

# Computational Intelligence in Games - Quick Reference -

Alexander Dockhorn

March 27, 2019

This short overview tries to provide you a guide for all equations and algorithms mentioned in the lecture Computational Intelligence in Games at the Otto von Guericke University Magdeburg. It can be quite hard to understand all the little details involved without being sure about the symbols used. We hope this guide helps you during your review of the course topics.

I currently do not recommend printing this guide, because it is the first version of this kind of overview. Even if I reviewed the contents multiple times, errors are still likely. Please always refer to the actual version in the Git-Repository.

In case you have any ideas how to improve this guide please let us know by writing an e-mail to: [dockhorn@ovgu.de](mailto:dockhorn@ovgu.de).

## Contents

<b>1</b>	<b>(Evolutionary) Game Theory</b>	<b>4</b>
1.1	Basics in Game Theory . . . . .	4
1.2	General 2-Player Games . . . . .	4
1.3	Nash Equilibria . . . . .	5
1.4	Evolutionary Game Theory . . . . .	6
1.5	Additional Information on Evolutionary Game Theory . . . . .	7
<b>2</b>	<b>Reinforcement Learning</b>	<b>8</b>
2.1	General Notation . . . . .	8
2.2	Value and Action-Value Functions . . . . .	10
2.3	Dynamic Programming . . . . .	11
2.4	Monte Carlo Method . . . . .	12
2.5	Temporal Difference Learning . . . . .	14
<b>3</b>	<b>Simulation Based Search Algorithms</b>	<b>15</b>
3.1	Flat Monte Carlo . . . . .	15
3.2	Flat Upper Confidence Bounds (UCB) . . . . .	15
3.3	Monte Carlo Tree Search (MCTS) . . . . .	16
3.4	Evolutionary Algorithms (EA) . . . . .	17
3.5	Rolling Horizon (RH) . . . . .	18
<b>4</b>	<b>Multi Objective Optimization</b>	<b>19</b>
4.1	General Notation . . . . .	19
4.2	NSGA-II Algorithm . . . . .	19

## List of Equations

1	Expected payoff of pure strategy $s_i$ , $E(s_i, y)$ . . . . .	5
2	Expected payoff of mixed strategy $x$ , $E(x, y)$ . . . . .	5
3	Fitness of strategy $s_i$ , $f_i(x)$ . . . . .	6
4	Average population fitness, $\Phi(x)$ . . . . .	6
5	Replicator Equation for infinite populations, $\dot{x}$ . . . . .	6
6	Weak Selection, $\dot{x}$ . . . . .	7
7	Probability for replacing a strategy using weak selection $\rho$ . . . . .	7
8	Comparison of using T,R,S,P or b,c for describing a payoff matrix . . . . .	7
9	Simple (undiscounted) Return, $G_t$ . . . . .	9
10	Discounted Return, $G_t$ . . . . .	9
11	Expected Reward, $r(s, a)$ . . . . .	9
12	Expected Reward, $r(s, a)$ . . . . .	9
13	Value Function, $v(s)$ . . . . .	10
14	Consistency Condition of the Value Function, $v(s)$ . . . . .	10
15	Action-Value Function, $q(s, a)$ . . . . .	10
16	Consistency Condition of the Action-Value Function, $q(s, a)$ . . . . .	10
17	Derivation of the Incremental Mean and Incremental Monte Carlo, $\mu_k, V(s)$	13
18	n-Step Return, $G_t^{t+n}(s)$ . . . . .	14

## List of Algorithms

1	Iterative Policy Evaluation . . . . .	11
2	Policy Iteration . . . . .	11
3	Value Iteration . . . . .	11
4	Monte Carlo Method $v(s)$ . . . . .	12
5	Monte Carlo Method $q(s, a)$ . . . . .	12
6	Constant- $\alpha$ Monte Carlo Method $v(s)$ . . . . .	13
7	Temporal Difference Learning . . . . .	14
8	One Step Q-Learning . . . . .	15
9	Flat Monte Carlo . . . . .	15
10	Flat UCB . . . . .	15
11	General Evolutionary Algorithm . . . . .	17
12	Rolling Horizon Random Search . . . . .	18
13	Rolling Horizon Evolutionary Algorithm (RHEA) . . . . .	18
14	Crowding Distance Computation . . . . .	19
15	NSGA-II Algorithm . . . . .	20

# 1 (Evolutionary) Game Theory

## 1.1 Basics in Game Theory

Symbol	Name	Description
$N$	number of agents	the number of agents in an N-player game
$S_i$	strategies of player $i$	each player has a set of strategies, which he can choose to play
$s_i$	chosen strategy of player $i$	
$m_i$	number of strategies in $S_i$	
$\pi_i$	payoff function of player $i$	provides a reward to agent $i$ after each agent chose his action
$A$	payoff matrix player 1 (or both)	represents the payoff matrix of the first player, or of both players in case the payoff is symmetric
$B$	payoff matrix player 2	

## 1.2 General 2-Player Games

Symbol	Name	Description
$C$	strategy "Cooperate"	often used in standard literature to describe cooperative behavior, which comes with additional cost for the playing agent
$D$	strategy "Defect"	often used in standard literature to describe defective behavior, which comes without a cost for the playing agent, while ignoring costs for the other player
$T$	temptation	Reward for defecting, when the other player is Cooperating
$R$	reward for mutual cooperation	reward in case both players are cooperating
$S$	suckers payoff	reward for the player who cooperated against a player who defected him
$P$	punishment	reward for the player when both players chose defect

### 1.3 Nash Equilibria

Symbol	Name	Description
$\mathcal{S}$	strategy profile	the chosen strategies per player
$x$	mixed strategy	probability distribution on all possible strategies in $S_i$
$E(s_i, y)$	expected payoff of a pure strategy	expected payoff of pure strategy $s_i$ against mixed strategy $y$ . See Equation (1)
$E(x, y)$	expected payoff of a mixed strategy	expected payoff of mixed strategy $x$ against mixed strategy $y$ . See Equation (2)

Expected payoff of pure strategy  $s_i$ ,  $E(s_i, y)$

$$A = \begin{bmatrix} \pi_{s_1, s'_1} & \pi_{s_1, s'_2} \\ \pi_{s_2, s'_1} & \pi_{s_2, s'_2} \end{bmatrix}, \quad E(s_i, y) = \sum_{j=1}^{|S_j|} \pi_{s_i, s'_j} y_j \quad (1)$$

Expected payoff of mixed strategy  $x$ ,  $E(x, y)$

$$A = \begin{bmatrix} \pi_{s_1, s'_1} & \pi_{s_1, s'_2} \\ \pi_{s_2, s'_1} & \pi_{s_2, s'_2} \end{bmatrix}, \quad E(x, y) = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} \pi_{s_1, s'_1} & \pi_{s_1, s'_2} \\ \pi_{s_2, s'_1} & \pi_{s_2, s'_2} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad (2)$$

**Definition 1 (Strategy profile)** A strategy profile is the set of strategies taken by agents to play in one game. Assume a strategy profile with  $n$  strategies  $S = (s_1, \dots, s_n)$ . Each player gets one of them. All strategies except  $s_i$  are denoted by  $s_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n)$

**Definition 2 (Best Response:)**  $s_i$  is a best response to  $s_{-i}$  if and only if

$$\pi_{s_i, s_{-i}} \geq \pi_{s', s_{-i}} \quad \forall s' \in S$$

**Definition 3 (Nash Equilibrium:)** The strategy profile  $S = (s_1, \dots, s_n)$  is a Nash Equilibrium if and only if for each player  $i$ , his strategy  $s_i$  is a best response to  $s_{-i}$

## 1.4 Evolutionary Game Theory

Symbol	Name	Description
$x$	frequency of strategies	similar to the mixed strategy, but here the population frequency is given by $x$
$s_i$	strategy with index $i$	index corresponds to the frequency vector $x$
$f_i(x)$	fitness of strategy $s_i$	See Equation (3)
$\Phi(x)$	average population fitness	See Equation (4)
$\dot{x}$	replicator equation of $x$	shows the time and fitness dependent development of each strategies frequency, is computed for every $s_i$ separately, See Equation (5)

**Evolutionary Stable Strategy:** A strategy  $\sigma$  is an evolutionary stable strategy if:

- $\pi_{\sigma,\sigma} > \pi_{\tau,\sigma}$  for all possible mutant strategies  $\tau$
- or,  $\pi_{\sigma,\sigma} = \pi_{\tau,\sigma}$  and  $\pi_{\sigma,\tau} > \pi_{\tau,\tau}$

**Fitness of strategy  $s_i$ ,**  $f_i(x)$

$$fitness(s_i) = f_i(x) = payoff(s_i) = \sum_{j=1}^n x_j \pi_{ij} \quad (3)$$

**Average population fitness,**  $\Phi(x)$

$$\Phi(x) = \sum_{j=1}^n x_j f_j(x) \quad (4)$$

**Replicator Equation for infinite populations,**  $\dot{x}$

$$\dot{x} = x_i [f_i(x) - \Phi(x)] \quad (5)$$

There seems to be something wrong with our current definition of replicator equations for finite populations. Please ignore this part for now! We will update the slides and this overview after a careful revision of this part.

## 1.5 Additional Information on Evolutionary Game Theory

Symbol	Name	Description
$\beta$	selection intensity	used to apply weak selection, See Equation (6)
$\rho$	fixation probability	probability that a mutant strategy can take over the whole population, See Equation (7)
$\Delta f$	fitness difference of two strategies	
$b$	benefit	benefit for playing with a cooperating player, See Equation (8)
$c$	cost	cost for playing cooperation

**Weak Selection,**  $\dot{x}$

$$f_A(i) = 1 - \beta + \beta \pi_A(i)$$

$$\beta = 0 \Rightarrow \text{neutral selection} \quad (6)$$

$$\beta \ll 1 \Rightarrow \text{weak selection}$$

$$\beta = 1 \Rightarrow \text{selection by payoff}$$

**Probability for replacing a strategy using weak selection**  $\rho$

$$\rho = \frac{1}{1 + \exp^{-\beta \Delta f}} \quad (7)$$

**Comparison of using T,R,S,P or b,c for describing a payoff matrix**

$$\begin{array}{cc}
 & \begin{array}{cc} C' & D' \end{array} \\
 \begin{array}{c} C \\ D \end{array} & \begin{bmatrix} b-c & -c \\ b & 0 \end{bmatrix}
 \end{array}
 \quad
 \begin{array}{cc}
 & \begin{array}{cc} C' & D' \end{array} \\
 \begin{array}{c} C \\ D \end{array} & \begin{bmatrix} R & S \\ T & P \end{bmatrix}
 \end{array}
 \quad (8)$$

## 2 Reinforcement Learning

### 2.1 General Notation

Symbol	Name	Description
$t$	timestep $0 \dots T$	
$T$	final timestep $T$	
$\mathcal{S}$	set of possible states	
$s_t$	state at time $t$	$s_t \in \mathcal{S}$ , the $t$ is dropped in case the equation is independent of the actual time
$\mathcal{A}(s_t)$	available actions in $s_t$	
$a_t$	action at time $t$	$a_t \in \mathcal{A}(s_t)$ , the $t$ is dropped in case the equation is independent of the actual time
$r_t$ or $R_t$	reward at time $t$	after the agent left state $s_t$ by executing action $a_t$ , he receives the reward $r_{t+1}$ , where $r_{t+1} \in \mathbb{R}$ , depending on the context we use upper-case R to not confuse the single reward with the probability distribution at the bottom of this table
$\pi$	policy function	a probability distribution on the actions depending on the state
$\pi(a s)$	policy	probability of choosing action $a$ if in state $s$
$G_t, G_s$	(expected) return	sum of rewards from $s$ or timestep $t$ till the end $T$ (undiscounted return, See Equation (9)) or a non-ending episode (discounted return, See Equation (10))
$\gamma$	discount rate	affects the influence of future rewards on the return calculation (discounted return, See Equation (10))
$p(s', r s, a)$	state and reward probability depending on $a$ and $s$	also known as the markov property, probability of transitioning to state $s'$ and receiving reward $r$ is only dependent on choosing action $a$ in state $s$ ,
$p(s' s, a)$	transition probability	probability of transitioning into $s'$ after the agent executed its action $a$ in state $s$
$r(s, a)$	expected reward	valid for all Markov Decision Processes, See Equation (11)
$r(s, a, s')$	expected reward after the follow-up state is known	valid for all Markov Decision Processes, See Equation (12)



**Simple (undiscounted) Return,**  $G_t$ 

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T \quad (9)$$

**Discounted Return,**  $G_t$ 

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^T \gamma^k R_{t+k+1}, \quad \gamma \in [0, 1] \quad (10)$$

**Expected Reward,**  $r(s, a)$ 

$$r(s, a) = \mathbb{E}[R_{t+1} | s_t = s, a_t = a] = \sum_{r \in \mathbb{R}} \sum_{s' \in \mathbb{S}} r(s, a, s') \cdot p(s' | s, a) \quad (11)$$

**Expected Reward,**  $r(s, a)$ 

$$r(s, a, s') = \mathbb{E}[R_{t+1} | s_t = s, a_t = a, s_{t+1} = s'] = \frac{\sum_{r \in \mathbb{R}} r p(s', r | s, a)}{p(s' | s, a)} \quad (12)$$

## 2.2 Value and Action-Value Functions

Symbol	Name	Description
$v_\pi(s)$	value function	rates the value/expected return of a state following policy $\pi$ , See Equation (13) and Equation (14)
$\pi^*$	optimal policy	policy with the best possible expected return for all states
$v^*$	optimal value function	value function of the optimal policy $\pi^*$
$q_\pi(s, a)$	action-value function	rates the value/expected return of choosing an action in a given state following policy $\pi$ , See Equation (15) and Equation (16)
$V(s)$	approximation of $v(s)$	used in multiple iterative algorithms
$Q(s, a)$	approximation of $q(s, a)$	used in multiple iterative algorithms

**Value Function,**  $v(s)$

$$v_\pi(s) = \mathbb{E}_\pi[G_t | s_t = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \middle| s_t = s \right] \quad (13)$$

**Consistency Condition of the Value Function,**  $v(s)$

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[G_t | s_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \end{aligned} \quad (14)$$

**Action-Value Function,**  $q(s, a)$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | s_t = s, a_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \middle| s_t = s, a_t = a \right] \quad (15)$$

**Consistency Condition of the Action-Value Function,**  $q(s, a)$

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \end{aligned} \quad (16)$$

## 2.3 Dynamic Programming

---

### Algorithm 1: Iterative Policy Evaluation

---

**Input:** policy  $\pi$  to be evaluated

Initialize an array  $V(s) = 0$ , for all  $s \in S$

**repeat**

$\Delta \leftarrow 0$

**foreach**  $s \in S$  **do**

$v \leftarrow V(s)$

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

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

**until**  $\Delta < \theta$  (*a small positive number*)

**Output:**  $V \approx v_\pi$

---



---

### Algorithm 2: Policy Iteration

---

**Input:** policy  $\pi$  to be evaluated

Initialize an array  $V(s) = 0$ , for all  $s \in S$

**repeat**

    apply Iterative Policy Evaluation given  $\pi$

$policy-stable \leftarrow true$

**foreach**  $s \in S$  **do**

$a \leftarrow \pi(s)$

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

**If**  $a \neq \pi(s)$  **then**  $policy-stable \leftarrow false$

**until**  $policy-stable$

**Output:**  $V \approx v^*$ ,  $\pi \approx \pi^*$

---



---

### Algorithm 3: Value Iteration

---

Initialize an array  $V(s) = 0$ , for all  $s \in S$

**repeat**

$\Delta \leftarrow 0$

**foreach**  $s \in S$  **do**

$v \leftarrow V(s)$

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

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

**until**  $\Delta \leftarrow 0$  (*a small positive number*)

**Output:** a deterministic policy  $\pi$  such that

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

---

## 2.4 Monte Carlo Method

---

**Algorithm 4:** Monte Carlo Method  $v(s)$ 


---

**Input:** policy  $\pi$  to be evaluated

Initialize an array  $V(s) = 0$ , for all  $s \in S$

Initialize an empty list  $Returns(s)$ , for all  $s \in S$

**while** *true* **do**

    Generate an episode using  $\pi$

**foreach** *state  $s$  appearing in the episode* **do**

$G \leftarrow$  collected return after the first occurrence of  $s$

$Returns(s).append(G)$

$V(s) \leftarrow average(Returns(s))$

**Output:**  $V \approx v_\pi$

---



---

**Algorithm 5:** Monte Carlo Method  $q(s, a)$ 


---

**Input:** policy  $\pi$  to be evaluated

Initialize an array  $Q(s, a) = 0$ , for all  $s \in S, a \in A$

Initialize an empty list  $Returns(s, a)$ , for all  $s \in S$  and  $a \in A$

**while** *true* **do**

    Choose  $s_0 \in S$  and  $a_0 \in A(s_0)$

    Generate an episode starting from  $s_0, a_0$  following  $\pi$

**foreach**  *$(s, a)$  in the episode* **do**

$G \leftarrow$  return following the first occurrence of  $(s, a)$

$Returns(s, a).append(G)$

$Q(s, a) \leftarrow \max_a Q(s, a)$

**foreach**  *$s$  in the episode* **do**

$\pi(s) \leftarrow \arg \max_a Q(s, a)$

**Output:**  $Q \approx q_\pi, \quad \pi \approx \pi^*$

---

**Derivation of the Incremental Mean and Incremental Monte Carlo,  $\mu_k, V(s)$** 

$$\begin{aligned}
\mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\
&= \frac{1}{k} \left( x_k + \sum_{j=1}^{k-1} x_j \right) \\
&= \frac{1}{k} (x_k + (k-1)\mu_{k-1}) \\
&= \mu_{k-1} \frac{1}{k} (x_k - \mu_{k-1})
\end{aligned} \tag{17}$$

$$\begin{aligned}
V(s) &\leftarrow \frac{1}{k} \sum_{i=1}^k G_s(i) \\
&\leftarrow \frac{1}{k} \left( G_s(k) + \sum_{i=1}^{k-1} G_s(i) \right) \\
&\leftarrow \frac{1}{k} (G_s(k) + (k-1) V(s)) \\
&\leftarrow V(s) + \frac{1}{k} (G_s(k) - V(s))
\end{aligned} \tag{18}$$

---

**Algorithm 6:** Constant- $\alpha$  Monte Carlo Method  $v(s)$ 

---

**Input:** policy  $\pi$  to be evaluatedInitialize an array  $V(s) = 0$ , for all  $s \in S$ **while** *true* **do**    Generate an episode using  $\pi$         **foreach** *state  $s$  appearing in the episode* **do**             $G_s \leftarrow$  collected return after the first occurrence of  $s$              $V(s) \leftarrow V(s) + \alpha[G_s - V(s)]$ **Output:**  $V \approx v_\pi$ 

---

## 2.5 Temporal Difference Learning

Symbol	Name	Description
$G_t^{t+n}$	n-step return	See Equation (19)
$h$	horizon = $t + n$	number of future reward steps used for return calculation, if $h > T$ the return is the sum of rewards till the end of the episode
$V_t$	approximation of $v$ at timestep $t$	See Equation (19)
$\Delta_t(s_t)$	error based adaptation	Weighted error used in temporal difference learning, See Equation (19)

**n-Step Return,**  $G_t^{t+n}(s)$

$$\begin{aligned}
 G_t^{t+n}(s) &= R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_h + \gamma^n V_t(s_{t+n}) \\
 \Delta_t(s_t) &= \alpha [G_t^{t+n}(V_t(s_{t+n})) - V_t(s_t)] \\
 V_{t+1}(s) &= V_t(s) + \Delta_t(s), \quad \forall s \in \mathcal{S}
 \end{aligned} \tag{19}$$

---

### Algorithm 7: Temporal Difference Learning

---

**Input:** policy  $\pi$  to be evaluated

Initialize an array  $V(s)$  arbitrarily, for all  $s \in \mathcal{S}$

**foreach** for each episode **do**

    Initialize  $S$

**foreach**  $s \in \mathcal{S}$  **do**

$A \leftarrow$  action given by  $\pi$  for  $S$

        Take action  $A$ ; observe reward  $R$ , and next state,  $S'$

$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

**Output:**  $V \approx v_\pi$

---

---

**Algorithm 8:** One Step Q-Learning

---

Initialize a matrix  $Q(s, a)$  arbitrarily, for all  $s \in S$ , and  $a \in A$   
**foreach** *terminal state*  $s$  **do**  $Q(s, \cdot) = 0$   
**foreach** *episode* **do**  
    **foreach** *state*  $s_t$  *in episode* **do**  
        Choose  $a_t$  from  $s_t$  using policy derived from  $Q$   
        take action  $a_t$ , observe reward  $r$ , and follow-up state  $s_{t+1}$   
         $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$   
         $s \leftarrow s_{t+1}$   
**Output:**  $Q(s_t, a_t)$

---

### 3 Simulation Based Search Algorithms

#### 3.1 Flat Monte Carlo

---

**Algorithm 9:** Flat Monte Carlo

---

**Input:** Initial state  $s_t$  and MDP of the environment  
Initialize an array  $Q(s_t, a) = 0$ , for all  $a \in A$   
  
**foreach**  $a \in A$  **do**  
    Simulate  $k$  Episodes using random policy  
     $Q(s_t, a) \leftarrow \text{mean/expected return of simulated episodes}$   
**Output:**  $\arg \max_a Q(s_t, a)$

---

#### 3.2 Flat Upper Confidence Bounds (UCB)

---

**Algorithm 10:** Flat UCB

---

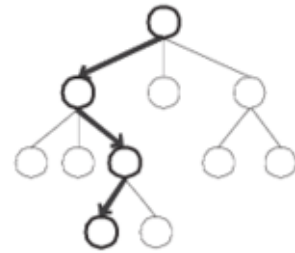
**Input:** any UCB variant  
Initialize arrays  $Q(s_t, a) = 0$ ,  $N(s_t) = 0$ ,  $N(s_t, a) = 0$ ,  $\forall a \in A, \forall a \in A$   
Initialize an empty list of  $Returns(s_t, a) \forall a \in A$   
  
**for**  $k$  *iterations* **do**  
    Select first action  $a$  of the simulation using UCB  
    Simulate episode using random policy  
    Update visit counts  $N(s_t) \leftarrow k$  and  $N(s_t, a) \leftarrow N(s_t, a) + 1$   
     $Returns(s_t, a).append(\text{return of episode } k)$   
     $Q(s_t, a) \leftarrow \text{mean}(Returns(s_t, a));$   
**Output:**  $\arg \max_a Q(s_t, a)$

---

### 3.3 Monte Carlo Tree Search (MCTS)

- **1. Tree Selection:**

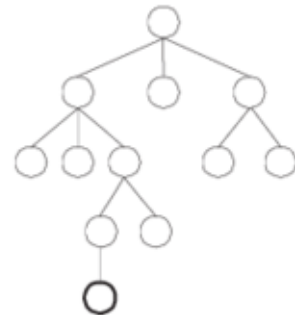
Following the **tree policy** (i.e. UCB1), navigate the tree until reaching a node with at least one child state that was not explored yet.



- **2. Expansion:**

Add a new node in the tree, as a child of the node reached in the tree selection step.

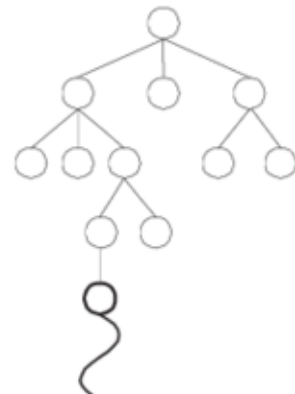
- this node resembles an action



- **3. Monte Carlo Simulation/Rollout:**

Following the **default policy**, advance the state until a terminal state or a pre-defined maximum number of steps.

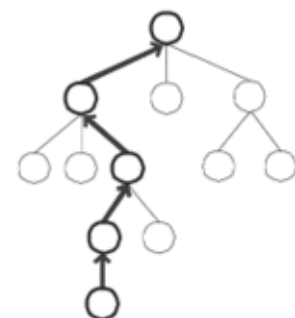
- calculate the return of this episode.



- **4. Back-propagation:**

Update the values of  $Q(s, a)$ ,  $N(s)$  and  $N(s, a)$  of the nodes visited in the tree during steps 1 and 2.

- Values of nodes visited during the simulation will not be updated.





### 3.4 Evolutionary Algorithms (EA)

Symbol	Name	Description
$i$	individual/solution	a possible solution for the problem at hand
$P$	population of individuals	set of possible solutions
$t$	time-step or generation	
$P_t, P(t)$	population at time $t$	
$\mu, N$	(number of) individuals in $P$	equal to $ P $
$M$	mating pool	individuals selected for genetic operators
$\lambda$	(number of) individuals in $M$	set of children
$f_i$	fitness of individual $i$	scoring the fitness of an individual based on an objective function
$p_i$	selection probability of individual $i$	probability of an individual being selected for the mating pool

---

**Algorithm 11:** General Evolutionary Algorithm

---

**Input:** Initial Population  $P(t)$

$t \leftarrow 0$

**repeat**

    Select individuals for genetic operators:  $M(t) \leftarrow \text{selection}(P(t))$

    Apply Crossover:  $M'(t) \leftarrow \text{crossover}(M(t))$

    Apply Mutation:  $M''(t) \leftarrow \text{mutation}(M'(t))$

    Update  $P(t)$ :  $P(t+1) \leftarrow \text{reproduction\_scheme}(P(t) \cup M''(t))$

$t \leftarrow t + 1$

**until** *termination condition*

**Output:** Final Population  $P(t)$

---

### 3.5 Rolling Horizon (RH)

---

**Algorithm 12:** Rolling Horizon Random Search
 

---

**Input:** current state  $s_t$ , horizon  $h$

$E \leftarrow$  simulate a random episode starting at  $s_t$  till timestep  $t + h$

**repeat**

$E' \leftarrow$  simulate a random episode starting at  $s_t$  till timestep  $t + h$

**if** *Return of  $E'$*  > *Return of  $E$*  **then**

$E \leftarrow E'$

**until** *time is over*

**Output:** First action of  $E$

---



---

**Algorithm 13:** Rolling Horizon Evolutionary Algorithm (RHEA)
 

---

**Input:** current state  $s_t$ , horizon  $h$

$P \leftarrow$  population of random episodes starting at  $s_t$  till timestep  $t + h$

**repeat**

$M \leftarrow$  select individuals from  $P$  for genetic operators

$M' \leftarrow$  apply crossover to  $M$

$M'' \leftarrow$  apply mutation to  $M'$

$P \leftarrow$  apply reproduction scheme to  $(P \cup M'')$

**until** *time is over*

**Output:** First action of the episode with the highest return in  $P$

---

## 4 Multi Objective Optimization

### 4.1 General Notation

Symbol	Name	Description
$S$	search space	the space of all possible solutions
$O$	objective space	$O = \{\vec{f}(x) \in \mathbb{R}^m   x \in S\}$
$\vec{f}(x)$	objective vector	all objective values for individual $x$
$f_i$	objective function	assigns an objective value to each solution
$f_i(x)$	$i$ -th objective value of $x$	
$r_i$	rank of solution $i$	see the lecture for ranking methods
$d_i$	(crowding) distance of $i$	specifically used in NSGA-II
$\Omega_i(x)$	cone domination	cone-dominating version of objective function $f_i$
$C(A, B)$	convergence metric	compare the convergence between two non-dominated set of points $A$ and $B$
$HV(A)$	hyper-volume of $A$	calculate the hyper-volume of a non-dominated set of solutions $A$
$MHV(i)$	marginal hyper volume of individual $i$	$MHV(i) = HV(entire\ Set) - HV(set\ without\ i)$

### 4.2 NSGA-II Algorithm

---

**Algorithm 14:** Crowding Distance Computation

---

**Input:** A single front  $F$  with  $N$  individuals

Initialize distances for all individuals:  $d_i = 0, \quad \forall i \in F$

**foreach** each objective  $m$  **do**

Sort individuals in  $F$  regarding their objective values  $f_m(i)$

Assign  $d_i \leftarrow \infty$  to the elements with index 1 and  $N$  in the sorted list [the elements with lowest or highest value in  $f_m(i)$ ]

$f_m^{min}, f_m^{max} \leftarrow$  minimum and maximum observed value for objective value  $f_m(i)$  for individuals in  $F$

**for** all individuals  $i$  with index  $i_{idx} = 2$  to  $N - 1$  **do**

$d_i \leftarrow d_i + \left| \frac{f_m(i_{idx+1}) - f_m(i_{idx-1})}{f_m^{max} - f_m^{min}} \right|$

**Output:** Crowding Distance  $d(i), \quad \forall i \in F$

---

---

**Algorithm 15:** NSGA-II Algorithm

---

**Input:** Individuals of Population  $P(t)$  and the Mating Pool  $M''(t)$ Fronts  $(F_1, F_2, \dots) \leftarrow$  non-dominated sorting of  $P(t) \cup M''(t)$ Set  $P(t+1) = \emptyset$ ,  $i = 1$ ,  $N = |P(t)|$ **repeat**

$$\begin{array}{|l} P(t+1) \leftarrow P(t+1) \cup F_i \\ i \leftarrow i + 1 \end{array}$$
**until**  $|P(t+1)| + |F_i| > N$ Sort all individuals in  $F_i$  by their crowding distance $P(t+1) \leftarrow P(t+1) \cup \{\text{first } N - |P(t+1)| \text{ elements in } F_i\}$ **Output:** Next Population  $P(t+1)$ 

---