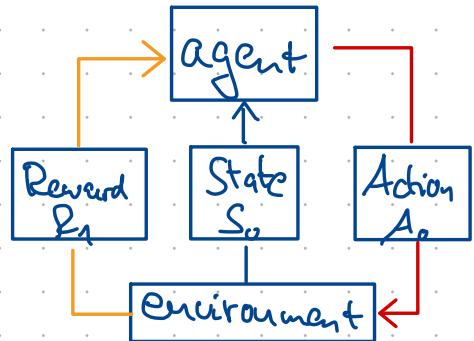




What is RL?

Reinforcement Learning is a computational approach to understanding and automating goal-directed learning and decision making. It is distinguished from other computational approaches by its emphasis on learning by an agent from direct interactions with its environment.



Goal: Maximize cumulative expected reward

example "walking": $r = \min(v_1, v_{max}) - 0.5 \cdot (v_1^2 + v_2^2) - 0.5 \cdot f + 0.02$

↓ ↓ ↓ ↓
fast straight elegant long

Discounted Reward is a trade-off between immediate reward and long-term goals:

$$G_t = \sum_{k=t}^{\infty} \gamma^k \cdot R_{t+k+1}, \quad \gamma \in [0, 1]$$

Other definitions:

A Policy is a mapping from perceived states to actions to be taken.

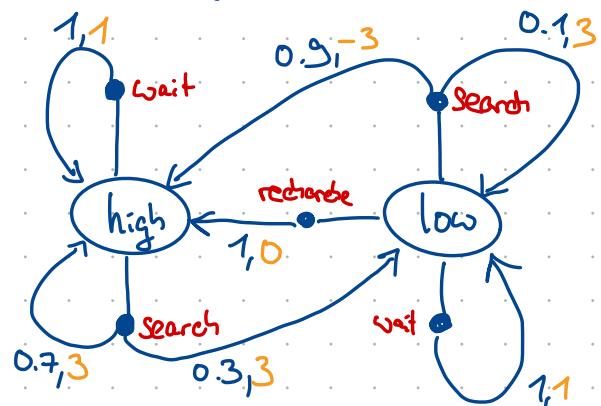
The value of a state is the total amount of reward an agent can expect to accumulate, starting from that state.

An episodic task has a temporal end, whereas a continuous task hasn't.

(Finite) Markov Decision Processes ("The problem")

MDPs are a formalization of sequential decision making, where the action does not only influence the immediate reward, but also subsequent states and, therefore, future rewards.

Example: transition diagram of a "garbage collector robot"



A finite MDP is defined by:

- a finite set of states S
- a finite set of actions A
- a finite set of rewards R
- the one-step dynamics of the environment
- a discount rate $\gamma \in [0, 1]$

The dynamics of a finite MDP is a probability distribution for random values $s' \in S$ and $r \in R$ occurring at time t , given particular values for the previous state and action:

$$p(s', r | s, a) = P\{S_{t+1} = s' | S_t = s, A_t = a\} \quad (\text{for all } s, s', a \text{ & } r)$$

→ $p: S \times R \times S \times A$ is an ordinary deterministic function of four arguments. $\sum_s \sum_r p = 1$

→ based on the dynamics, one can compute facts of interest about the environment, such as:

- the state-transition probabilities: $p(s' | s, a) = P\{S_{t+1} = s' | S_t = s, A_t = a\} = \sum_{r \in R} p(s', r | s, a)$
- the expected rewards for state-action pairs: $r(s, a) = E[R_{t+1} | S_t = s, A_t = a] = \sum_{r \in R} r \sum_{s' \in S} p(s', r | s, a)$
- the expected rewards for state-action-next-state triplets

Policies and Value Functions ("The Solution")

Gridworld Example:

...
-1
2	0	5

A policy is a mapping from states to probabilities of selecting each possible action.

→ $\pi(a|s)$ is the probability for choosing $A_t = a$ when in state $S_t = s$ under the policy π .

A state-value function is an estimate of the expected reward when following a policy, starting from that state.

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

The Bellman Expectation Equation estimates the value of any state as the sum of the intermediate reward and the total discounted reward under policy π :

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

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma \cdot V_\pi(s')] \quad \text{for all } s \in S$$

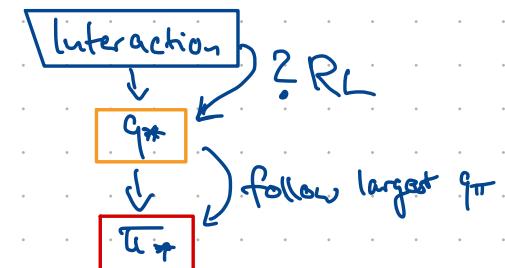
≡ multiply the sum of the reward r and discounted value $\gamma \cdot V_\pi$ of the next state s' by its corresponding probability $\pi \cdot p$.

Summing over all possibilities yields the expected values.

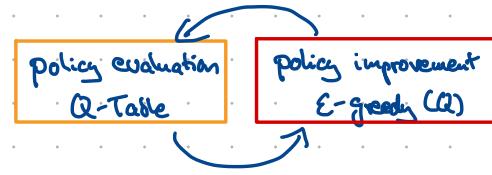
The action-value function yields the return if the agent starts in s , takes action a and then follows policy π :

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

An Optimal policy π^* is a policy where $V_{\pi^*}(s) \geq V_\pi(s)$ for all $s \in S$.



Monte Carlo Methods



Prediction Problem: Given a policy, how might the agent estimate the value function for that policy?

→ With MC prediction methods, the goal is to collect state-action pairs for many episodes and track their reward.

Taking the mean reward of state-action pairs over all episodes yields the estimated Q-Table.

For multiple visits of state-action pairs within an episode, there is the first-visit-MC (unbiased, better) or every-visit-MC.

→ Procedure: Input: policy π , num-episodes, Output: value-function $Q(s,a) \approx q_\pi$ for large num-episodes

Generate an episode and calculate the cumulative reward for the first occurrence of a state-action pair

$$Q(s,a) \leftarrow \text{returns_sum}(s,a) / N(s,a) \text{ for all } s \in S, a \in A$$

Control Problem: Estimate the optimal policy.

→ With MC control methods, the policy is updated based on the generated Q-Table.

This is done with an ϵ -greedy update: $\pi' \leftarrow \epsilon\text{-greedy}(Q)$ where ϵ sets the exploration-exploitation trade-off.

$$\pi'(a|s) \leftarrow \begin{cases} 1 - \epsilon + \epsilon / |A(s)| & \text{if } a \text{ maximizes } Q(a,s) \\ \epsilon / |A(s)| & \text{else} \end{cases}$$

($\epsilon / |A(s)|$ is the equiprobable random policy including the optimal action)

An optimal solution is found with the "Greedy in the Limit with Infinite Exploration" (G_LE) algorithm.

In order to update more frequently, one can implement incremental mean or constant alpha approach.

→ Constant alpha example:
$$Q(s,a) \leftarrow Q(s,a) + \alpha (G_t - Q(s,a))$$
, where $\alpha \in [0,1]$

$\Rightarrow \delta_t$ ("error" between estimated reward and actual reward)

Temporal Difference Methods

Review of Monte Carlo Methods:

- In order to learn, the agent has to terminate an episode.
- For most RL tasks, this is either slow (chess) or not desired (car \rightarrow crash).
- Constant- α update equation:
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (G_t - Q(S_t, A_t))$$

TD Control Methods update the value function after each time step.

- Sarsa is an on-policy TD method. It uses the current reward and value of the next state-action pair to estimate Q :

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \cdot Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$
- Sarsamax is an off-policy TD method. It uses the current reward and greedy value of the next state to estimate Q :

$$(Q\text{-learning}) \quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \cdot \max_{a \in A} Q(S_{t+1}, a) - Q(S_t, A_t))$$
- Expected Sarsa is an on-policy TD method. It uses the current reward and probability distribution for the next action to estimate Q :

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

- All Sarsa methods are guaranteed to converge for small α and ϵ under the SLE conditions (In practice, also w/o SLE)
- In Praxis:
 - It shows that Optimism when initializing Q yields better results.
 - On-policy methods have a better online learning performance. Expected Sarsa often works best.

RL in Continuous Spaces

The previous (model-free) methods work with finite MDPs, where state-action pairs can be represented in a discrete table.

In continuous spaces, states and/or actions are infinite: $s \in \mathbb{R}^n$, $a \in \mathbb{R}^m$. There are two ways of applying RL methods:

- Discretization (only works for small underlying spaces \rightarrow "curse of dimensionality")

- \rightarrow The infinite state space gets mapped to a discrete representation, e.g. by uniformly binning the space.
- \rightarrow Tile-Coding overlays multiple grids to the space, resulting in a bit-vector representation.

The estimated value function is the average of the weights of the activated tiles.

- \rightarrow More advanced discretization methods are adaptive tile-coding or coarse coding with a radial-basis function.

- Function Approximation



\rightarrow the state is represented as a feature vector $x(s) = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$.

\rightarrow the mapping from s to $x(s)$ can make use of kernel functions to derive non-linear features.

\rightarrow the scalar value output is the dot product of feature vector and weights: $\hat{v}(s, w) = x(s)^T \cdot w$

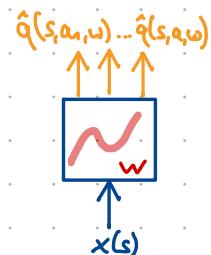
\rightarrow the weights can be learned using gradient descent: $J(w) = \mathbb{E}_{\pi} [(V_{\pi}(s) - x(s)^T \cdot w)^2] \rightarrow \Delta w = \alpha (V_{\pi}(s) - x(s)^T \cdot w) \cdot x(s)$

\rightarrow for discrete action-values, the value function of each action $\hat{q}(s, a_i, w)$ can be approximated with a weight matrix.

\rightarrow non-linear value functions can be approximated using non-linear activation functions:

$$\hat{v}(s, w) = f(x(s)^T \cdot w) \rightarrow \Delta w = \alpha (V_{\pi}(s) - \hat{v}(s, w)) \cdot \nabla_w \hat{v}(s, w), \nabla_w \hat{v}(s, w) \text{ is the } \underline{\text{partial derivative!}}$$

\Rightarrow This is the basis of using neural networks for RL!



Deep Q-Networks

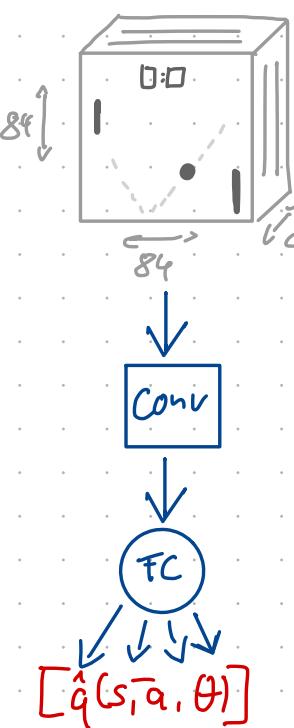
Review: It is unfeasible to discretize large state-spaces in a Q-table, i.e. the SLE requirement cannot be fulfilled.

Therefore, (non-linear) function approximation is used to generalize the value-functions from sparse samples.

Neural networks are powerful non-linear function approximators, but two pitfalls in combination of RL need to be solved:

- there is a high correlation between state-spaces and actions that lead to oscillations and divergence in online learning
- the target needs to be decoupled from the action, because Q-learning "updates a guess with a guess".

The Deep-Q-Network algorithm avoids these pitfalls with experience replay and fixed Q-targets:



- In experience replay, the agent samples $\langle S_t, A_t, R_t, S_{t+1} \rangle$ tuples and saves them in a replay buffer
 - randomly sampling from the experience buffer when learning tackles state-action-correlation and increases data efficiency, because the samples can be used multiple times for updates.
- With fixed Q-targets, the TD-target is calculated from a copy of the network and only updated after C learning steps. TD error

$$\Delta \omega = \lambda \cdot (R + \gamma \cdot \max_a \hat{q}(S', a, \omega^-) - \hat{q}(S, A, \omega)) \cdot \nabla_\omega \cdot \hat{q}(S, A, \omega)$$

TD target current estimation
- algorithm: initialize replay memory \mathcal{D} to capacity N and target value function \hat{Q} with $\theta^- = \theta$
 - for episode = 1 to M
 - for $t = 1, T$: take \times samples with ϵ -greedy policy and store them in replay buffer.
 - Sample in mini-batches and perform gradient descent on θ , set θ^- to θ after C steps

Improving DQNs

- Double DQNs: The TD target depends on the action that maximizes the $\hat{q}(s, a, \omega)$ function for state s .

This estimation tends to be too high (overestimation of Q-values), especially in the early stages of training.

It can be avoided by choosing $\max_a q(s, a, \omega)$ with network ω , and evaluating $q(s', a, \omega')$ with the fixed target network ω' .

$$\text{TD target} = R + \gamma \cdot \hat{q}(S', \arg \max_a \hat{q}(S', a, \omega'), \omega')$$

Select a Evaluate \hat{q}

- Prioritized Experience Replay: Some experiences are more important than others \rightarrow Sample experience based on TD error $\delta_{t,i}$!
(+ small min. priority ϵ)

The priority-random-trade-off is controlled with hyperparameter α .

The update rule needs to be adapted, because the underlying distribution of the experience sampling causes a bias.

$$p_t = \delta_t + \epsilon, \quad p_i = \frac{p_i^\alpha}{\sum_n p_n^\alpha}, \quad \Delta \omega = \alpha \cdot \left(\frac{1}{n} \cdot \frac{1}{p_i} \right)^b \cdot \delta_i \cdot \nabla \omega \cdot \hat{q}(S, A_i, \omega)$$

- Dueling Networks: In many RL tasks, the choice of action is not significant in many states.

Dueling networks address this by using a network with "two streams" that estimate the state value $V''(s)$ and the "advantage" $A''(s, a)$ separately, where $A(s, a) = Q(a, s) - V(s)$.

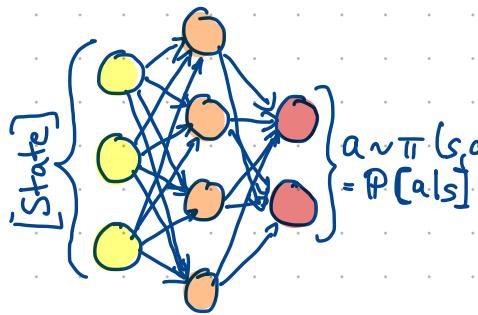
For a stable gradient computation, the output $Q(s, a, \theta)$ of a dueling network is calculated as follows:

$$Q(s, a, \theta, \alpha, \beta) = V(s, \theta, \beta) + (A(s, a, \theta, \alpha) - \frac{1}{|\mathcal{A}|} \cdot \sum_{a'} A(s, a', \theta, \alpha))$$

- Rainbow DQN: A "Rainbow" DQN combines the three improvements above as well as multi-step learning,

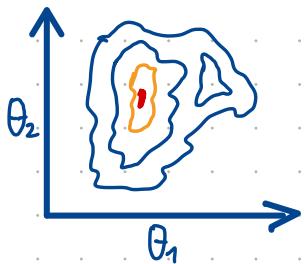
distributional RL and Noisy Net in a single network. This lead to Super-human performance.

Policy Based Methods



- In value-based methods, we approximate a value function $\hat{q}(s, a)$ to determine the best policy π^* .
- Instead, we can directly learn the **policy mapping** $\pi: s \rightarrow a$ (deterministic) or the **action distribution** $\pi(s, a)$ (stochastic).
- In discrete action spaces, the output is a probability for action, in continuous spaces it is a value e.g. [-1, 1].
- Benefits of policy-based methods are:
 - **Simplicity**: P.B.M. directly get to the problem at hand
 - **Stochastic policies**: P.B.M. can learn true stochastic policies (unlike E-greedy) \rightarrow good for "aliasing" states.
 - **Continuous action spaces**: P.B.M. tend to work well here.

Hill climbing



- Hill climbing is an iterative algorithm that slightly perturbs the weights Θ_{new} with gaussian noise and evaluates $\Theta_{best} \leftarrow \Theta_{new}$.
- It is a set of black box algorithms that require no back-propagation.
- Improvements of the vanilla algorithm are:
 - Steepest-Hill-Ascent: A small number of policies are evaluated and the best among them is chosen.
 - Simulated annealing: The search radius is reduced gradually.
 - Adaptive noise scaling: Decrease the search radius when a new best policy is found, else increase.
- These algorithms can also be used with Evolutionary Strategies where the next weights are determined by a weighted-distribution of the "previous" generation. This allows efficient, massive parallel computing and tends to work well with RL tasks with "sparse reward".

Policy Gradient Methods

- Policy Gradient Methods try to estimate the weights of an optimal policy through gradient ascent: $\theta \leftarrow \theta + \alpha \cdot \hat{g}$
- The ultimate goal is to make (state, action) pairs more likely when the reward is positive, and less likely otherwise
- In order to work with episodic and continuous tasks, we define a trajectory τ as: $(S_0, a_0, s_1, a_1, \dots, s_H, a_H, s_{H+1})$, H : "horizon"
- The return of a trajectory is denoted as $R(\tau)$, the corresponding expected return $U(\theta)$ is: $U(\theta) := \sum_{\tau} P(\tau; \theta) \cdot R(\tau)$
- The algorithm we use to optimize the policy is called REINFORCE:

1) Use policy π_θ to collect m trajectories $\{\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(m)}\}$ with horizon H . ("policy rollout")

2) Use the trajectories to estimate the gradient $\nabla_\theta U(\theta)$:

$$\nabla_\theta \cdot U(\theta) \approx \hat{g} := \frac{1}{m} \cdot \sum_{i=1}^m \sum_{t=0}^{H-1} \nabla_\theta \cdot \log \pi_\theta(a_t^{(i)} | s_t^{(i)}) \cdot R(\tau^{(i)})$$

"likelihood ratio(policy) gradient trick":
 $P(x; \theta) \cdot \frac{\nabla_\theta P(x; \theta)}{P(x; \theta)} = P \cdot \nabla_\theta \log P$

3) Use the estimated gradient to update the weights of the policy:

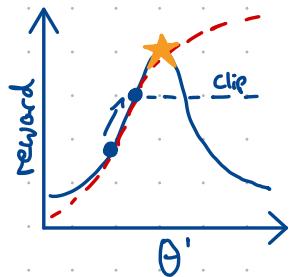
$$\theta \leftarrow \theta + \alpha \cdot \hat{g}$$

(For continuous action spaces, we can try to output a parameterized probability distribution and sample an action at $\sim \mathcal{N}(\mu, \sigma^2)$.

Modifications to make this work are described below.)

Proximal Policy Optimization

- PPO is a state-of-the-art RL algorithm for complex environments and continuous action spaces
- Drawbacks of the REINFORCE algorithm that motivates PPO are: 1) data inefficiency (trajectories are thrown away) 2) noisy gradients 3) no credit assignment
- Noise reduction is reached by sampling multiple trajectories and averaging their gradients and by applying "batch normalization" on the rewards.
- Credit assignment (also noise reduction) is done by changing the coefficient from total rewards $R(t)$ to future rewards $\bar{R}(t)$ $\rightarrow g = \sum R_t^{\text{true}} \cdot \nabla_{\theta} \log \pi(a_t | s_t)$
- In order to re-use old trajectories, we can compensate the change in the underlying distribution by importance sampling $\rightarrow \sum_t P(x; \theta) \cdot \frac{P(x; \theta')}{P(x; \theta)} \cdot f(x)$
 "old π " re-weighting factor



- PPO consists of two key principles:

- The Surrogate function: We use importance sampling to estimate the gradient of θ' based on trajectories that were sampled from an old policy $P(x; \theta)$. We further assume that the old and new policies are similar to each other (proximity). Hence we can define a Surrogate function as optimization objective for gradient ascent:

$$g = \nabla_{\theta} \cdot L_{\text{sur}}(\theta', \theta), \text{ where } L_{\text{sur}} = \sum_t \frac{\pi_{\theta'}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \cdot \bar{R}_t^{\text{true}}$$

- Clipping policy updates: The gradient estimation approximation is only valid if the policies are similar enough. To avoid over-training (and possibly landing in a very bad policy), we clip the surrogate objective if the ratio of the policies exceed $1 \pm \epsilon$.

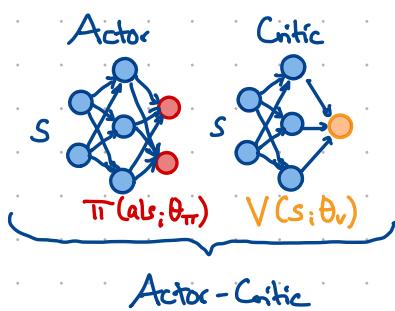
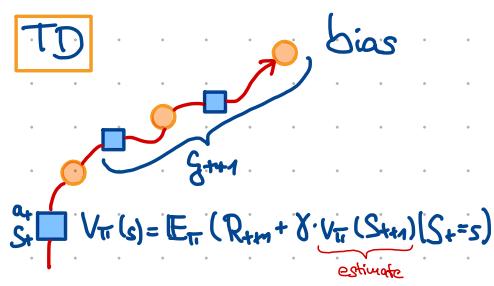
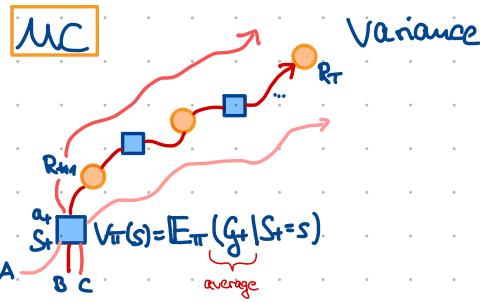
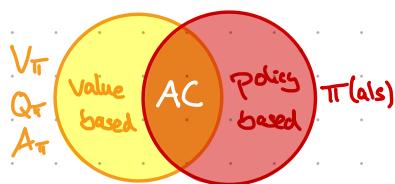
$$L_{\text{sur}}^{\text{clip}} = \sum_t \min \left\{ \frac{\pi_{\theta'}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \cdot \bar{R}_t^{\text{true}}, \text{clip} \left(\frac{\pi_{\theta'}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \right) \cdot \bar{R}_t^{\text{true}} \right\}$$

⇒ PPO algorithm: First collect trajectories based on θ and set $\theta' = \theta$.

↑
repeat: Compute the gradient $\nabla_{\theta} \cdot L_{\text{sur}}^{\text{clip}}(\theta', \theta)$ and update $\theta' \leftarrow \theta + \alpha \cdot \nabla_{\theta} \cdot L_{\text{sur}}^{\text{clip}}(\theta', \theta)$.

After a few iterations, set $\theta' = \theta$ and repeat from the beginning ("policy rollout").

Actor-Critic-Methods



$$\theta_{\text{trn}} = \theta_t + \alpha \cdot (\hat{g} - \hat{V}(s)) \cdot \nabla \log \pi(\text{als})$$

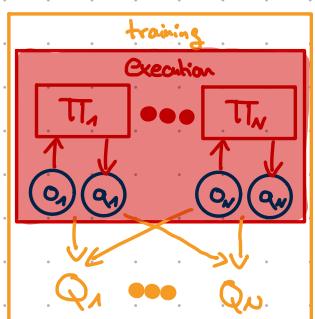
baseline

- AC-methods try to use value-based techniques to reduce the variance of policy-based methods.
- In policy-based methods, the agent learns to "act", but actions are reinforced based on the outcome (win/loss).
 - good actions are discouraged when the total reward is negative, so many trajectories are needed! (MC)
- In value-based methods, the agent learns to "estimate" a state-value, e.g. via TD-methods.
 - combining both approaches will eventually speed up learning and reduce the variance!
- Bias / Variance in RL algorithms come from the approaches to estimate expected total reward:
 - Monte-Carlo-Methods average the reward of multiple trajectories transitioning through state s ⇒ Variance
 - Temporal-Difference-Methods use an estimate of the next state-value to estimate $V_\pi(s_t)$ ⇒ bias
- A basic Actor-Critic Algorithm works as follows:
 - Use the actor for interaction with the environment to get a (s, a, r, s') -tuple.
 - Use the TD-estimate of that experience tuple $r + \gamma \cdot V(s'; \theta_V)$ to train the critic.
 - Use the critic to calculate the advantage $A(s, a) = r + \gamma \cdot V(s'; \theta_V) - V(s; \theta_V)$ to train the actor ("baseline").
- Using n-step bootstrapping is a trade-off between bias ($n=1$) and variance ($n=\infty$).
 - In general advantage estimation (GAE) all n-step bootstraps are combined: $\sum_{t=1}^T g_t \cdot (1-\lambda) \cdot \gamma^t$
- AC-methods are on-policy (=stable) and replace the replay buffer with parallel training.
 - A3C: asynchronous advantage actor-critic ; A2C: (synchronous) advantage actor-critic

Multi-Agent RL

- Multi-Agent scenarios are very common in RL applications and can either be collaborative, competitive or mixed.
- Learning is very unstable with vanilla RL-methods, because for each individual agent, the environment is not constant because of the other agents policies.

→ This problem is solved in **MADDPG** that is using a "shared" critic with a centralized value function that considers all agents policies.



- A Markov Game is a multi-agent extension of MDPs with partially observed states.

$$A = A_1 \dots A_N, O = O_1 \dots O_N, \pi_{\theta_i}: O_i \times A_i \mapsto [0, 1], T: S \times A_1 \times \dots \times A_N \mapsto S \text{ (transition depends on } S \text{ + all actions)}$$

- The gradient of a single agent in MADDPG is calculated as follows using the deterministic policy gradient approach.

$$\nabla_{\theta_i} J(\mu_i) = \mathbb{E}_{x, a, \text{and}} \left[\nabla_{\theta_i} \mu_i(a_i | o_i) \nabla_a Q_i^{\mu}(x, a_1, \dots, a_N) | a_i = \mu_i(o_i) \right]$$

where:

- $J(\mu)$ is the objective function, i.e. the total expected reward of agent i following policy μ ,

- $\mu(a_i | o_i)$ is the deterministic policy of agent i , given observation O_i ,

- $Q_i^{\mu}(x, a_1, \dots, a_N)$ is the centralized value function considering all observations x and all agents actions a .

- Instead of "knowing" all agents policies, they can also be learned with an entropy regularizer.
- In competitive scenarios, agents tend to overfit to specific strategies. To avoid such behavior, it is more robust to learn K ensemble sub-policies and then choosing uniformly random from the ensemble.