- "world" / environment where our agent acts, is characterized by it's state $s_t$ at time $t$, states evolve in discrete time steps $s_t \to s_{t+1}$
- behaviour of the "world" is fully characterized by the transition probability given current state $s_t = s$ and action $a_t = a$, whats the probability to put to state $s_{t+1} = s'$ and receive reward $R_{t+1} = r$?

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

  important: this is stationary, doesn't depend in which time step we reached $s$ (Markov property)

- behavior of the agent is determined by its policy $\pi(a \mid s)$: probability to choose amongst possible action $a$ (hopefully good ones) in state $s$
- if world and/or agent behave deterministically, $p(s', r \mid s, a)$ and/or $\pi(a \mid s)$ reduce to delta-distributions

- value function of a state $s$ for given policy $\pi$ expected reward total reward when the game starts in state $s$ and agent follows the policy $\pi$:

$$V_\pi(s) = E_{p, \pi} \left[ \sum_{t'=t+1}^{\infty} \gamma^{t'-t-t} R_t \mid s_t = s \right]$$

- self-consistence of values leads to Bellmann equations

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

- the optimal value function chooses best action: $\boxed{V^*(s) = \max_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma V^*(s') \right]}$

– determine optimal value by <u>value iteration</u>:

  ⓪   initial guess $V^{(0)}(s)$     (e.g. $V^{(0)}(s) = 0$)

  ①   for $\tau = 1, 2, \ldots$    (or until convergence)

$$\forall s: \quad V^{(\tau)}(s) = \max_a \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma V^{(\tau-1)}(s')\right]$$

converges to a fixed point ( global optimum ? )

– once we know the optimal value of each state, we can derive optimal policy:

$$\boxed{\pi^*(s) = \arg\max_a \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma V^*(s')\right]}$$
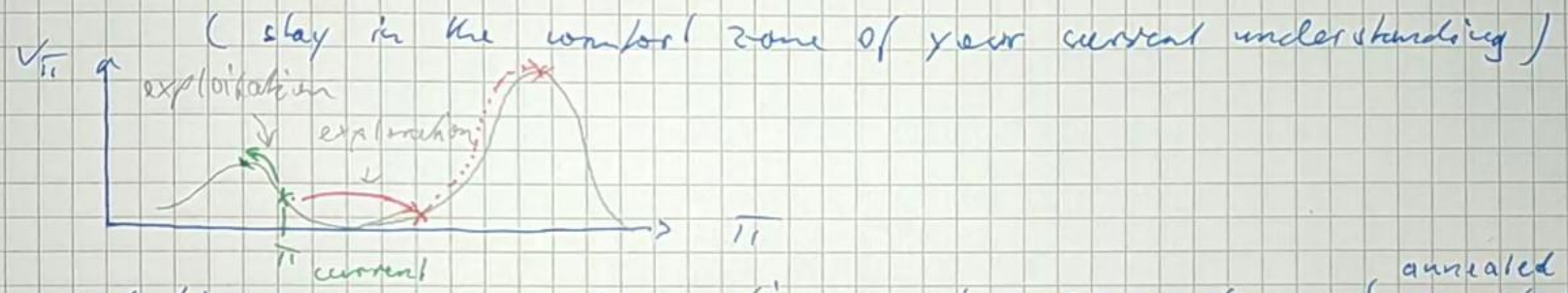
## <u>Model-free RL methods</u>

- problem so far: to determine $V^*(s)$ and $\pi^*(s)$, we need to know the transition probabilities $p(s', r \mid s, a)$

- usually, we don't know them $\Rightarrow$ <u>model-based RL</u> $\triangleq$ first learn $p(s', r \mid s, a)$ ("world model") and derive $V^*(s)$, $\pi^*(s)$ by value iteration
  <u>adv</u>: understand the "world"     <u>disadv</u>: very difficult

- <u>model-free RL</u> learns good policy without fully understanding the "world"
  <u>adv</u>: easier     <u>disadv</u>: black-box solutions that are hard to interpret

- meta-algorithm: <u>repeat many times</u>: ⓐ play according to current policy
                                        ⓑ collect reward data
                                        ⓒ use data to improve policy

⇒ simultaneous learning of policy and (implicitly) world behavior

<u>exploration - exploitation trade-off</u>

- exploration : learn more about the world
- exploitation : fine-tune the current policy and maximize reward
  ( stay in the comfort zone of your current understanding )



- common solution : <u>ε-greedy policy</u> : hyperparameter ε ( ~~annealed~~ annealed during training )
  - with probability $(1-ε)$ : execute the best action according to your current policy → exploit
  - —"— ε : execute random action → explore

  → guarantees that after infinitely many trials, every combination of state s and action a is tried infinitely often $\hat{=}$ have seen all world

<u>Monte-Carlo (MC) value estimation</u>

- choose action according to ε-greedy policy and estimate values by averaging actual rewards : - playing the game once $\hat{=}$ "episode" ⇒ episode index $τ = 1, 2, ...$
  - doing a move in game $τ$ : index $t = 1, 2, ... T_τ$
    $T_τ$ : duration of episode $τ$

return of step $t$ in game $τ$ = total reward of the rest of the game $τ$ : $G_{τt} = \sum_{t'=t+1}^{T_τ} γ^{t'-t-1} R_{τt'}$

- Determine value $V_\pi(s)$ as the average over all games where we encountered state $s$

$$V_\pi(s) \approx \frac{1}{N_s} \sum_{\tau, t \, : \, S_{\tau t} = s} G_{\tau t}$$

         $\underset{\text{we saw state } s}{\swarrow \text{how often}}$

can be computed incrementally:    update $V_\pi(s)$ whenever we come to state $s$:

$$N_s \leftarrow N_s + 1$$

new value $\nearrow$      $\searrow$ old value

$$V_\pi(s) \leftarrow V_\pi(s) + \frac{1}{N_s} \underbrace{( G_{\tau t} - V_\pi(s))}_{\text{correction of current guess } V_\pi(s)}$$

generalize to arbitrary learning rate $\alpha$

$$\boxed{V_\pi(s) \leftarrow V_\pi(s) + \alpha \, ( G_{\tau t} - V_\pi(s))}$$

advantage: this also works when the policy $\pi$ is not constant over training
(we need this, because $\pi$ should improve!)

disadvantage: slow

Temporal Difference (TD) value estimation

- computing $G_{\tau t}$ is inconvenient, because one must wait to the end of episode $\tau$ before $G_{\tau t}$ is determined, immediate updates after each move are better

$\Rightarrow$ approximation.    $G_{\tau t} \approx \underset{\uparrow}{R_{\tau_t t+1}} + \gamma \, \underset{\uparrow}{V_\pi(s_{\tau, t+1})}$

                      actual reward        current guess for value of the next

$$\boxed{V_\pi(s_{\tau t}) \leftarrow V_\pi(s_{\tau t}) + \alpha \left[ R_{\tau, t+1} + \gamma \, V_\pi(s_{\tau, t+1}) - V_\pi(s_t) \right]}$$

- more accurate variant : <u>n-step TD value estimation</u>

  approximate $G_{\gamma t}$ with real rewards for the next $n$ moves (plain TD: $n=1$)

  $$G_{\gamma t} \approx R_{\tau, t+1} + \gamma R_{\tau, t+2} + \dots + \gamma^{n-1} R_{\tau, t+n} + \gamma^n V_{\overline{\pi}}(s_{\tau, t+n+1})$$

  update rule $\left[ V_{\overline{\pi}}(s_{\tau, t}) \leftarrow V_{\overline{\pi}}(s_{\tau, t}) + \alpha \left[ \sum_{t'=t+1}^{t+n} \gamma^{t'-t-1} R_{\tau, t'} + \gamma^n V_{\overline{\pi}}(s_{\tau, t+n+1}) - V_{\overline{\pi}}(s_{\tau, t}) \right] \right]$

  main advantage: actually observed return $R_{\tau, t}$ is used in $n$ updates

  $\Rightarrow$ make better use of our data

- once we have determined the value of the current policy, we can use this to improve the policy

  $$\pi'(s) \leftarrow \underset{a}{\text{arg max}} \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma V_{\overline{\pi}}(s) \right]$$

  <u>but</u>: we need to know the transition probabilities

  <u>Q</u> : can we also do this in a model-free setting?

## Q - learning

- define "action-value function" $Q_\pi(s,a)$  ("Q-function")

- remember Bellman eq. for value function: $V_\pi(s) = \sum_a \pi(a|s) \underbrace{\sum_{s',r} p(s',r|s,a)\left[r + \gamma V_\pi(s')\right]}$

$$= Q_\pi(s,a)$$

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

$$V_\pi(s) = \sum_a \pi(a|s) Q_\pi(s,a)$$

interpretation: value of state $s$, when we continue with action $a$, and then

follow policy $\pi$ for the rest of the game ($=$ ignore policy in first move)

$\Rightarrow$ optimal policy chooses the action that maximizes action-value

$$\pi^*(s) = \arg\max_a Q^*(s,a)$$

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

$$V_{\pi*}(s') = \max_a Q^*(s',a)$$

Bellman equations for Q-functions

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

same holds for sub-optimal policies

$\Rightarrow$ once we know $Q_\pi(s,a)$ for old policy, we get improved policy

without knowing transition prob.  $\boxed{\forall s: \pi'(s) = \arg\max_a Q_\pi(s,a)}$

(guaranteed: $\pi'$ is not worse than $\pi$)

## Q-learning algorithm

⓪ initialize $Q^{(0)}(s,a)$ arbitrary ( good guess → faster convergence )

but make sure that $~~~~~~~~~~~~~~~~Q^{(0)}(s,a) = 0~$ for all terminal states

initial policy $\pi^{(0)}(s) = \arg\max_a Q^{(0)}(s,a)$

① for episodes $\tau = 1, 2, \ldots,$   (until convergence)

ⓐ play episode according to $\varepsilon$-greedy policy :

$$a_{\tau,t} \sim \begin{cases} \arg\max_a Q^{(\tau-1)}(s_{\tau t}, a) & \text{with prob. } 1-\varepsilon \\ \text{uniform}(a) & \text{with prob. } \varepsilon \end{cases}$$

[ alternative exploration – exploitation trade-off : softmax policy

$$a_{\tau,t} \sim \frac{e^{Q^{(\tau-1)}(s_{\tau,t}, a)/\beta}}{\sum_{a'} e^{Q^{(\tau-1)}(s_{\tau,t}, a')/\beta}}$$

softmax function
for vector $v$

$$\frac{\exp(v_i/\beta)}{\sum_{i'} \exp(v_{i'}/\beta)}$$

role of temperature $\beta$ :    $\beta \to 0$ : always choose
$a = \arg\max_{a'} Q(s,a')$

anneal $\beta$ as
$\tau$ increases

$\beta \to \infty$ : $a \sim \text{uniform}(a')$

⟹ game outcome $s_{\tau,0}, a_{\tau 0}, r_{\tau 1}, s_{\tau 1}, a_{\tau 1}, r_{\tau 2}, \ldots, r_{\tau T_\tau}$

ⓑ update $Q$-function : for $t = 1, \ldots T_\tau$

$$Q^{(\tau+1)}(s_{\tau, t-1}, a_{\tau, t-1}) \leftarrow Q^{(\tau)}(s_{\tau, t-1}, a_{\tau, t-1}) + \alpha\left[r_{\tau t} + \gamma V^{(\tau)}(s_{\tau, t}) - Q^{(\tau)}(s_{\tau, t-1}, a_{\tau, t-1})\right]$$

two variants according to how we define the value $V^{(\tau)}(s_{\tau,t})$ of the rest of the game:

$$V^{(\tau)}(s_{\tau t}) = \arg\max_a Q^{(\tau)}(s_{\tau,t}, a) \qquad \text{"Q-learning"}$$

$$V^{(\tau)}(s_{\tau t}) = Q^{(\tau)}(s_{\tau t}, a_{\tau t}) \qquad \text{"SARSA algorithm"}$$

②  update policy:  $\boxed{\forall s: \quad \hat{\pi}(s) = \arg\max_a Q^{(\tau_{end})}(s, a)}$   final policy

[variant:   $n$-step  Q-learning

$$Q(s_{\tau,t-1}, a_{\tau,t-1}) \leftarrow Q(s_{\tau,t-1}, a_{\tau,t-1}) + \alpha\left[\sum_{t'=t}^{t+n-1} \gamma^{t'-t} r_{t'} + \gamma^n V(s_{\tau,t+n}) - Q(s_{\tau,t-1}, a_{\tau,t-1})\right]$$

advantages of Q-learning:   — can update policy without knowing $p(s', r \mid s, a)$
                            "model-free method"

   — supports off-policy learning  :  we can mix game outcomes
      obtained with different policies in updating Q
      ( at the beginning of training , policy is bad , but we can still use
        game outcomes in the updates later on )
   ⇒ learning with experience buffer