

Reinforcement learning

- Episodic vs Continuity tasks →

* Episodic → Interaction ends after some time step T .

$s_0, a_0, \dots, r_t, s_T$ will define ending point.

* Continuous → Interaction continuous without limit.

$s_0, a_0, \dots, r_t, s_t, \dots$ Neverending task.

- Reward →

* Agent should maximize reward over time, not immediate.

$$G_t = R_{t+1} + \dots + R_{t+1000} + \dots \leftarrow \text{Return.}$$

* Discounted return →

- Discounting rewards over time, the more immediate are more likely to happen.

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

- Markov Decision Process → MDPs

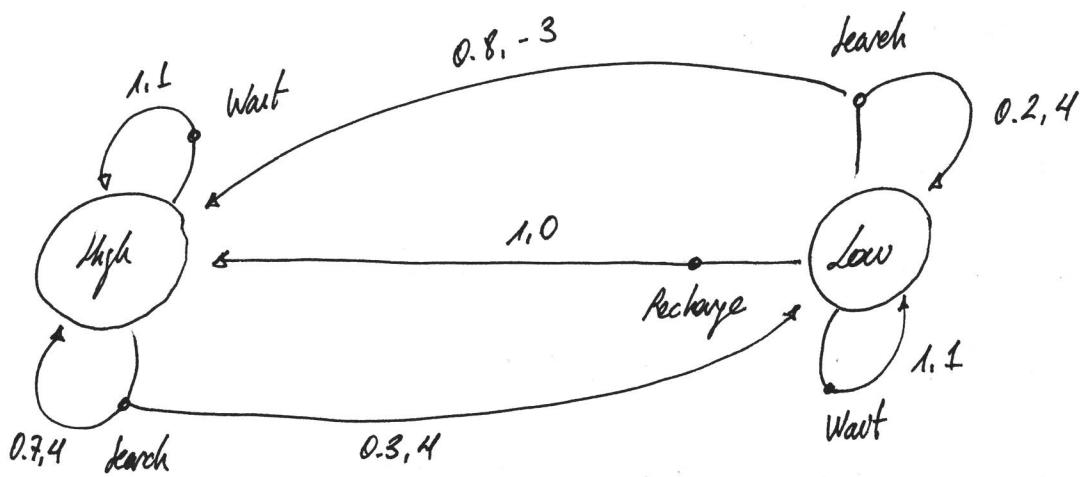
$S \equiv$ All non-terminal states

$S^* \equiv$ All states, including terminal (Episodic task).

$A \equiv$ Action space of all possible actions.

$A(s) \equiv$ set of actions available in state $s \in S$

$$A = \begin{bmatrix} \text{Search} \\ \text{Recharge} \\ \text{Wait} \end{bmatrix} \quad S = \begin{bmatrix} \text{High} \\ \text{Low} \end{bmatrix}$$



* Environment doesn't consider any previous steps/actions, it has access →

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

One-step dynamics

Function $p: S \times R \times S \times A \rightarrow [0, 1]$ ordinary deterministic function.

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1, \text{ for all } s \in S, a \in A(s).$$

* Total transition probability →

$$p(s' | s, a) = P(S_t = s_{t+1} | S_t = s, A_t = a) = \sum_{r \in R} p(s', r | s, a)$$

* Expected reward for state-action path →

$$r: S \times A \rightarrow R$$

$$r(s, a) = E[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_{r \in R} r \sum_{s' \in S} \frac{p(r | s, a)}{p(s' | s, a)}$$

* Expected reward for state-action-next-state triple →

$$r: S \times A \times S \rightarrow R$$

$$r(s, a, s') = E[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in R} r \frac{p(s' | r | s, a)}{p(s' | s, a)}$$

- Policies \rightarrow

- * mapping from set of states to set of possible actions per stage.
- * Agent will need to learn for the state which are the appropriate actions.
- * Deterministic policy \rightarrow

$$\pi: S \rightarrow A$$

$$\pi(\text{low}) = \text{Recharge}$$

$$\pi(\text{high}) = \text{Search}$$

- * Stochastic policy \rightarrow

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

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

- Probability that agent takes action a while in state s .

$$\pi(\text{recharge} | \text{low}) = 0.5$$

$$\pi(\text{search} | \text{high}) = 0.9$$

$$\pi(\text{wait} | \text{low}) = 0.4$$

$$\pi(\text{wait} | \text{high}) = 0.1$$

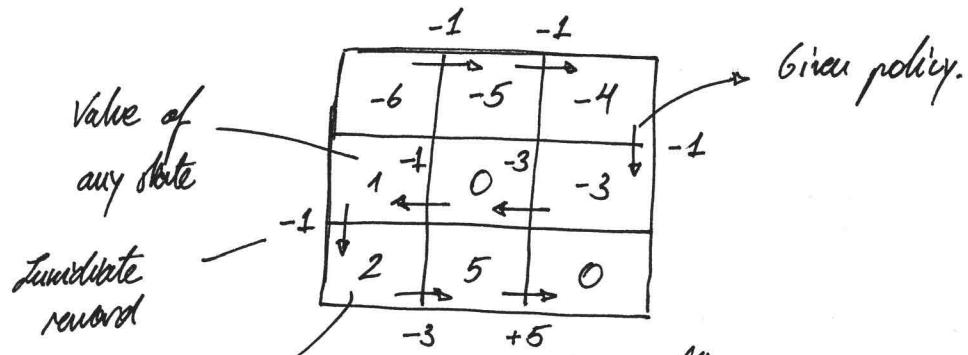
$$\pi(\text{search} | \text{low}) = 0.1$$

- State - Value function →

- * For each state, state-value function yields expected return, if the agent started in that state, and follow the policy for all the steps.

$E_{\pi}[\cdot] = \text{Expected value of a random variable.}$

$$V_{\pi}(s) = E_{\pi} [G_t | s_t = s] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | s_t = s \right], \text{ for all } s \in S$$

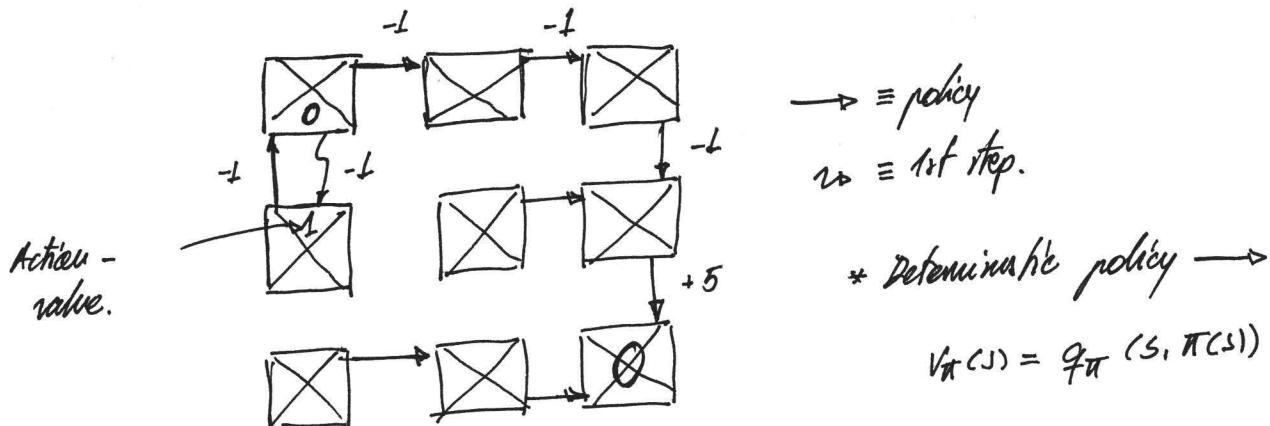


discounted value of state follows.

- Action - value function →

- * For each state, action ; the action-value function yields expected return, if the agent started in that state, does the action, and then follows the policy for all next steps.

$$q_{\pi}(s, a) = E_{\pi} [G_t | s_t = s, A_t = a] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | s_t = s, A_t = a \right], \text{ for all } s \in S, a \in A(s)$$



- Bellman Expectations Equations \rightarrow

* Value - state \rightarrow

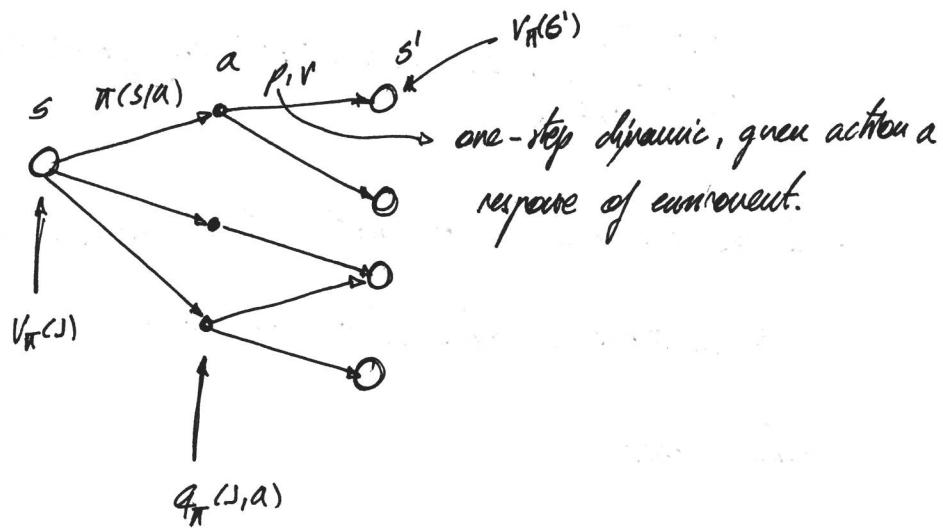
$$V_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] =$$

$$= \sum_{a \in A(s)} \pi(a|s) \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma V_{\pi}(s'))$$

* Action - Value \rightarrow

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, a_{t+1}) | s_t = s, A_t = a] =$$

$$= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma \sum_{a' \in A(s')} \pi(a'|s') q_{\pi}(s'|a'))$$



- Bellman Optimality Equations \rightarrow

* State - value \rightarrow

$$V_t(s) = \max_{a \in A(s)} E[R_{t+1} + \gamma V_{t+1}(s_{t+1}) | s_t = s] =$$

$$= \max_{a \in A(s)} \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma V_{t+1}(s'))$$

- Expresses the value of the state under the optimal policy, in terms of the successor states.

* Action - value \rightarrow

$$q_{\pi^*}(s, a) = E[R_{t+1} + \gamma \max_{a' \in A(s_{t+1})} q_{\pi^*}(s_{t+1}, a') | s_t = s, A_t = a] =$$

$$= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma \max_{a' \in A(s')} q_{\pi^*}(s', a'))$$

- Expresses the value of the state-action pair s, a under an optimal policy in terms of values of the successor state-action pairs under optimal policy.

* Optimal policy \rightarrow

π^* satisfies $\pi^* > \pi$ for all π .

$\pi^* \geq \pi$ if and only if $V_{\pi^*}(s) \geq V_{\pi}(s)$ for all $s \in S$.

- Denoting Bell Equations \rightarrow

* Probability space \rightarrow

Mathematical construct that models a real-world process consisting of states that occur randomly.

Elements:

Ω = sample space, set of possible outcomes.

\mathcal{F} = set of events, each event is a set containing 0 or more outcomes.

P = Assignment probabilities to the events.

* Law of total expectation \rightarrow

X = random variable of probability space (Ω, \mathcal{F}, P)

Y = random variable of probability space (Ω, \mathcal{F}, P)

$$E[X] = E[E[X|Y]]$$

- If the set of values is finite for X, Y . It can be proved as \rightarrow

$$E[X] = \sum_i E[X|A_i] P(A_i)$$

$$E[E[X|Y]] = E \left[\sum_x x P(X=x|Y) \right] = \sum_y \sum_x x P(X=x|Y) \cdot P(Y=y)$$

$$= \sum_y \sum_x x P(X=x, Y=y) = \sum_x x \sum_y P(X=x|Y=y) = \sum_x x P(X=x)$$

= $E[X]$ If they are finite, we can switch numbers.

* Shatty point →

law of total expectation

$$\begin{aligned}
 q_{\pi}(s, a) &= \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a] = \left[E[X] = \sum_i E[X|A_i] P(A_i) \right] = \\
 &= \sum_{s' \in S} \sum_{r \in R} P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a) \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r] \\
 &= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r] = \\
 &\quad \left[\mathbb{E}_{\pi} [G_t | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r] = \mathbb{E}_{\pi} [G_t | S_{t+1} = s', R_{t+1} = r] \right] \\
 &= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) \mathbb{E}_{\pi} [G_t | S_{t+1} = s', R_{t+1} = r] = \\
 &= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_{t+1} = s', R_{t+1} = r] = \\
 &= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) \left[r + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s'] \right] = \\
 &\quad \begin{array}{l} \text{--- } G_{t+1} \text{ doesn't depend on } R_{t+1}, \\ \text{--- } \text{dropped.} \end{array} \\
 &= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) \left[r + \gamma V_{\pi}(s') \right]
 \end{aligned}$$

$$V_{\pi}(s) = \sum_{a \in A(s)} \pi(a|s) q_{\pi}(s, a)$$

→ With these two
State-value, action-value
functions.

* optimal equations \rightarrow

$$q_*(s, a) = \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma v_*(s'))$$

$$v_*(s) = \max_{a \in A(s)} q_*(s, a)$$

$$\pi_*(s) = \arg \max_{a \in A(s)} q_*(s, a)$$

for all $s \in S$.

{ with these two,

optimal state-value, action-value
functions can be derived.

Dynamic Programming

- Policy evaluation →

- * Assumption that the environment is finite MDP. S, A, R are finite and its dynamics are given by $p(s', r | s, a)$ where $s \in S, a \in A(s), r \in R$, and $s' \in S^+$
- * DP can be applied to continuous state and value spaces, but an exact solution is not always guaranteed.

* Main idea →

The value functions are used to organize and structure the search for good policies.

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

- The existence and uniqueness for $V_{\pi}(s)$ is guaranteed as long as $\gamma < 1$ and eventual termination is guaranteed from all states under policy π .

- If environmental dynamics are known $p(s', r | s, a) \rightarrow$
151 simultaneous linear equations in 151 unknown
To compute this, it is really expensive.

- Iterative concept \rightarrow

1. Assign first value to state-values \rightarrow

$$V_0(s_i) = x \quad \leftarrow \text{Policy at iteration } 0.$$

2. For each iteration, update state-value function \rightarrow

$$V_{K+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r + \gamma V_K(s')]$$

for all $s \in S$.

$V_{K+1}(s)$ \equiv State-value function under policy $K+1$, approximation to π with iter.

3. If $K \rightarrow \infty$, $V_K = V_\pi$. $\{V_K\}$ sequence will converge to $V_\pi \Rightarrow$
Iterative policy evaluation.

- To compute this \rightarrow

* Two arrays, $V_K(s)$ and $V_{K+1}(s)$

* One array, update values in the same array.

Second option converges faster, since it uses the updated values as soon as they are available.

- Algorithm \rightarrow

Inputs: MDP, policy π , max iteration number θ .

Output: State-value function $\equiv V_\pi(s)$

Initialize $V(s)$ for all $s \in S^+$ arbitrarily, except $V(\text{terminal}) = 0$.

Repeat until $\Delta < \theta$:

$$\Delta \leftarrow 0$$

for $s \in S$:

$$V \leftarrow V(s)$$

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

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

- Estimation Action Values →

$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma V_{\pi}(s_{t+1}) / s_t = s, A_t = a]$$
$$= \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

Input: State-value function V

Output: Action-value function Q

For $s \in S$:

For $a \in A(s)$:

$$q(s, a) \leftarrow \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma V(s'))$$

End

End

Return q .

- Policy Improvement →

* Now to find a better policy? Consider state s , action a and then follow policy π' .

$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma V_{\pi}(s') / s_t = s, A_t = a]$$

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

* If $q_{\pi}(s, \pi'(a)) \geq V_{\pi}(s) \rightarrow$

$$V_{\pi}(s) \leq q_{\pi}(s, \pi'(s)) = \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_{\pi}(s') / s_t = s, A_t = \pi'(s)] =$$

deterministic
policy. →

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_{\pi}(s') / s_t = s] \leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s)) / s_t = s] =$$
$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma \mathbb{E}_{\pi'} [R_{t+2} + \gamma V_{\pi}(s_{t+2})] / s_t = s] = \text{threshold until all policies}$$

followed by π'

Rewards under
policy π'

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(s_{t+2}) / s_t = s] \leq \dots = V_{\pi'}(s)$$

* New policy \rightarrow

$$\pi'(s) = \arg \max_a q_{\pi}(s, a)$$

- * If $V_{\pi}(s) = V_{\pi'}(s)$, that mean the policy reaches optimal.
- * This case was for deterministic policy, it can be shown for stochastic also.

* Algorithm \rightarrow

Input: MDP, value function V

Output: Policy π'

For $s \in S$ do

 For $a \in A(s)$ do

$$Q(s, a) \leftarrow \sum_{s' \in S} p(s', r|s, a) (r + \gamma V(s'))$$

End

$$\pi'(s) \leftarrow \arg \max_{a \in A(s)} Q(s, a)$$

End

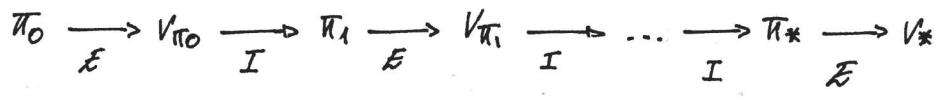
Return π'

- * If some state has more than one action with same $Q(s, a) \rightarrow$

$$\pi'(a|s) = 0 \text{ if } a \notin \arg \max_{a' \in A(s)} Q(s, a')$$

Other, foll prob.

- Policy Iteration →



$E \equiv$ Iterative policy evaluation.

$I \equiv$ Policy improvement

* Algorithm →

Input: MDP, small R^+ .

Output: policy $\pi \approx \pi_*$

Because we have a finite MDP,
it has a finite number of policies.
This process will converge to the optimal.

Initialization policy → $\pi(a|s) = 1/|A(s)|$ for all $s \in S$ and $a \in A(s)$.

Repeat

$V(s) \leftarrow$ Policy evaluation.

$\pi' \leftarrow$ Policy improvement.

If $\pi = \pi'$ then break

$\pi \leftarrow \pi'$

return π

- Truncated Policy Iteration →

* on policy valuation → Instead of θ , replace with fix number of iterations.

Optimal Action-Value

	a_1	a_2	a_3
s_1	1	② -3	
s_2	-2	1	③
s_3	⑤	4	-1

Incorrect estimate Action-Value

	a_1	a_2	a_3
s_1	0	④ -2	
s_2	-3	-2	⑥
s_3	②	1	0

same policy optimal.

$$\pi_*(s_1) = a_2$$

$$\pi_*(s_2) = a_3$$

$$\pi_*(s_3) = a_1$$

$$A = \{a_1, a_2, a_3\}; S \in \{s_1, s_2, s_3\}$$

I can start with the estimate of
I reach same optimal policy.

- Value Iteration \rightarrow save only one line

Policy Iteration \rightarrow for $s \in S$:
$$V(s) \leftarrow \sum_{s'} \sum_r p(s', r | s, \pi(s)) (r + \gamma V(s'))$$

Policy Improvement \rightarrow for $s \in S$:

for $a \in A(s)$:

This can be combined $\left(\begin{array}{l} Q(s, a) \leftarrow \sum_{s'} \sum_r p(s', r | s, a) (r + \gamma V(s')) \\ \pi'(s) \leftarrow \arg \max_{a \in A(s)} Q(s, a) \end{array} \right)$

Policy Iteration \rightarrow for $s \in S$:

$$V(s) \leftarrow \sum_{s'} \sum_r p(s', r | s, \pi(s)) (r + \gamma V(s'))$$

Policy Improvement \rightarrow for $s \in S$:

$$\pi'(s) \leftarrow \arg \max_{a \in A(s)} \frac{f(a)}{\sum_{s'} \sum_r p(s', r | s, a) (r + \gamma V(s'))}$$

Policy Iteration \rightarrow for $s \in S$:

$$V(s) \leftarrow \frac{f(\pi(s))}{\sum_{s'} \sum_r p(s', r | s, \pi(s)) (r + \gamma V(s'))}$$

$$V(s) \leftarrow \arg \max_{a \in A(s)} \frac{\sum_{s'} \sum_r p(s', r | s, a) (r + \gamma V(s'))}{f(a)}$$

* Algorithm →

Input: MDP, small positive number θ .

Output: policy π^*

$$V(s) \leftarrow 0 \text{ for all } s \in S^+$$

Repeat until $\Delta < \theta$:

$$\Delta \leftarrow 0$$

loop over
states to
converge to
the state-value
function.

for $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_{a \in A(s)} \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma V(s'))$$

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

for $s \in S$:

$$\pi(s) = \arg\max_{a \in A(s)} \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) (r + \gamma V(s'))$$

Return π

↑
Obtain policy correspondingly to the final
value-function.

Monte Carlo Methods

- Restatement of goal →
 - * For any episode $s_0, a_0, \dots, r_n, s_n, \dots, s_T$; the agent's goal is to maximize the reward →
- Find best policy π to max $E_{\pi} \left[\sum_{t=1}^T R_t \right]$
- Prediction problem →
 - * Here the environment dynamics are unknown. Given policy π , determine V_{π} (or q_{π}) by interacting with the environment (Monte Carlo).
 - * Off-policy →
 - One policy to evaluate, another to interact with the environment.
 - Generates episodes following policy δ , where $\delta \neq \pi$. Then use episodes to evaluate π .
 - * On-policy →
 - Agent interacts with the environment following policy π . From the generated episodes, estimate π .

- On-policy \rightarrow

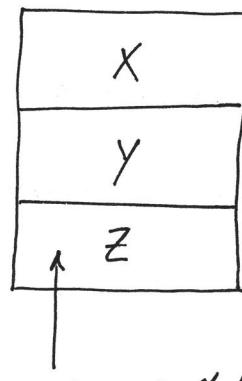
* MC predictions \rightarrow

- Ancestry through the sampled returns for each state would give the expected return, calculated by sampling as long as we have enough visits to the state.

- State values \rightarrow

$$S^+ = \{x, y, z\}$$

$A = \{\uparrow, \downarrow\}$; Actions are not deterministic



Terminal state.

- Evaluate policy through episodes generated following it \rightarrow

$$\pi(x) = \uparrow, \quad \pi(y) = \downarrow$$

Episodes

1. $x, \uparrow, -2, y, \downarrow, 0, y, \downarrow, 3, z.$

2. $y, \downarrow, 2, y, \downarrow, 1, y, \downarrow, 0, z.$

3. $y, \downarrow, 1, x, \uparrow, -3, y, \downarrow, 3, z.$

$V_{\pi}(x) = E_{\pi} [G_t | s_t = x] \equiv$ Expected return after state x .

$$V_{\pi}(x) = \frac{\text{Returns}}{\text{Occurrences.}}$$

* First-Visit MC Method →

- Only consider the first visit per episode.

$$V_{\pi}(x) = \frac{1+0}{2}$$

* Every-Visit MC Method →

- Consider all the visits in the same episode.

$$V_{\pi}(x) = \frac{\overline{3+3} + \overline{3+1+0} + \overline{1+3}}{7} = \frac{14}{7} = 2.$$

* Properties →

- Both methods converge to $V_{\pi}^{(1)}$ as # visits converge to infinite.
- First-visit: each return is an independent event, identically distributed with finite variance of $V_{\pi}(s)$. By "law of large numbers" the sequence of averages of these estimates converges to the expected value.
- First-visit unbiased, Every-visit biased.
- Initially, MSE lower of every visit, with episodes occurred first-visit.

* Algorithm \rightarrow

Input: Policy π , num_episodes.

Output: Value-function v ($\approx v_\pi$ if num_episodes large enough).

Initialise $N(s) = 0$ for all $s \in S$.

Initialise return_nu(s) = 0 for all $s \in S$.

for $i=1$ to num_episodes:

 | Generate episode following π .

 | for $t=0$ to $T-1$:

 | if s_t is a final unit (with return b_t) then

 | | $N(s_t) \leftarrow N(s_t) + 1$

 | | $return_nu(s_t) \leftarrow return_nu(s_t) + b_t$

 | end

 | end

end

$v(s) \leftarrow return_nu(s)/N(s)$ for all $s \in S$.

- Action-value function \rightarrow

- * If a model is available, state-value functions are enough to determine a policy \rightarrow just look to the next step that maximizes reward.
- * Without the model, we need to explicitly estimate the value of each action, in order for the values to be useful in suggesting a policy.
- * Determine policy looking at each state-action pair \rightarrow

Episode

1. $x, \uparrow, -2, y, \downarrow, 0, y, \downarrow, 3, z.$

2. $y, \uparrow, 2, y, \downarrow, 1, y, \downarrow, 0, z.$

3. $y, \downarrow, 1, x, \uparrow, -3, y, \downarrow, 3, z.$

- * MC prediction is not used with a deterministic policy \rightarrow

$$\begin{array}{l} \pi(x) = \uparrow \\ \pi(y) = \downarrow \end{array} \quad \left| \quad \rightarrow q_{\pi}(x, \downarrow) = ? ; q_{\pi}(y, \uparrow) = ? \right.$$

These won't ever be used.

- * To compare alternatives we need to estimate the value of all the actions for each state, not the ones we currently favor.

* Instead we stochastic. A number of episodes, needs to pour converge to 00.
 the number of visits will approximate to the prob. distribution forced
 on M.C.

* First-visit, Every-visit.

* Algorithm →

Input: Policy π , max-episodes.

Output: Action-value function $Q(s, a)$ ($\approx q_\pi$ if max-episodes large).

Initialise $N(s, a) = 0$ for all $s \in S, a \in A(s)$.

Initialise return- $mu(s, a) = 0$ for all $s \in S, a \in A(s)$.

For $i=1$ to max-episodes:

 Generate episode s_0, A_0, \dots, s_T using π .

 For $t=0$ to $T-1$:

 If (s_t, a_t) is a first visit (with return G_t) then

$N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$

 return- $mu(s_t, a_t) \leftarrow return-mu(s_t, a_t) + G_t$

 End

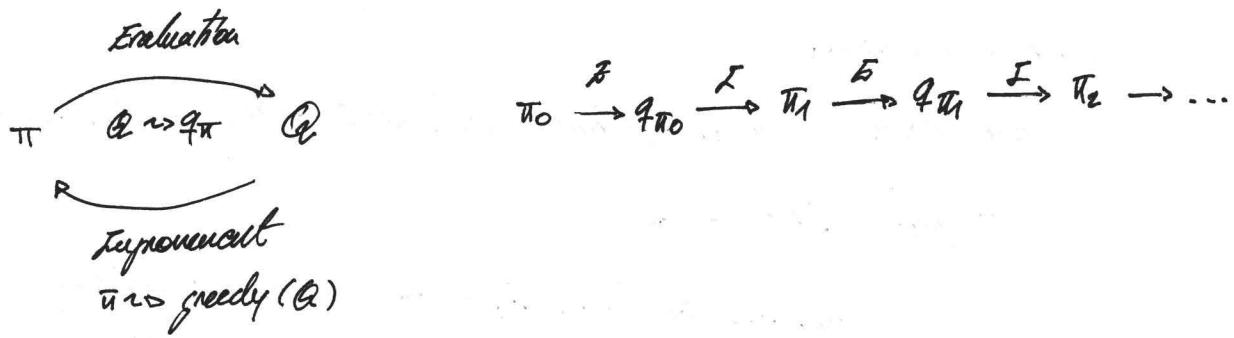
 End

End

$Q(s, a) \leftarrow return-mu(s, a) / N(s, a)$ for all $s \in S, a \in A(s)$.

return Q .

- Value Iteration



* Policy improvement \rightarrow greedy policy $\pi(s) = \arg \max_a q(s, a)$

$$q_{\pi_{K+1}}(s, \pi_{K+1}(s)) = q_{\pi_{K+1}}(s, \arg \max_a q_{\pi_K}(s, a)) = \max_a q_{\pi_K}(s, a)$$

$$\geq q_{\pi_K}(s, \pi_K(s)) \geq V_{\pi_K}(s)$$

Each π_{K+1} is uniformly better than π_K , or as good.

* This is based in two assumptions \rightarrow converges to optimal policy.

1. Episodes have exploratory starts, so you can get samples for all state, actions instead of greedy greedy, which might just give you a local optimum.

2. Policy evaluation is done over an infinite number of episodes.

* Policy evaluation problem \rightarrow

- As in dynamic programming, two approaches, do policy evaluation to a number of iterations until we converge up to a level of approximation for the action and state-value functions.

- Do as in DP value-iteration, after each episode evaluate policy

* Incremental mean \rightarrow

- Run n episodes.

- Keep track of reward following each state-value (action-value) not per episode \rightarrow

$x_1, x_2, x_3, \dots, x_n \rightarrow$ reward per episode
for a state.

- AC situation \rightarrow

$$\text{action-value} = \mu_n = \frac{\sum_{j=1}^n x_j}{n}$$

- Incremental mean, update estimate after each unit \rightarrow

$$\mu_1 = x_1$$

$$\mu_2 = \frac{x_1 + x_2}{2}$$

$$\mu_K = \frac{x_1 + x_2 + \dots + x_K}{K}$$

$$\mu_n = \frac{\sum_{j=1}^n x_j}{n}$$

$$\mu_K = \frac{1}{K} \sum_{j=1}^K x_j = \frac{1}{K} (x_K + \mu_{K-1}(K-1)) \rightarrow$$

$$\mu_K = \mu_{K-1} + \frac{1}{K} (x_K - \mu_{K-1})$$

\hookrightarrow proportional update.

\hookrightarrow previous mean

- Algorithm \rightarrow

$$\mu \leftarrow 0$$

$$k \leftarrow 0$$

while $k < n$:

$$n \leftarrow k+1$$

$$\mu \leftarrow \mu + \frac{1}{n} (x_n - \mu)$$

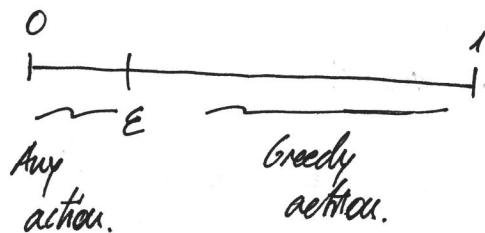
* Overall, you run policy evaluation & improvement so on. This concept can be used to cut all the iterations for the evaluation and update immediate.

* Epsilon-Greedy Policy \rightarrow

$$\pi(a|s) \leftarrow \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A(s)|} & \text{if } a = \arg \max_{a' \in A(s)} Q(s, a') \\ \frac{\epsilon}{|A(s)|} & \text{Otherwise.} \end{cases}$$

$|A(s)| \equiv$ Number of possible actions from state $s \in S$.

- ϵ $\hat{=}$ More likely to pick the non-greedy action.



* Exploration / Exploitation dilemma →

- One potential solution is to implement in a way that gradually reduces ϵ .
- Start favoring exploration, along with the agent knowing the environment dynamics. Limit exploration →
 - * Start with explicable random policy. → $\epsilon = 1$.
 - * As time, agent becomes more greedy.
- * Greedy in the limit with infinite exploration = GLIE →
 - To ensure convergence to optimal policy π^*
 1. Condition → every state-action is visited infinitely many times.
 2. Condition → Policy converges to a policy that is greedy with respect to action-value estimate Q .
 - To satisfy these conditions →
 - * $\epsilon_i \equiv \epsilon$ at i th time step.
 - * $\epsilon_i > 0$ for all time steps i .
 - * ϵ_i decays to 0 in the limit as the time step i approaches infinite → $\lim_{i \rightarrow \infty} \epsilon_i = 0$.
 - E.g. → $\epsilon_i = 1/i$

- MC Control: GLIE Algorithm \rightarrow

Input: positive integer num-episodes.

Output: policy π ($\approx \pi_*$ if num-episodes is large enough).

Initialise $Q(s_t, a) = 0$; $N(s_t, a) = 0$ for all $s \in S, a \in A(s)$.

for $i = 1$ to num-episodes:

$$\epsilon \leftarrow 1/i$$

$$\pi \leftarrow \epsilon\text{-greedy } Q$$

generate episode using π .

you don't even know Q at the beginn. act on fly following policy to estimate Q

\hookrightarrow After calculate policy

for $t=0$ to $T-1$:

if (s_t, a_t) is a first visit (with return g_t)

$$N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (g_t - Q(s_t, a_t))$$

end

end

and

return π .

- MC control: constant-alpha \rightarrow

"forgetful mean"

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$$

Intermediate return
Expected value of action-value

* learning on $Q(S_t, A_t)$
with trace decays matter

$G_t = \text{total return}$

Update by this difference

* Instead use α factor \rightarrow

"runny mean"

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (G_t - Q(S_t, A_t))$$

* With this, most later rewards are discounted

* setting α values \rightarrow

- $\alpha=0 \rightarrow$ never updated.

$\} \rightarrow \alpha \in (0, 1)$

- $\alpha=1 \rightarrow$ immediate return

- smaller values of α encourages the agent to consider a longer history, increasing more recent.

- If α is too close to one it may prevent for the agent to find the π_α .

- MC control: GVE K-constant \rightarrow

Input: mu-episodes, κ

Output: policy π

Initialize Q arbitrarily (e.g. $Q(s,a) = 0$ for all $s \in S, a \in A(s)$).

for $i=1$ to mu-episodes:

$$\epsilon \leftarrow 1/i$$

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

generate episode following π .

for $t=0$ to $T-t$:

if (s_t, a_t) is a first-visit:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \kappa (b_t - Q(s_t, a_t))$$

end

end

end

return π .

Temporal-Difference Methods

- Temporal-Difference prediction \rightarrow

- * Given policy π , determine v_π (or q_π).

- * Focus on update step \rightarrow

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (G_t - Q(s_t, a_t))$$

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t))$$

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

$$= \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(s_{t+1}) | s_t = s]$$

Bellman equation.

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

- * Touches the Bellman equation expectation and the MC estimation.

- * Now, the state-value doesn't depend on the final state and reward, but just on the next one.

- * We can update immediately.

- * In practice, TD converges faster than MC.

- * TD is guaranteed to converge as long as the step-size is small enough.

* Summary DP, MC and TD →

$$r_{\pi}(s) = \mathbb{E}_{\pi} [G_t | s_t = s] \quad (1)$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | s_t = s]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s]$$

- MC estimates function (1), it is an estimation because the expectation is unknown and it is calculated from samples.

- Dynamic programming is an estimate because π_{θ} is unknown and it is approximated using $V(\cdot)$

- Temporal Difference is a difference estimate for both reasons, it samples the expected values and it uses the current state V instead of π_{θ} .

* TD error →

$$\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

- learning difference between estimated $V(s_t)$ and the one given by iteration, reward.

* AC error is the accumulated TD error, if the $V(t)$ are not updated until the end of the episode.

$$\begin{aligned}
 \underbrace{G_t - V(s_t)}_{\text{AC Error}} &= R_{t+1} + \gamma G_{t+1} - V(s_t) + \gamma V(j_{t+1}) - \gamma V(j_{t+1}) = \\
 &= R_{t+1} + \gamma V(j_{t+1}) + \gamma(G_{t+1} - V(j_{t+1})) = \\
 &\quad \underbrace{-V(s_t)}_{\text{AC Error}} \\
 &= s_t + \gamma(G_{t+1} - V(j_{t+1})) = s_t + \gamma j_{t+1} + \gamma^2(G_{t+2} - V(j_{t+2})) \\
 &= s_t + \gamma s_{t+1} + \gamma^2 s_{t+2} + \dots + \gamma^{T-t} \underbrace{(G_t - V(s_t))}_{0-0} = \\
 &= \sum_{k=t}^{T-1} \gamma^{k-t} s_k
 \end{aligned}$$

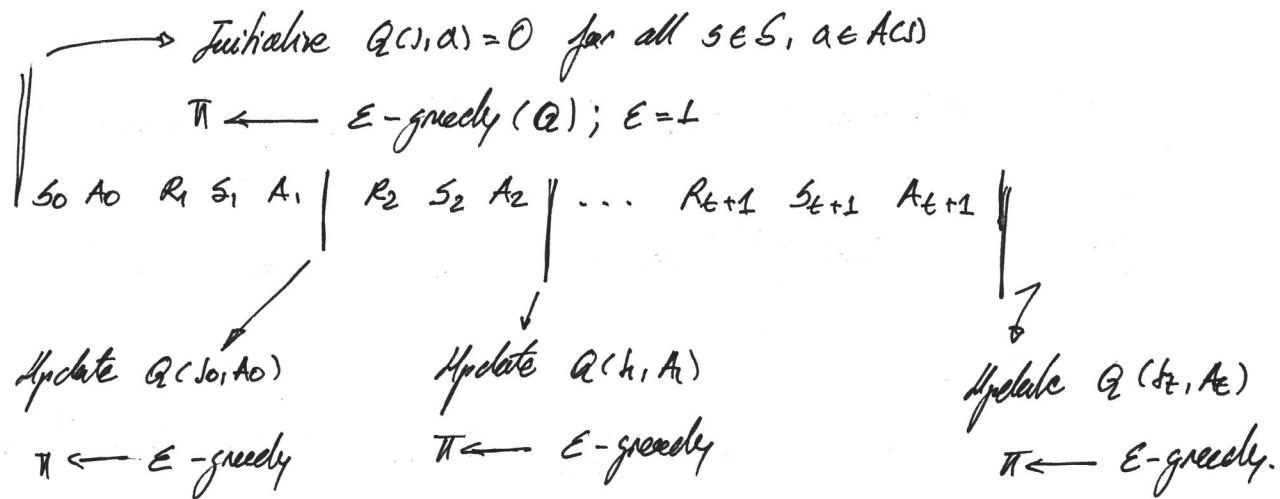
* the identity is not exact if $V(t)$ is updated during the episode.

But if the step-size α is small enough it may hold an approximation.

$$Q(s_0, a_0) \leftarrow Q(s_0, a_0) + \alpha (R_1 + \gamma(Q(s_1, a_1) - Q(s_0, a_0)))$$

- Sarsa \rightarrow

* Determine π_* \rightarrow



- Guaranteed to converge to optimal action-value function as long as α is small enough, and GLIE conditions are met \rightarrow

1. $\epsilon_i > 0$ for all time-step i .

2. $\lim_{i \rightarrow \infty} \epsilon_i = 0$.

- It will be fast as long as we run epochs long enough \rightarrow

$$\pi_*(s) = \arg \max_{a \in A(s)} q_*(s, a)$$

* Sarsa Algorithm \rightarrow

Input: π , α -epoches, α

Output: $Q(s, a)$ if α -epoches is large enough.

Initialize Q arbitrarily (e.g. $Q(s, a) = 0$ for all $s \in S, a \in A(s)$).

including $Q(\text{terminal-state}, \cdot) = 0$.

For $i = 1$ to α -epoches:

$E \leftarrow 1/E$

Observe s_0

choose action A_0 from policy π using Q (e.g. ϵ -greedy).

$t \leftarrow 0$.

Repeat

take action A_t and observe s_{t+1}, R_{t+1}

choose action A_{t+1} using policy derived from Q (e.g. ϵ -greedy).

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha (R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(s_t, A_t))$$

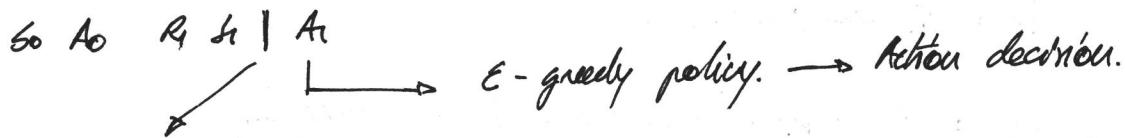
$t \leftarrow t+1$

until s_t is terminal.

end

return Q

- Q-greedy (ϵ -max) \rightarrow



$$Q(s_0, a_0) \leftarrow Q(s_0, a_0) + \gamma (R_t + \gamma \max_{a \in A} Q(s_t, a) - Q(s_0, a_0))$$

* For the update amine greedy policy, it approximates to the optimal policy.

* Algorithm \rightarrow

Input: π , ϵ -greedy, α

Output: Q ($\geq q_\pi$ if ϵ -greedy large enough).

Initialize Q arbitrarily.

For $i=1$ to num_episodes :

$$\epsilon \leftarrow 1/i$$

Observe s_0

$$t \leftarrow 0$$

Repeat

choose a_t using policy derived from Q (ϵ -greedy).

Take a_t and observe s_{t+1}, R_{t+1}

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

$$t \leftarrow t+1$$

until s_{t+1} is terminal

end

return Q .

- Expected Sarsa →

- * guaranteed to converge to the q_{π} under the same conditions as Sarsa.

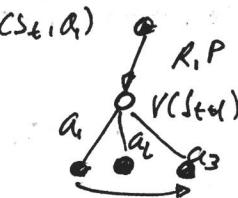
Illustrates the revenue from the random selection of A_{t+1}

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \sum_{a \in A} \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right)$$

* Algorithm →

Input: π , num-epochs, α

Output: $Q(\cdot | S_t)$ is num-epochs large enough).



Initialize Q arbitrarily.

For $i = 1$ to num-epochs:

$$\epsilon \leftarrow 1/i$$

Observe S_0

$$t \leftarrow 0$$

repeat

 Choose action A_t , using policy derived from $Q(\epsilon\text{-greedy})$

 Take action A_t and observe R_{t+1}, S_{t+1}

$$Q(S_t, A_t) \leftarrow$$

$$Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \sum_{a \in A(S_{t+1})} \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right)$$

$$t \leftarrow t+1$$

until t_{end} is reached

end

return Q

- Analyze performance \rightarrow

* All these algorithms converge \Rightarrow These conditions are met \rightarrow

1. The value of E decays according to GUE.

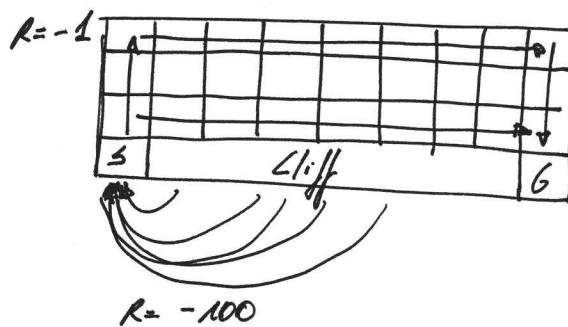
2. κ is sufficiently small.

* Differences \rightarrow

- Sarsa, Expected sarsa \rightarrow on-policy TD. \Rightarrow Better on-line performance than off-policy.

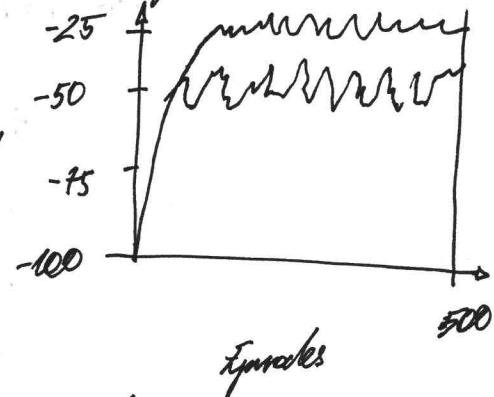
- max_{π} \rightarrow off-policy

- Expected sarsa better performance than sarsa.



safe path
optimal path

law of rewards during episode.



Results are from a single run but smoothed by averaging reward runs from 10 successive episodes.

- Sarsa and Q-learning are $\epsilon = 0.1 \rightarrow$

1. Q-learning achieves worse on-line performance but learns the optimal policy.

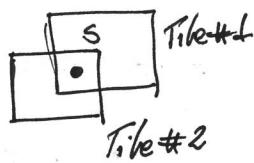
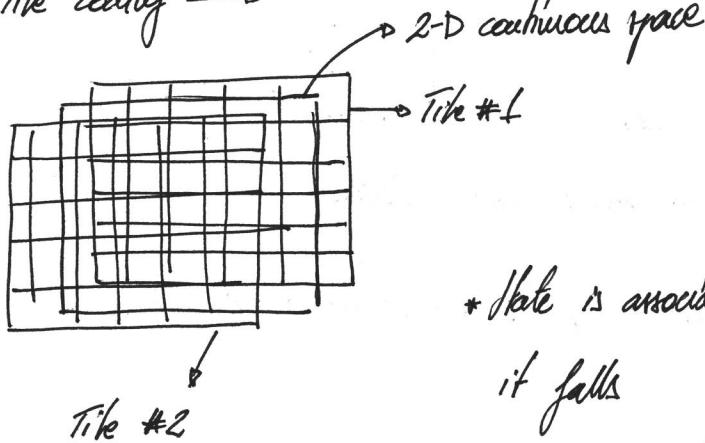
2. Sarsa achieves better on-line performance, but learns a suboptimal "safe" policy.

Reinforcement learning in continuous spaces

* Discretization →

- Approximate continuous function through uniform/non-uniform discretization.

- Tile coding →



* State is associated to each cell of the tile in which it falls

* Has different tiles with certain effect. →
Random selection of cells, r.v. effect.

* Algorithm → S, A, P, R, γ, w, n

for $i=1$ to m :

Initialize tile i with n/m tiles

for $j=1$ to n/m :

Initialize tile j with zero weight.

repeat

$s \leftarrow$ random state from S

$$\Delta V(s) \leftarrow \max_a [R(s,a) + \gamma V(P(s,a)) - V(s)]$$

for $i=1$ to m :

$w \leftarrow$ weight of active-tile $c_{i,i}$

$$w \leftarrow w + \frac{\kappa}{m} \Delta V(s)$$

until time expires

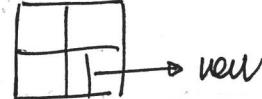
$$v(s) = \sum_{i=1}^n \delta_i(s) w_i \rightarrow \text{weight for each file.}$$

* Adaptive tile coding →

- Start with given number of cells and divide when appropriate

- Example →

* split when we are not
leaving with our current representation.



* Stop with number of splits or number of iterations.

* Function Approximation →

- Represent state-value and action-value functions →

$$\hat{v}(s, w) \approx v_\pi(s)$$

$$\hat{q}(s, a, w) \approx q_\pi(s, a)$$

| → parameter w that shapes the function.

- State-value Approximation →

$$\hat{v}(s, w) = x(s)^T w = \sum_{j=1}^n x_j(s) w_j$$



$$x(s) = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$$

feature vector

* Value function → $\hat{v}(s, w) = x(s)^T w$

* Minimize Error → $J(w) = \frac{1}{2} \sum_{s \in S} [(v_\pi(s) - x(s)^T w)^2]$

* Error Gradient →

$$\nabla_w J(w) = -2(v_\pi(s) - x(s)^T w) x(s)$$

* Update rule →

$$\Delta w = -\alpha \frac{1}{2} \nabla_w J(w) = \alpha (v_\pi(s) - x(s)^T w) x(s)$$

- Action-value Approximation →

$$\hat{q}(s, a, w) = x(s)^T w_1$$

$$\hat{q}(s, a_1, w) = x(s)^T w_1 \rightarrow (x_1(s, a), \dots, x_n(s, a)) \begin{pmatrix} w_{11} & \dots & w_{1m} \\ \vdots & & \vdots \\ w_{n1} & \dots & w_{nm} \end{pmatrix}$$

$$= (\hat{q}(s, a_1, w), \dots, \hat{q}(s, a_n, w))$$

* Allows to control two actions at the same time

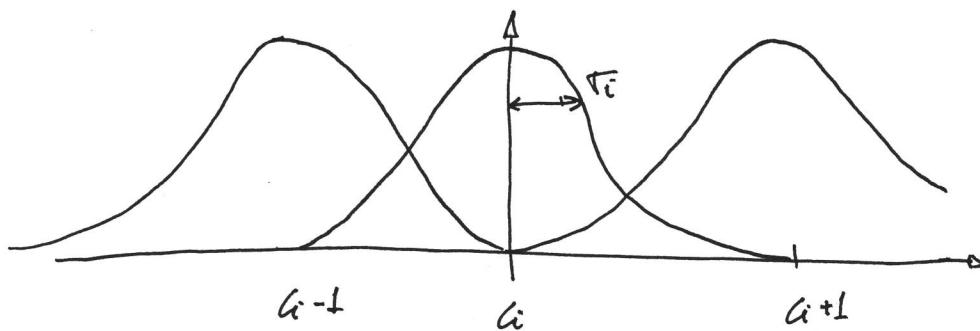
- Kernel features →

$$x(s) = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix} \rightarrow \left. \begin{array}{l} x_1(s) = s \\ x_2(s) = s^2 \\ x_3(s) = s^3 \end{array} \right\} \text{Kernel functions.}$$

* Radial Basis Function (RBF) →

- Associate each feature with a kernel function

$$\phi_i(s) = \exp\left(-\frac{\|s - c_i\|^2}{2r_i^2}\right)$$



- Non-linear function Approximation →

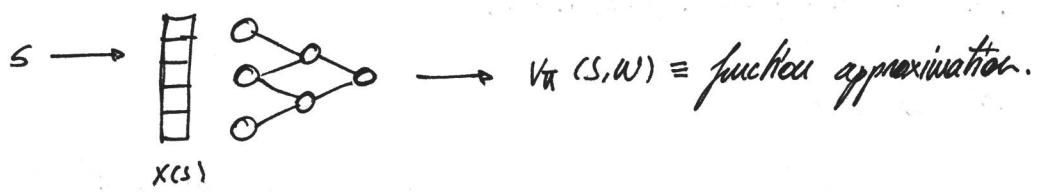
$$\hat{v}(s, w) = f(x(s)^T w)$$

$$\Delta w = \alpha (u_T(s) - \hat{v}(s, w)) P_w \hat{v}(s, w)$$

← DNN.

Deep Reinforcement Learning

- * Neural Networks as Value functions \rightarrow



$s \rightarrow V_{\pi}(s) \rightarrow V_{\pi}(s) \in \mathbb{R} = \text{state value function under optimal policy.}$

- backpropagation and gradient descent to update params $w \rightarrow$

$$\text{error} \rightarrow (V_{\pi}(s) - \hat{V}(s, w))^2$$

$$\Delta w = \alpha (V_{\pi}(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)$$

$$\Delta w = \alpha (q_{\pi}(s, a) - \hat{q}_{\pi}(s, a, w)) \nabla_w \hat{q}_{\pi}(s, a, w)$$

- Now, question is how to approximate the optimal value functions.

- * Sante Carlo learning \rightarrow

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t))$$

$$\text{Actual Return} \Rightarrow G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$\Delta w = \alpha \frac{(V_{\pi}(s_t) - \hat{V}(s_t, w)) \nabla_w}{G_t}$$

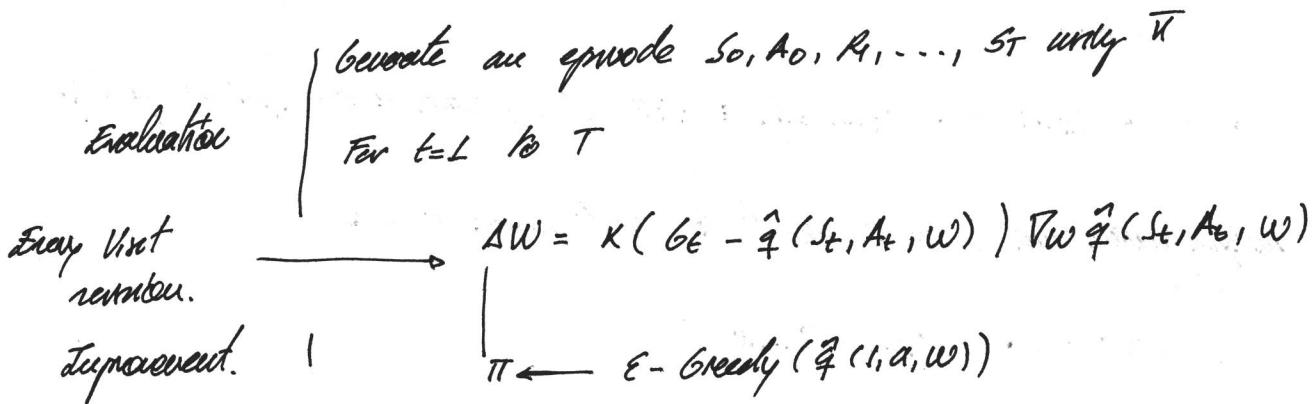
- RL is guaranteed to converge on a local minimum optima, in case of function approximation, it will converge to a global optima.

* Function approximation →

Initialize W with random values.

Policy $\pi \leftarrow \epsilon\text{-greedy } (\hat{q}(s, a, w))$

Repeat till converge:



* Temporal Difference learning →

$$V(s_t) \leftarrow V(s_t) + \alpha (r_{t+1} + \gamma \underline{V(s_{t+1})} - V(s_t))$$

\hookrightarrow Estimated return.

$$\underline{AW = \alpha [r_{t+1} + \gamma \hat{V}(s_{t+1}, w) - \hat{V}(s_t, w)] D_w \hat{V}(s_t, w)}$$

^{TD-TIME}

- TD(0) coupled with function approximation →

Initialize w randomly, $\pi \leftarrow \epsilon\text{-Greedy } (\hat{q}(s, a, w))$

Repeat many episodes:

Initial State s

while s non terminal

A \leftarrow state s under π .

R, S' \leftarrow env(s, A)

A' \leftarrow state s' under π

Update: $AW = \alpha [R + \gamma \hat{q}(s', a', w) - \hat{q}(s, a, w)] D_w \hat{q}(s, a, w)$

s \leftarrow s', A \leftarrow A'

* Q-learning →

- SARSA & Q-learning are TD approaches, they might not converge to a global optima with non linear function approximation.

- Update step → $\Delta W = \alpha (R + \gamma \max_a \hat{q}(s', a, w) - \hat{q}(s, a, w))$

- Differences →

SARSA

On-policy

Good on-line
performance

Q-values affected
by exploration

Q-learning

Off-policy

Bad On-line
performance

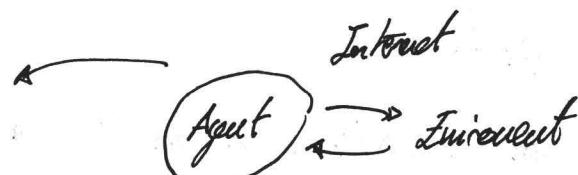
Q-values unaffected by
exploration.

- Off-policy → decouples actions taken and learning value functions.

* Experience Replay →

$(s_t, a_t, r_{t+1}, s_{t+1})$
$(s_{t+1}, a_{t+1}, r_{t+2}, s_{t+2})$
$(s_{t+2}, a_{t+2}, r_{t+3}, s_{t+3})$

Replay buffer



sample batch randomly
to learn.

- Every action a_t affect next state s_{t+1} , which means that each experience tuples can be highly correlated.

- Sampling at random from the buffer, breaks the correlations and prevents action values from diverging.
- Function approximation \rightarrow may be able to approximate, evaluate unseen scenarios, while abstraction often leaves 'holes' from previous learning cases.
- Experience replay also helps with rare state, actions \rightarrow
 - * If same action, state is repeated then it will be hard to learn another option action.
 - * Keeping experience and learning later will help with it.
- * Fixed Q-TARGETS \rightarrow
 - Helps to address one type of correlation: consecutive experience tuples.

$$\Delta w = \alpha \left[\frac{R + \gamma \max_a \hat{q}(s, a, w) - \hat{q}(s, a, w)}{\text{TD Target}} - \frac{\nabla_w \hat{q}(s, a, w)}{\text{Current Value}} \right] \nabla_w \hat{q}(s, a, w)$$

$\hat{q}_\pi(s, a)$

- * Both we calculated based on w \rightarrow correlations between target and parameters are not changing.
- * Fix the parameter in the target \rightarrow
 w^- = copy of w that we don't change during certain # of steps.
 decoupled.

$$\Delta w = \alpha \left[R + \gamma \max_a \hat{q}(s, a, w^-) - \hat{q}(s, a, w) \right] \nabla_w \hat{q}(s, a, w)$$

* Algorithm →

Initalise replay memory D with capacity N .

Initalise \hat{q} with all random.

Initalise target action-value weights $w^- \leftarrow w$

For episode $e=1$ to M :

Inital repeat from x_1 (frames)

Inital State $s \leftarrow \phi(x_1)$

For time-step = 1 to T :

$a \leftarrow$ State s using $\bar{a} \leftarrow \epsilon\text{-Greedy } (\hat{q}(s, a, w))$

$r, x_{t+1} \leftarrow \text{env}(s, a)$

$s' \leftarrow \phi(s, a, r, x_{t+1})$

Memory $D \leftarrow$ Experience tuple (s, a, r, s')

$s \leftarrow s'$

Sample

Minibatch $\leftarrow (s_j, a_j, r_j, s_{j+1}); j \in \text{range } N$

Learn

let target $y_j = r_j + \max_a \hat{q}(s_{j+1}, a, w^-)$

Update $\Delta w = \alpha [y_j - \hat{q}(s_j, a_j, w)] \nabla_w \hat{q}(s_j, a_j, w)$

Every C steps, reset $w^- \leftarrow w$

* DQN *Implementation* →

- Overestimation of Q values →

$$\Delta w = \alpha [R + \gamma \max_a \hat{q}(s', a, w) - \hat{q}(s, a, w)] \nabla_w \hat{q}(s, a, w)$$



$$R + \gamma \hat{q}(s', \arg \max_a \hat{q}(s', a, w), w)$$



↓
State



Q-Value State Action with max Q-Value.

- * Take especially on first steps, because Q values are still evaluting and we might have enough samples to figure it out.

- * It has been proven that overestimator values.

- * Double Q-learning →

best action.

$$R + \gamma \hat{q}(s', \arg \max_a \hat{q}(s', a, w), w')$$

↓ Evaluate action.

- Two function approximations that not agree on the best action.

If they don't agree, Q-value returned not that high.

- This prevents individual high rewards by chance that may have ramifications in the long run.

- Requirements →

Fixed Q
Targets.

- * choose randomly one to update at each step, using the other for evaluation.

- Prioritized Experience Replay →

* TD Error → $\delta_t = R_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a, w) - \hat{q}(s_t, a_t)$

* Priority → $p_t = |\delta_t| + c \rightarrow$ memory entry $\langle s_t, a_t, r_{t+1}, s_{t+1}, p_t \rangle$

* Sampling probability → $p_{\text{采}} = \frac{p_t^\alpha}{\sum_k p_k^\alpha}$

- A low δ_t doesn't mean that we meet the optimal, it may be that the estimation is off → +c

- This purely sampling may produce some overfit if we don't weight due to introduce more uniform prob →

* $\alpha=0 \rightarrow$ Uniform probability.

* $\alpha=1 \rightarrow$ Priority probability.

* Modified update rule →

$$\Delta w = \alpha \left(\frac{1}{N} \cdot \frac{1}{p_{\text{采}}} \right)^\delta \delta_t \nabla_w \hat{q}(s_t, a_t, w)$$

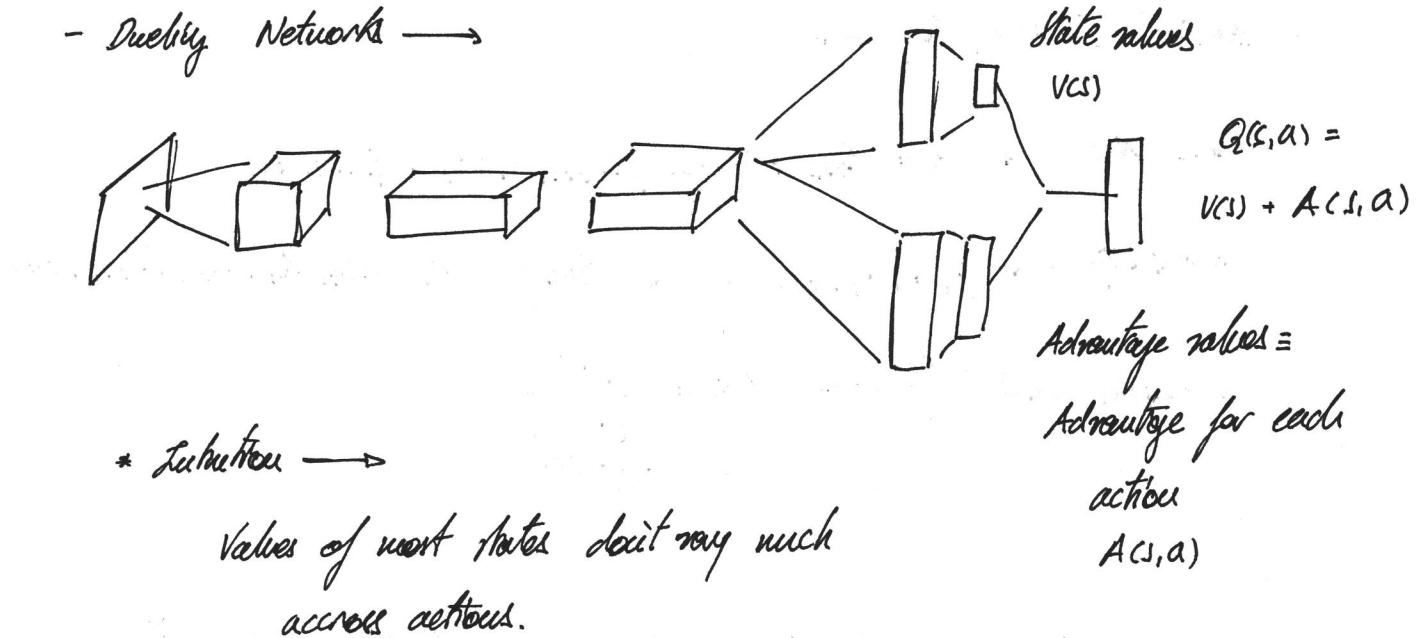
- before it was an expectation from a sampling of all my experiences, now, it must match our distribution.

- Importance sampling weight →

δ = hyperparameter. → $\delta=0$ Uniform

N = buffer size. $\delta=1$ Pure priority distribution.

- Dueling Networks \rightarrow



* Intuition \rightarrow

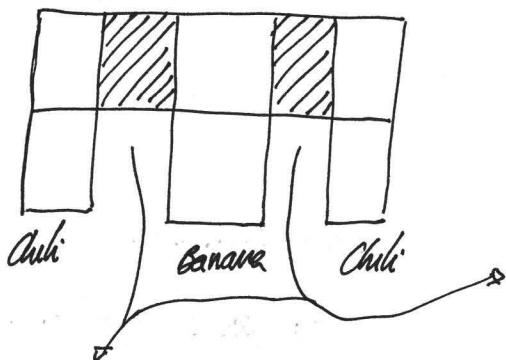
Values of most states don't vary much
across actions.

Policy Based Methods

- Reasons \rightarrow

Deterministic $\pi: s \rightarrow a$; Stochastic $a \sim \pi(s, a) = \text{PLAIS}$

- * Policy based methods can learn true stochastic policies \rightarrow
 - Randomness is derived from the environment, not as a ϵ -greedy policy which is an approximation.
- * Altered States \rightarrow
 - Two or more states that we perceive to be identical but are actually different.



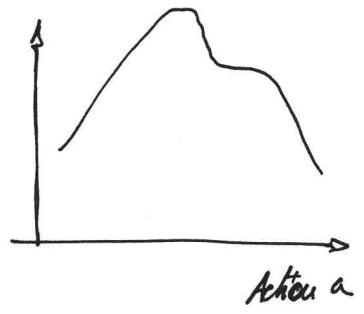
- * George and banana problem.
- * state defined by walls and floor type.
- * In the "11" cells actions to take are to be based on previous experience, states are equal.

- * Actions will be the same since features are equal
- * This will make George to go back and forth in one on the points, it won't get out until ϵ -greedy provider may act \rightarrow inefficient.
- * Continuous Action spaces \rightarrow

- It is now an optimization problem to find maximum.

- If the function is highly action space dimensional this is not a trivial problem.

Value $\hat{q}(s, a)$



- Policy function approximation \rightarrow

$$a \sim \pi(s, a, \theta) = P[a|s, \theta]$$

* Linear function with softmax policy \rightarrow

- For discrete set of actions.
- We can sample across actions.

$$f(s, a, \theta) = x(s, a)^T \theta$$

$$\pi(a|s, \theta) = \frac{e^{f(s, a, \theta)}}{\sum_{a'} e^{f(s, a', \theta)}}$$

* Linear function with gaussian policy \rightarrow

$$f(s, a, \theta) = x(s, a)^T \theta$$

$$\pi(a|s, \theta) = N(\mu, \sigma^2)$$

$$\mu = f(s, a, \theta)$$

σ^2 = fixed or parameterized

- Continuous functions.
- Sample action from distribution.

- objective function →

* Evaluate the expected value of the rewards under that policy.

$$J(\theta) = \mathbb{E}_{\pi} [R(z)] ; z = \text{trajectory, complete or partial episode.}$$

* State values →

- Start state value → $J_V(\theta) = \mathbb{E}_{\pi} [G_1] = \mathbb{E}_{\pi} [V(s_1)]$

- Average state value →

$$J_{\bar{V}}(\theta) = \mathbb{E}_{\pi} [V(s)] = \sum_s d(s) V(s)$$

Better values
per occurrence
of state.

$$d(s) = \frac{N(s)}{\sum_{s'} N(s')} ; \text{ trajectory probability}$$

* Action values →

$$J_Q(\theta) = \mathbb{E}_{\pi} [Q(s, a)] = \sum_s d(s) \sum_a \pi(s, a, \theta) q(s, a)$$

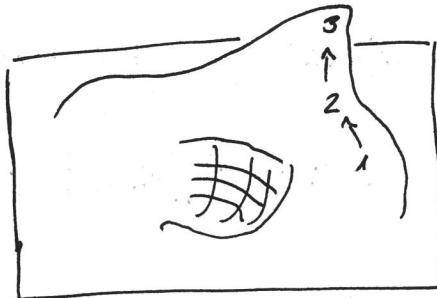
* Rewards →

- Work without state, action value functions.

$$J_{\bar{r}}(\theta) = \mathbb{E}_{\pi} [r] = \sum_s d(s) \sum_a \pi(s, a, \theta) R_s^a$$

- Stochastic Policy Search →

- * Start with arbitrary policy
- * Add noise to policy at each step, choose the next proximity.



- * Adaptive noise → large noise at the begining to avoid local maxima, decrease noise as you get closer to max.

- Policy Gradient →

- * Having a policy π_θ and objective function $J(\theta)$ →
compute gradient $\nabla_\theta J(\theta)$ and update $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

- * If the function is not differentiable, not possible to find gradient directly.

- * Gradient Estimation →

For $k \in [1, n]$: $\xrightarrow{\text{each dimension}}$
 $\xrightarrow{\text{Step}} \xrightarrow{\text{unit vector in } k^{\text{th}} \text{ dimension}}$

$$\frac{\partial J(\theta)}{\partial \theta_k} = \frac{J(\theta + \epsilon e_k) - J(\theta)}{\epsilon}$$

$$\nabla_\theta J(\theta) \approx \left\langle \frac{\partial J(\theta)}{\partial \theta_1}, \frac{\partial J(\theta)}{\partial \theta_2}, \dots, \frac{\partial J(\theta)}{\partial \theta_n} \right\rangle$$

Computationally sufficient.

* Compute gradient analytically \rightarrow

- Need access to the underlying policy.

$$J(\theta) = \mathbb{E}_{\pi} [R(z)] ; R(z) \text{ score function.}$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\pi} [R(z)] = \mathbb{E}_{\pi} [\nabla_{\theta} (\log \pi(z, a, \theta)) R(z)]$$

$$\Delta \theta = \alpha \nabla_{\theta} (\log \pi(z, a, \theta)) R(z)$$

Need to compute gradient from policy, not score function.

- Derivation \rightarrow

$$\nabla_{\theta} \mathbb{E}_x [f(x)] ; x \sim P[x|\theta]$$

$$\nabla_{\theta} \mathbb{E}_x [f(x)] = \nabla_{\theta} \sum_x P(x|\theta) f(x) = \sum_x \nabla_{\theta} P(x|\theta) f(x) =$$

Likelyhood Ratio trick

$$= \sum_x P(x|\theta) \frac{\nabla_{\theta} P(x|\theta)}{P(x|\theta)} f(x) = \sum_x P(x|\theta) \nabla_{\theta} (\log P(x|\theta)) f(x) =$$

$$= \mathbb{E}_x [\nabla_{\theta} (\log P(x|\theta)) f(x)]$$

- Monte Carlo Policy Gradient \rightarrow

Incentive & Anti-incentive.

For each episode $z = s_0, a_0, r_1, s_1, \dots, s_T$

For $t=1 : T-1 :$

$$\Delta \theta = \alpha \nabla_{\theta} (\log \pi(s_t, a_t, \theta)) g_t$$

$$\theta = \theta + \Delta \theta$$

$\left. \begin{array}{l} \\ \end{array} \right| \leftarrow \text{Algorithm Reinforce.}$

- Untrained Policy Gradient →

- Constrained Policy Gradient \rightarrow

- * With policy gradient, we focus on the overall performance based on the reward, disregarding state, action \rightarrow This could lead to two different policies A, B with similar rewards and really different behaviors.

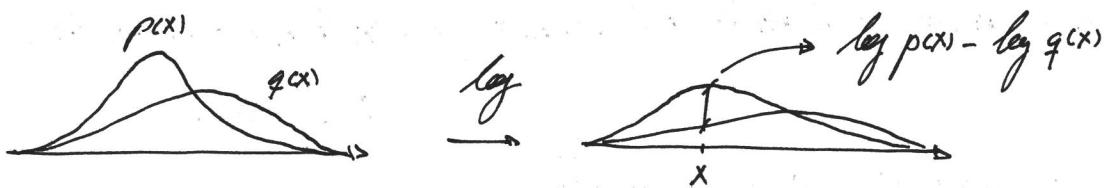
$$\mathcal{J}(\theta) = \mathbb{E}_{\pi} [R(z)] - \beta D[\pi(\cdot, \cdot, \theta) - \pi(\cdot, \cdot, \theta')]$$

- * Difference between policies $\rightarrow D[\pi(\cdot, \cdot, \theta) - \pi(\cdot, \cdot, \theta')] \leq \delta$

- Stabilizes the algorithm
- Prevents to have big changes in behavior

- * KL-Divergence \rightarrow Policy as probability function, how different?

$$D_{KL}(p \parallel q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx = \int_{-\infty}^{\infty} p(x) [\log p(x) - \log q(x)] dx$$



$$D_{KL}(p \parallel q) \neq D_{KL}(q \parallel p)$$

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi} [R(z)] - \beta D(\pi(\cdot, \cdot, \theta), \pi(\cdot, \cdot, \theta'))$$

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi} [R(z)] - \beta D_{KL}(\pi(\cdot, \cdot, \theta) \parallel \pi(\cdot, \cdot, \theta'))$$

Actor - Critic Methods

- Loss function \rightarrow

- * We have a way to improve our policy, but we need a way to keep track how we are doing.
- * Non episodic environment, meaning that we can compute as it interacts with the environment.

$$\Delta \theta = \alpha \nabla_{\theta} (\log \pi(s_t, a_t, \theta)) R(z)$$

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



Policy
update

$$\Delta \theta = \alpha \nabla_{\theta} (\log \pi(s_t, a_t, \theta)) Q(s_t, a_t)$$

Value
Update

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \beta \underbrace{(R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))}_{\text{TD Difference.}}$$

- * β = learning rate for value updates.

- * This two can run in parallel and they don't need an entire episode to be updated.

- Two function approximation \rightarrow

* Q-learning with function approximation

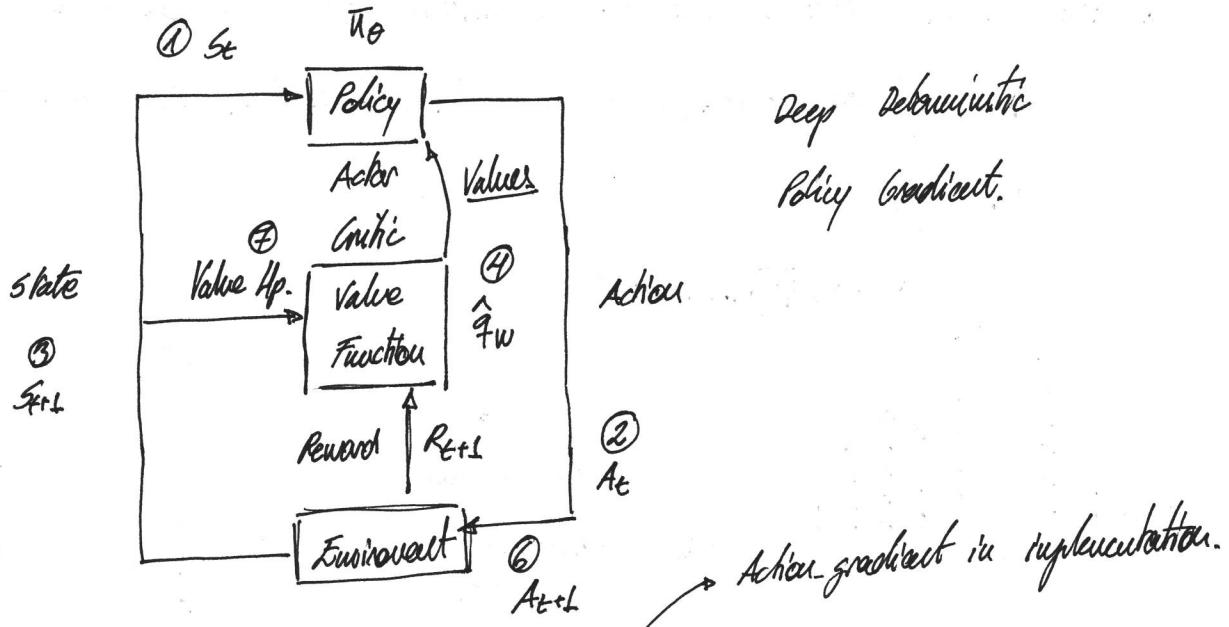
$$\Delta \theta = \alpha \nabla_{\theta} (\log \pi(s_t, a_t, \theta)) \hat{q}(s_t, a_t, w) \quad \text{Policy}$$

$$\Delta w = \beta (R_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w) - \hat{q}(s_t, a_t, w)) \nabla_w \hat{q}(s_t, a_t, w) \quad \text{Value}$$

$\pi(s, a, \theta) \rightarrow$ Action ; probability of taking an action in state.

$\hat{q}(s, a, w) \rightarrow$ Critic ; expected return in that state-action.

③ Update policy



$$⑤ \text{ Policy Update} \rightarrow \Delta \theta = \alpha \nabla_{\theta} (\log \pi(s_t, a_t, \theta)) \hat{q}(s_t, a_t, w)$$

$$⑥ \text{ Value Update} \rightarrow \Delta w = \beta \left[R_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w) - \hat{q}(s_t, a_t, w) \right] \nabla_w \hat{q}(s_t, a_t, w)$$

Target.

- Advantage function \rightarrow

* Policy gradient $\rightarrow \Delta\theta = \alpha \nabla_{\theta} J(\theta)$

* Gradient \rightarrow

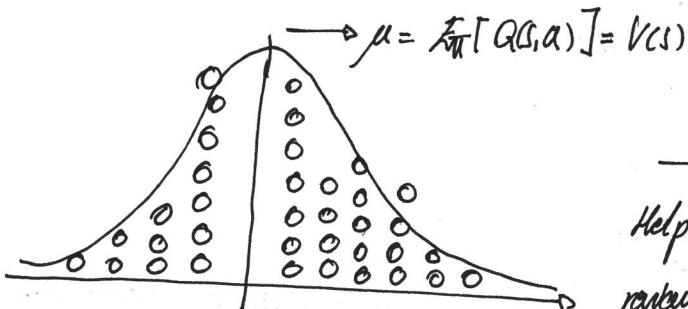
$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi} [\nabla_{\theta} (\log \pi(s, a, \theta)) R(z)] = \mathbb{E}_{\pi} [\nabla_{\theta} (\log \pi(s, a, \theta)) \hat{q}(s, a, w)]$$

Approximate the expectation by taking small steps.

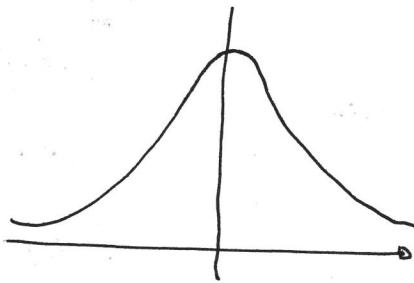
$$\Delta\theta = \alpha \nabla_{\theta} (\log \pi(s, a, \theta)) \hat{q}(s, a, w)$$

* Need to ensure small step \rightarrow working with stochastic values, they can vary a lot.

* With the expectation, there is also a variance of the samples \rightarrow given by the action value function estimation, since each time step noisy.



Helps reduce
variance
of estimate.



$$A(s, a) = Q(s, a) - V(s)$$

* Incentive \rightarrow

How much reward beyond the expected we get with the action.

- Actor-Critic with Advantage \rightarrow

* Gradient $\rightarrow \nabla_{\theta} J(\theta) = E_{\pi} [\pi_{\theta}(\log \pi(s, a, \theta)) (\hat{q}(s, a, w) - \hat{v}(s, a, w))]$

* Policy update \rightarrow

$$\Delta \theta = \alpha \nabla_{\theta} [\log \pi(s, a, \theta)) (\hat{q}(s, a, w) - \hat{v}(s, a, w))]$$

* Use the TD error of as estimation of the advantage function \rightarrow

$$TD \text{ error } \delta = r + \gamma \hat{v}(s', w) - \hat{v}(s, w) \approx \hat{q}(s, a, w) - \hat{v}(s, w)$$

$$= R_{t+1} + \gamma \hat{v}(s_{t+1}, w) - \hat{v}(s_t, w)$$

