# Different than the rest of machine learning

Predictive machine learning



Data     Model     Predictions

Reinforcement learning

Environment

- Markov Decision Processes (MDP)

- The Bellman equation

- Q-networks

- Policy gradients

# Temporal Difference (TD) Learning



Learn at each time step



Each network gets a cost
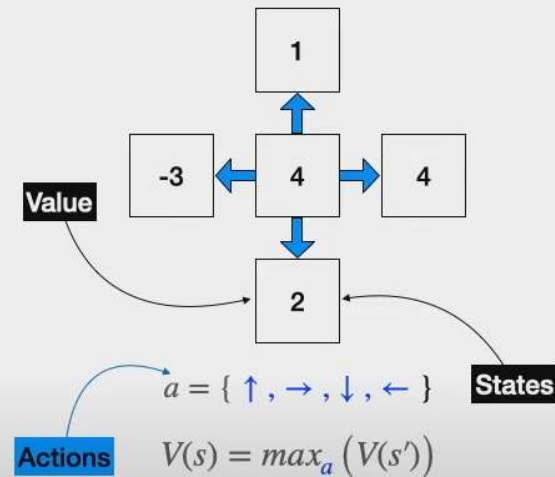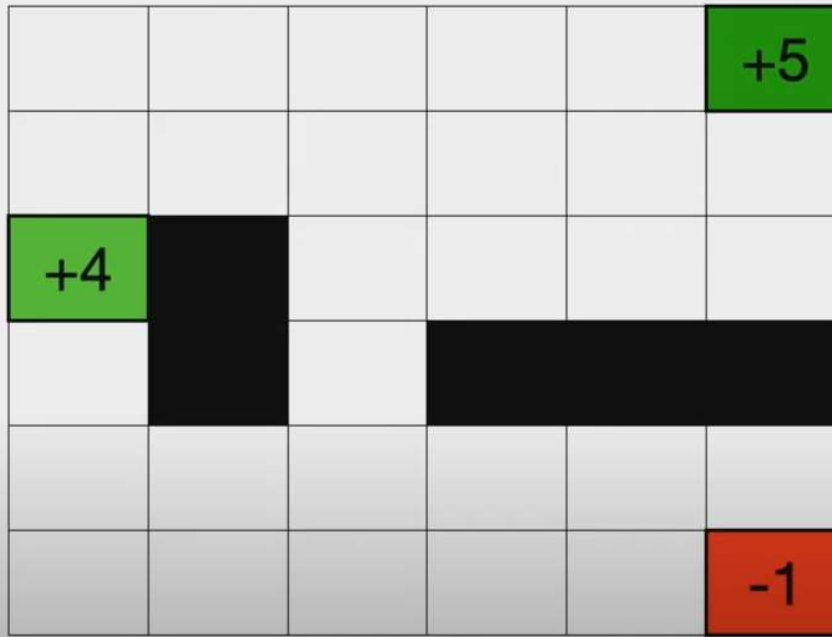
$$\delta = R_t + \gamma V(S_{t+1}) - V(S_t)$$

$$\boxed{\delta^2}$$

$$\boxed{\delta \ln \pi(A_t | S_t)}$$

Critic loss

Actor loss

# Values, states, actions, and policy

# Value and policy



Value

-3 ← 4 → 4

1

2

$a = \{\ \uparrow, \rightarrow, \downarrow, \leftarrow\ \}$  States

Actions  $V(s) = max_a\left(V(s')\right)$

**Bellman equation**
**Value of state as a maximum of the value of its neighbors for the neighbors states obtained by applying any possible action**
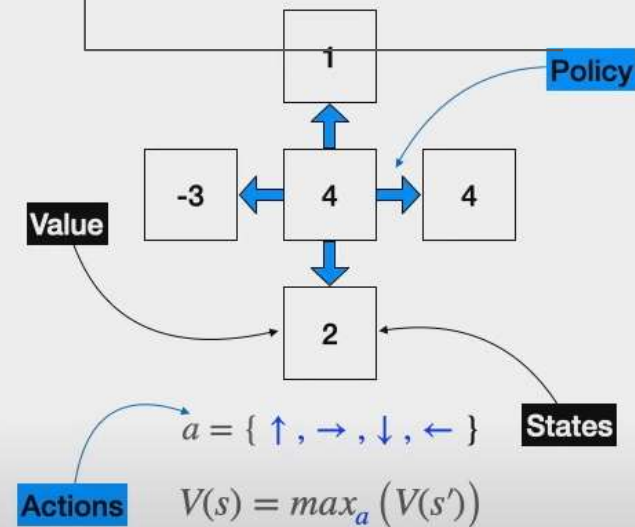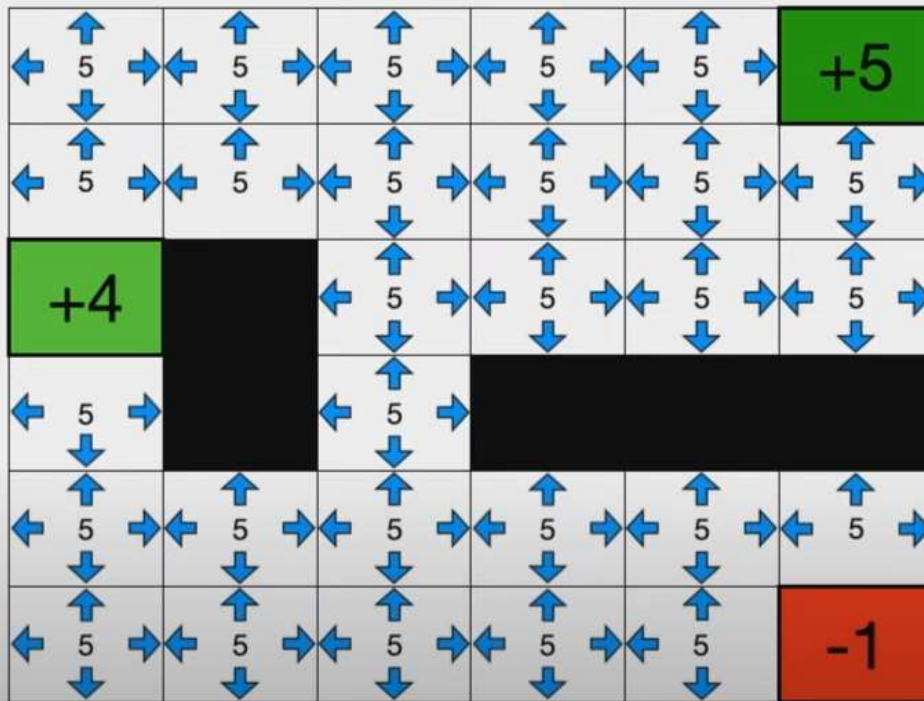
**Best possible decision**
**⇔ POLICY**

# Value and policy



Policy

Value

States

Actions

$$a = \{\, \uparrow\, ,\, \rightarrow\, ,\, \downarrow\, ,\, \leftarrow\, \}$$

$$V(s) = max_a\left(V(s')\right)$$

# Markov decision processes (MDP)

Example of a simple MDP with three states (green circles) and two actions (orange circles), with two rewards (orange arrows).
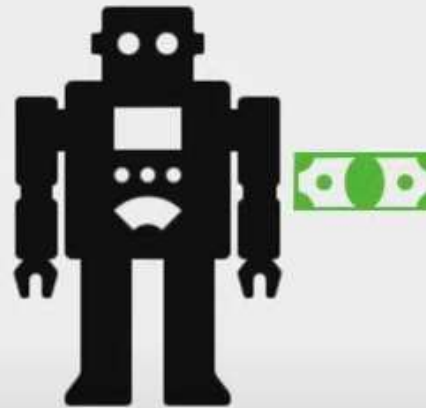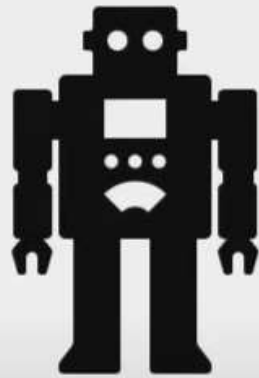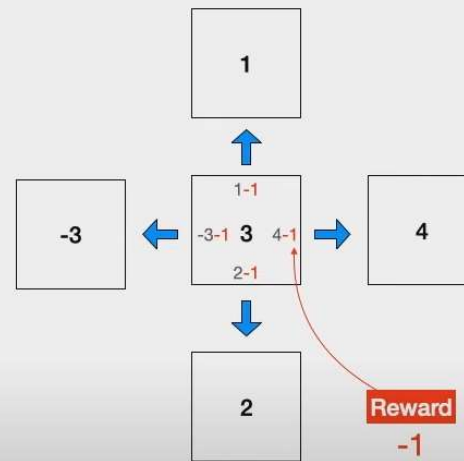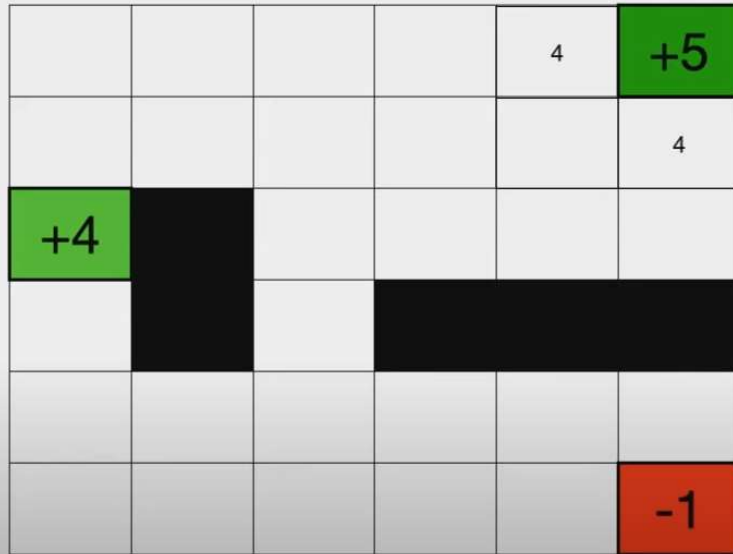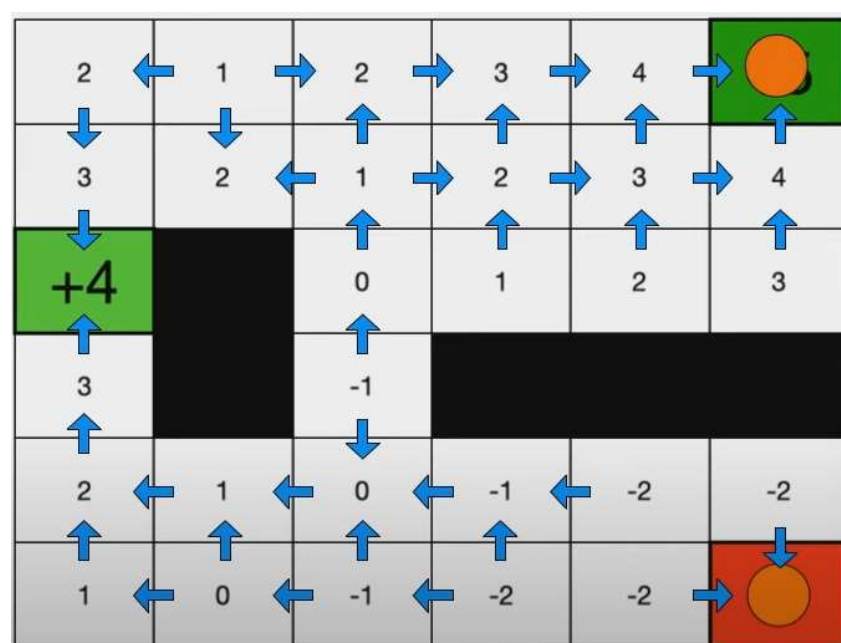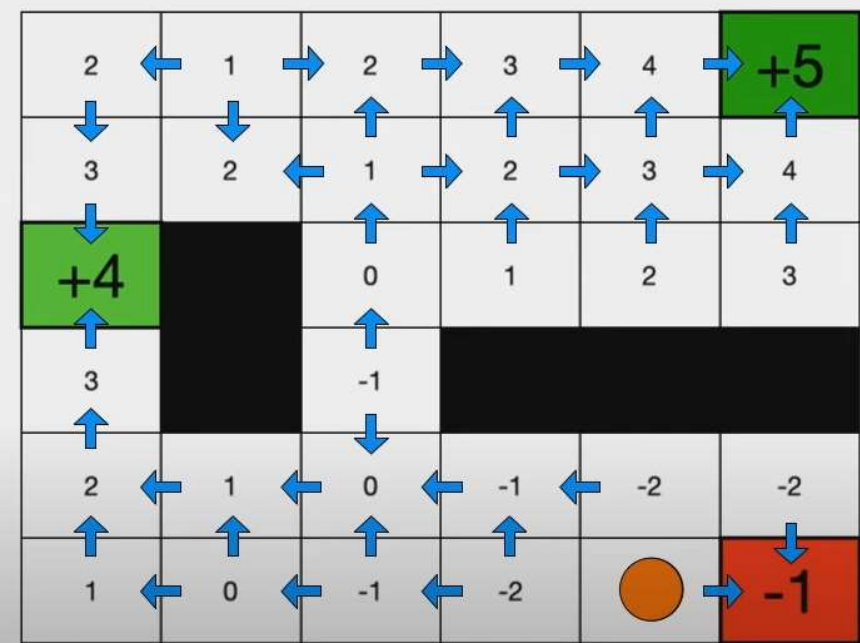
# Reward

# Reward



$$V(s) = max_a \left( R(s,a) + V(s') \right)$$

$$a = \{ \uparrow , \rightarrow , \downarrow , \leftarrow \}$$
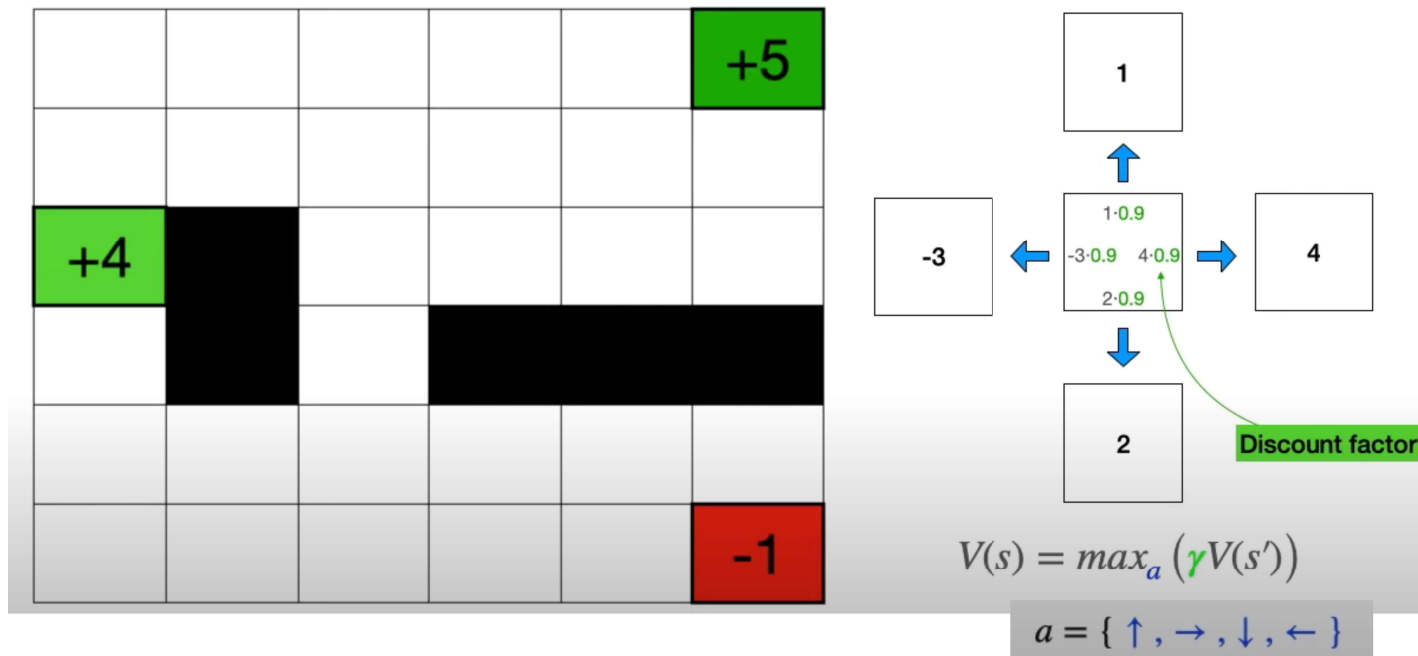
**Bellman equation  1st change**
each value of state is the maximum of points you can obtain if you walk
in the best possible way, & that best possible way is given by policy
(point to the neighbor that has the highest value notice)
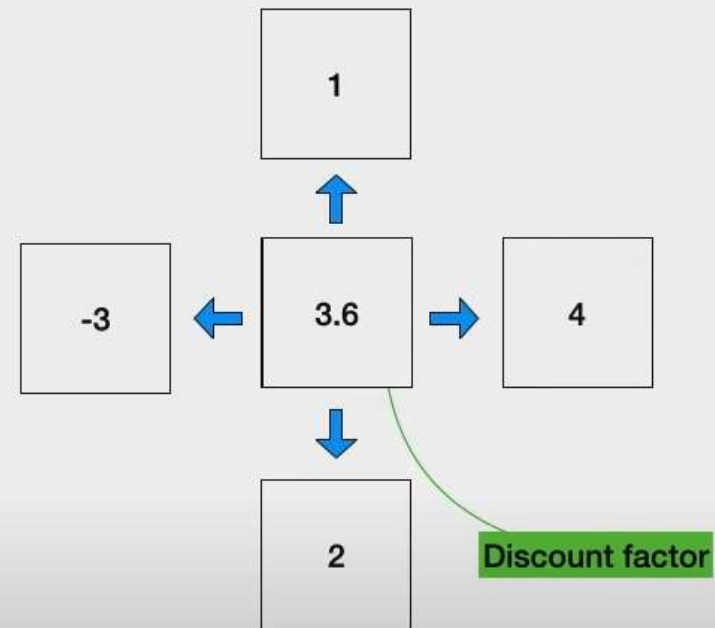
# Reward

# Discount factor

# Discount factor



$$V(s) = max_a \left( \gamma V(s') \right)$$
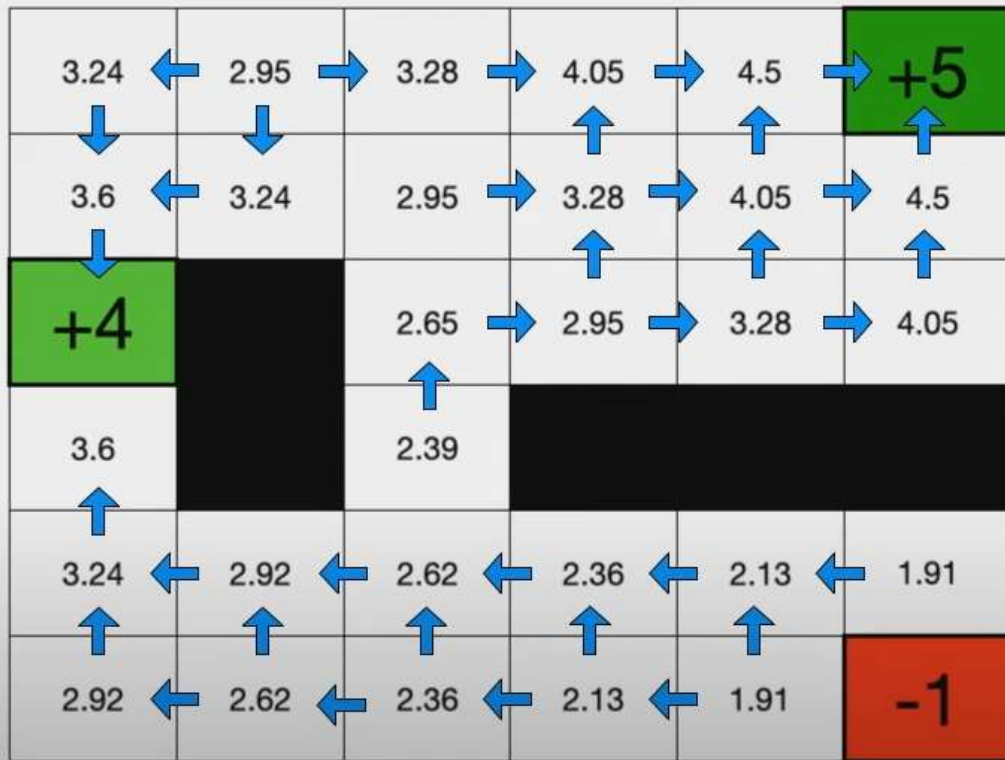
$$a = \{ \uparrow , \rightarrow , \downarrow , \leftarrow \}$$

**Bellman equation  2nd change**
Value of state as a maximum over all its neighbors where all the states i can obtain by applying the actions of a DISCOUNT FACTOR gamma times the value of that new state
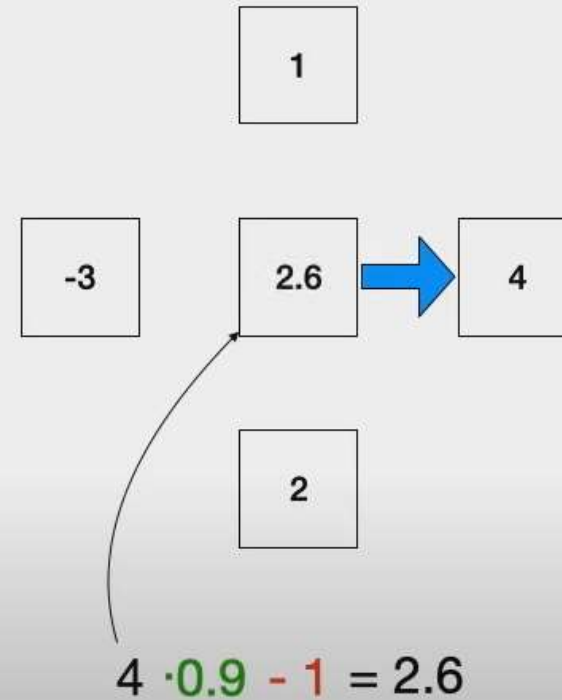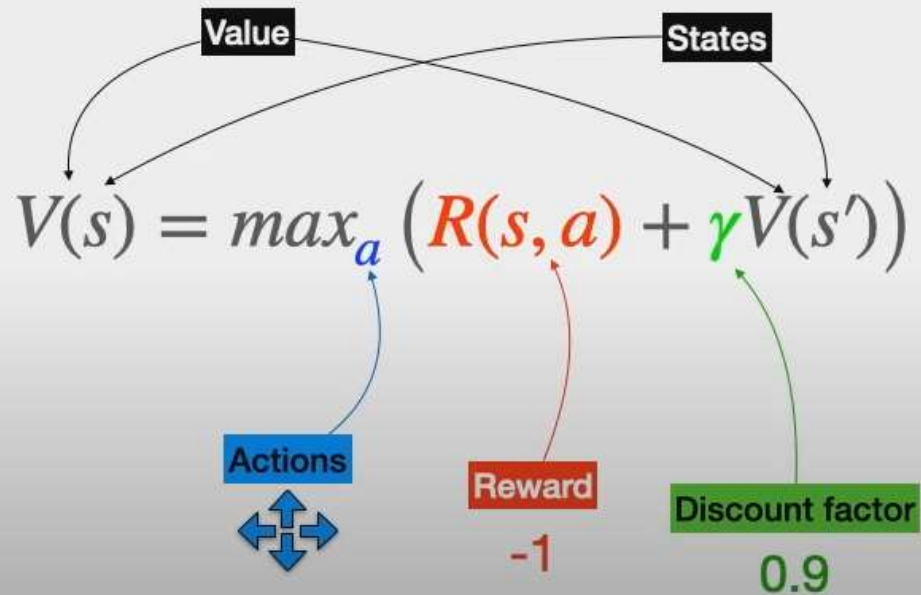
# Discount factor

| | | | | | |
|---|---|---|---|---|---|
| 3.24 ← | 2.95 → | 3.28 → | 4.05 → | 4.5 → | **+5** |
| ↓ | ↓ | | ↑ | ↑ | ↑ |
| 3.6 ← | 3.24 | 2.95 → | 3.28 → | 4.05 → | 4.5 |
| ↓ | | | ↑ | ↑ | ↑ |
| **+4** | ■ | 2.65 → | 2.95 → | 3.28 → | 4.05 |
| ↓ | | ↑ | | | |
| 3.6 | | 2.39 | ■ | ■ | ■ |
| ↑ | | | | | |
| 3.24 ← | 2.92 ← | 2.62 ← | 2.36 ← | 2.13 ← | 1.91 |
| ↑ | ↑ | ↑ | ↑ | ↑ | |
| 2.92 ← | 2.62 ← | 2.36 ← | 2.13 ← | 1.91 ← | **-1** |

| | | |
|---|---|---|
| | 1 | |
| | ↑ | |
| -3 ← | 3.6 → | 4 |
| | ↓ | |
| | 2 | |

Discount factor

$$V(s) = max_a \left( \gamma V(s') \right)$$

$$a = \{ \uparrow, \rightarrow, \downarrow, \leftarrow \}$$

# Bellman equation

# Bellman equation

$$V(s) = max_a \left( R(s,a) + V(s') \right) \qquad V(s) = max_a \left( \gamma V(s') \right)$$

Value

States

$$V(s) = max_a \left( R(s,a) + \gamma V(s') \right)$$

Actions

Reward
-1

Discount factor
0.9

| | | |
|---|---|---|
| | 1 | |
| -3 | 2.6 | 4 |
| | 2 | |

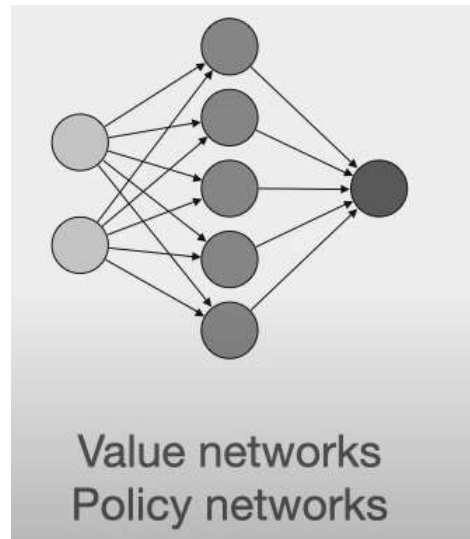$4 \cdot 0.9 - 1 = 2.6$

# Solving the Bellman equation

**Bellman equation allways satisfied**

**Best possible decision ⇔ POLICY**

If we have a much bigger example : the algo says : you have to visit all the boxes, not only once but several times until the values start converging & that's impossible to do when you have a very very large universe with millions of states & actions… : it is very expensive to go over all of them, so what can we do ?
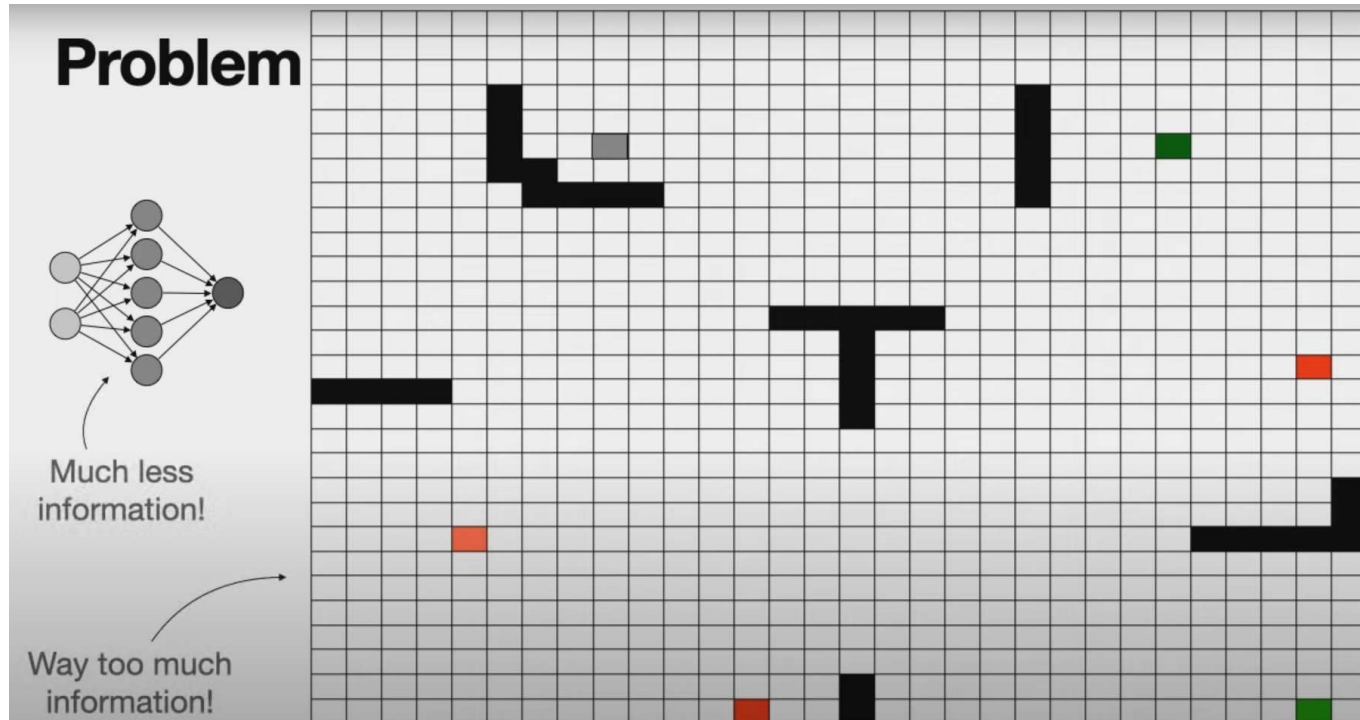
THIS IS WHERE **NEURAL NETWORK** COME TO OUR RESCUE

NN will come in 2 ways : one of them for calculating the value which is called a **value network (Q NETWORK)** & the other one for calculating the **policy** (**Policy NETWORK**)
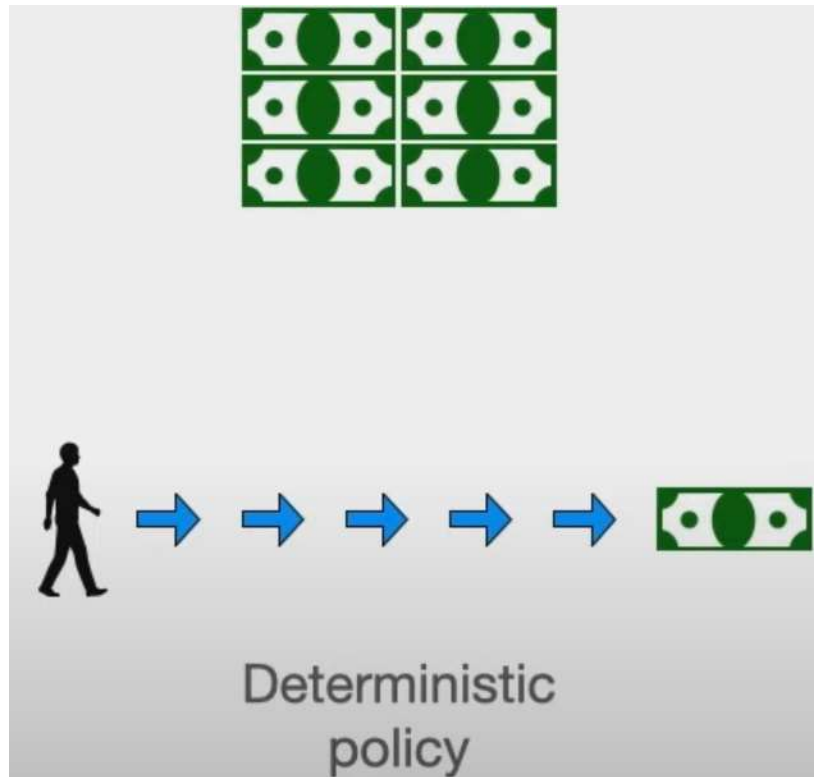


Value networks
Policy networks

In this way we won't have to go over every state several times to learn the value of the policy we simply have to let our agent wander and the NN will be smart enough to pick up information from the places that this agent manages to visit & propagate them to the entire grid.
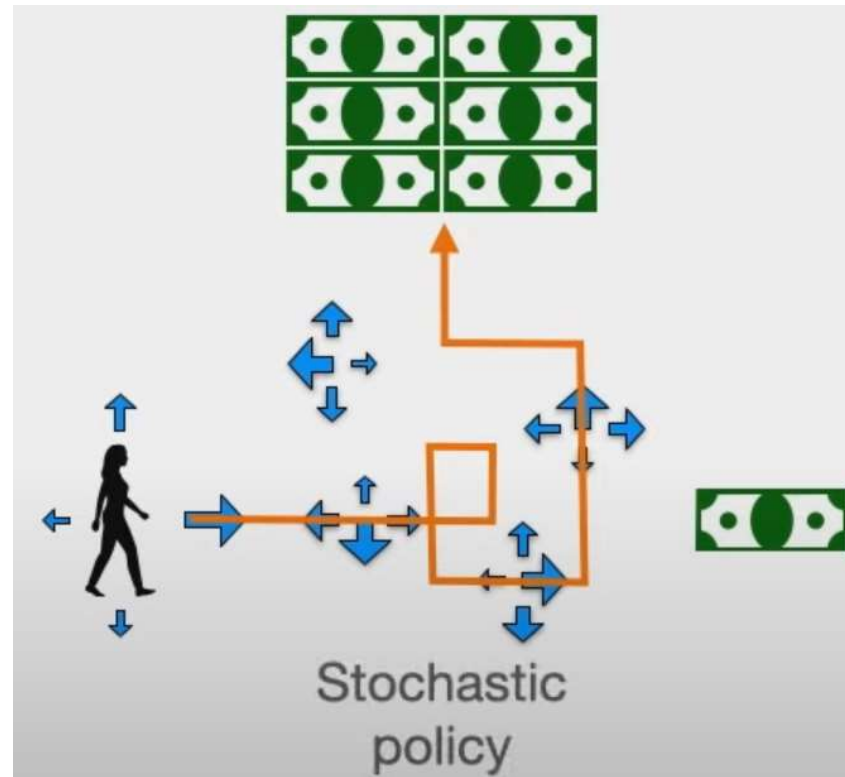
This not only good for saving time but also for saving space
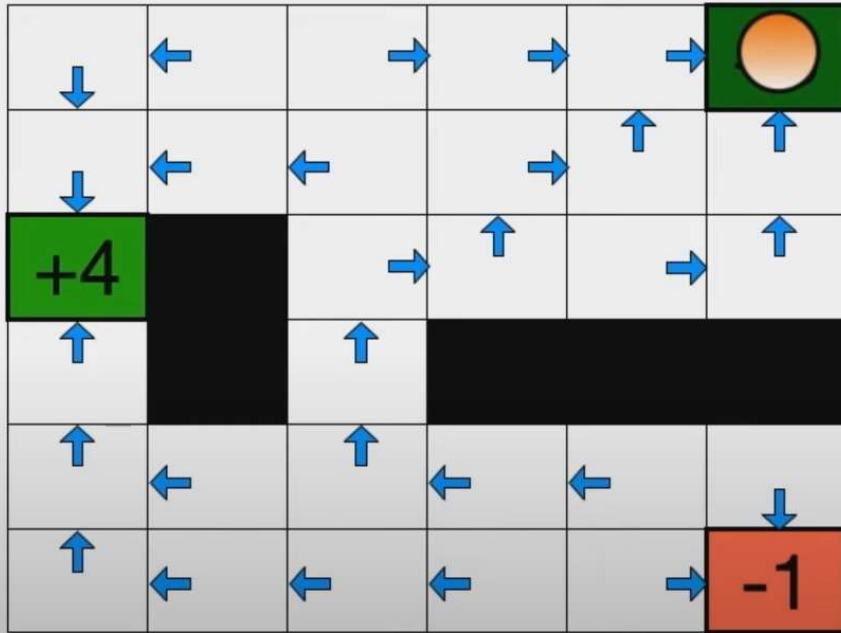
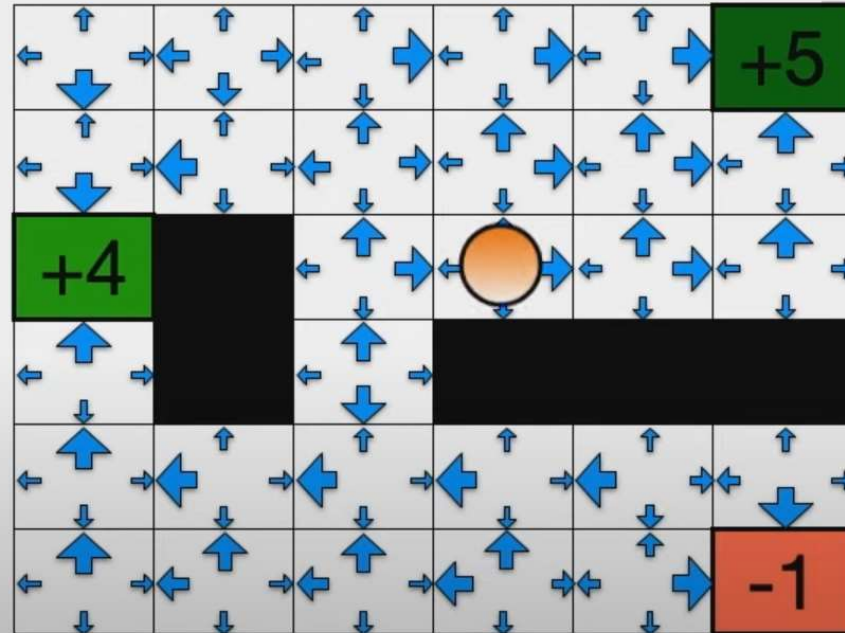# Deterministic vs stochastic policies

**EXPLOITATION**

**EXPLORATION**
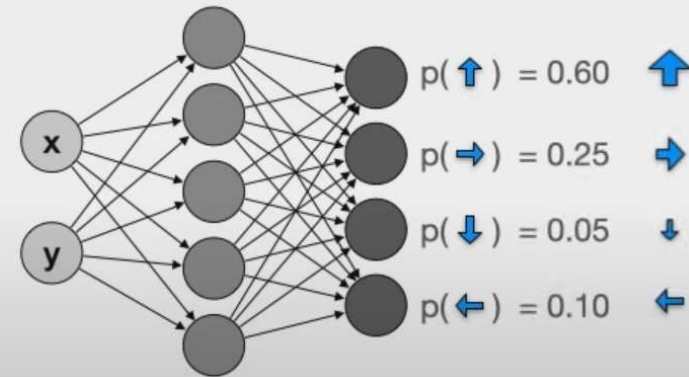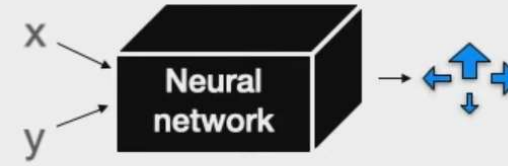
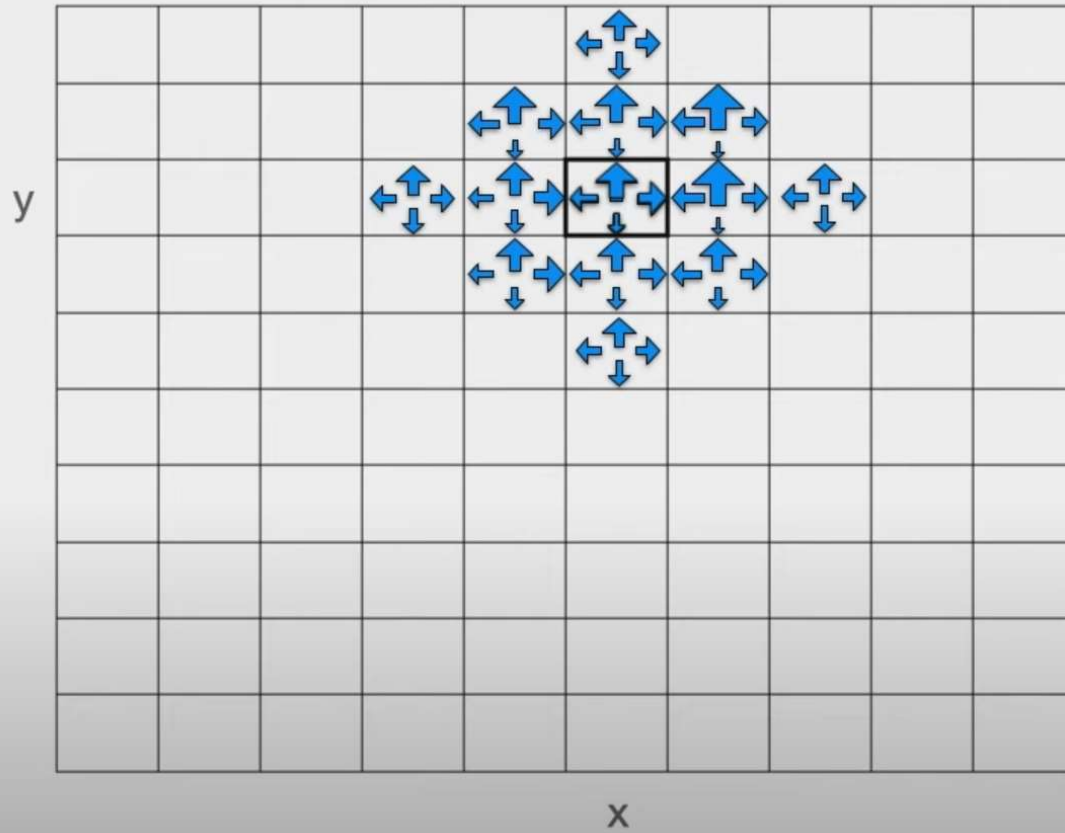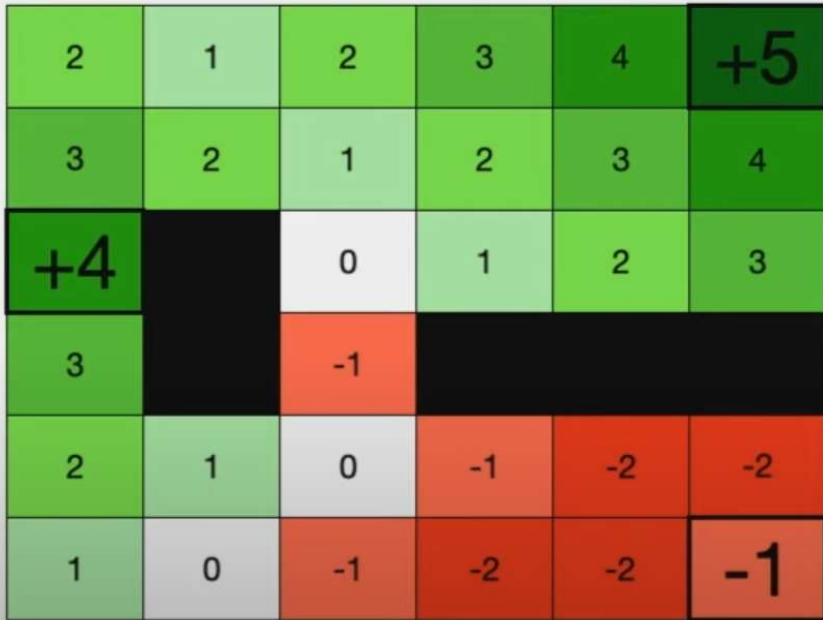# Deterministic vs stochastic



Deterministic

Stochastic

# Neural networks

# Policy network

# Two neural networks
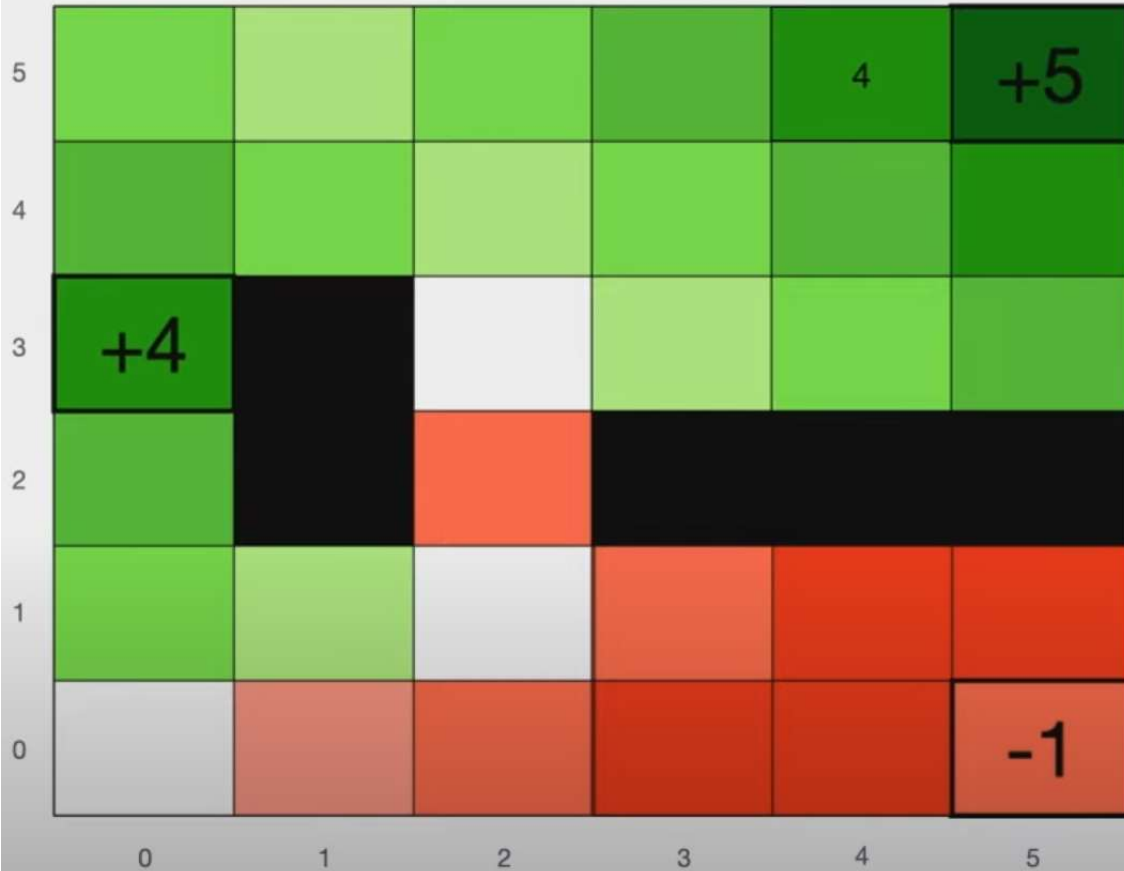
## Value neural network

| | | | | | |
|---|---|---|---|---|---|
| 2 | 1 | 2 | 3 | 4 | **+5** |
| 3 | 2 | 1 | 2 | 3 | 4 |
| **+4** | | 0 | 1 | 2 | 3 |
| 3 | | -1 | | | |
| 2 | 1 | 0 | -1 | -2 | -2 |
| 1 | 0 | -1 | -2 | -2 | **-1** |

## Policy neural network

# Value neural networks



$$V(s) = max_a \left( R(s,a) + \gamma V(s) \right)$$

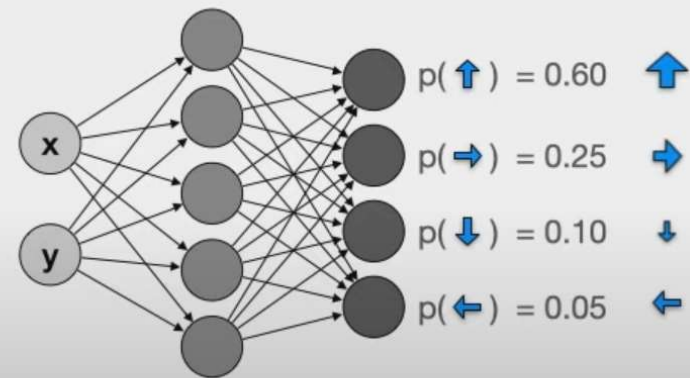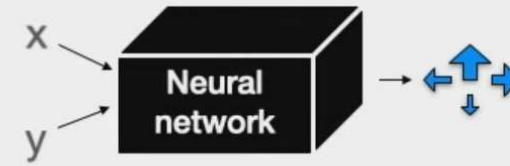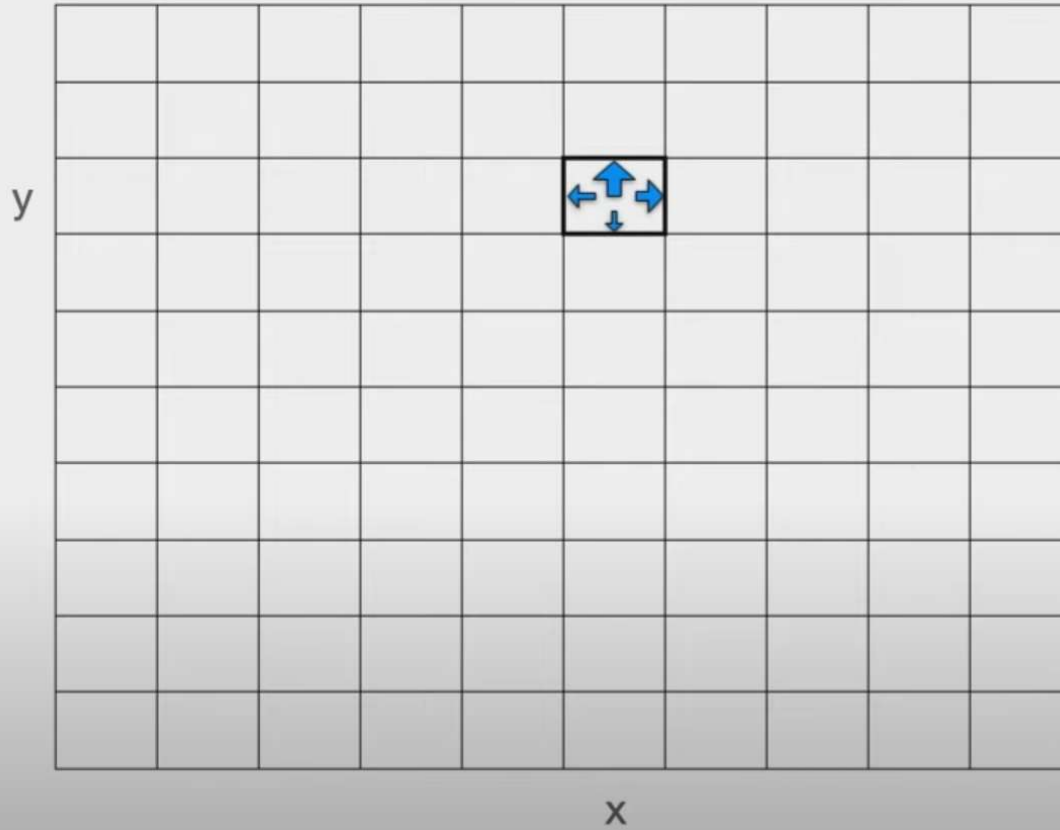| | | | |
|---|---|---|---|
| 3 | | 1.3 | |
| 2 | -2.7 | 0.45 | 3.2 |
| 1 | | 4.9 | |
| | 2 | 3 | 4 |

2 → Neural network →
3 →

Error function

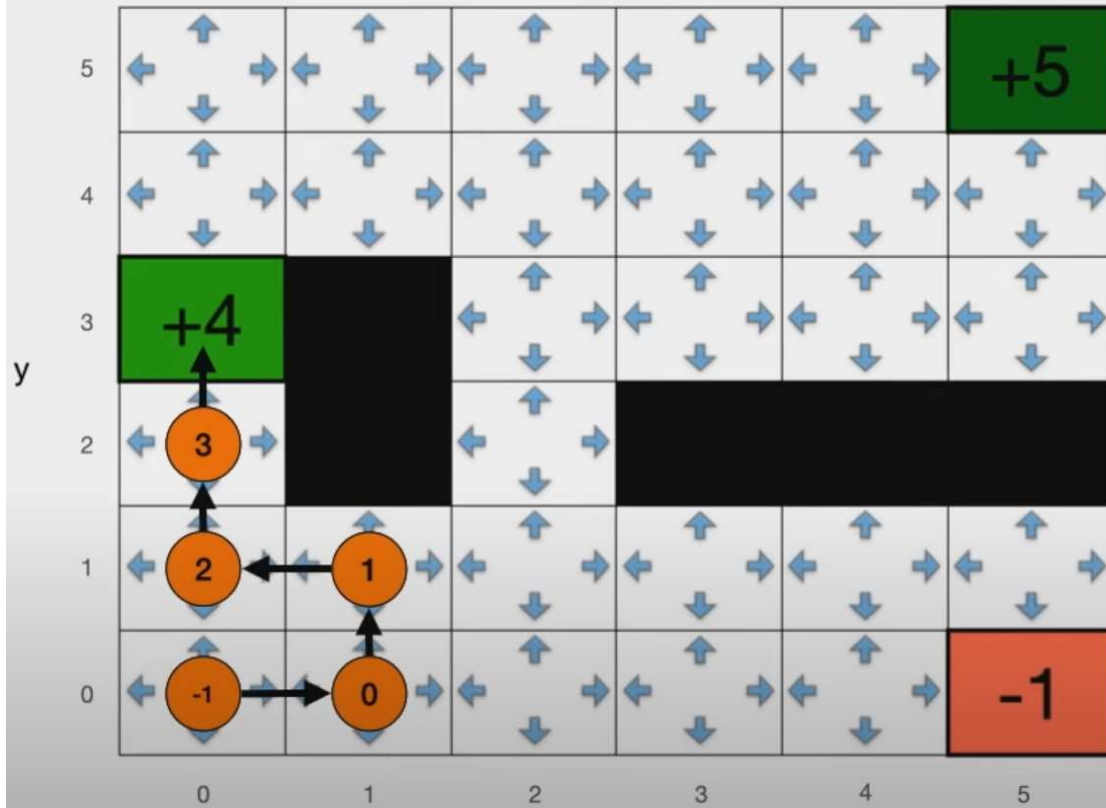$$\left(V(s) - max_a\left(R(s,a) + \gamma V(s)\right)\right)^2$$

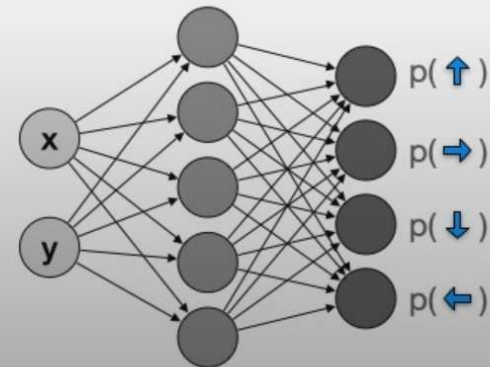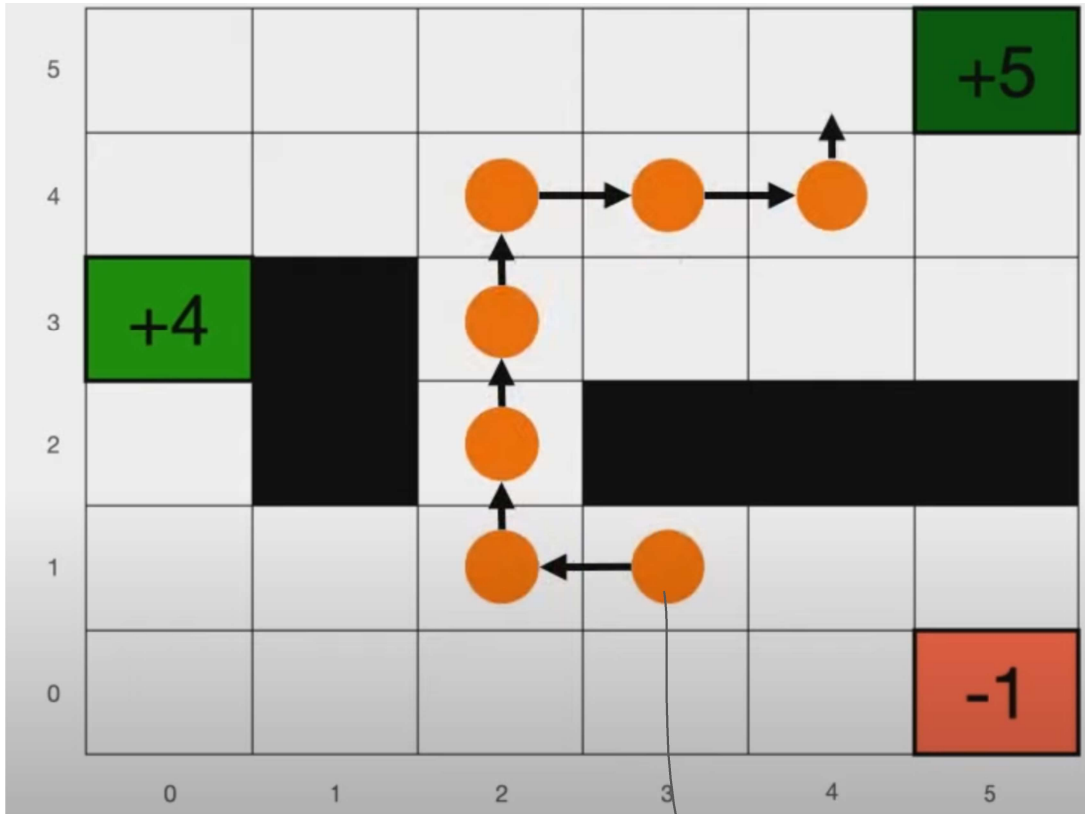$$(0.2 - \max\{2.2, 0.3, -3.7, 3.9\})^2$$

# Policy neural network

# How to train it?



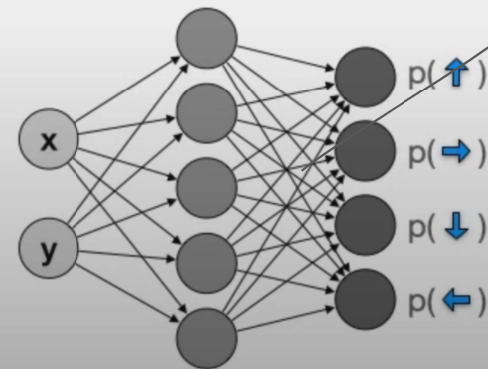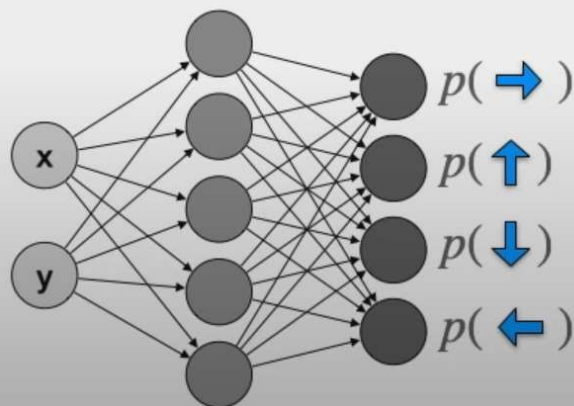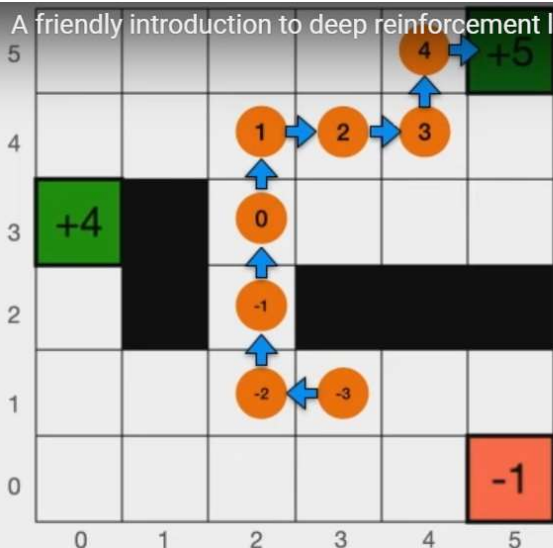| Gain | x | y | Direction | Change |
|------|---|---|-----------|--------|
| 3 | 0 | 2 | ⬆ | increase |
| 2 | 0 | 1 | ⬆ | increase |
| 1 | 1 | 1 | ⬅ | increase |
| 0 | 1 | 0 | ⬆ | stay |
| -1 | 0 | 0 | ➡ | decrease |

| Gain | x | y | Direction | Change |
|------|---|---|-----------|--------|
| 3 | 0 | 2 | ⬆ | increase |
| 2 | 0 | 1 | ⬆ | increase |
| 1 | 1 | 1 | ⬅ | increase |
| 0 | 1 | 0 | ⬆ | stay |
| -1 | 0 | 0 | ➡ | decrease |

we take another walk

That is the table we give to NN

# Training the policy neural network

# Training the policy network



| Gain | x | y | action | Probability | Increase log loss | change |
|---|---|---|---|---|---|---|
| 4 | 4 | 5 | ➡️ | $p(\rightarrow) = 0.3$ | 4 $\ln(p(\rightarrow))$ | increase |
| 3 | 4 | 4 | ⬆️ | $p(\uparrow) = 0.9$ | 3 $\ln(p(\uparrow))$ | increase |
| 2 | 3 | 4 | ➡️ | $p(\rightarrow) = 0.1$ | 2 $\ln(p(\rightarrow))$ | increase |
| 1 | 2 | 4 | ➡️ | $p(\rightarrow) = 0.2$ | 1 $\ln(p(\rightarrow))$ | increase |
| 0 | 2 | 3 | ⬆️ | $p(\uparrow) = 0.5$ | 0 $\ln(p(\uparrow))$ | stay |
| -1 | 2 | 2 | ⬆️ | $p(\uparrow) = 0.4$ | -1 $\ln(p(\uparrow))$ | decrease |
| -2 | 2 | 1 | ⬆️ | $p(\uparrow) = 0.3$ | -2 $\ln(p(\uparrow))$ | decrease |
| -3 | 3 | 1 | ⬅️ | $p(\leftarrow) = 0.7$ | -3 $\ln(p(\leftarrow))$ | decrease |