

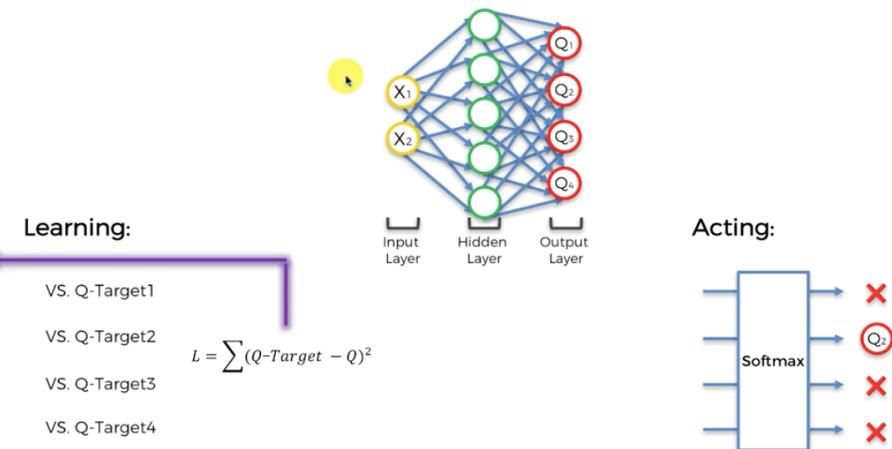
Action selection policies

Action Selection Policies

Artificial Intelligence A-Z

© SuperDataScience
udemy

Action Selection Policies

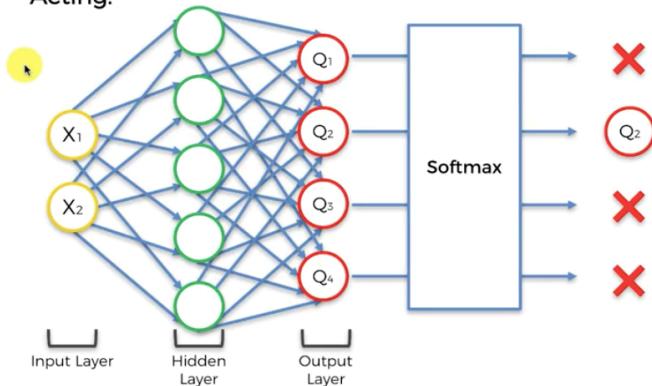


Artificial Intelligence A-Z

© SuperDataScience
udemy

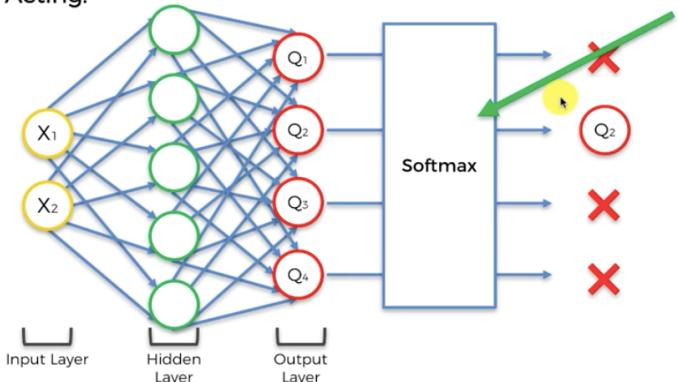
Action Selection Policies

Acting:



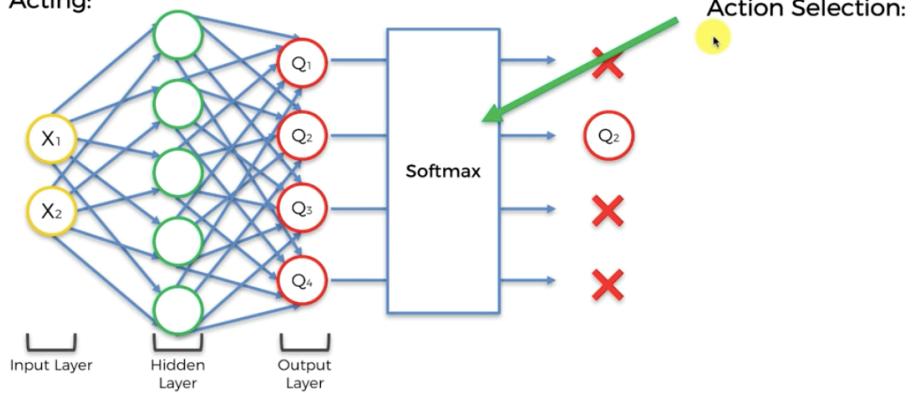
Action Selection Policies

Acting:



Action Selection Policies

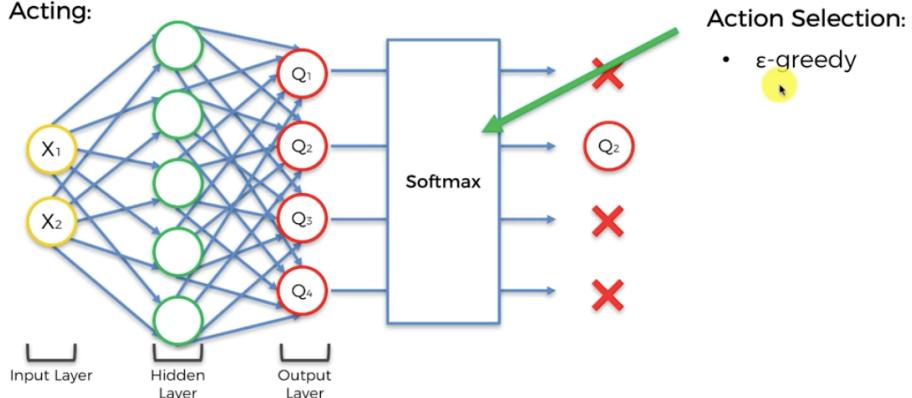
Acting:



Action Selection:

Action Selection Policies

Acting:

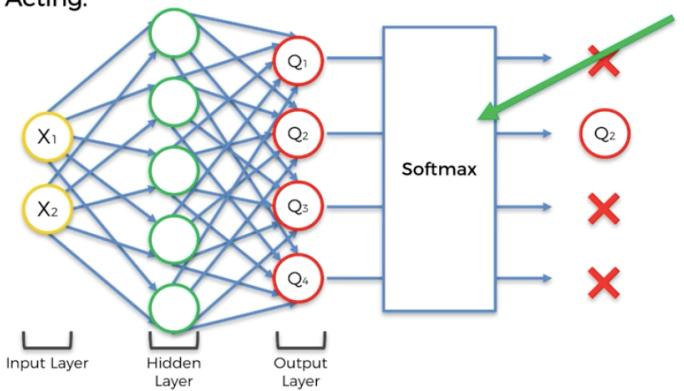


Action Selection:

- ϵ -greedy

Action Selection Policies

Acting:

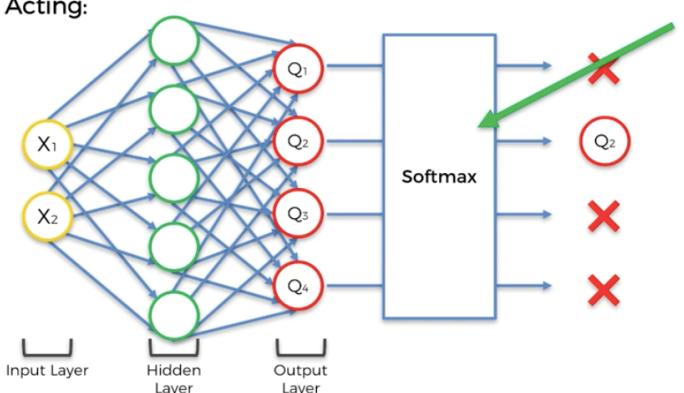


Action Selection:

- ϵ -greedy
- ϵ -soft ($1 - \epsilon$)

Action Selection Policies

Acting:

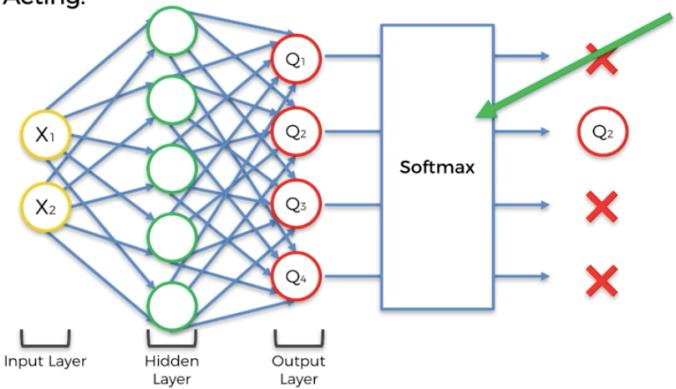


Action Selection:

- ϵ -greedy
- ϵ -soft ($1 - \epsilon$)
- Softmax

Action Selection Policies

Acting:



Action Selection:

- ϵ -greedy
- ϵ -soft ($1-\epsilon$)
- Softmax

Exploration
vs
Exploitation

The reason for having different action selection: Exploration and Exploitation.

Sometimes the environment forces the agent to explore different actions.

Exploration and exploitation is really important since sometimes the agent get stuck in the local maximum. For example, the agent might find based on its initial exploration that a certain action is the best action even though it's not and there is a bias toward this action and the agent thinks this is a very good action and it keeps getting good reward and it keeps doing that but there might be a better action than this. Here is when action selection going to help us.

ϵ – greedy: This will select the one with the highest Q value except for the epsilon percent of the time. For instance, if the epsilon is equal to 0.1 or 10% then 10% of the time, action is going to be selected at random and 90% it is going to be selected based on the best action based on the highest Q value.

That's why it's called greedy because we greedily choose good action except for a little epsilon. The lower the epsilon, the greedier algorithm is.

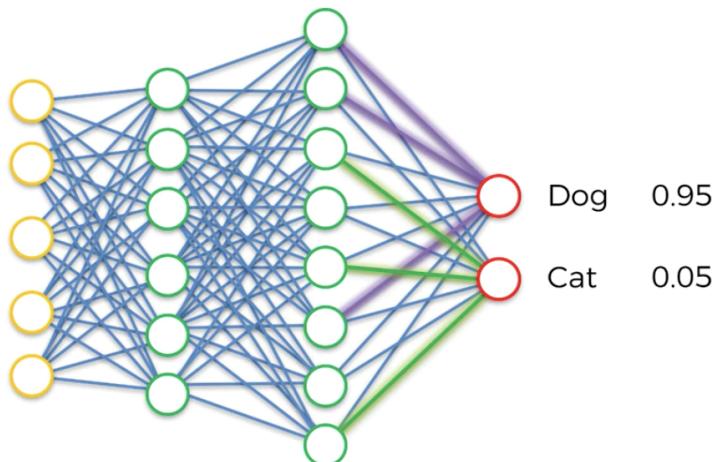
ϵ – soft ($1 - \epsilon$): This is exactly opposite of the ϵ – greedy that it will select at random except epsilon percent of the time.

Action Selection Policies



.....

Flattening



Artificial Intelligence A-Z

© SuperDataScience

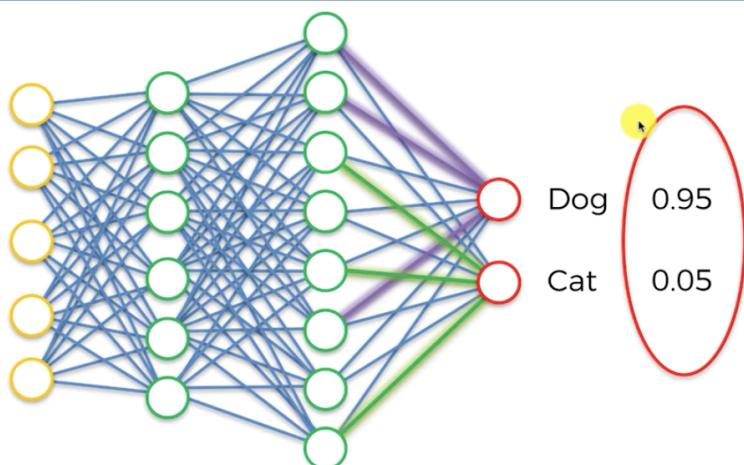
udemy

Action Selection Policies



.....

Flattening



Artificial Intelligence A-Z

© SuperDataScience

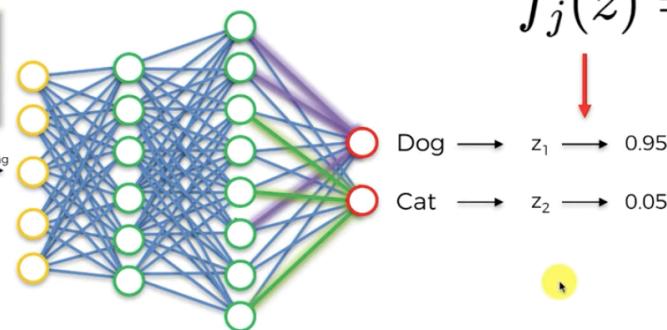
udemy

Action Selection Policies



.....

Flattening

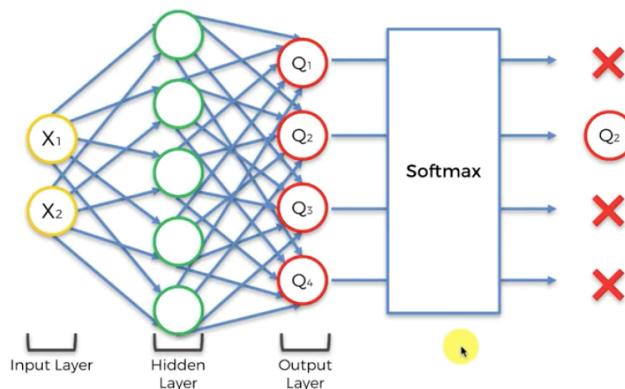


$$f_j(z) = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

This is the softmax function that we are applying. After applying we can add our number in the total of number 1.

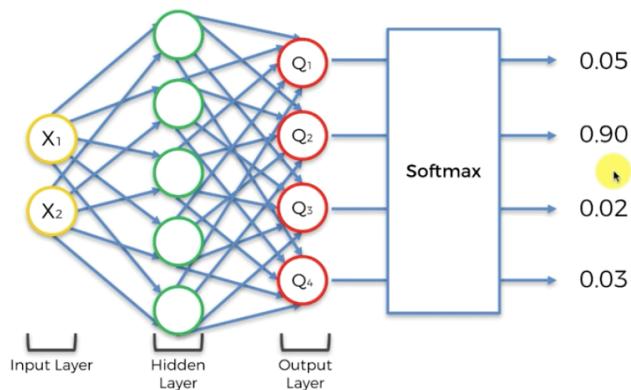
Action Selection Policies

Acting:



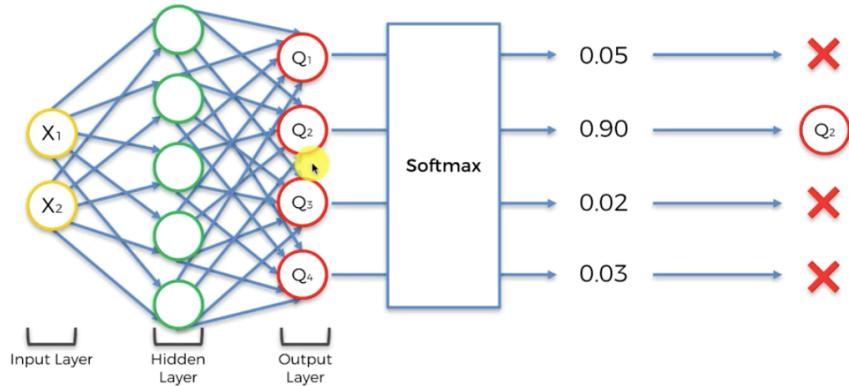
Action Selection Policies

Acting:



Action Selection Policies

Acting:

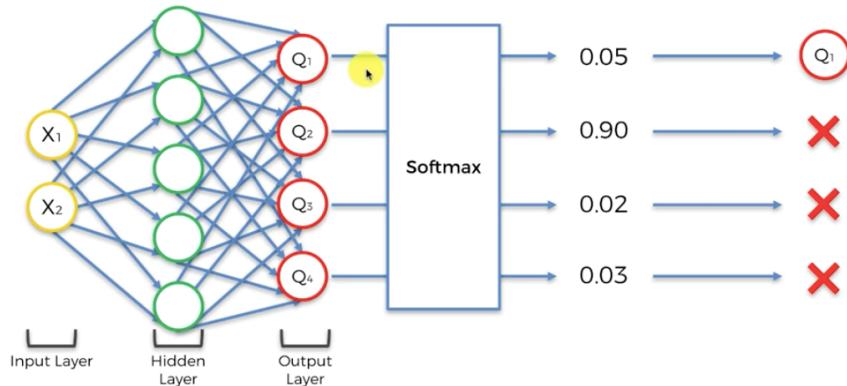


In here the best action is going to have the highest probability because it has the highest Q value. In here we are going to take Q_2 90% of the times but still 5% of times we are going to take Q_1 and so on.

After each update, these values are going to be change.

Action Selection Policies

Acting:



Artificial Intelligence A-Z

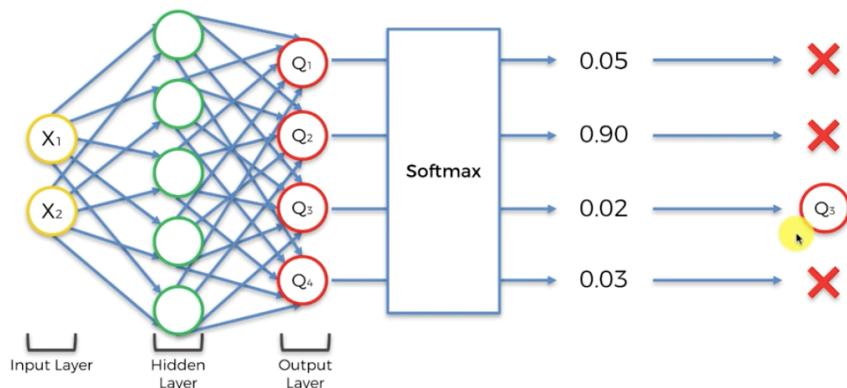
© SuperDataScience

udemy

5% of time we choose Q_1 .

Action Selection Policies

Acting:



Artificial Intelligence A-Z

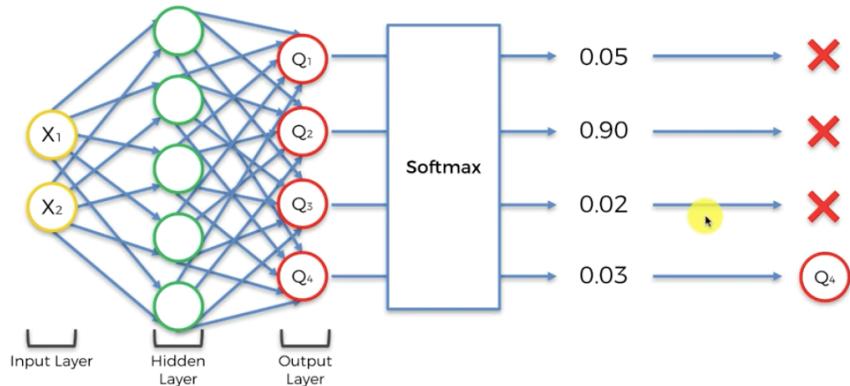
© SuperDataScience

udemy

2% of time we choose Q_3 .

Action Selection Policies

Acting:



Artificial Intelligence A-Z

© SuperDataScience

udemy

3% of time we choose Q_4 .

Additional Reading

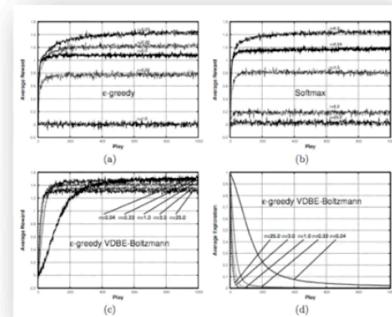
Additional Reading:

Adaptive ϵ -greedy Exploration in Reinforcement Learning Based on Value Differences

Michel Tölk (2010)

Link:

<http://tokic.com/www/tokicm/publikationen/papers/AdaptiveEpsilonGreedyExploration.pdf>



Artificial Intelligence A-Z

© SuperDataScience

udemy