the return obtained by each seq. of actions an agent could choose.

$$V_{\pi}(A) = \sum_a \pi(a|s) \text{ for an action } a \quad [r + 0.7 V_{\pi}(s')] \quad \text{there's only one state & reward.}$$

$$V_{\pi}(A) = \frac{1}{4} (5 + 0.7 V_{\pi}(B)) \text{ MOVE RIGHT}$$

$$+ \frac{1}{4} (0 + 0.7 V_{\pi}(C)) \text{ MOVE DOWN}$$

$$+ \frac{1}{4} (0 + 0.7 V_{\pi}(A)) \text{ MOVE LEFT}$$

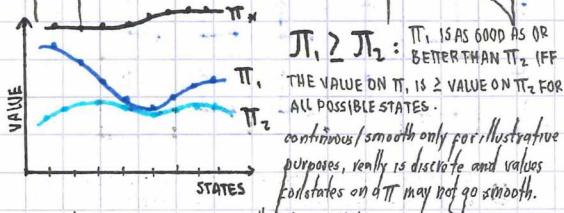
$$+ \frac{1}{4} (0.7 V_{\pi}(A)) \text{ MOVE UP} \quad \frac{1}{2} 0.7 V_{\pi}(A)$$

MANY UNKNOWN, BUILD AN EQUATION SYSTEM TO SOLVE.

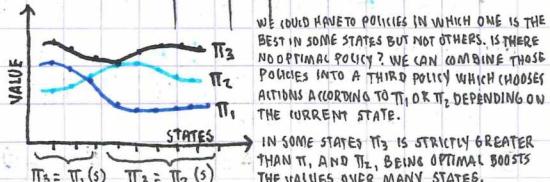
Here, the Bellman Equation reduced an unmanageable infinite sum over possible futures, to a simple linear algebra problem. This may be possible for MDPs of moderate size, in complex problems may end up with 10^{45} (chess) linear equations. However, algorithms based on Bellman Equations can scale up well for large problems.

Optimal Policies /

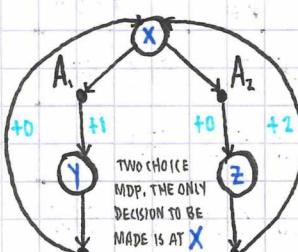
We have talked about a policy as something that is given. It specifies how an agent behaves, given this way of behaving, we then aim to find a value function. But RL does not only seek to evaluate specific policies, but we want to obtain a policy that obtains as much reward as possible in long run.



π^* : optimal policy, good as or better than all the other policies, will have the highest value on every state. There will always exist at least one optimal policy, there can be more than one.



In some states π_3 is strictly greater than π_1 and π_2 , being optimal boosts the values over many states.



* $\pi_1(X)$ offers an immediate reward of +1.

* $\pi_2(X)$ offers a better +2 reward after delay.

* there are only two deterministic policies.

$$\pi_1(X) = A_1, \quad \pi_2(X) = A_2$$

Optimal policy? ... depends on gamma γ

$$\gamma = 0, V_{\pi_1}(X) = 1, V_{\pi_2}(X) = 0, V_{\pi_1}!$$

$$\gamma = 0.9, V_{\pi_1}(X) = 1 + 0.9 \times 0 + 0.9^2 \times 1 + \dots$$

$$5.3 \approx \frac{1}{1 - 0.9^2} \leftarrow \sum_{k=0}^{\infty} 0.9^{2k} \text{ (exclude odd steps)}$$

$$V_{\pi_2}(X) = 0 + 0.9 \times 2 + 0.9^2 \times 0 + \dots V_{\pi_2}!$$

$$2 \sum_{k=0}^{\infty} 0.9^{2k+1} \approx 2 \times 0.9 \times \frac{1}{1 - 0.9^2} \approx 9.5$$

We can only directly solve small MDPs...

⇒ 2 deterministic policies ⇒ brute-force search

General MDP ⇒ 4¹⁸ deterministic policies ⇒ X unfeasible to do brute-force search.

The number of possible policies is equal to the number of possible actions to the number of possible states. Luckily, there is another way to organize our search of the policy space. The solution comes in the form of another set of bellman equations called the bellman optimality equations.

* an optimal policy is defined as a policy with the highest possible value function in all states.

* at least one optimal policy always exists, but there may be more than 1.

* the exponential number of possible deterministic policies makes searching for the optimal policy by brute force intractable.

Optimal Value Functions /

REMEMBER

$$\pi_1 \geq \pi_2 \text{ if and only if } V_{\pi_1}(s) \geq V_{\pi_2}(s) \text{ for all } s \in \mathcal{S}$$

An optimal policy is one that is **as good as** or **better than** any other policy, thus has the **greatest value possible in every state**.

THIS HOLDS FOR EVERY STATE IN OUR STATE SPACE

$$V_{\pi^*}(s) = \mathbb{E}_{\pi^*}[G_t | S_t = s] = \max_{\pi} V_{\pi}(s) \quad \forall s \in \mathcal{S}$$

The value of optimal policy π^* = maximum value over all policies for all states.

OPTIMAL POLICIES ALSO SHARE SAME OPTIMAL ACTION-VALUE FUNCTION

BELLMAN EQUATION FOR STATE-VALUE FUNCTION \Rightarrow BELLMAN OPTIMALITY EQN FOR V^* .

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s' \in \mathcal{V}} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

THIS EQUATION HOLDS FOR ANY POLICY, EVEN THE OPTIMAL POLICY. WE CAN SUBSTITUTE THE OPTIMAL POLICY IN IT AND SIMPLIFY:

$$V^*(s) = \max_a \sum_{s'} p(s', r | s, a) [r + \gamma V^*(s')]$$

because if π is an optimal policy, we can rewrite the equation in a special form which does not reference the policy itself.

There always exist an optimal deterministic policy, one that selects an optimal action in every state, such a deterministic optimal policy will assign probability 1 for an action that achieves the highest value, and probability 0 for all other actions. we can express this by replacing $\sum_a \pi(a|s)$ with a **max over a**.

BELLMAN EQUATION FOR ACTION-VALUE FUNCTION \Rightarrow BELLMAN OPTIMALITY EQN FOR q^* .

$$q_{\pi}(s, a) = \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a'|s') q_{\pi}(s', a')]$$

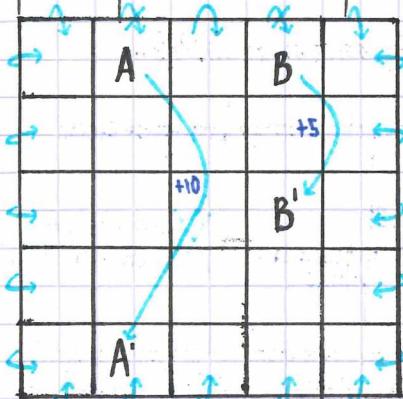
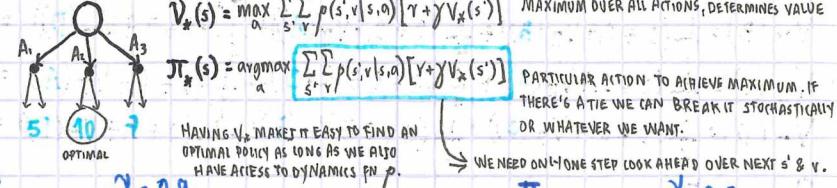
THIS EQUATION HOLDS FOR ANY POLICY, INCLUDING THE OPTIMAL POLICY. WE CAN REMOVE REFERENCES TO POLICY AND SIMPLIFY:

$$q^*(s, a) = \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \max_{a'} q^*(s', a')]$$

The Bellman Optimality Equations relate the value of a state or state-action pair, to its possible successors under **any optimal policy**.

Using Optimal Value Functions to Get Optimal Policies /

The optimal value functions are not the ultimate goal, but we want to find the optimal policy itself, but given an optimal value function it's possible to find the associated optimal policy (one of them).



* all actions on state A transition to state A' with a reward of +10

* all actions on state B transition to state B' with a reward of +5

* all 'wall bounces' transition to the same state with a reward of -1

* any other state transition has no reward.

* the discount factor is $\gamma = 0.9$

	A	B	..	
A				
+10		B'		

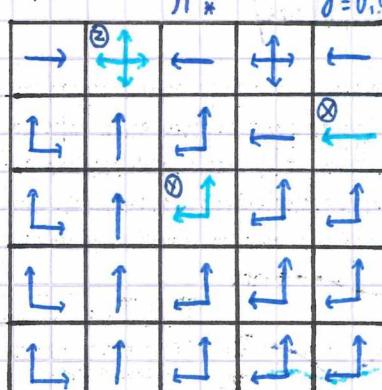
values for all states

* associated optimal values for each step (requires a lot of work to get this table)

* values at bottom not negative or too low, we're not following anymore a random-uniform policy, but an optimal policy gives these values.

* the optimal policy won't choose to bump into walls.

* the optimal value of state A is much higher than the immediate reward of +10.



*: A ONE STEP LOOK AHEAD (CONSIDERS EACH ACTION AND THE POTENTIAL NEXT STATES AND REWARDS. THIS CASE IS SIMPLE BECAUSE EACH ACTION TAKES US DETERMINISTICALLY TO ONE STATE AND ONE REWARD. $(s, a \rightarrow \text{one } s', r)$)

↑: $0 + 0.9 \times 17.5 \quad \leftarrow: 0 + 0.9 \times 17.8 \quad \uparrow: 14.0 \quad \leftarrow: 16.0$

↓: $0 + 0.9 \times 14.4 \quad \rightarrow: -1 + 0.9 \times 16.0 \quad \downarrow: 13.0 \quad \rightarrow: 13.4$

FOR MAXIMIZING ↑ ACTION THE RESULT IS 16.0, WHICH IS EQUAL TO $V^*(s)$.

*: IN THIS STATE, TWO DIFFERENT ACTIONS ↑, ↓ GIVE SAME OPTIMAL VALUE $(0 + 0.9 \times 17.8) = 17.8$. IN THIS STATE THERE ARE TWO DIFFERENT OPTIMAL ACTIONS, AND AN OPTIMAL POLICY IS FREE TO CHOOSE EITHER WITH SOME PROBABILITY.

*: REGARDLESS OF ACTION IN STATE A WE TRANSITION TO A' WITH +10

THE OPTIMAL ACTION IN EACH STATE MUST BE SELECTED BY EVERY OPTIMAL POLICY.

Policies and Value Functions Summary

policies tell an agent how to behave in their environment

DETERMINISTIC POLICIES / state $\xrightarrow{\pi}$ action maps each state to an action.

STOCHASTIC POLICIES / state $\xrightarrow{\pi}$ actions $\pi(s) = a$ maps each state to a distribution over all possible actions.

by definition, a policy depends only on the **current state**, it doesn't depend on me or previous states. This best thought as a restriction on the state, not on the agent.

the goal: to find a policy that obtains as much reward possible
Bellman equations help us evaluate policies, but they do not achieve the ultimate goal of finding an optimally performing policy.

An **optimal policy** achieves the highest val possible in all states.

OPTIMAL STATE-VALUE FN / $V_{\pi^*}(s) = \max_{\pi} V_{\pi}(s) \forall s \in S.$

OPTIMAL ACTION-VALUE FN / $q_{\pi^*}(s, a) = \max_{\pi} q_{\pi}(s, a) \forall s \in S, a \in A.$

value functions capture the future total reward/return under a particular policy.

STATE-VALUE FUNCTION / $V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$ the expected return from current state under a policy.

ACTION-VALUE FUNCTION / $q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$ the expected return from current state if the agent first selects action a and then follows policy π after.

Value functions simplify things by aggregating many possible **future returns** into a single value.

Bellman eqns/ define a relationship between value of a state/state-action pairs and its possible successors.

BELLMAN EQN STATE-VALUE FN / $V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s' \in V} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$ value of current state as a sum of the values of all the successor states & rewards.

BELLMAN EQN ACTION-VALUE FN / $q_{\pi}(s, a) = \sum_{s' \in V} p(s', r | s, a) [r + \gamma \sum_a \pi(a' | s') q_{\pi}(a', s')]$ value of state-action pair as a sum over the values of all possible next state-action pairs and rewards.

the Bellman Equations can be directly solved (e.g. with a **linear system**) to find the value functions. however, the optimal Bellman Equations cannot be solved linearly, due to the use of "max" in them.

LIKE ALL VALUE FN THESE OPTIMAL VN FN HAVE BELLMAN EQN

$V_{\pi^*}(s) = \max_a \sum_{s' \in V} p(s', r | s, a) [r + \gamma V_{\pi^*}(s')]$ these Bellman Equations do not reference a specific policy, this amounts to replacing the policies in the Bellman Equations with a max over all actions.

$q_{\pi^*}(s, a) = \sum_{s' \in V} p(s', r | s, a) [r + \gamma \max_a (q_{\pi^*}(s', a'))]$

Practice Quiz/

* A policy is a function which maps from **states** to **probability distributions** over **actions**.

THIS STATEMENT IS GENERAL BECAUSE IT INCLUDES STOCHASTIC POLICIES, BUT ALSO DETERMINISTIC SINCE THE DISTRIBUTION CAN HAVE A SINGLE ACTION WITH PROBABILITY 1.

* the term **backup** most closely resembles **update**.

* it is true that at least one **deterministic optimal policy** exists in every MDP.

IF WE HAVE A POLICY π_1 THAT PERFORMS BETTER IN SOME STATES, AND A POLICY π_2 THAT PERFORMS BETTER IN OTHERS, A SPliced POLICY π_3 CAN BE FORMED SO THAT DEPENDING ON THE STATE, THE ACTION FROM EITHER π_1 OR π_2 IS TAKEN.

* adding a constant to all rewards **changes** the set of **optimal policies** in episodic tasks

THIS IS BECAUSE CAN MAKE EPISODES LONGER (REDUCED PENALTY) OR ACTIONS LESS ATTRACTIVE DEPENDING ON THE SIGN OF THE CONSTANT.

* adding a constant to all rewards **doesn't change** the **optimal policies** in continuing tasks

SINCE THE TASK IS CONTINUING, THE AGENT WILL ACCUMULATE THE SAME AMOUNT OF EXTRA REWARD W/E BEHAVIOUR.

* the **optimal state-value function** is unique in every finite MDP.

the Bellman optimality equation is actually a system of equations, one for each state, so if there are N states, then there are N equations with N unknowns. If the dynamics of the environment are known, then in principle one can solve this system of equations for the optimal value function using any one of the variety of methods for solving systems of non-linear equations.

All optimal policies share the same optimal state-value function.

* Select the eqn that relates V_{π^*} to q_{π^*} , assume π RANDOM UNIFORM.

- $V_{\pi^*}(s) = \sum_{a, r, s'} \pi(a|s) p(s', r | s, a) [r + q_{\pi^*}(s')]$
- $V_{\pi^*}(s) = \sum_{a, r, s'} \pi(a|s) p(s', r | s, a) [r + \gamma q_{\pi^*}(s')]$
- $V_{\pi^*}(s) = \sum_{a, r, s'} \pi(a|s) p(s', r | s, a) q_{\pi^*}(s')$
- $d) V_{\pi^*}(s) = \max_a q_{\pi^*}(s, a)$

* Select the eqn that maps q_{π^*} to V_{π^*} , using 4-arg fn p .

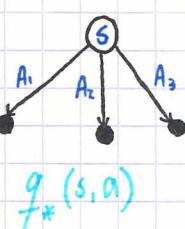
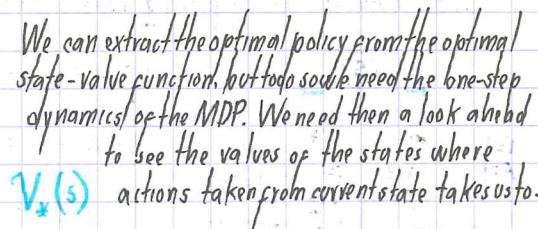
- $q_{\pi^*}(s, a) = \sum_{s', r} p(s', r | s, a) [r + V_{\pi^*}(s')]$
- $b) q_{\pi^*}(s, a) = \sum_{s', r} p(s', r | s, a) \gamma [r + V_{\pi^*}(s')]$
- $c) q_{\pi^*}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi^*}(s')]$

* Give an eqn for some π_{π^*} in terms of V_{π^*} & 4-arg p .

- $\pi_{\pi^*}(a|s) = 1 \text{ if } V_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + V_{\pi^*}(s')] \text{ else } 0$
- $b) \pi_{\pi^*}(a|s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi^*}(s')]$
- $c) \pi_{\pi^*}(a|s) = \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi^*}(s')]$
- $d) \pi_{\pi^*}(a|s) = 1 \text{ if } V_{\pi^*}(s) = \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi^*}(s')] \text{ else } 0$

* Write a policy π_{π^*} in terms of q_{π^*}

- $\pi_{\pi^*}(a|s) = q_{\pi^*}(s, a)$
- $b) \pi_{\pi^*}(a|s) = \max_a q_{\pi^*}(s, a)$
- $c) \pi_{\pi^*}(a|s) = 1 \text{ if } a = \arg\max_{a'} q_{\pi^*}(s, a') \text{ else } 0$



We can get the optimal policy with much less work if we had the optimal action-value function. We simply select the action with the highest value in each state.

Prediction and Control

LESSON 01 / POLICY EVALUATION (PREDICTION)

- * understand the distinction between policy evaluation & control.
- * explain the setting of dynamic programming, also limitations.
- * outline the iterative policy evaluation algorithm for estimating state values under a given policy.
- * apply iterative policy evaluation to compute value functions.

LESSON 02 / POLICY ITERATION (CONTROL)

- * understand the policy improvement theorem.
- * use a value function for a policy to produce a better policy.
- * outline the policy iteration algorithm for finding optimal policy.
- * understand "the dance of policy and value".
- * apply policy iteration to compute optimal policies + value functions.

LESSON 03 / GENERALIZED POLICY ITERATION

- * understand the framework of generalized policy iteration.
- * outline value iteration, an important example of GPI.
- * understand distinction of (a) synchronous programming.
- * describe brute-force for searching an optimal policy.
- * describe monte-carlo for searching a value function.
- * understand the advantage of DP and "bootstrapping" over these alternatives for searching the optimal policy.

WEEK 04 - PAGES 73-88 OF RL BOOK → LECTURES

Policy Improvement

GIVEN π , WE CAN FIND THE OPTIMAL POLICY BY CHOOSING THE GREEDY ACTION, WHICH MAXIMIZES THE BELLMAN OPTIMALITY EQN IN EACH STATE.

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')]$$

Instead of optimal value function select a greedy action with respect to an arbitrary policy V_π .

* the new policy π' must be different than π . If greedification doesn't improve it, then the policy was already greedy respect to V_π .

$$\text{iff } \pi'(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')] \quad \forall s \in S$$

V_π obeys Bellman Optimality EQN, π is already optimal policy.

* the new policy π' is a strict improvement on π unless already optimal.

POLICY IMPROVEMENT THEOREM

$$q_\pi(s, \pi'(s)) \geq q_\pi(s, \pi(s)) \quad \forall s \in S \quad \text{then } \pi' \text{ must be better}$$

↓
an action but taken with diff. policy π'

$\{\pi'$ IS AT LEAST AS GOOD OR BETTER THAN π
IF $\forall s \in S$ THE VALUE OF ACTION SELECTED BY π' IS \geq THE VALUE OF ACTION SELECTED BY π .
STRICTLY BETTER IF VALUE IS $>$ IN AT LEAST ONE STATE.

The policy improvement theorem only guarantees a new policy is an improvement on the original, can't expect to find optimal policy easily.

Policy Evaluation vs Control

Evaluation is the task of determining the value function for a specific policy.

Control is the task of finding a policy to obtain as much reward as possible, finding a policy that maximizes the value function.

Control is the ultimate goal of RL, but the task of Evaluation is necessary 1st step. It's hard to improve our policy if we don't have a way to assess how good it is. Dynamic Programming techniques can be used to solve both these tasks, if we have access to the dynamics function p ($p(s',r|s,a)$).

Iterative Policy Evaluation

Dynamic Programming are obtained by turning Bellman EQN, into update rules.

THE BELLMAN EQUATION FOR STATE-VALUE GIVES US A RECURSIVE EXPRESSION

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')]$$

INSTEAD OF AN EQUATION THAT HOLDS FOR THE TRUE VALUE FUNCTION, USE IT AS AN UPDATE RULE. THIS PROCEDURE ITERATIVELY APPLIED WILL REFINE THE ESTIMATE OF THE VALUE FUNCTION.

arbitrary initialization for approximate value function.
 V_0 : each iteration produces a better approximation using update rule.
 Each iteration applies to $V_{s \in S}$, which is called a sweep. Applying update repeatedly leads to a better approx to the s-v fn V_π .

for any choice of V_0 , the update leaves the approx unchanged, $V(s) = V(s) \forall s \in S$, then $V_k = V_\pi$ and we have found the value function, this is because

$\lim_{k \rightarrow \infty} V_k = V_\pi$. V_π is the unique solution to the Bellman Equation, the only way for V_k to be unchanged after update is if it already obeys Bellman EQN.

We store two arrays, each with an entry for every state. V stores the current approximate value function, the other V' the updated values. Compute new values from old one step at a time, without changing old. At the end of a full sweep compare both arrays, if different, update V' into V & iterate.
 IN PLACE UPDATE IS POSSIBLE & WILL CONVERGE FASTER (BECAUSE SOME STATES WILL USE UPDATED STATE VALUES) AND CONSUME LESS MEMORY, BUT TWO ARRAY IS EASIER TO EXPLAIN.

$R = -1$ forming

$\gamma = 1$ episodic

$\theta = 0.001$

25% 25% π random policy

25% 25% terminal states (NEVER UPDATE IT)

V_0	V_1	V_2	\dots	V_π
0 0 0 0	0 -1 -1 -1	0 -1.7 -2 -2	-14 -18 -20 -20	0 -14 -20 -22
0 0 0 0	-1 -1 -1 -1	-1.7 -2 -2 -2	-18 -20 -20 -18	-14 -18 -20 -20
0 0 0 0	-1 -1 -1 -1	-2 -2 -2 -1.7	-20 -20 -18 -18	-14 -18 -20 -20
0 0 0 0	-1 -1 -1 0	-2 -2 -1.7 0	-22 -20 -14 0	-14 -18 -20 -22

1st sweep 2nd sweep final sweep converge

ITERATIVE POLICY EVALUATION FOR ESTIMATING $V \approx V_\pi$

$$V \leftarrow \bar{V}, V' \leftarrow \bar{V}$$

loop:

$$\Delta \leftarrow 0$$

foreach $s \in S$:

$$V'(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |V'(s) - V(s)|)$$

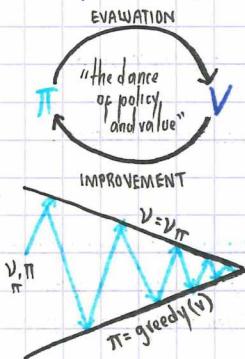
$$V \leftarrow V'$$

$$\text{until } \Delta < \theta \quad (\text{a small positive number})$$

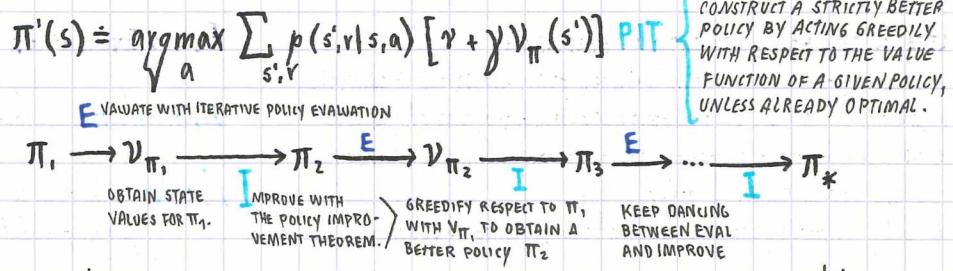
Output: $V \approx V_\pi$.

Policy Iteration, Control/

Policy Iteration/ The value function computed for a given policy can be used to find a better policy. We can find the optimal policy by iteratively evaluating & improving a sequence of policies.



dimensionality ~ geometry of space of policies & value functions is more complicated, but the same intuition holds.



Each policy generated this way is deterministic. There are a finite number of deterministic policies, so this iterative improvement must eventually reach an optimal policy and optimal π_* . Every policy will be greedy with respect to previous value function but not its own if it's not an optimal policy, but if it's an optimal policy then it's greedy with respect to its own value function. Meaning should terminate the algorithm since policy won't improve anymore.

POLICY ITERATION (USING ITER. POL. EVAL) FOR ESTIMATING $\pi \approx \pi_*$

1. INITIALIZATION:

$v(s) \in \mathbb{R}$ and $\pi(s) \in \Delta(s)$ arbitrarily $\forall s \in S$

2. POLICY EVALUATION: make v reflect the value of π

loop: $\Delta \leftarrow 0$

foreach $s \in S$:

$$v \leftarrow v(s)$$

$$v(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma v(s')] \quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation).

3. POLICY IMPROVEMENT: make a greedy π' respect to v .

policy-stable \leftarrow true

foreach $s \in S$:

old-action $\leftarrow \pi(s)$

$$\pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v(s')]$$

If old-action $\neq \pi(s)$, then policy-stable \leftarrow false.

if policy-stable = true, stop and return $V \approx V_*$, $\pi \approx \pi_*$, else go to 2.

Generalized Policy Iteration/

Regular policy iteration is a fairly rigid procedure. Alternates between evaluating the current policy and improving it. The generalized policy iteration framework allows much more freedom, while maintaining the optimality guarantees of converging to an optimal policy approx.

VALUE ITERATION

Policy Iteration Algorithm runs every step all the way to completion. We can imagine relaxing this, each evaluation step brings the estimate a little closer to the value of the current policy, and then a little more greedy but not totally, but should still converge to π_* . In Value Iteration, we still sweep over all states and greedify with respect to value function, but we don't run policy evaluation to completion, we perform a single sweep over all states instead of sweeping until some $\Delta < \theta$.

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')] \quad \text{UPDATE DOES NOT REFERENCE ANY SPECIFIC POLICY, HENCE THE NAME "VALUE ITERATION".}$$

We can recover the optimal policy from this value function with argmax $\pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')] \dots \pi \approx \pi_*$

SYNCHRONOUS DP/

Methods that perform systematic sweeps over all states are called synchronous. This can be bad if the state space is large, but for small very precise & straight forward.

MONTE CARLO METHOD updates states randomly, needs a lot returns to average.
POLICY ITERATION polynomial/time in $|S|$ and $|A|$, very fast. \Rightarrow USES BOOTSTRAPPING/DP STATES
BRUTE FORCE SEARCH exponential/time in $|S| \cdot |A|^{|S|}$ deterministic pols. \Rightarrow ALL COMBINATIONS OF ACTIONS

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

ASYNCHRONOUS DP/

Updates values of states in any order, not systematic, may update some states more, but guarantee the convergence must update every state some time, can update states near those that have changed val.

MONTE CARLO METHOD updates states randomly, needs a lot returns to average.
POLICY ITERATION polynomial/time in $|S|$ and $|A|$, very fast. \Rightarrow USES BOOTSTRAPPING/DP STATES
BRUTE FORCE SEARCH exponential/time in $|S| \cdot |A|^{|S|}$ deterministic pols. \Rightarrow ALL COMBINATIONS OF ACTIONS

T	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	T

T	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	T

T	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	T

T	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	T

Practice Quiz/

* the value of any state under an optimal policy is greater than or equal to the value of a state under a non-optimal policy.

* It is false that if a policy is greedy with respect to the value function of the uniform random policy then it's guaranteed to be optimal policy.

* It is true that if a policy is greedy with respect to its own value function then it is the optimal policy.

* If π' is greedy respect to v_π , then $v_\pi \geq v_{\pi'}$.

* If $v_\pi = v_{\pi'}$, then π and π' may not be equal, they could have different tie-breaking mechanisms.

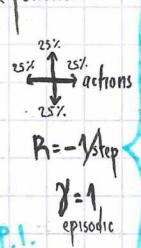
* VALUE ITERATION and POLICY ITERATION are special cases of G.P.I.

* it is false that all GENERALIZED POLICY ITER. algorithms are synchronous.

* as described, VALUE ITERATION & POLICY ITERATION are synchronous.

* ASYNCHRONOUS methods scale better than SYNCHRONOUS methods.

* DP algos are planning methods because they use a model to impr. policy.



$$R = -1 \text{ step}$$

$$\gamma = 1 \text{ episodic}$$

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	11
12	13	14	T

T	1	2	3
4	5	6	7
8	9	10	