Studend science session 20.05.2021

# Topic: Reward value effects on Q-learning agents

*prepared: msc.eng. Martin Kaloev*

# Abstract

R.U "Angel Kunchev"

- One of the great advances in the field of artificial intelligence is the creation of general artificial intelligence (AGI). Or artificial intelligence allowing an agent to perform tasks of human complexity. This is due to the development of neural networks that implement in practice the concepts described by adaptive programming.
- Adaptive programming is a set of algorithms and methods describing the relationships between actions and states.
- Artificial intelligence of this type plays an important role in the modern world, in the areas like targeted advertising; digital assisted trade; robotics; transport and others.
- What all these agents have in common is the ability to learn from observed processes and interact with users.
- This article discusses the work of such an agent.The focus of the studie is the relationship between the change in the policy used by the agent and the change in the value of the rewards received from an action.
- The agent's task is to solve a two-dimensional spatial problem using a greedy epsilon policy. The algorithms used by the agent are an algorithm for determine the quality of actions, a sarsa algorithm and algorithms for evaluation quality of the policy.

# Related work

- Articles on a similar topic can be found with the following titles:
- An Analysis of Q-Learning Algorithms with Strategies of Reward Function
- https://www.researchgate.net/publication/50247491_An_Analysis_of_Q-Learning_Algorithms_with_Strategies_of_Reward_Function
- Q-Learning Algorithms: A Comprehensive Classification and Applications
- https://www.researchgate.net/publication/335805245_Q-Learning_Algorithms_A_Comprehensive_Classification_and_Applications
- Playing Atari with Deep Reinforcement Learning
- https://arxiv.org/pdf/1312.5602v1.pdf
- Low-rank State-action Value-function Approximation
- https://arxiv.org/abs/2104.08805v1

# Algorithms list

R.U "Angel Kunchev"

- The list of used algorithms:
- *-Sarsa:*
- *-Quality of action equation by R.E.Bellman*
- *-Value of action equation by R.E.Bellman*
- *-Policy based on action equation.*
- Those equations belongs to the family of discrete mathematics.
- Discrete mathematics is the study of mathematical structures that are countable or otherwise distinct and separable. Examples of structures that are discrete are combinations, graphs, and logical statements.
- Please note the slide named: **Model and equation**
- Manuals for application of algorithms:
- **Reinforcement Learning:**
- **An Introduction**
- by:
- *Richard S. Sutton and Andrew G. Barto*
- A Bradford Book
- The MIT Press
- Cambridge, Massachusetts
- London, England
- *at: http://incompleteideas.net/book/first/ebook/node41.html*

Sarsa:

$$Q(s_t, a_t) = Q(s_t, a_t) + a^*(r_t + y^*Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$

Q_l belman curent quality for action

$$Q(s_t, a_t) = Q(s_t, a_t) + a^*\left(r_t + y^*\max_a Q(s_{t+1}, a) - Q(s_t, a_t)\right)$$

which can also be written as

$$Q(s_t, a_t) = (1 - \alpha)^*Q(s_t, a_t) + \alpha^*\left(r_t + y^*\max_a Q(s_{t+1}, a)\right)$$

Value:

$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1})|S_t = s]$$

Quality

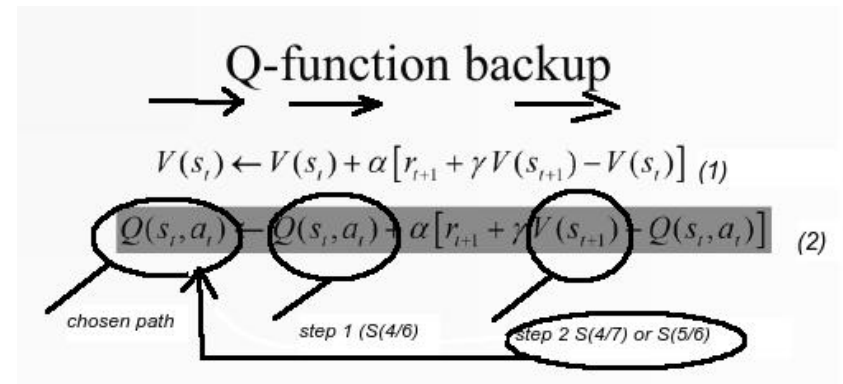$$q_\pi(s, a) = \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$$

Policy:

$$\text{TDerror}(s) = V^\pi(s) - \sum_{s'} T(s, \pi(s), s')[r(s, \pi(s), s') + \gamma V^\pi(s')]$$

# Model and Equation

Q-function backup

$$V(s_t) \leftarrow V(s_t) + \alpha \left[ r_{t+1} + \gamma V(s_{t+1}) - V(s_t) \right] \quad (1)$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma V(s_{t+1}) - Q(s_t, a_t) \right] \quad (2)$$

chosen path    step 1 (S(4/6))    step 2 S(4/7) or S(5/6)

# Model and Equation Explained

- In the slide: Model and equation you can see an equation that can be interpreted as:
- The quality of a certain action (Q) is determined by the change from initial state (S) to another state (St).
- Q has two attributes (S-state) and (A or V -value of action). Actions of higher quality have priority in choosing a policy or strategy.

- In this particular model, S (4/6) can choose to go into several new states. The state S (4/7) has a higher quality compared to the state S (5/6). From the agent's point of view, it's a better policy to go around a wall instead of crashing into it.

# Test recreation

- To repeat the tests performed, please follow the steps below:
- 1) Use a virtual tool of your choice to navigate the World Wide Web to reach current github address
- 2) Download a copy of the software solution used in the tests. (Note **the github navigation** slides)
- 3) Use a virtual tool of your choice to navigate the World Wide Web to reach current jypiterlab address
- 4) Use the menus to access the online version of the software lab (Note **the jypiterlab navigation** slides)
- 5) Load a file named main _ ** _ main in the online version of jypiterlab and run the simulation. (Note the slides **jypiterlab: how to use it**)

# The github navigation

# The jypiterlab navigation

# jypiterlab: How to use it

# Test Explained

- The virtual tool main _ ** _ main creates a simulation of an agent solving a randomly generated maze in an optimal way. The user can change the value (1) of the actions for each step and the penalties on the assessment of action in violation of the rules of the maze.

- This allows different agent behavior to be observed at different values set by the user.

- In addition, the functionality in cell 7 can be replaced by a Bellman algorithm with a Sarsa algorithm (2).

- 1*see Slide Model and equation explained again
- 2*see sarsa_cell_update sniped

# Results

- Tests are performed exhausting all possible combinations. After the tests are completed, a statistical analysis is performed as to which factors play the biggest role in the violation of the rules of the labyrinth by the agent. The results of the tests show that the difference between the value of the detour compared to the value of the penalty in the event of a breach plays the biggest role. The detour path means the path connecting two points that are closest to the breakthrough point.

# Discussions and Optimizations

- A possibility to optimize the virtual agent is the addition of a library signaling a possible breakthrough. An example of such a library is shown in alarm_lib.py provided in the project repository.

$$fw = \begin{cases} true, [\sum_{n=1}^{Vp} n * Rp] > [\sum_{n=1}^{Vw} n * Rw] \\ false, [\sum_{n=1}^{Vp} n * Rp] < [\sum_{n=1}^{Vw} n * Rw] \end{cases}$$

# Where are we?

- Let's take a few minutes to discuss the role of Q-learning as the basis of many of the technologies used in the modern world.

# Conclusion

- The article analyzes the change in the behavior of an agent when changing the values of actions.

- Answers to questions about why rewards often have a negative value instead of a positive one. The moments in which the agent violates the rules of the given task are considered and analyzed. A solution for the specific problems of a problem derived in an equation and a software library is proposed.

# Thank you for your attention