

OUTLINE

- 1) THEORY
- 2) IMPLEMENT THE GAME (ENVIRONMENT)
- 3) IMPLEMENT THE AGENT
- 4) IMPLEMENT THE MODEL

REINFORCEMENT LEARNING

Is an area of machine learning concerned with how software agents (our computer player) ought to take actions in an environment (our game in this case) in order to maximize the notion of cumulative reward.

In other words, RL is teaching a software agent how to behave in an environment by telling how good it's doing.

DEEP Q LEARNING

This approach extends reinforcement learning by using a deep neural network to predict the actions.

OVERVIEW

AGENT

- Game
- model

Store both
in our agent

Action:

- [1, 0, 0] → straight
- [0, 1, 0] → right turn
- [0, 0, 1] → left turn

Training: (loop)

- state = get_state(game)
- action = get_move(state)
 - model.predict()
- reward, game_over, score = game.play_step(action)
- new_state = get_state(game)
- remember
- model.train()

State: (11 Boolean values)

[danger straight, danger right,
danger left,
direction left, dir. right,
direction up, dir. down,
food left, " right,
" up, " down]

GAME (Pygame)

- play_step(action)
 - reward, game_over, score

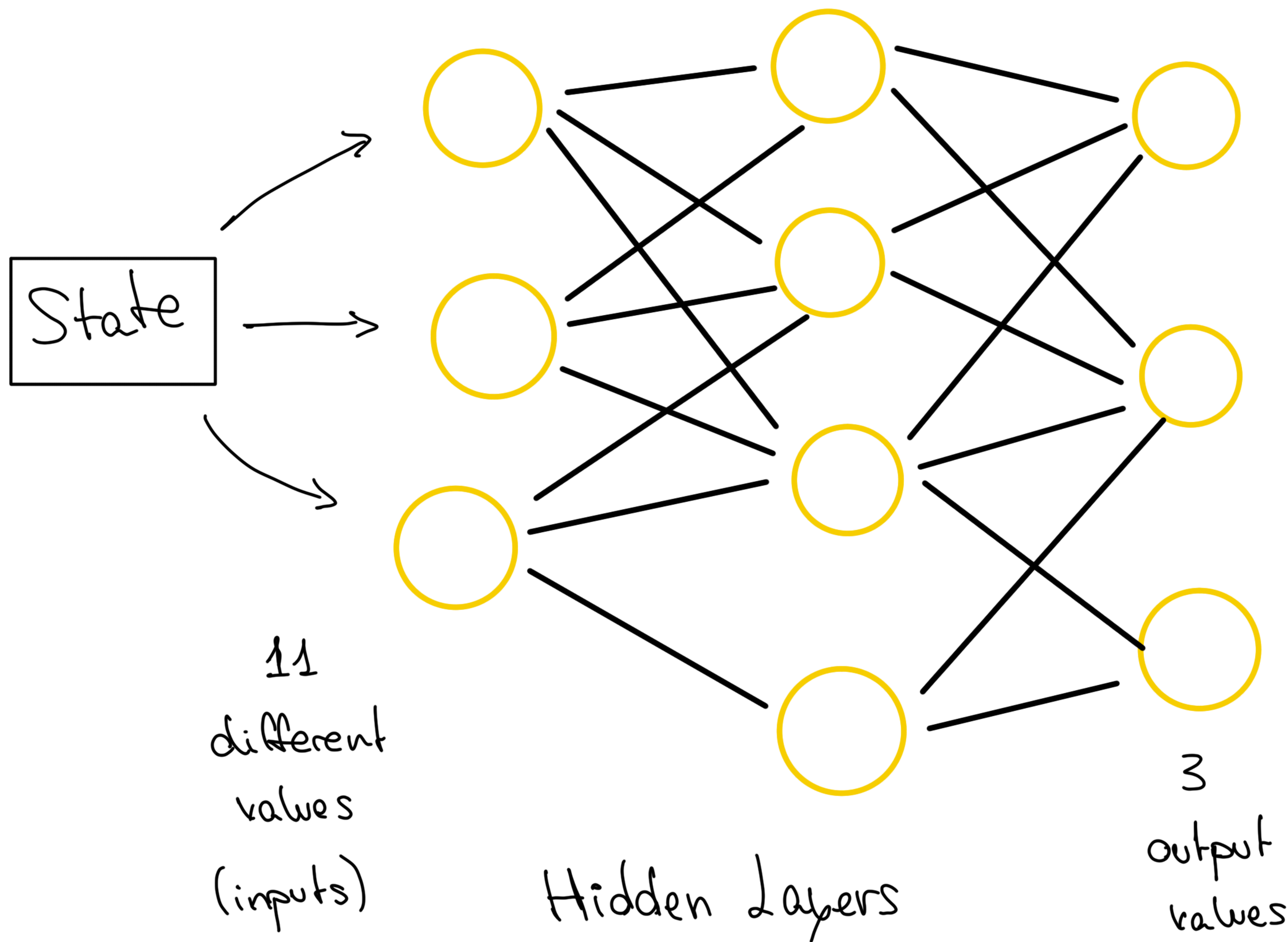
Reward:

- eat food: +10
- game over: -10
- else: 0

MODEL (PyTorch)

- Linear_QNet (DQN)
 - model.predict(state)
 - action

MODEL



⇒ ACTION

e.g. $[5.0, 2.7, 0.1]$

max → $[1, 0, 0]$

(Deep) Q Learning

We want to improve! \rightarrow Q value = Quality of action

0. Init Q Value (= init model)

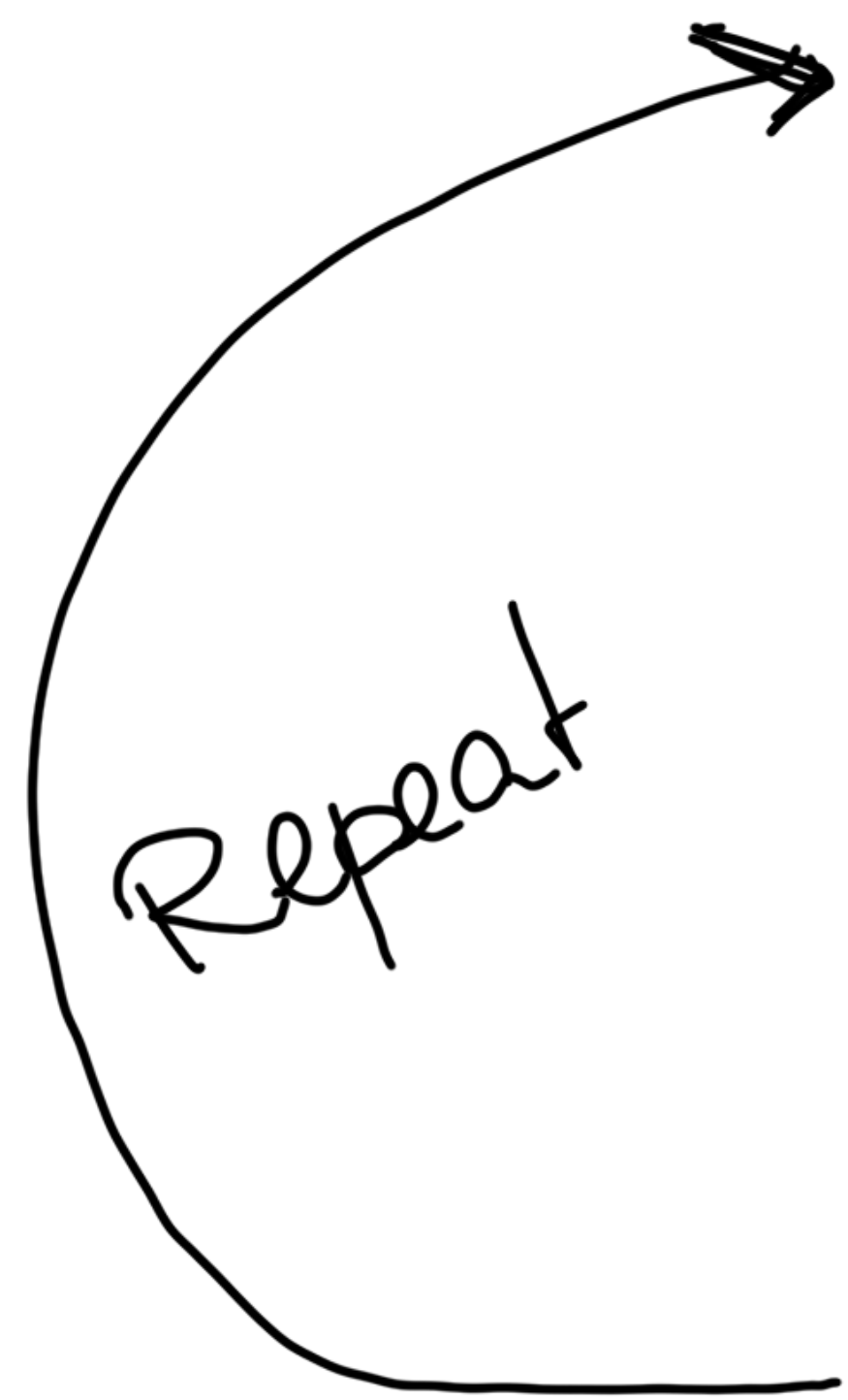
1. Choose action (model.predict(state))

2. Perform action

3. Measure reward

4. Update Q value (+ train model)

Repeat



Bellman Equation (for calculating the new Q value)

$$\underbrace{\text{New } Q(s, a)}_{\text{New Q value for that state and that action}} = \underbrace{Q(s, a)}_{\text{Current Q value}} + \underbrace{\alpha}_{\text{Learning rate}} \left[\underbrace{R(s, a)}_{\text{Reward for taking that action at that state}} + \underbrace{\gamma \max_{a'} Q'(s', a')}_{\substack{\text{Max expected future} \\ \text{reward given the new } s' \\ \text{and all possible actions} \\ \text{at that new state.}}} - Q(s, a) \right]$$

state action

Simplified:

(Old) $Q = \text{model.predict}(\text{state}_0)$

$Q_{\text{new}} = R + \gamma \cdot \max(Q(\text{state}_1))$
 ↑
 Reward

Loss Function:

$$\text{Loss} = (Q_{\text{new}} - Q)^2$$

↳ mean squared error