

Approximate Q-Learning

Eirik Fagerhaug

Norwegian University of Science and Technology

Slides Adapted from/based on:

- Slides by Chris G. Willocks
- Grokking DRL (Miguel Morales)
- AI: A Modern Approach (Russel and Norvig)

Reading Material

Russel and Norvig;

- Chapter 23.4.(1,2,3)
- Chapter 22.1.1

Wiki (Assignments)



- Recap

- Generalization in Reinforcement Learning

- Artificial Neural Networks

- Bellman Update for ANNs

- Deep Q-Network

Q-Function

(Action-utility function)

The Q-function as a Bellman Equation:

$$Q(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q(s', a')]$$

Optimal policy w/r the Q-function:

$$\pi^*(s) = \arg \max_a Q(s, a)$$

Frozen Lake Environment



Start			
0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	Goal +1 15

Expected Utility

$$\max_a Q(s, a)$$

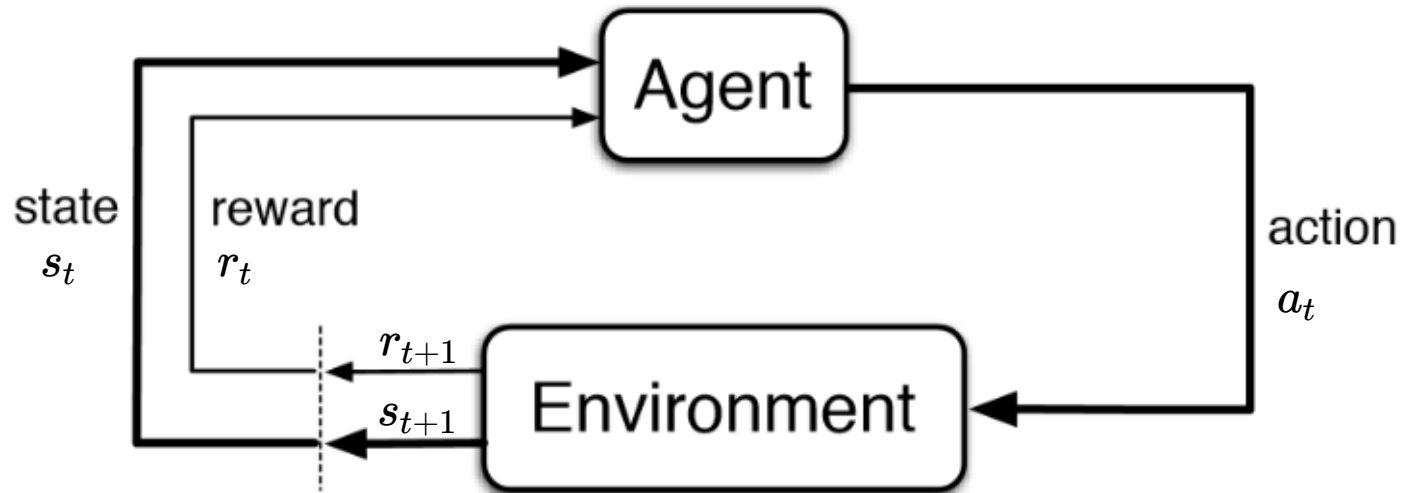
Start 0.54 0	0.50 1	0.47 2	0.46 3
0.56 4	5	0.36 6	7
0.59 8	0.64 9	0.62 10	11
12	0.74 13	0.86 14	Goal 15 +1

Q-Values

Start

0.52 0.54 0.53 0.53	0.50 0.34 0.32 0.33	0.47 0.44 0.42 0.43	0.46 0.31 0.30 0.31
0.36 0.56 0.37 0.38		0.16 0.36 0.36 0.20	
0.59 0.38 0.40 0.41	0.40 0.44 0.45 0.64	0.33 0.62 0.40 0.50	
	0.50 0.46 0.74 0.53	0.78 0.73 0.82 0.86	Goal

Agent-Environment interaction



- At each step t the agent:
 - Receives state s_t (and reward r_t)
 - Executes action a_t
- The environment:
 - Receives action a_t
 - Emits state s_{t+1} (and reward r_{t+1})

Bellman Update

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

Q-tables for large state- and action-spaces

For a game with 200x200 pixels, where we use the pixels as the state, we have:

$$|S| = ((255)^3)^{200 \times 200} \approx 6.6 \cdot 10^{288784}$$

(There are between 10^{78} and 10^{82} atoms in the observable universe)

- Recap
- • Generalization in Reinforcement Learning
 - Artificial Neural Networks
 - Bellman Update for ANNs
 - Deep Q-Network

Generalization in Reinforcement Learning

Problem with large state spaces:

- Too large to save in a Q-table
- Cannot possibly visit all possible states

Solution:

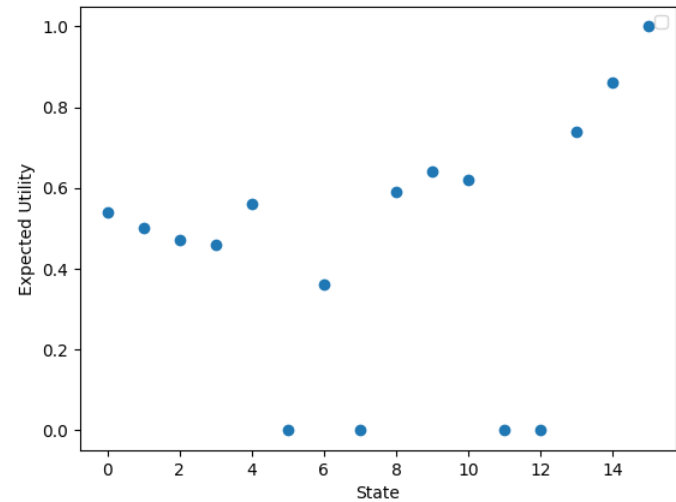
Function approximation; a process of constructing a compact approximation of the true utility- or Q- function.

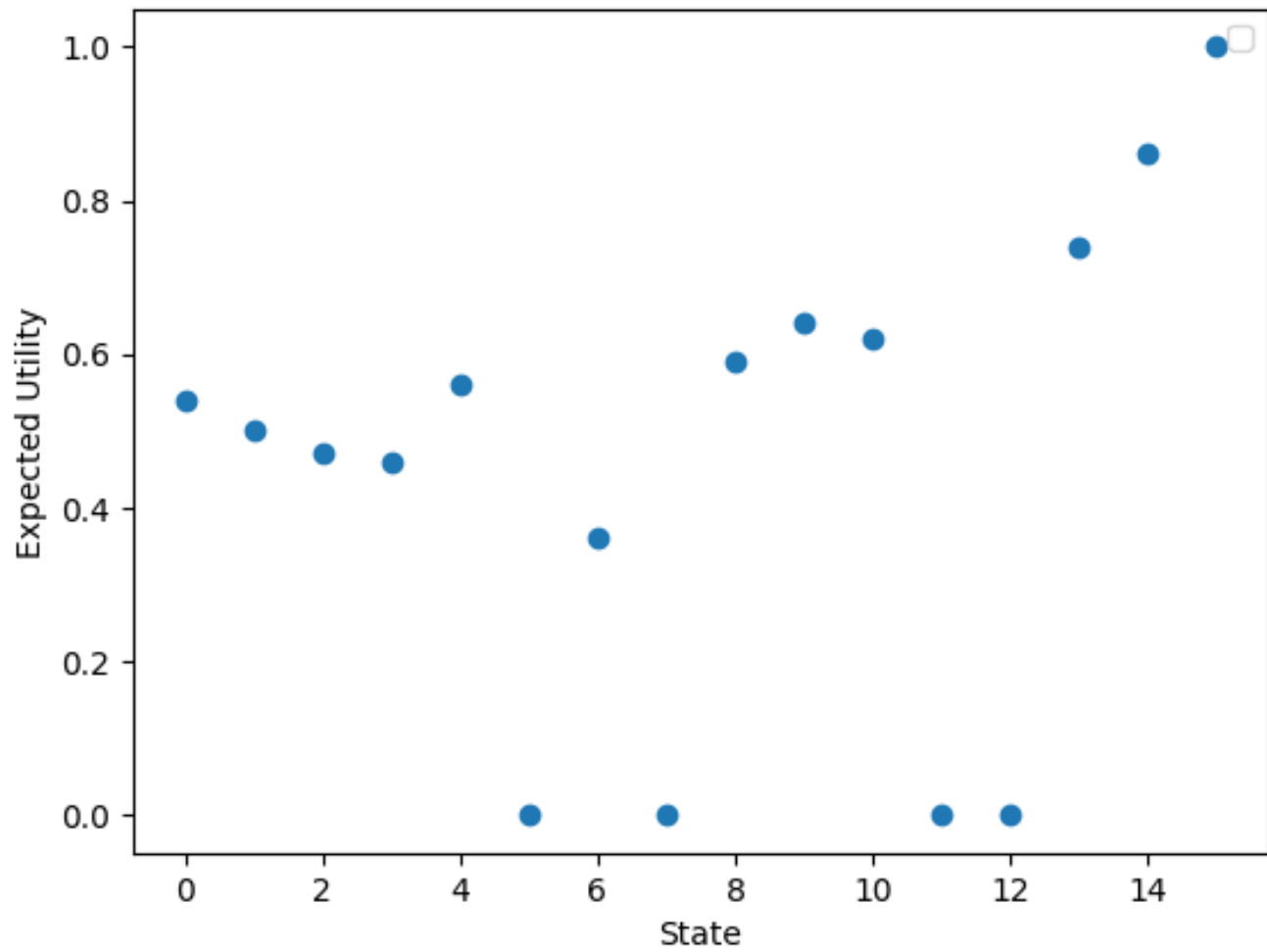
Example; A weighted linear combination of features

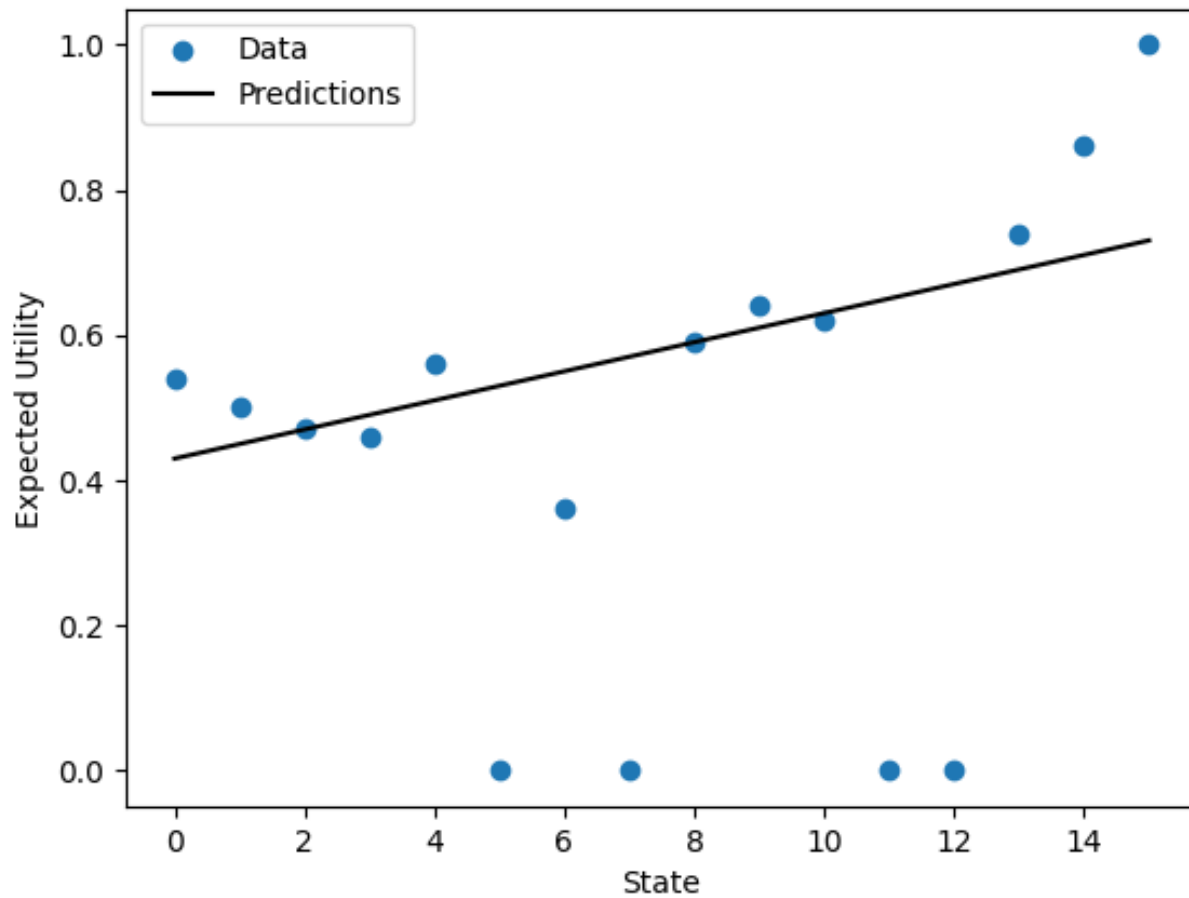
$$\hat{U}_w(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

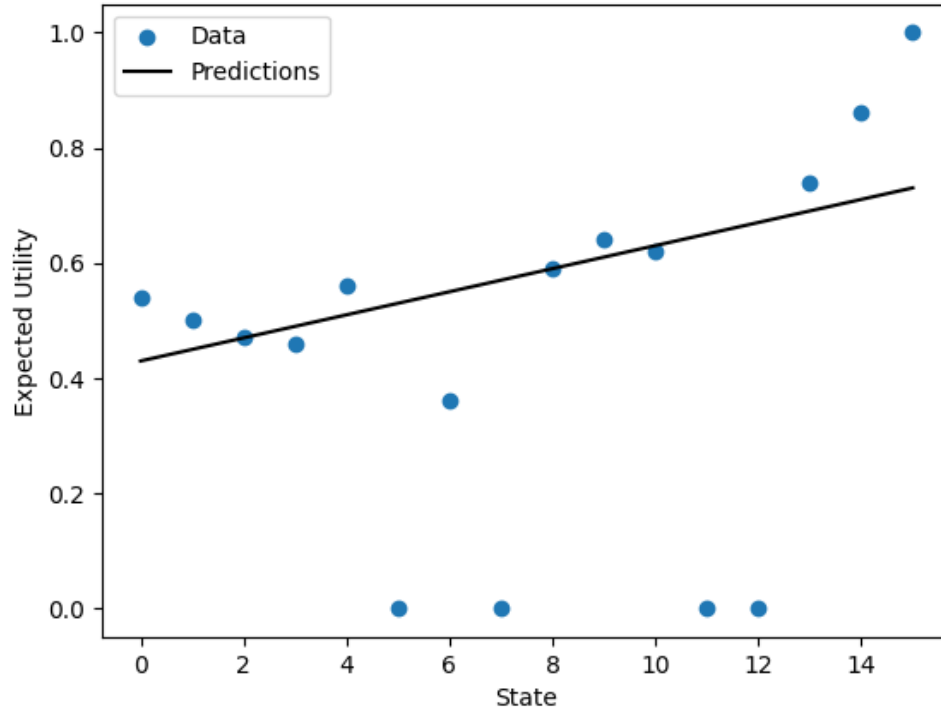
$$\max_a Q(s, a)$$

Start 0.54 0	0.50 1	0.47 2	0.46 3
0.56 4	5	0.36 6	7
0.59 8	0.64 9	0.62 10	11
12	0.74 13	0.86 14	Goal 15 +1









$$\hat{U}_w(x) = w_0 + w_1 x$$

$$\hat{U}_w(x) = 0.43 + 0.02x$$

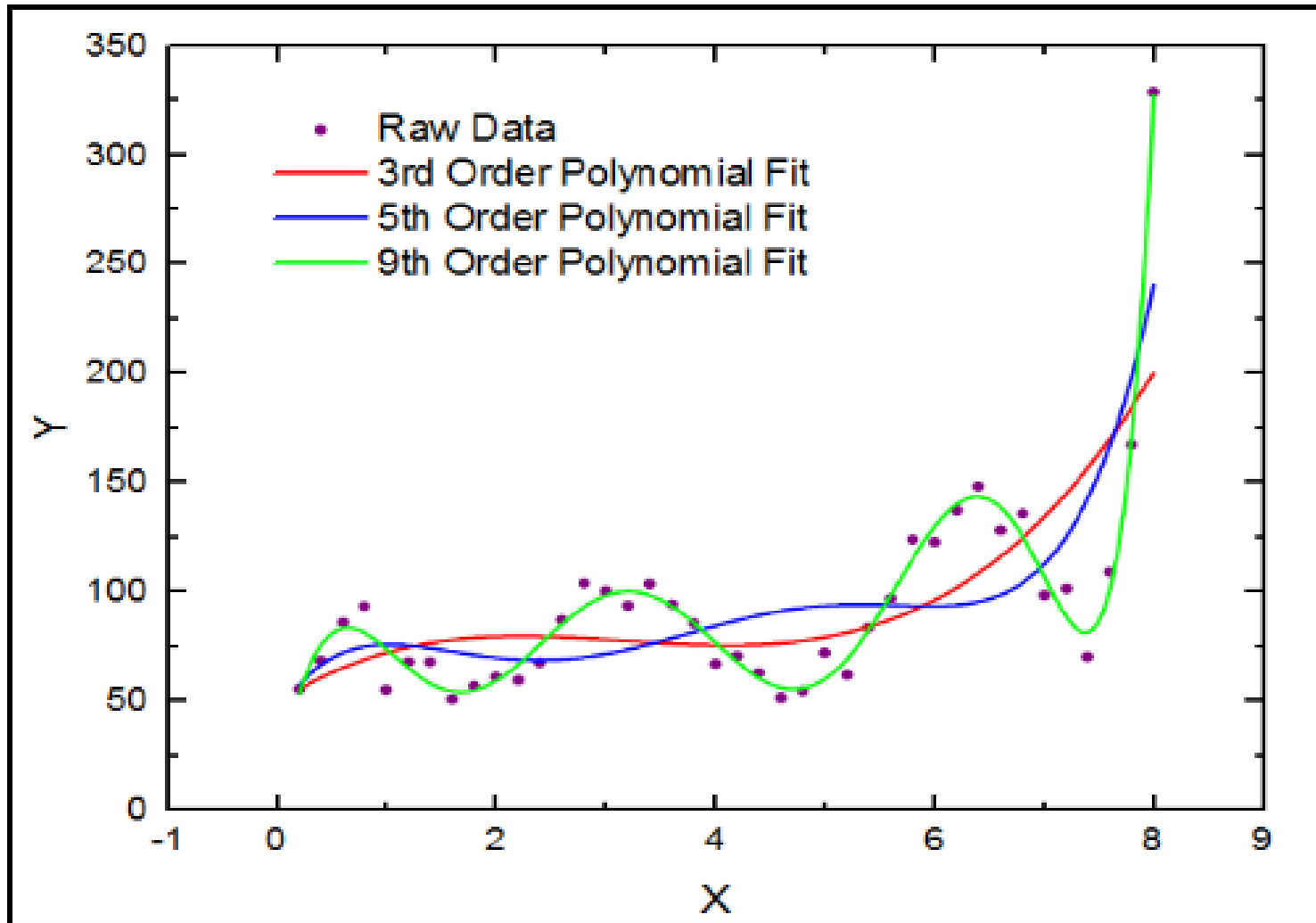
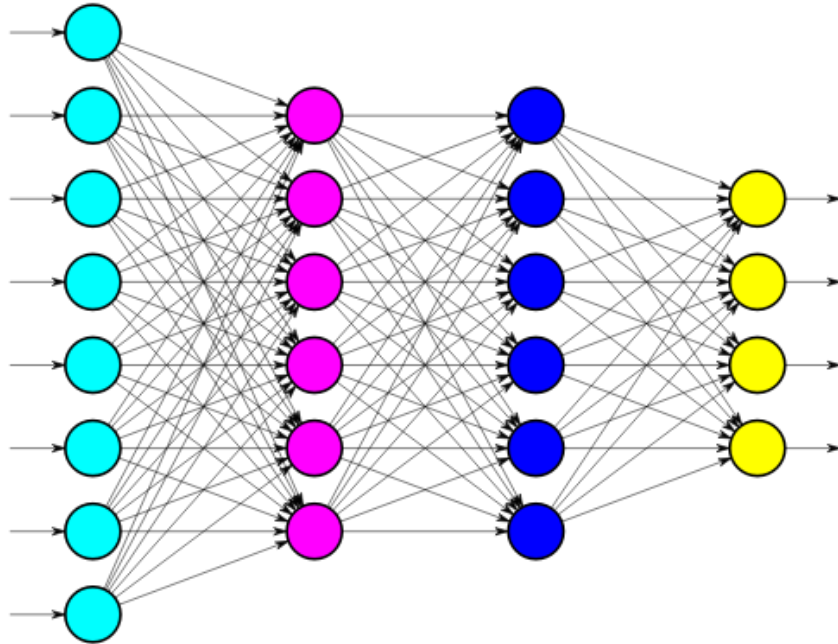
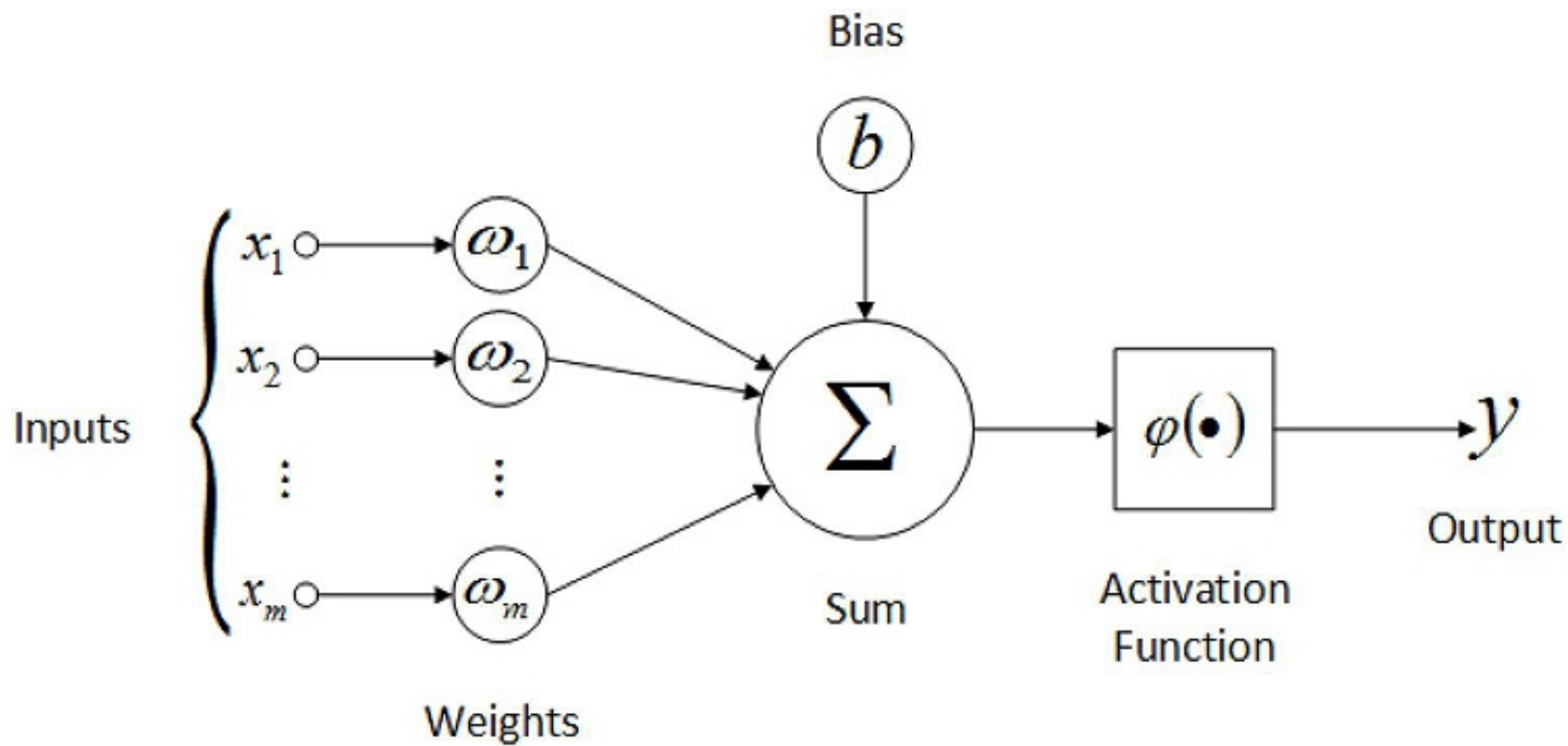


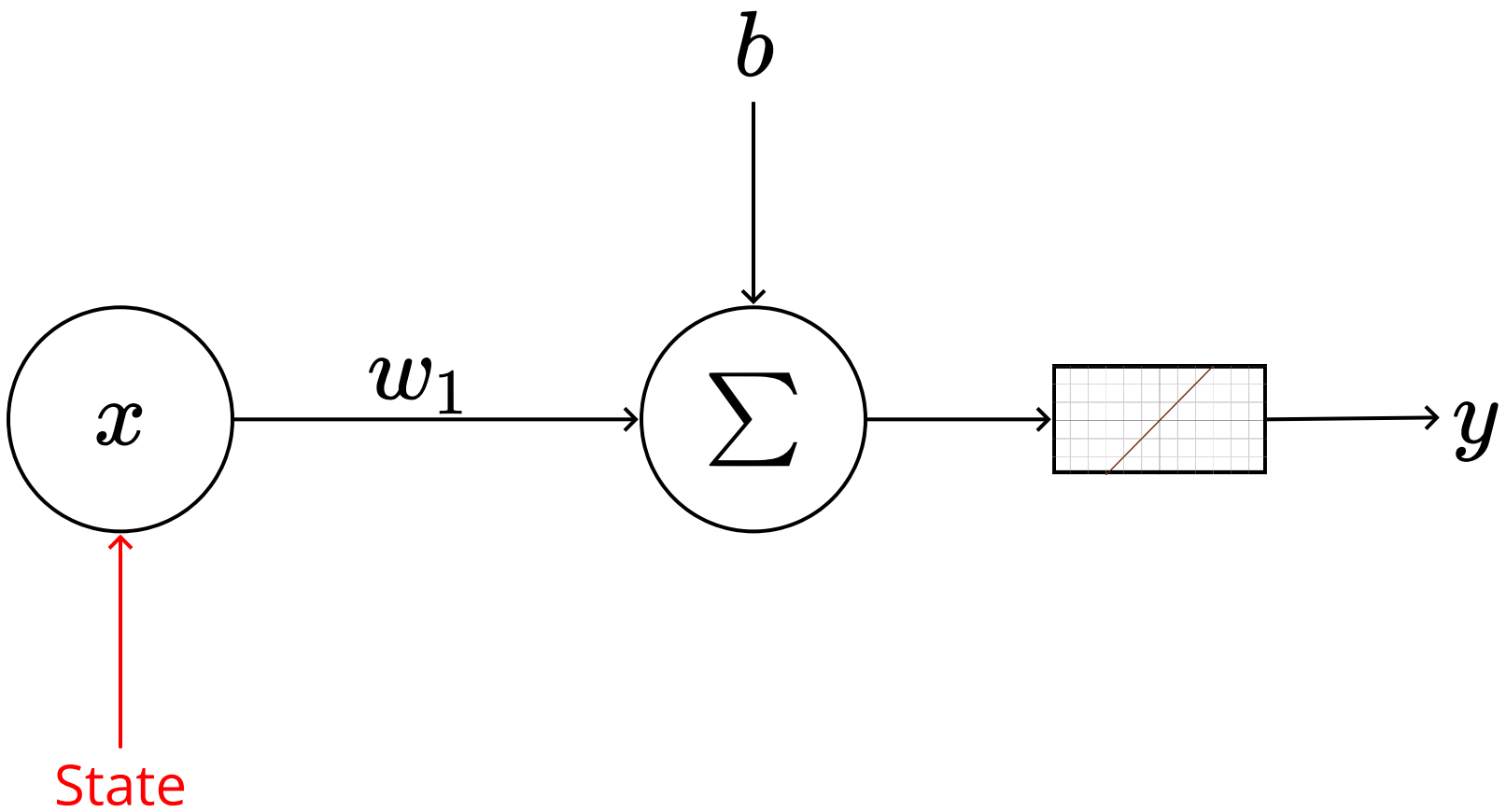
image from [WikiMedia](#) by [Shobha](#)

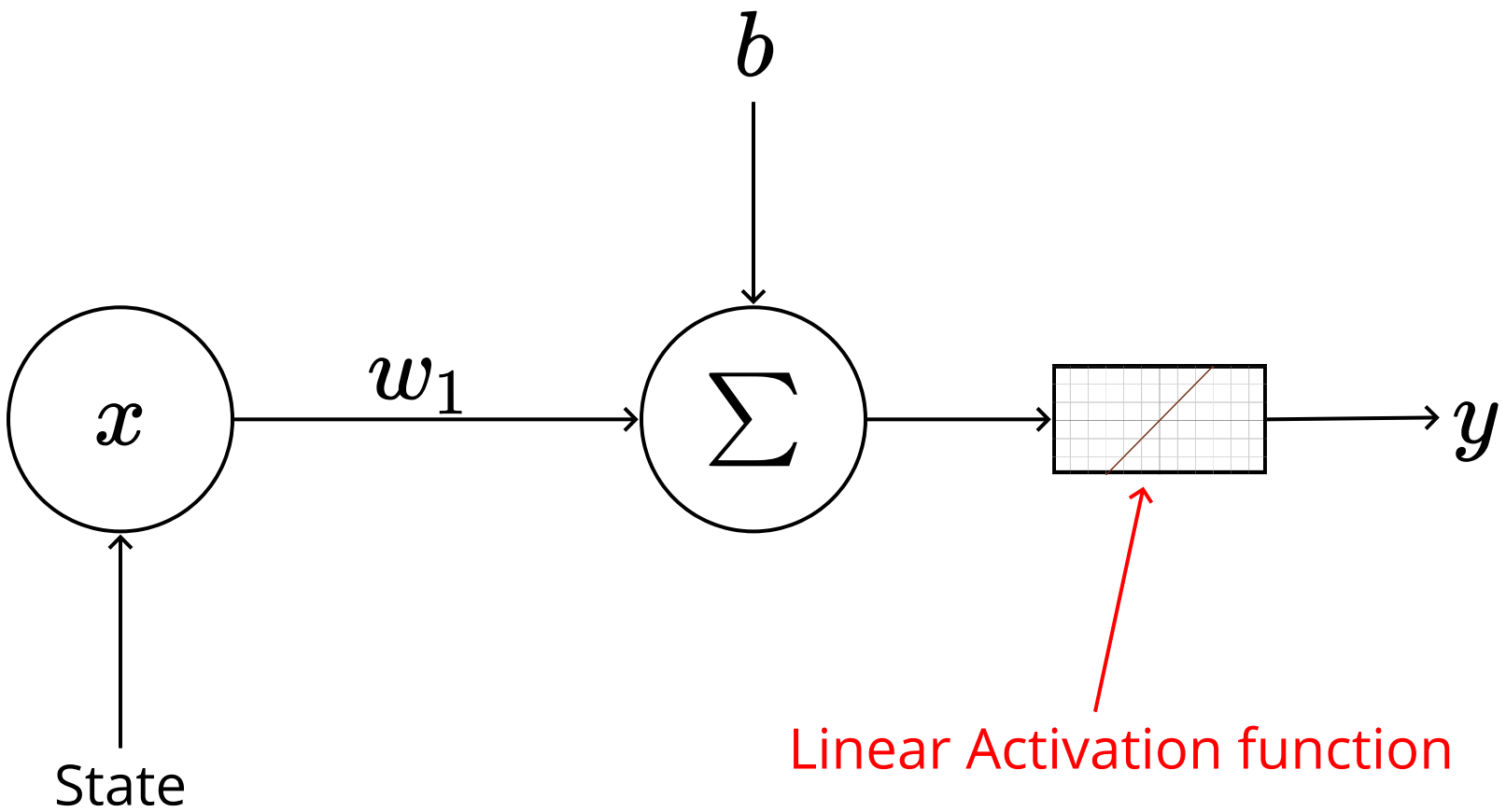
- Recap
- Generalization in Reinforcement Learning
- • Artificial Neural Networks
 - Bellman Update for ANNs
 - Deep Q-Network

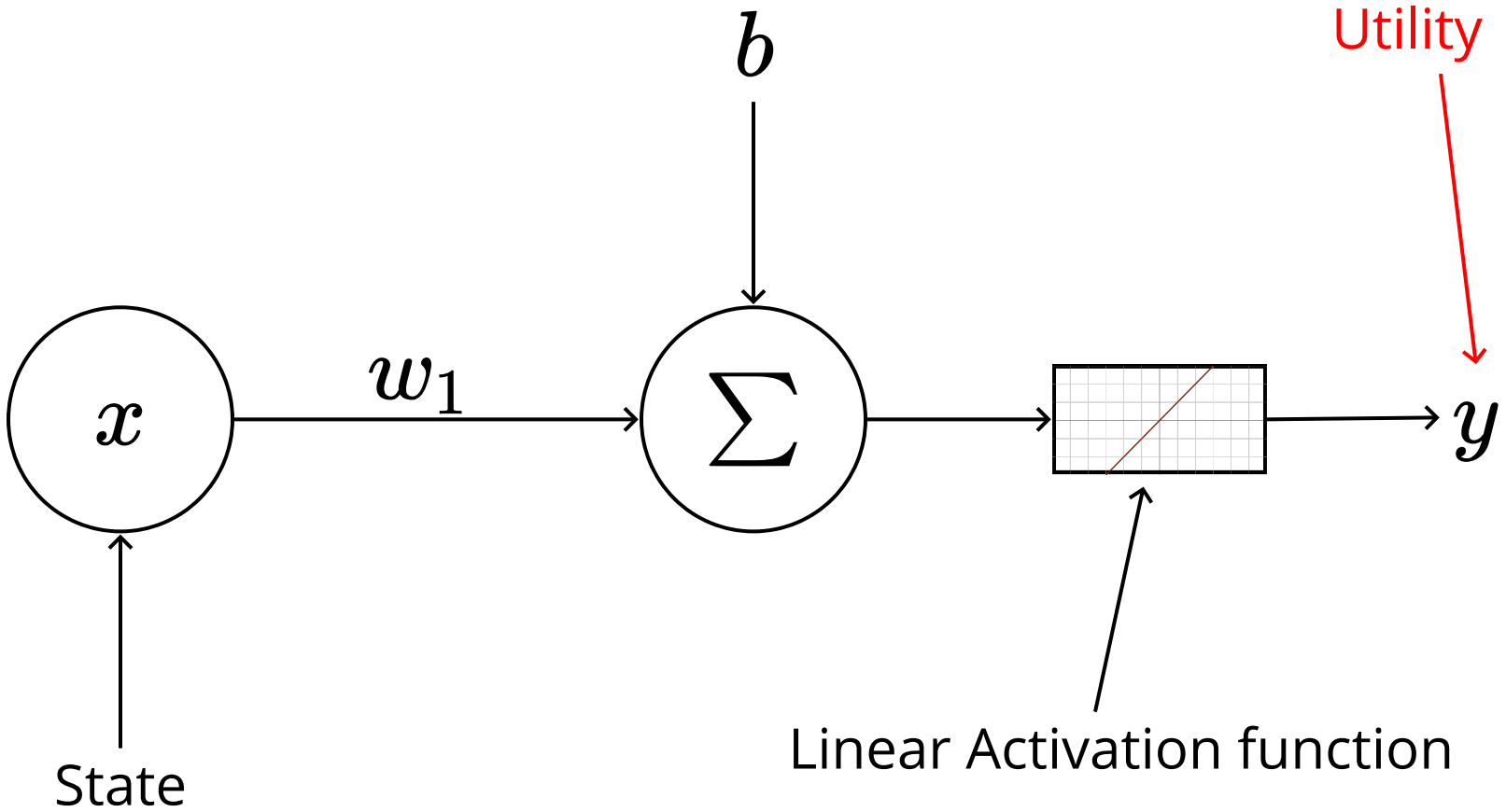
Artificial Neural Network

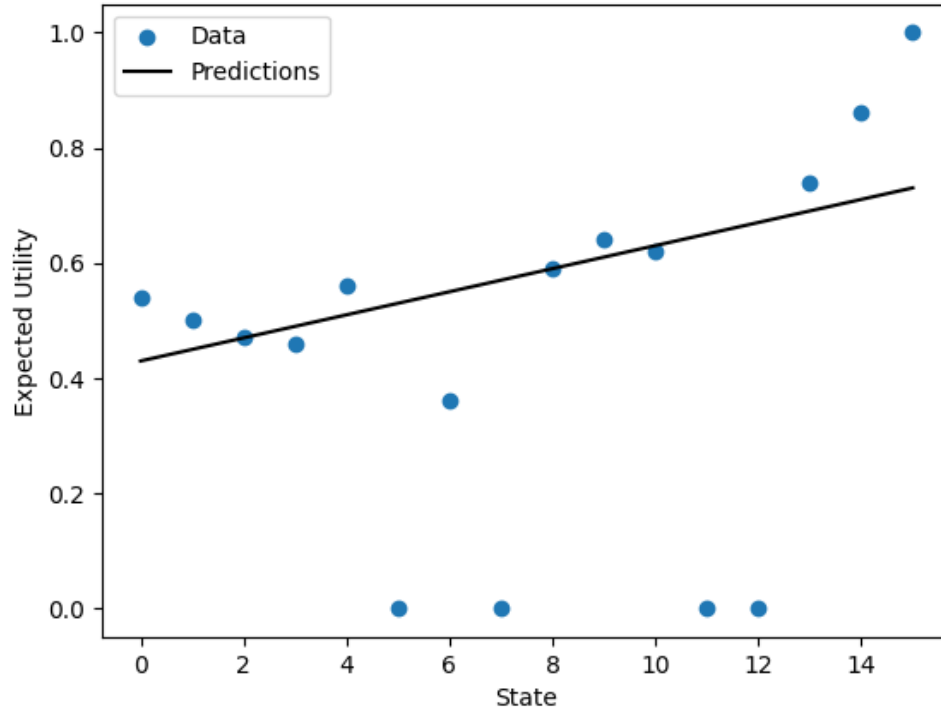






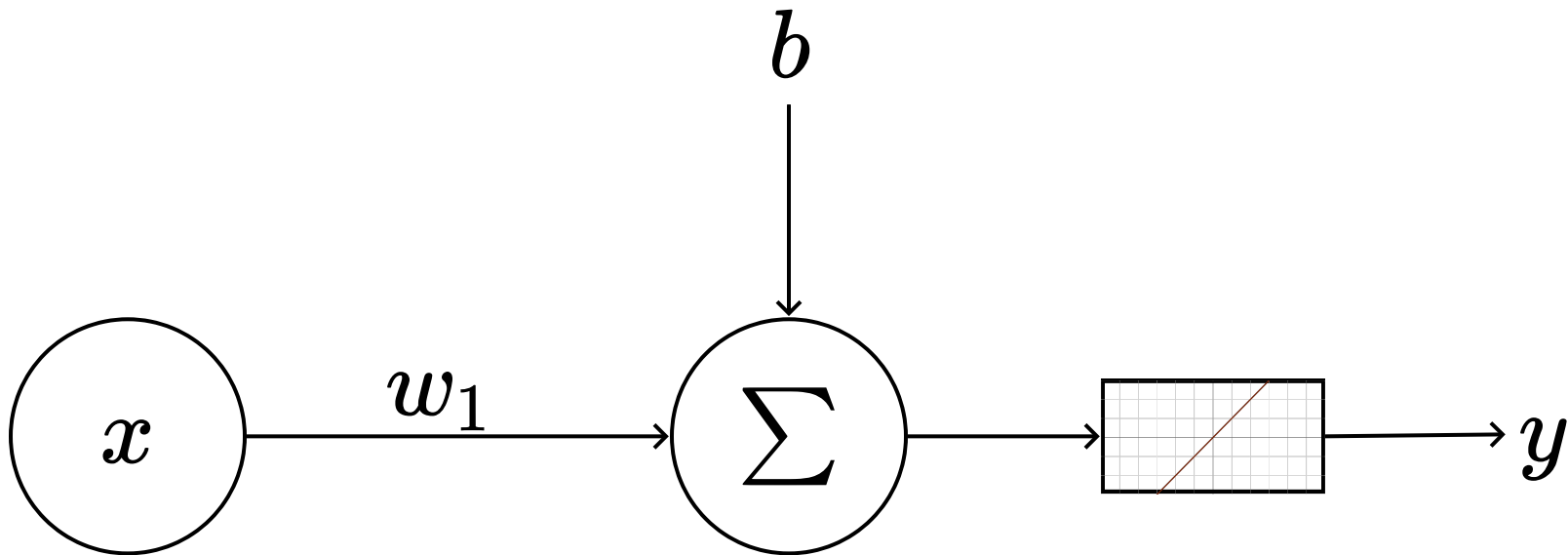






$$\hat{U}_w(x) = w_0 + w_1 x$$

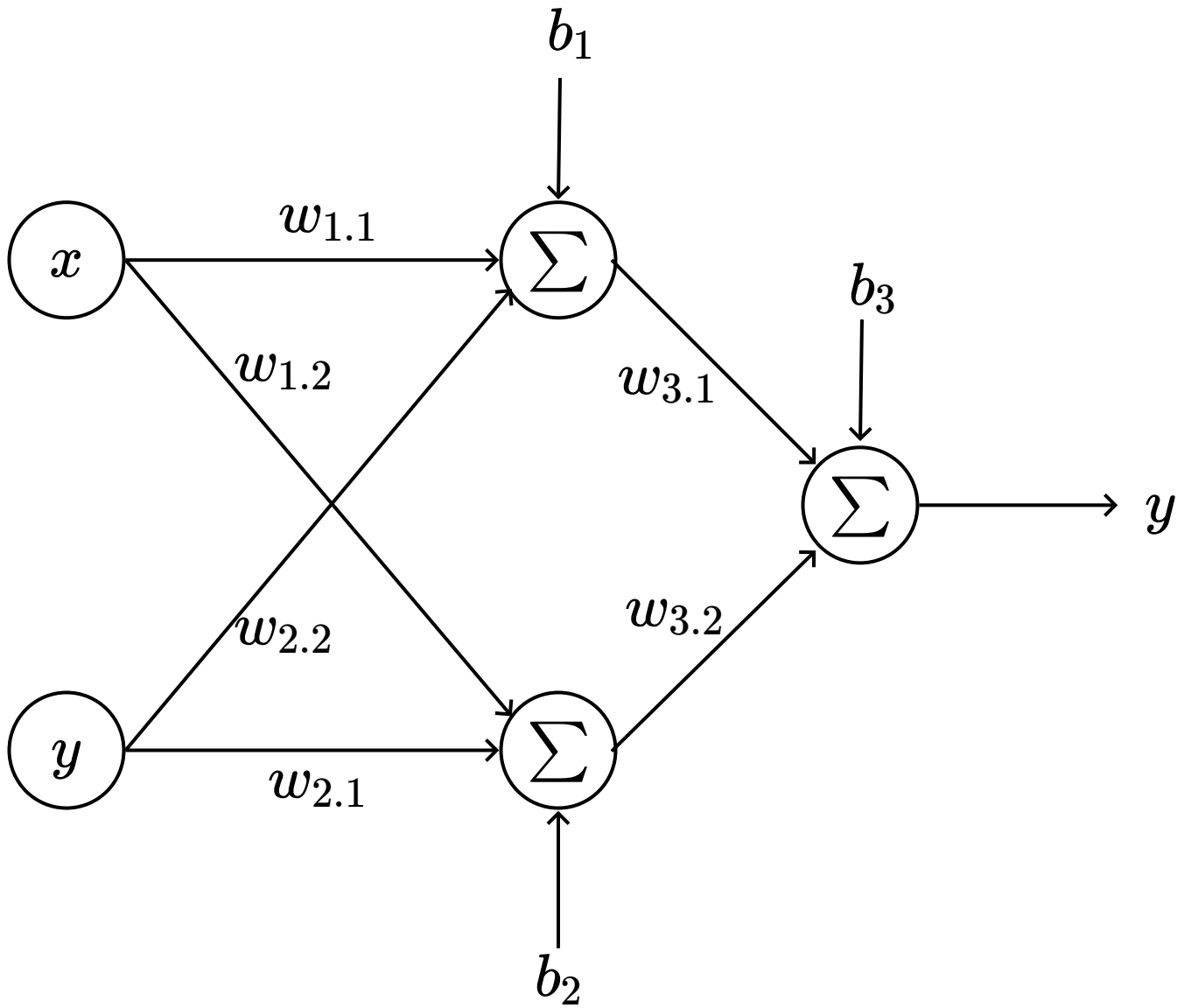
$$\hat{U}_w(x) = 0.43 + 0.02x$$

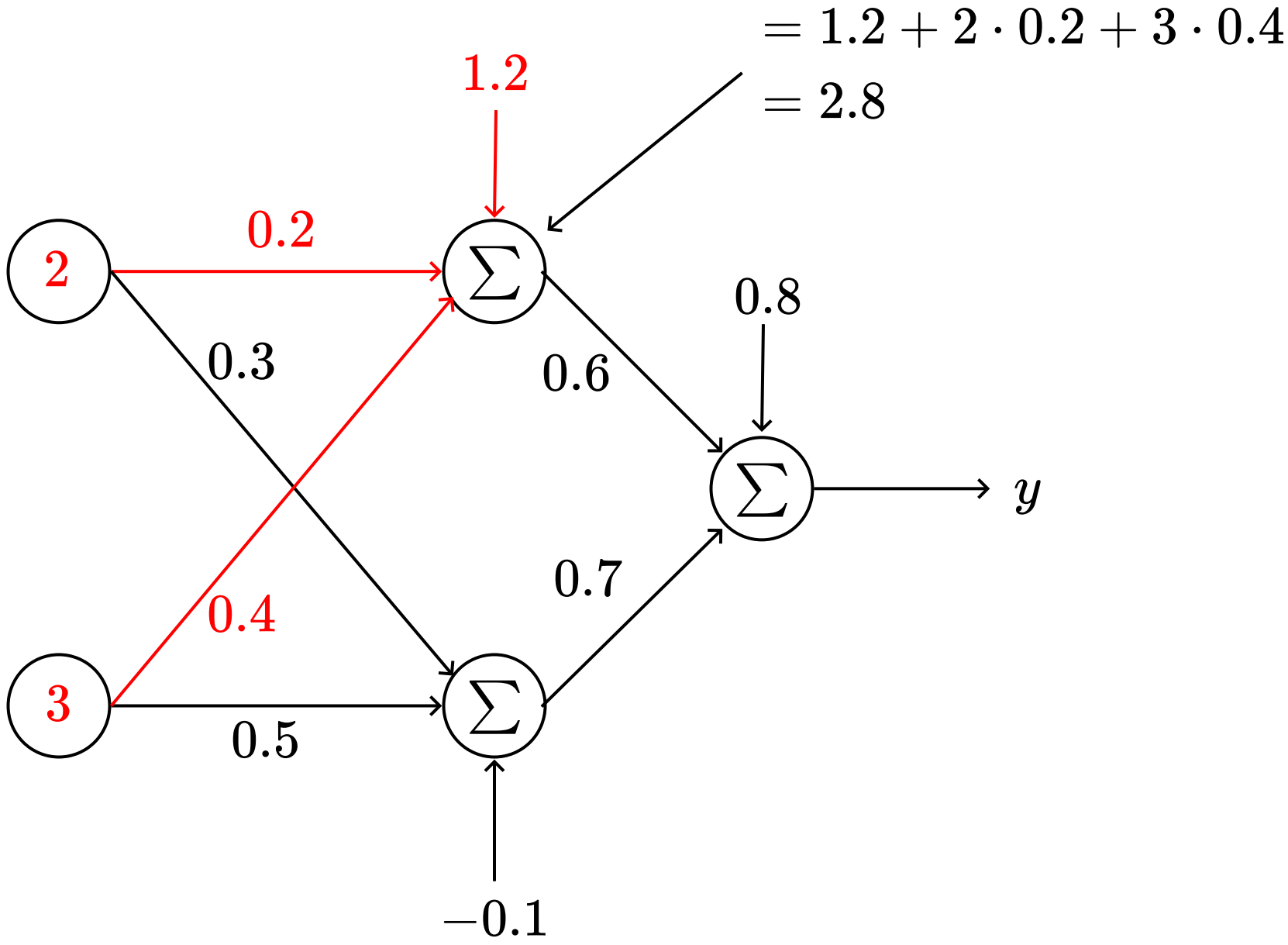


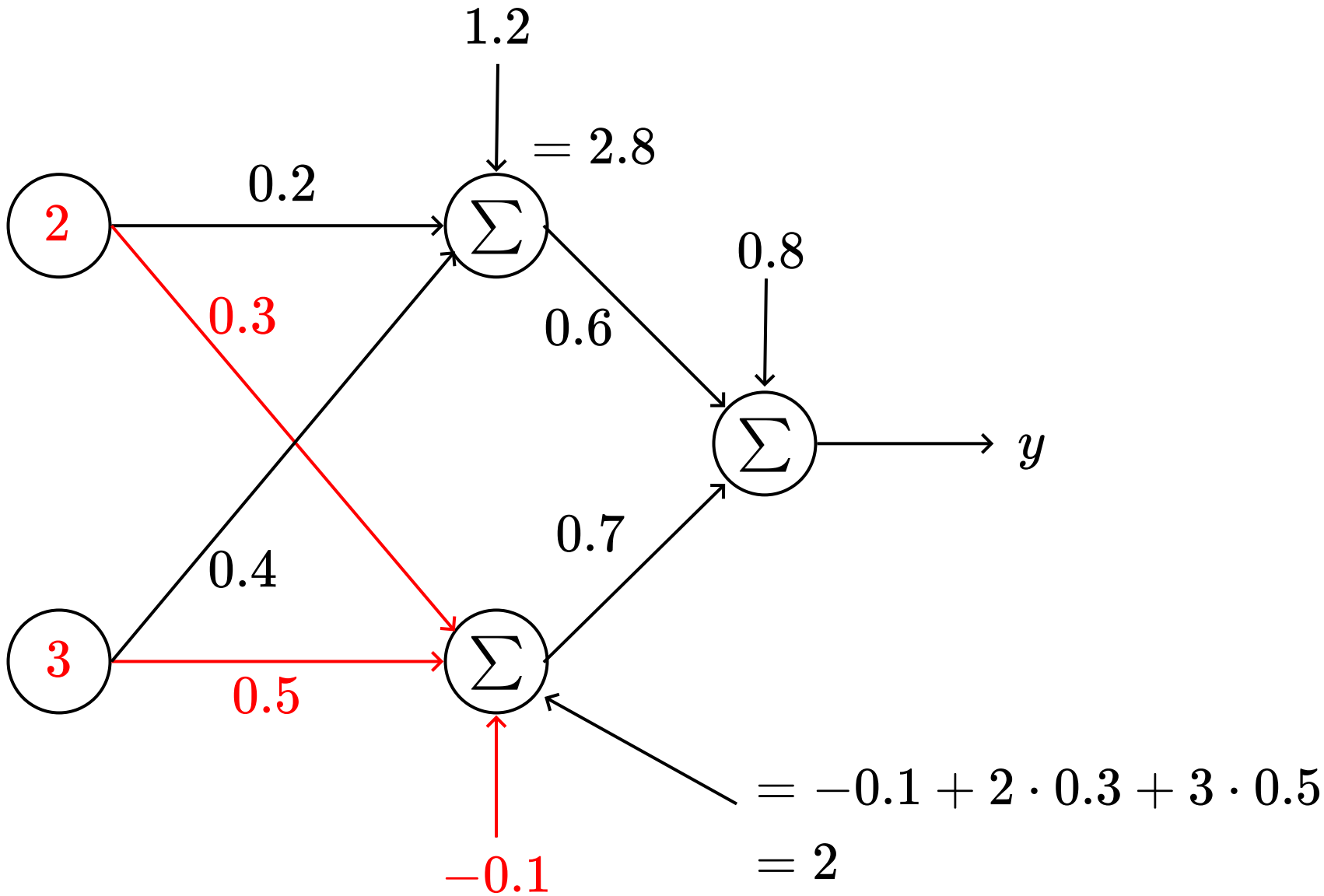
$$\hat{U}_w(x) = w_0 + w_1 x$$

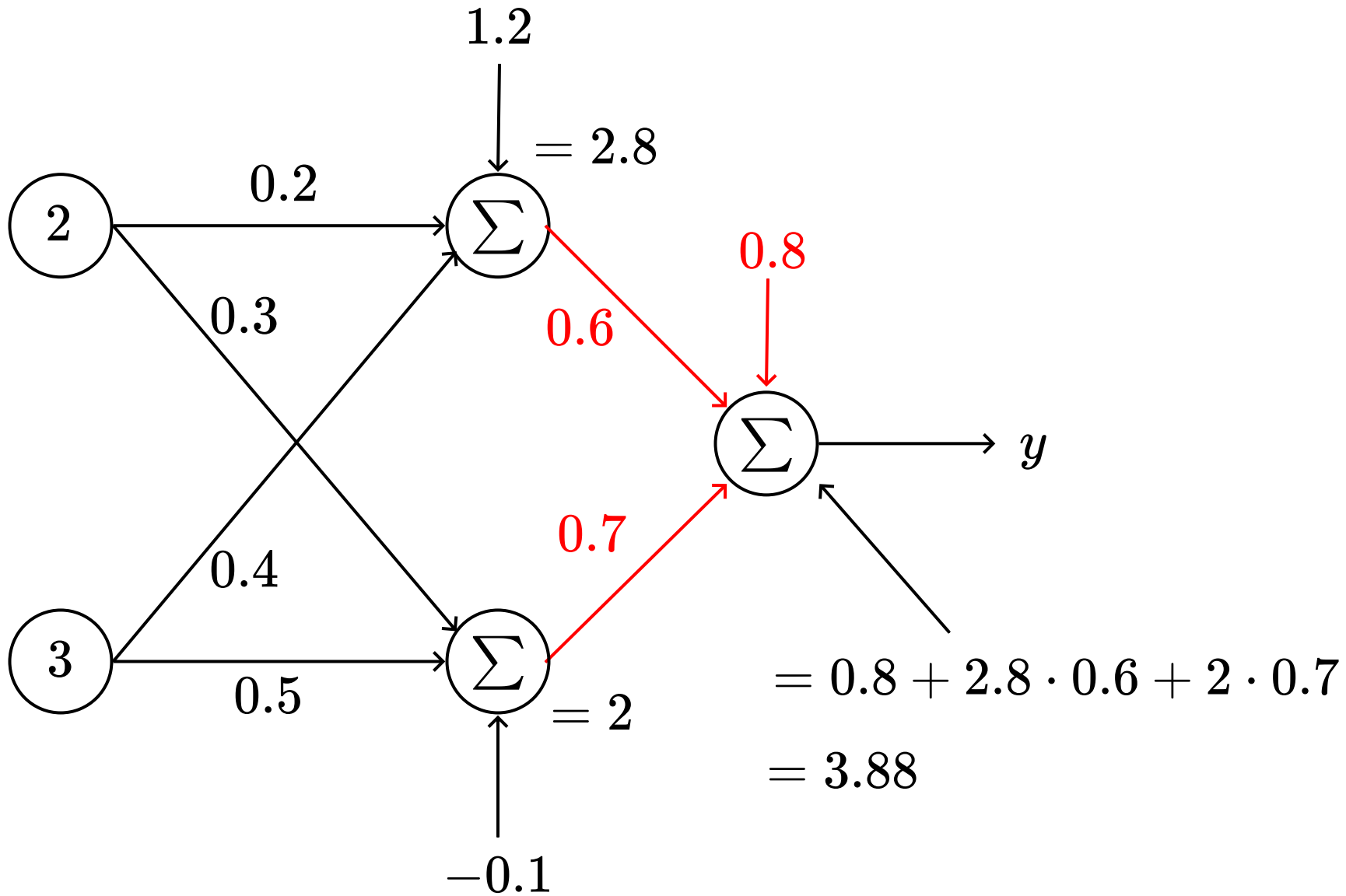
$$\longrightarrow b = 0.43, w_1 = 0.02$$

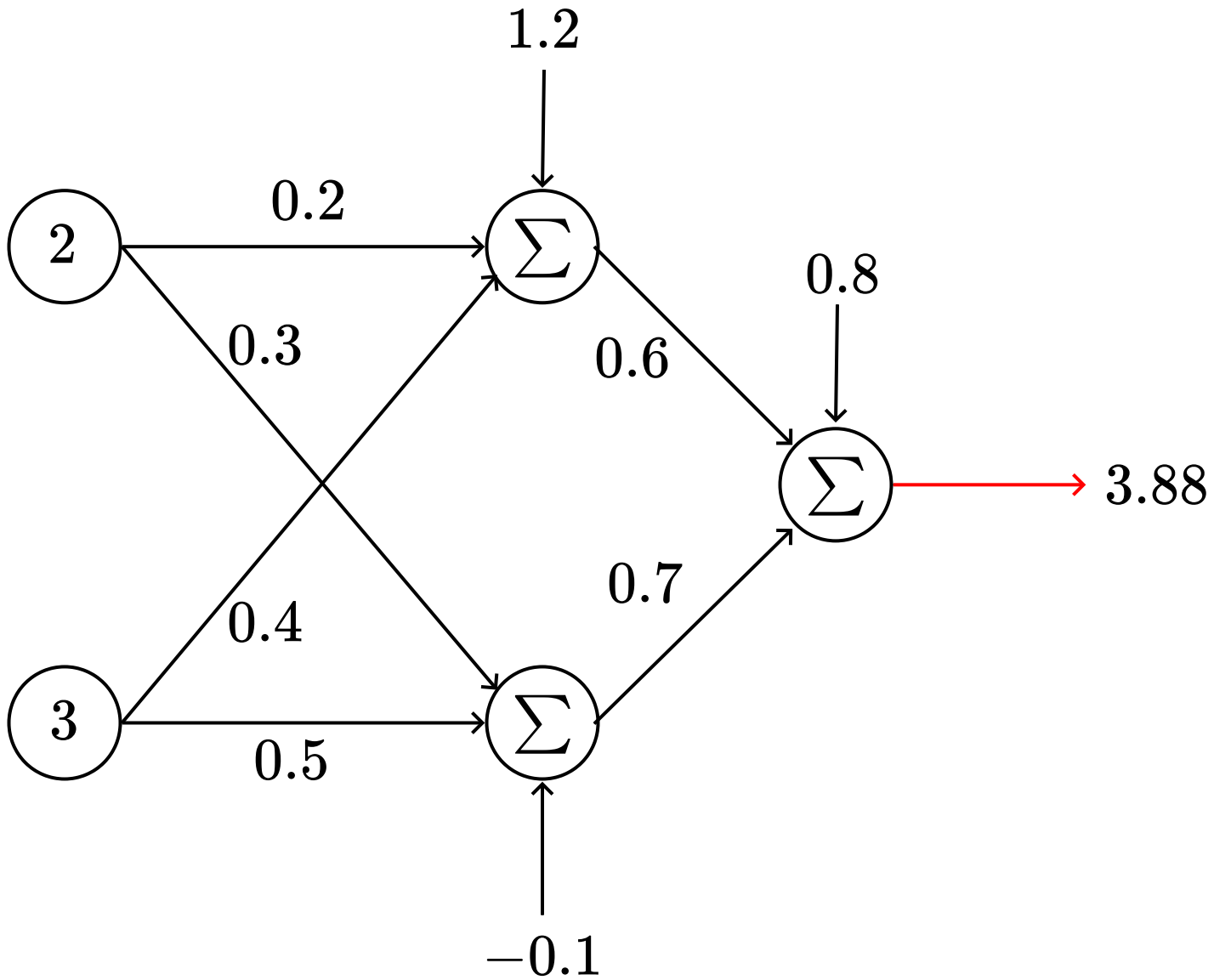
$$\hat{U}_w(x) = 0.43 + 0.02x$$

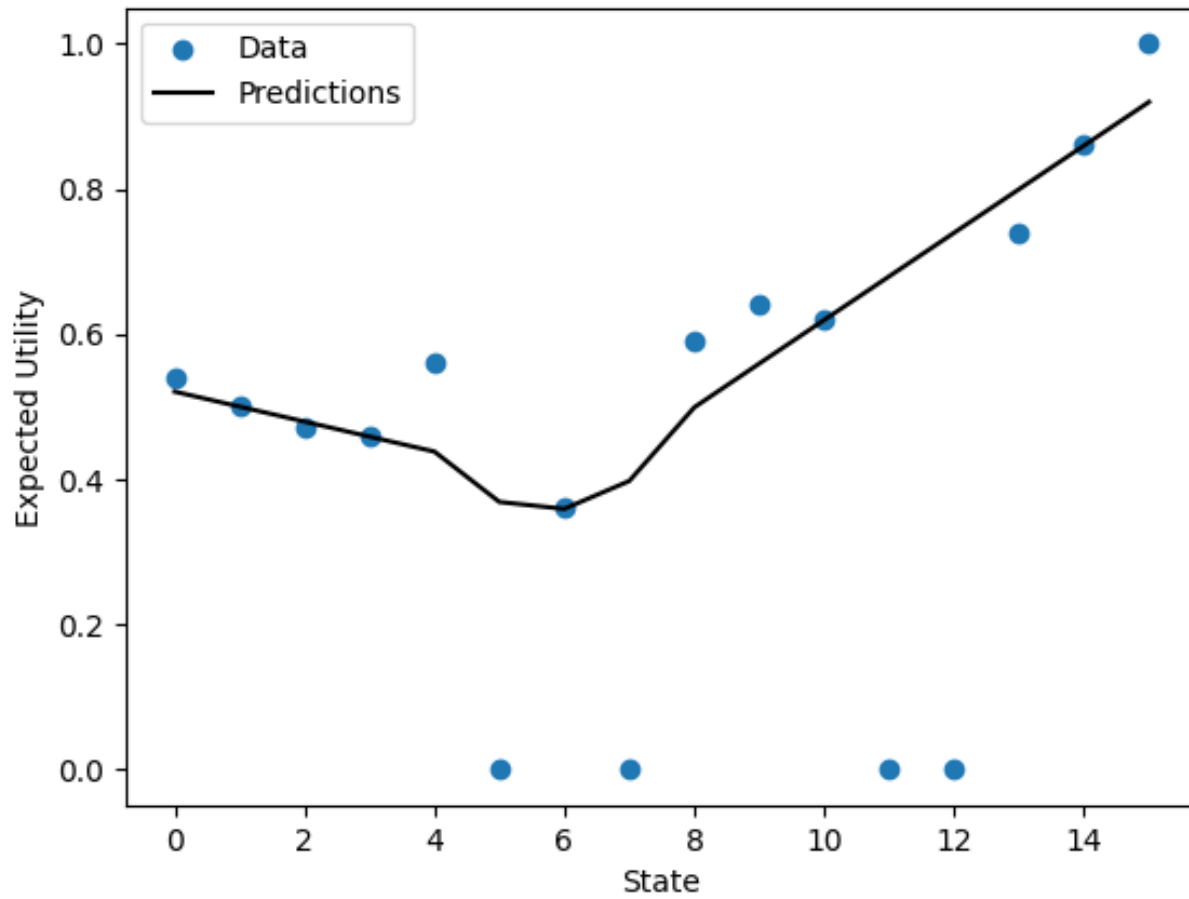


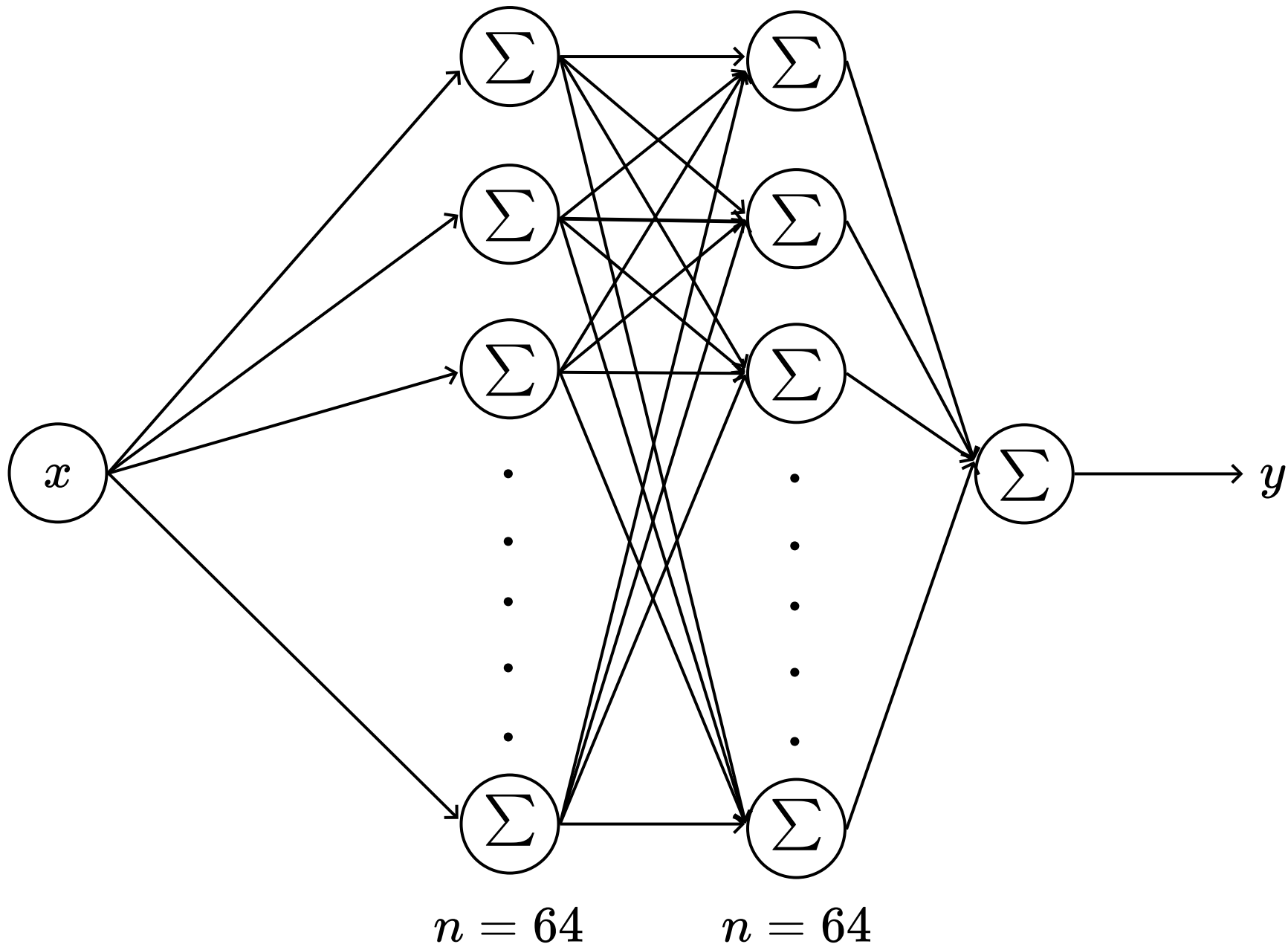


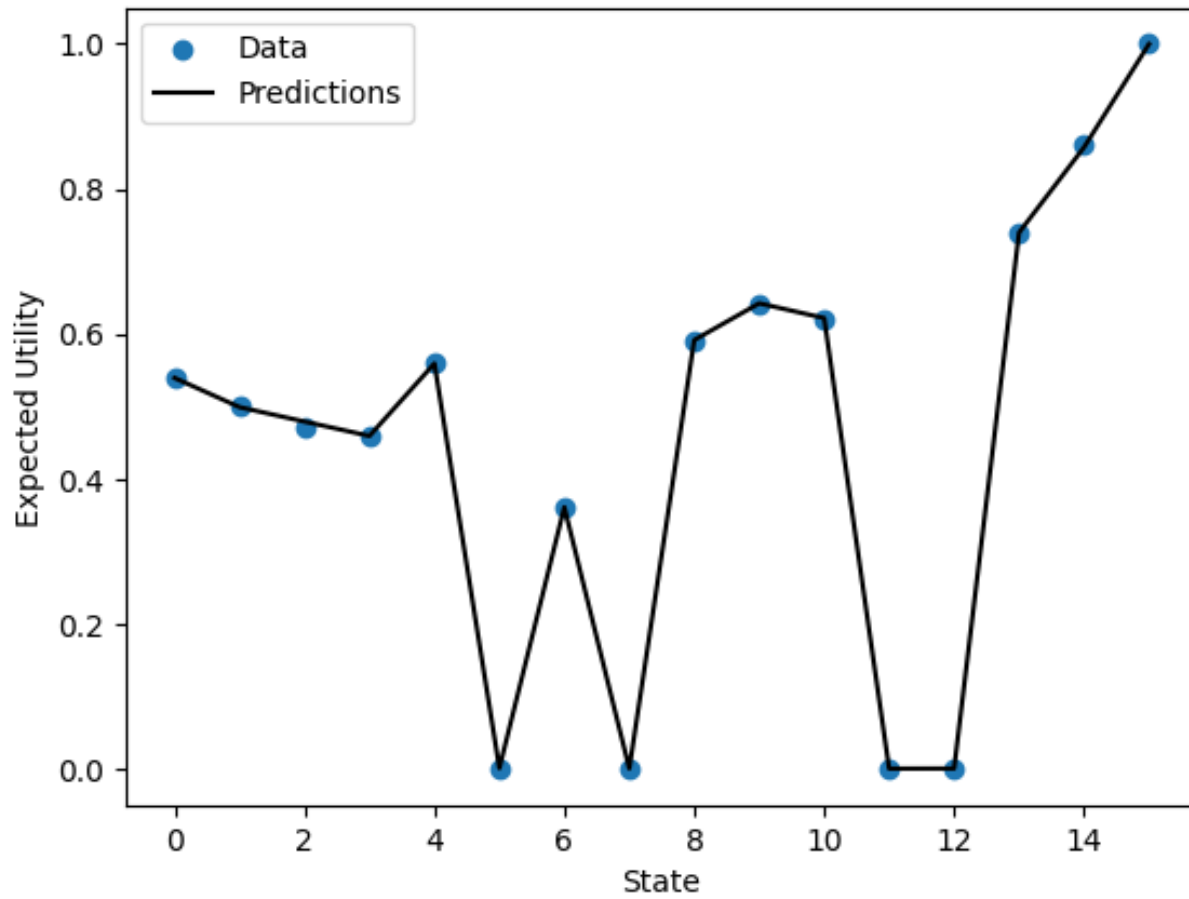


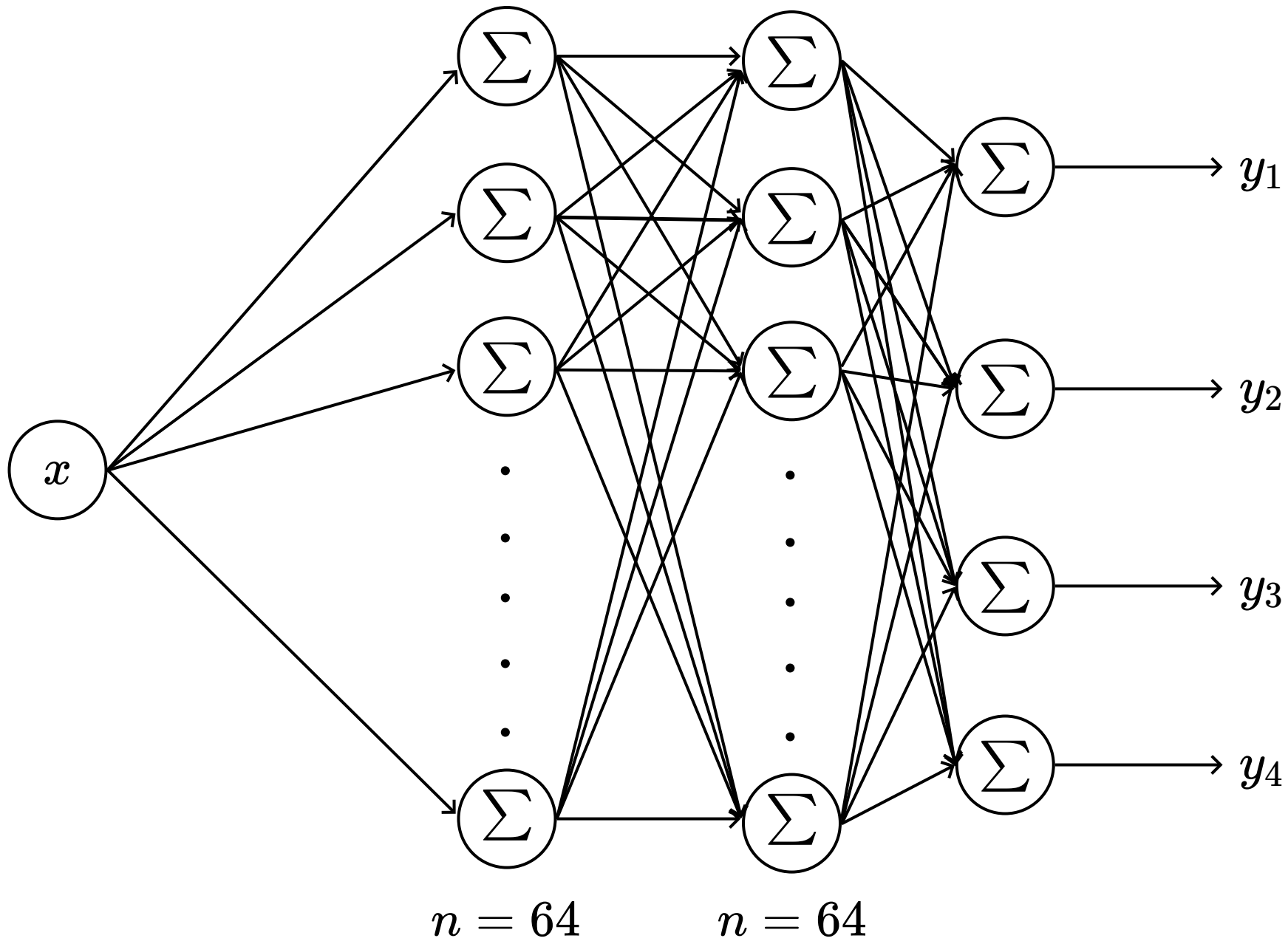


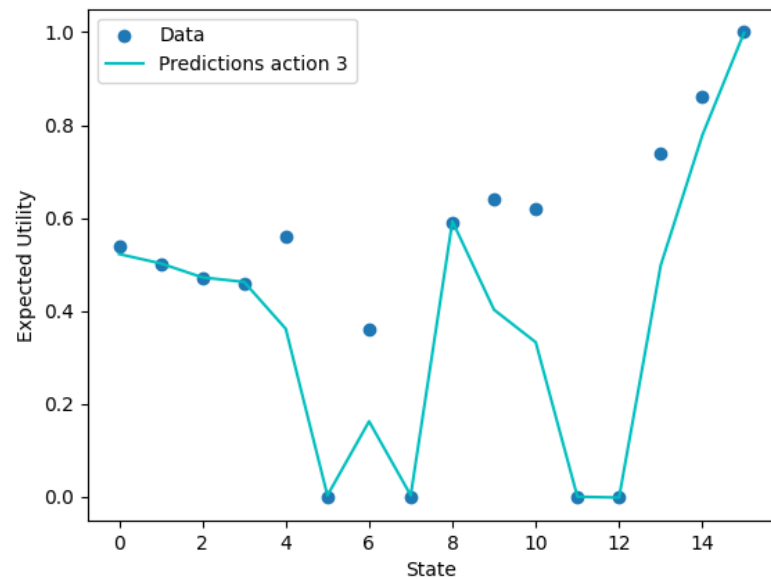
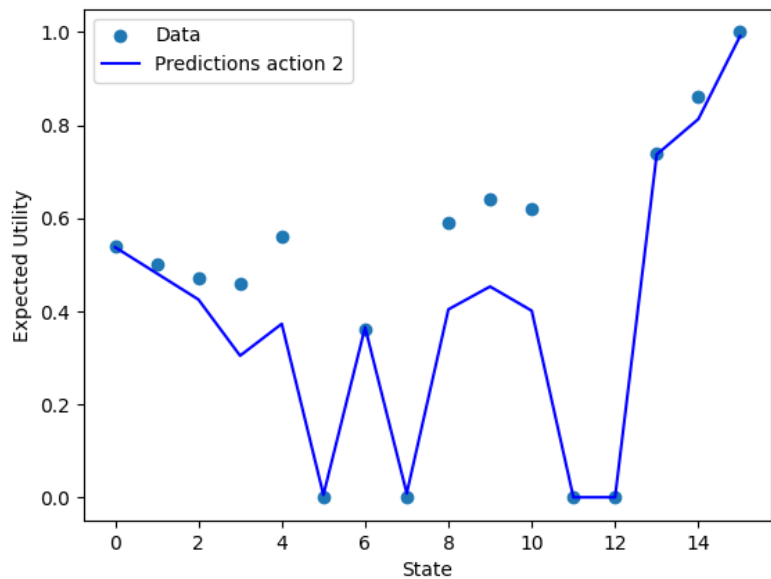
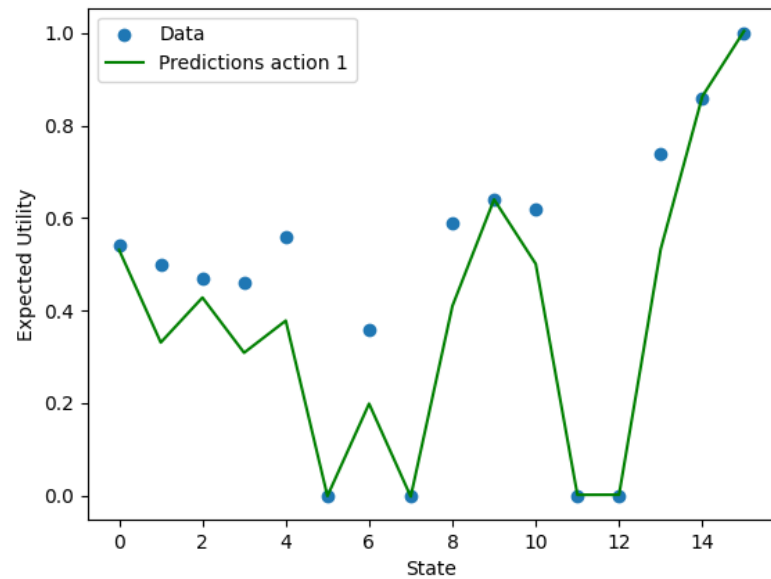
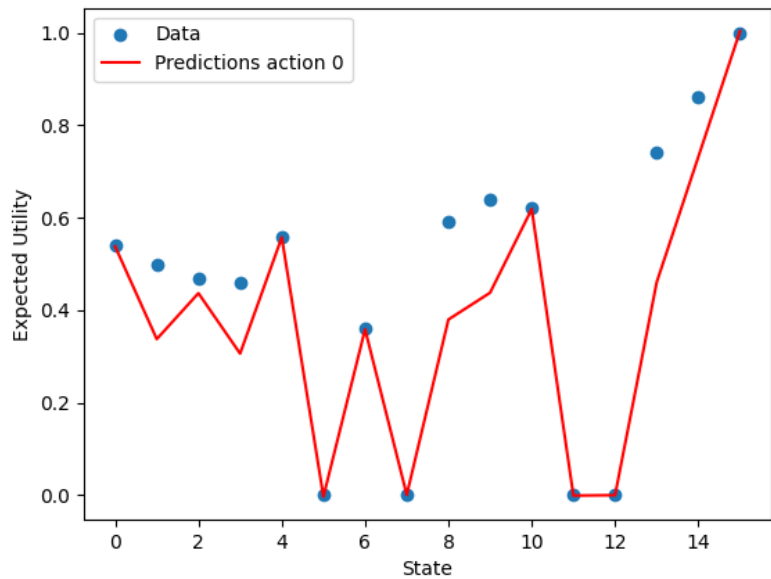












- Recap
- Generalization in Reinforcement Learning
- Artificial Neural Networks
- • Bellman Update for ANNs
- Deep Q-Network

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$



$$w \leftarrow w + \alpha [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)) \nabla_w Q(s, a, w)]$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[\underbrace{R(s, a, s') + \gamma \max_{a'} Q(s', a')} - Q(s, a) \right]$$



$$w \leftarrow w + \alpha \left[\underbrace{(R(s, a, s') + \gamma \max_{a'} Q(s', a', w)) - Q(s, a, w)} \nabla_w Q(s, a, w) \right]$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$



$$w \leftarrow w + \alpha [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)) \underline{\nabla_w Q(s, a, w)}]$$

$$\underline{Q(s, a)} \leftarrow Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$



$$\underline{w} \leftarrow w + \alpha [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)) \nabla_w Q(s, a, w)]$$

- Recap
 - Generalization in Reinforcement Learning
 - Artificial Neural Networks
 - Bellman Update for ANNs
- • Deep Q-Network

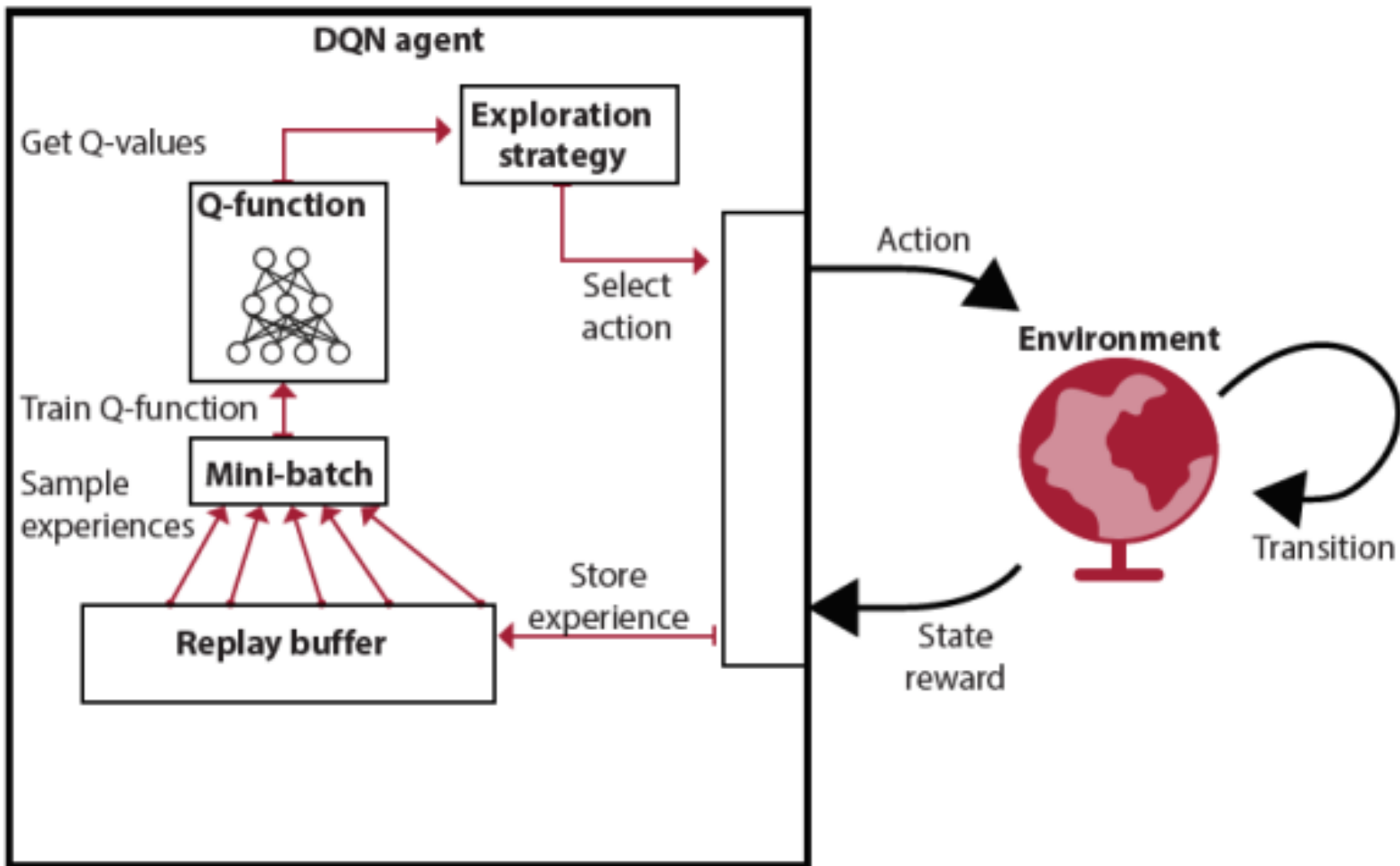


Image from Grokking Deep Reinforcement Learning

Experience Replay

Incremental methods have several problems:

- They are not **sample-efficient**.
- They do not work well with popular (and computationally efficient) ML algorithms.
- They may forget rare experiences that would be useful later on (**catastrophic forgetting**).

Experience replay and prioritized experience replay address these issues by storing experiences and reusing them. The experiences can then be sampled uniformly or by prioritization.

Deep Q-Network

```
1 initialize replay memory D
2 initialize action-value function Q with random weights
3 observe initial state s
4 repeat
5     select an action a
6         with probability  $\epsilon$  select a random action
7         otherwise select  $a = \operatorname{argmax}_a Q(s, a)$ 
8     carry out action a
9     observe reward r and new state s'
10    store experience  $\langle s, a, r, s' \rangle$  in replay memory D
11
12    sample random transitions  $\langle ss, aa, rr, ss' \rangle$  from replay memory D
13    calculate target for each minibatch transition
14        if ss' is terminal state then  $tt = rr$ 
15        otherwise  $tt = rr + \gamma \max_a Q(ss', aa)$ 
16    train the Q network using  $(tt - Q(ss, aa))^2$  as loss
17
18     $s = s'$ 
19 until terminated
```

- Recap
- Generalization in Reinforcement Learning
- Artificial Neural Networks
- Bellman Update for ANNs
- Deep Q-Network

Questions?

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$



$$w \leftarrow w + \alpha [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)) \nabla_w Q(s, a, w)]$$

$$w \leftarrow w - \frac{1}{2} \alpha \nabla_w [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w))^2]$$

$$\nabla_{\theta_i} L_i(\theta_i) \leftarrow \mathbb{E}_{s,a,r,s'} [(r + \gamma \max_{a'} Q(s', a'; \theta_i) - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i)]$$

$$\Delta_w = -\frac{1}{2} \alpha \nabla_w [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w))^2]$$

$$\Delta_w = \alpha [(R(s, a, s') + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)) \nabla_w Q(s, a, w)]$$