

## Reinforcement Learning Project

# Evolutionary Function Approximation for Reinforcement Learning

Yunian Pan

ECE department

May. 10th 2019



# Evolutionary Function Approximation for Reinforcement Learning

## Outline:

- ▶ Some Background
- ▶ Methods and Results
- ▶ Experiments
- ▶ Conclusion

In a word, DQN(Deep Q-Learning) + NEAT(Neuroevolution of augmenting topologies)

---

## Algorithm 1: Q-Learning

---

**Input:**  $S$  : set of all states;  $A$  : set of all actions;  $\sigma$  : standard deviation of initial weights;  $c$  : output scale;  $\alpha$  : learning rate;  $\gamma$  : discount factor;  $\lambda$  : eligibility decay rate;  $\epsilon_{td}$  : exploration rate;  $e$  : total number of episodes;

**Initialize:**  $N \leftarrow \text{INIT} - \text{NET}(S, A, \sigma)$  ;

**for**  $i \leftarrow 1$  to  $e$  **do**

$s, s' \leftarrow \text{null}, \text{INIT-STATE}(S)$ ;

**while** *Terminal-state*( $s$ ) **do**

$Q[] \leftarrow c \times \text{EVAL} - \text{NET}(N, s')$ ;

**With-prob**( $\epsilon_{td}$ )  $a' \leftarrow \text{RANDOM}(A)$ ;

**else:**  $a' \leftarrow \arg \max Q[j]$ ;

**if**  $s \neq \text{null}$  **then**

$\text{BACKPROP}(N, s, a, (r + \gamma \max_j Q[j])/c, \alpha, \gamma, \lambda)$  ;

**else**

$s, a \leftarrow s', a'$ ;

$r, s' \leftarrow \text{TAKE} - \text{ACTION}(a')$

**end**

**end**

**end**

---

# Genetic Algorithm

## Genetic Algorithms

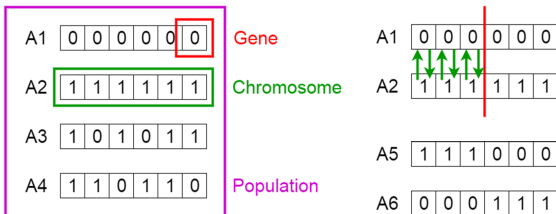


Figure 1: Illustration from Google image

# Genetic Algorithm

Standard evolution framework:

- (1) Initialize population
- (2) Evolve from  $1, \dots, n_{\text{generation}}$ :
  - a) Select from population according to fitness
  - b) Generate offspring through crossover and mutation
  - c) Replace population with offspring.

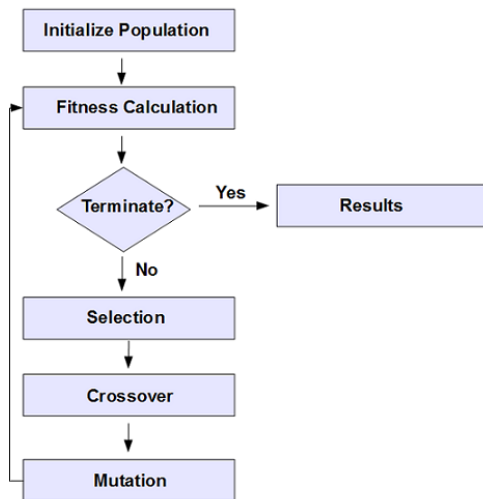
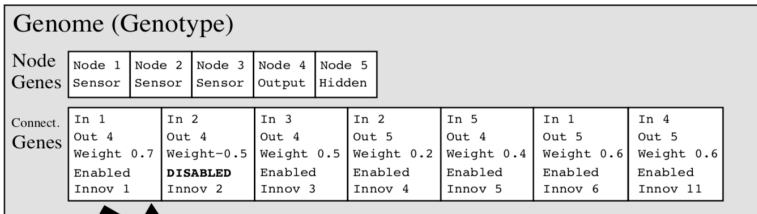


Figure 2: flow chart from Google image

# Network encoding



Network (Phenotype)

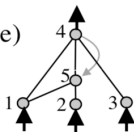


Figure 3: from Stanley, Kenneth O., and Risto Miikkulainen. (2002): 99-127.

# Crossover

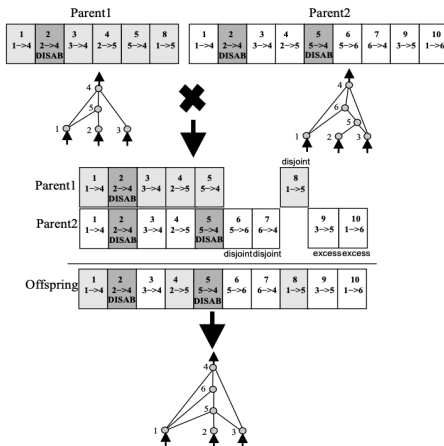


Figure 4: from Stanley, Kenneth O., and Risto Miikkulainen. (2002): 99-127.



## mutation

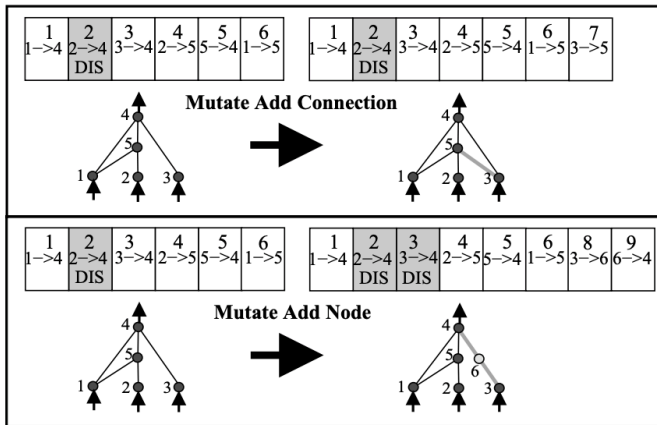


Figure 5: from Stanley, Kenneth O., and Risto Miikkulainen. (2002): 99-127.

## Algorithm 2: NEAT( $S, A, p, m_n, m_l, g, e$ )

**Input:**  $S$ : set of all states,  $A$ : set of all actions,  $p$ : population size,  $m_n$ : rate of adding node,  $m_l$ : rate of adding link,  $g$ : generations,  $e$ : episodes;

**Initialize:**  $P[] \leftarrow \text{Init-Populations}(S, A, p)$  ;

**for**  $i \leftarrow 1$  to  $g$  **do**

**for**  $j \leftarrow 1$  to  $e$  **do**

$N, s, s' \leftarrow \text{Random}(P), \text{null}, \text{INIT-STATE}(S)$ ;

**while**  $\text{Terminal-state?}(s)$  **do**

$Q[] \leftarrow \text{EVAL-NET}(N, s')$ ;

$a' \leftarrow \arg \max Q[j]$ ;  $s, a \leftarrow s', a'$ ;  $r, s' \leftarrow \text{TAKE-ACTION}(a')$ ;

$N.\text{fitness} \leftarrow N.\text{fitness} + r$

**end**

$N.\text{episodes} \leftarrow N.\text{episodes} + 1$

**end**

$P' \leftarrow \text{new array of size } p$ ;

**for**  $j \leftarrow 1$  to  $p$  **do**

$P'[j] \leftarrow \text{Breed-Net}(P[j])$ ;

**with-probability**  $m_n$ :  $\text{ADD-Node-Mutation}(P'[j])$ ;

**with-probability**  $m_l$ :  $\text{ADD-link-Mutation}(P'[j])$

**end**

$P[] \leftarrow P'[]$

**end**

---

## Algorithm 3: NEAT+Q( $S, A, p, m_n, m_l, g, e$ )

---

**Input:**  $S, A, p, m_n, m_l, g, e, \alpha, \lambda, \gamma, \epsilon$

**Initialize:**  $P[] \leftarrow \text{Init-Populations}(S, A, p)$  ;

**for**  $i \leftarrow 1$  to  $g$  **do**

**for**  $j \leftarrow 1$  to  $e$  **do**

$N, s, s' \leftarrow \text{Random}(P), \text{null}, \text{INIT-STATE}(S)$ ;

**while**  $\text{Terminal-state?}(s)$  **do**

$Q[] \leftarrow \text{EVAL-NET}(N, s')$ ;

**With-prob**( $\epsilon_{td}$ )  $a' \leftarrow \text{RANDOM}(A)$ ;

**else:**  $a' \leftarrow \arg \max Q[j]$ ;

**if**  $s \neq \text{null}$  **then**

$\text{BACKPROP}(N, s, a, (r + \gamma \max_j Q[j])/c, \alpha, \gamma, \lambda)$

**end**

$s, a \leftarrow s', a'$ ;

$r, s' \leftarrow \text{TAKE-ACTION}(a')$ ;

$N.\text{fitness} \leftarrow N.\text{fitness} + r$

**end**

$N.\text{episodes} \leftarrow N.\text{episodes} + 1$

**end**

    Crossover and mutation

**end**

---

## Boltzman Selection

Exploration probability:

$$\Pr(\cdot|s) = \frac{e^{Q(s,\cdot)/\tau}}{\sum_{a \in A} e^{Q(s,a)/\tau}}$$

Can we use it in network choosing? Yes.

$$\Pr(\cdot) = \frac{e^{S(\cdot)/\tau}}{\sum_{q \in P} e^{S(q)/\tau}}$$

---

**Algorithm 4:** Boltzman Selection( $P, \tau$ )

---

**Input:**  $P$ : population,  $\tau$ : softmax temperature**if**  $\exists N \in P \mid N.episodes = 0$  **then**  
| return  $N$ **else**|  $total \leftarrow \sum_{N \in P} e^{N.average/\tau};$ | **for**  $N \in P$  **do**| | with-prob ( $\frac{e^{N.average/\tau}}{total}$ ) return  $N$  else  
| |  $total \leftarrow total - e^{N.average/\tau}$ | **end****end**

---

## Some Comparison

- ▶ Online v.s. Offline
- ▶ Darwinian v.s. Lamarckian
- ▶ Annealing v.s. Without Annealing

## Experiments from the paper

- ▶ mountain car

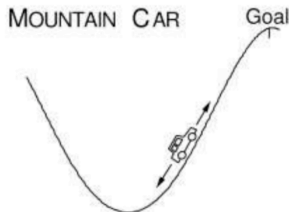
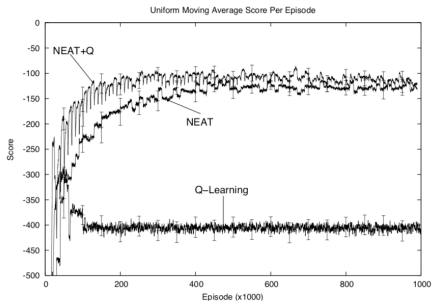


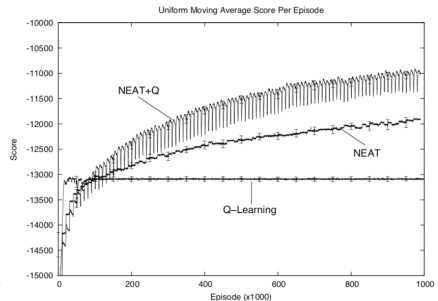
Figure 6: from Sutton and Barto (1998)

- ▶ Server job scheduling

# Results



(a) Mountain Car

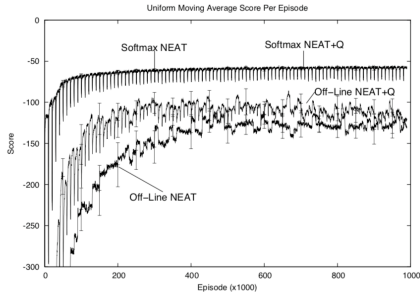


(b) Server Job Scheduling

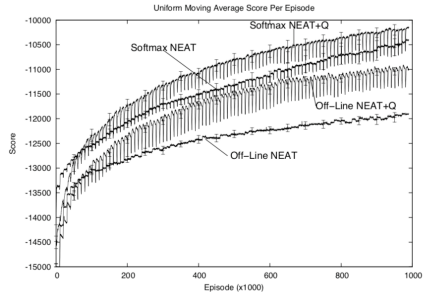
Figure 7: from Whiteson, Shimon, and Peter Stone. (2006): 877-917.  
Figure 6



# Results

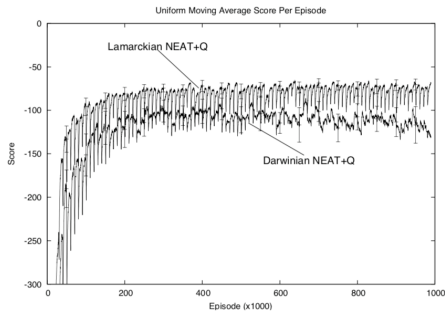


(a) Mountain Car

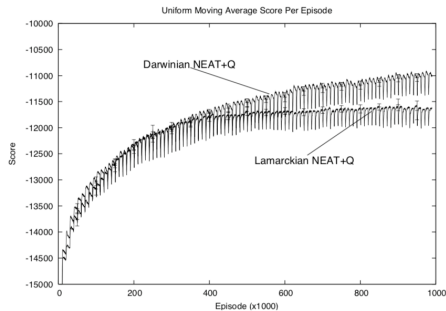


(b) Server Job Scheduling

Figure 8: from Whiteson, Shimon, and Peter Stone. (2006): 877-917.  
Figure 7

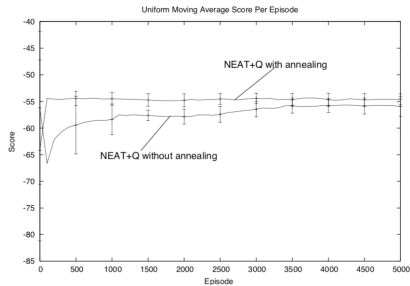


(a) Mountain Car

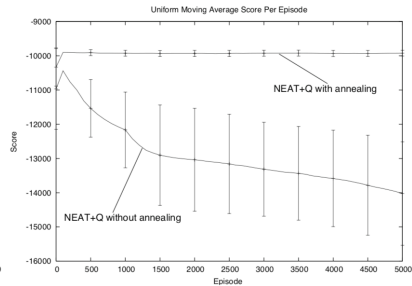


(b) Server Job Scheduling

Figure 9: from Whiteson, Shimon, and Peter Stone. (2006): 877-917.  
Figure 10



(a) Mountain Car



(b) Server Job Scheduling

Figure 10: from Whiteson, Shimon, and Peter Stone. (2006): 877-917. Figure 11

## Experiments

Let's try it!

- (a) Cartpole Balancing using NEAT and Q
- (b) test on Atari game Pacman on both NEAT and Q

## Conclusion

- ▶ NEAT outperform Q-learning in episodes, Generally, (NEAT+Q can perform better!)
- ▶ NEAT explore the function representation automatically.
- ▶ Some other Methods(DDQN and Duel QN, PG) can be combined with NEAT
- ▶ Non-stationary environment is challenging.

**Thank You!**