



A.D. 1308

**unipg**

DIPARTIMENTO  
DI INGEGNERIA

UNIVERSITÀ DEGLI STUDI DI PERUGIA

Dipartimento di Ingegneria

Laurea Magistrale in  
Ingegneria Informatica e Robotica

**Un fantastico titolo  
per la mia tesi di laurea!**

Relatore:

**Prof. Gabriele Costante**

Candidato:

**Paolo Speziali**

---

Anno Accademico 2023/2024



# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
<b>2</b>	<b>State of the Art</b>	<b>9</b>
<b>3</b>	<b>Prior Knowledge</b>	<b>11</b>
3.1	Deep Learning . . . . .	12
3.2	Reinforcement Learning . . . . .	14
3.2.1	Reinforcement Learning Algorithms . . . . .	17
3.3	Generative Adversarial Networks . . . . .	19
3.4	Causal Inference . . . . .	20
<b>4</b>	<b>Methodology</b>	<b>21</b>
<b>5</b>	<b>Experiments</b>	<b>23</b>
<b>6</b>	<b>Conclusion</b>	<b>25</b>



# List of Figures

3.1	A feed-forward neural network with two hidden layers . . . . .	12
3.2	A simple reinforcement learning task . . . . .	14
3.3	The example grid as a Markov Process graph, where the nodes represent the states and the edges represent the state transitions. . . . .	17
3.4	Partial taxonomy of algorithms in modern RL. . . . .	18



# Chapter 1

## Introduction





## Chapter 2

# State of the Art



## Chapter 3

# Prior Knowledge

## 3.1 Deep Learning

In the field of machine learning, the first models you are often introduced to are those for regression and classification that utilize linear combinations of fixed basis functions. According to [Bishop, 2008], these models possess useful analytical and computational properties, but their practical applicability is constrained since their capacity is limited to linear functions and they cannot understand the interaction between any two input variables. This leads to problems such as the **curse of dimensionality**, where the number of possible interactions between variables grows exponentially with the number of input variables. To address large-scale problems, it is essential to adapt the basis functions to the data, as demonstrated by support vector machines (SVMs). Alternatively, one can fix the number of basis functions in advance while allowing them to be adaptive, utilizing parametric forms for the basis functions with parameter values adjusted during training.

The most successful example of this approach in pattern recognition is the **feed-forward neural network**, also known as the multilayer perceptron (MLP).

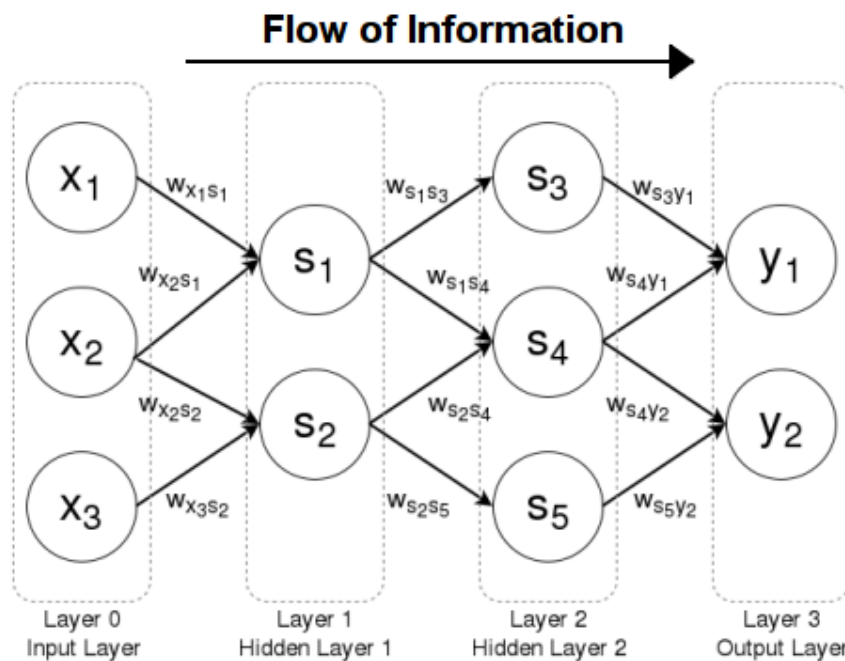


Figure 3.1: A feed-forward neural network with two hidden layers

Source: Brilliant

As explained in [Goodfellow et al., 2016], the goal of a feedforward network is to approximate a specific function  $f^*$ . For example, in a classification task,  $y = f^*(x)$  maps an input  $x$  to a category  $y$ . A feedforward network defines a mapping  $y = f(x; \theta)$  and learns the parameters  $\theta$  that yield the best approximation of this function.

These models are called **feedforward** because information flows unidirectionally

from the input  $x$ , through intermediate computations defining  $f$ , and ultimately to the output  $y$ , without feedback loops.

The training data directly specifies the output layer's behavior but not the intermediate layers, which are called **hidden layers** because their desired outputs are not provided by the training data. The learning algorithm must determine how to best utilize these hidden layers to approximate  $f^*$ . Every layer of the network computes a non-linear transformation of the previous layer's activations, this way a complex function can be approximated by composing simpler functions, one for each layer. Layers are composed of a set of **units**, where each unit is a node that computes a non-linear function of the weighted sum of its inputs and is only connected to units in the previous and the following layer.

In order to train the network and update the weights, MLPs use the **backpropagation** technique, which computes the gradient of the loss function with respect to the weights of the network. The weights are then updated using an optimization algorithm such as **stochastic gradient descent** (SGD) or one of its adaptive variants like **Adam**, whose hyperparameters are tuned according to the task during training.

A feedforward network is called a **deep neural network** if it has more than one hidden layer, the branch of machine learning that studies deep neural networks is called **deep learning**.

They are called networks due to their structure, which involves composing multiple functions together. Typically, this composition is represented by a directed acyclic graph. For instance, functions  $f^{(1)}$ ,  $f^{(2)}$ , and  $f^{(3)}$  might be connected in a chain to form  $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$ .

Functions  $f^{(1)}$  and  $f^{(2)}$  must be non linear, otherwise the composition would collapse into a single linear function. These non-linear functions are called **activation functions**. The last function  $f^{(3)}$  is typically a linear function that maps the output of the final hidden layer to the output layer. This is the same as applying a linear model to a transformed input  $\phi(x)$ , where  $\phi$  is a nonlinear transformation. The question then becomes how to choose the mapping  $\phi$ .

The strategy of deep learning is to learn  $\phi$ : in this approach, we use a model

$$y = f(x; \theta, w) = \phi(x; \theta)^\top w \quad (3.1)$$

Here, we have parameters  $\theta$  that are used to learn  $\phi$  from a broad class of functions, and parameters  $w$  that map  $\phi(x)$  to the desired output. This approach allows for greater flexibility: specifically, by using a broad family of functions  $\phi(x; \theta)$ , the human designer only needs to select the appropriate general function family rather than finding precisely an exact function.

The **universal approximation theorem** [Hornik et al., 1989] tells us that regardless of what function we are trying to learn, we know that a feedforward network with enough units will be able to *represent* this function, however we are not guaranteed that the training algorithm will be able to *learn* it.

## 3.2 Reinforcement Learning

“**Reinforcement learning** is learning [...] how to map situations to actions so as to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate reward but also the next situation and, through that, all subsequent rewards. These two characteristics, **trial-and-error search** and **delayed reward**, are the two most important distinguishing features of reinforcement learning.” [Sutton and Barto, 1998]

Let’s explore the basic concepts of reinforcement learning with a simple example as described in [Zhao, 2024]. Consider a grid world scenario as shown in Figure 3.2 where a robot, referred to as the **agent**, moves between cells, occupying one cell at a time. The white cells are accessible, while the orange cells are forbidden. The goal is for the agent to reach a target cell.

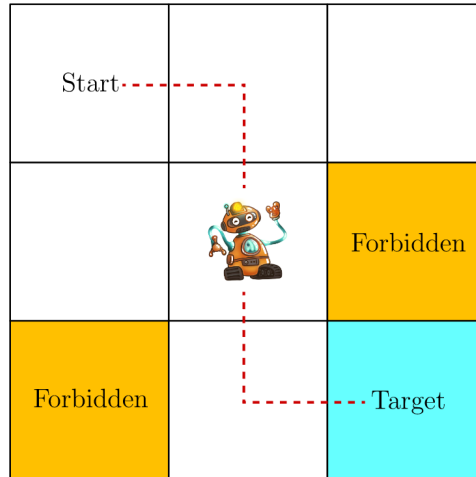


Figure 3.2: A simple reinforcement learning task

Source: [Zhao, 2024]

To achieve this, the agent needs a **policy** that guides it to the target efficiently, avoiding forbidden cells and unnecessary detours. If the grid layout is known, planning is simple. However, without prior knowledge, the agent must explore and learn through trial and error.

The agent’s status in the grid is defined by its **state**  $s_i \in \mathcal{S}$ , which represents its location relative to the **environment**. In the examples with nine cells, the state space is  $\mathcal{S} = \{s_1, s_2, \dots, s_9\}$ .

From each state, the agent can perform five **actions**: move up, right, down, left or stay put, denoted as  $a_1, a_2, \dots, a_5$ , this set of actions is the action space  $\mathcal{A} = \{a_1, \dots, a_5\}$ . The available actions can vary by state, for instance, at  $s_1$ , actions  $a_1$  (up) and  $a_4$  (left) would collide with the grid boundary, so the action space is  $\mathcal{A}(s_1) = \{a_2, a_3, a_5\}$ .

When taking an action, the agent may move from one state to another, a process known as **state transition**. For example, if the agent is at state  $s_1$  and selects action  $a_2$  (moving rightward), it transitions to state  $s_2$ . This process can be represented as:

$$s_1 \xrightarrow{a_2} s_2$$

The state transition process is defined for each state and its associated actions. Mathematically, state transitions are described by conditional probabilities. For example, for  $s_1$  and  $a_2$ , the conditional probability distribution is:

$$\begin{aligned} p(s_1 \mid s_1, a_2) &= 0, \\ p(s_2 \mid s_1, a_2) &= 1, \\ p(s_3 \mid s_1, a_2) &= 0, \\ p(s_4 \mid s_1, a_2) &= 0, \\ p(s_5 \mid s_1, a_2) &= 0, \end{aligned}$$

indicating that taking  $a_2$  at  $s_1$  guarantees the agent moves to  $s_2$ , with a probability of one, and zero probability for other states.

This is a deterministic state transition, but state transitions can also be stochastic, requiring conditional probability distributions. For instance, if random wind gusts affect the grid, taking action  $a_2$  at  $s_1$  might blow the agent to  $s_5$  instead of  $s_2$ , resulting in  $p(s_5 \mid s_1, a_2) > 0$ .

A **policy** is a function that maps states to actions, indicating the agent's behavior in the environment. In other words, a policy tells the agent what action to take at each state.

Mathematically, policies can be described by conditional probabilities. For example, the policy for  $s_1$  is:

$$\begin{aligned} \pi(a_1 \mid s_1) &= 0, \\ \pi(a_2 \mid s_1) &= 1, \\ \pi(a_3 \mid s_1) &= 0, \\ \pi(a_4 \mid s_1) &= 0, \\ \pi(a_5 \mid s_1) &= 0, \end{aligned}$$

indicating that the probability of taking action  $a_2$  at state  $s_1$  is one, and the probabilities of taking other actions are zero.

The above policy is deterministic, but policies may generally be stochastic. For instance, let's assume that at state  $s_1$  the agent may take actions to move either rightward or downward, each with a probability of 0.5. In this case, the policy for  $s_1$  is:

$$\begin{aligned} \pi(a_1 \mid s_1) &= 0, \\ \pi(a_2 \mid s_1) &= 0.5, \\ \pi(a_3 \mid s_1) &= 0.5, \\ \pi(a_4 \mid s_1) &= 0, \\ \pi(a_5 \mid s_1) &= 0. \end{aligned}$$

After executing an action at a state, the agent receives a reward  $r$  as feedback from the environment. The **reward** is a function of the state and action which predicts immediate reward and is denoted as:

$$R(s_t = s, a_t = a) = \mathbb{E}[r_t | s_t = s, a_t = a] \quad (3.2)$$

and it can be positive, negative, or zero. Different rewards influence the policy the agent learns. Generally, a positive reward encourages the agent to take the corresponding action, while a negative reward discourages it.

However we can't find good policies by simply selecting actions with the greatest immediate rewards since they do not consider long-term outcomes. To determine a good policy, we must consider the total reward obtained in the long run and an action with the highest immediate reward may not lead to the greatest total reward.

A **trajectory** is a state-action-reward chain. For example, the agent in our example may follow this trajectory:

$$s_1 \xrightarrow{a_2, r=0} s_2 \xrightarrow{a_3, r=0} s_5 \xrightarrow{a_3, r=0} s_8 \xrightarrow{a_2, r=1} s_9.$$

The **return** of this trajectory is the sum of all rewards collected along it, in the example above:

$$\text{return} = 0 + 0 + 0 + 1 = 1$$

Returns, also called total rewards or cumulative rewards, are used to evaluate policies. Returns can also be defined for infinitely long trajectories, which may diverge. Therefore, we introduce the concept of **discounted return** for infinitely long trajectories. The discounted return is the sum of the rewards from  $t$  to  $T$  (final step):

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (3.3)$$

where  $\gamma \in [0, 1]$  is the **discount factor**. If  $\gamma$  is close to 0, the agent emphasizes near-future rewards, resulting in a short-sighted policy. If  $\gamma$  is close to 1, the agent emphasizes far-future rewards.

When  $T$  is finite, we call the task episodic and each sequence up to the **terminal state** is an **episode**. Otherwise, we refer to **continuing tasks**.

The **Markov Decision Processes** (MDPs), a general framework for describing stochastic dynamical systems, allows us to formally presents these concepts. The **Markov property** refers to the **memoryless** property of a stochastic process. Mathematically, it means that

$$\begin{aligned} p(s_{t+1} \mid s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0) &= p(s_{t+1} \mid s_t, a_t), \\ p(r_{t+1} \mid s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0) &= p(r_{t+1} \mid s_t, a_t), \end{aligned} \quad (3.4)$$

where  $t$  represents the current time step and  $t + 1$  represents the next time step. This indicates that the next state or reward depends only on the current state and action and is independent of the previous ones.



An MDP is defined by a tuple  $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$ , let's break down the components:

- $\mathcal{S}$ : finite set of Markov states  $s$
- $\mathcal{A}$ : finite set of actions  $a$
- $P$ : state transition model for each action, a probability matrix that specifies

$$p(s_t + 1 = s' | s_t = s, a_t = a)$$

- $R$ : reward function

$$R(s_t = s, a_t = a) = \mathbb{E}[r_t | s_t = s, a_t = a]$$

- $\gamma$ : discount factor  $0 \leq \gamma \leq 1$

When the policy in a MDP is fixed, it reduces to a **Markov Process** (MP), this transformation simplifies the MDP by eliminating the decision-making aspect. A Markov process is referred to as a **Markov Chain** if it operates in discrete time and the number of states is either finite or countable.

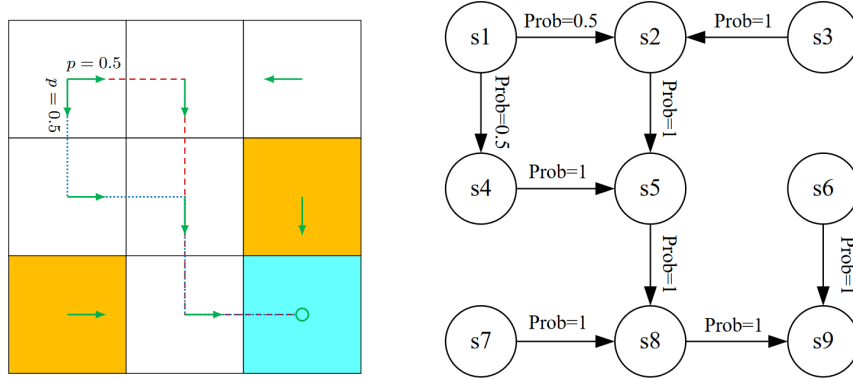


Figure 3.3: The example grid as a Markov Process graph, where the nodes represent the states and the edges represent the state transitions.

Source: [Zhao, 2024]

### 3.2.1 Reinforcement Learning Algorithms

In [OpenAI, 2018], the landscape of algorithms in modern reinforcement learning is explored.

The algorithms used in reinforcement learning almost always rely on **value functions**. The **value** is the expected return if the agent starts in that state or state-action pair and then follows a specific policy indefinitely.

There are four primary value functions to consider:

- **On-Policy Value Function,  $V^\pi(s)$ :** This function represents the expected return if the agent starts in state  $s$  and acts according to the policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s] \quad (3.5)$$

- **On-Policy Action-Value Function,  $Q^\pi(s, a)$ :** This function gives the expected return if the agent starts in state  $s$ , takes an arbitrary action  $a$ , and then always follows the policy  $\pi$ :

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s, a_0 = a] \quad (3.6)$$

- **Optimal Value Function,  $V^*(s)$ :** This function provides the expected return if the agent starts in state  $s$  and follows the optimal policy for the environment:

$$V^*(s) = \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s] \quad (3.7)$$

- **Optimal Action-Value Function,  $Q^*(s, a)$ :** This function represents the expected return if the agent starts in state  $s$ , takes an arbitrary action  $a$ , and then follows the optimal policy for the environment:

$$Q^*(s, a) = \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) \mid s_0 = s, a_0 = a] \quad (3.8)$$

Figure 3.4 presents a partial taxonomy of the modern reinforcement learning algorithms. One key distinction is whether the agent utilizes (or learns) a model of the environment. A model, in this context, refers to a function that predicts state transitions and rewards.

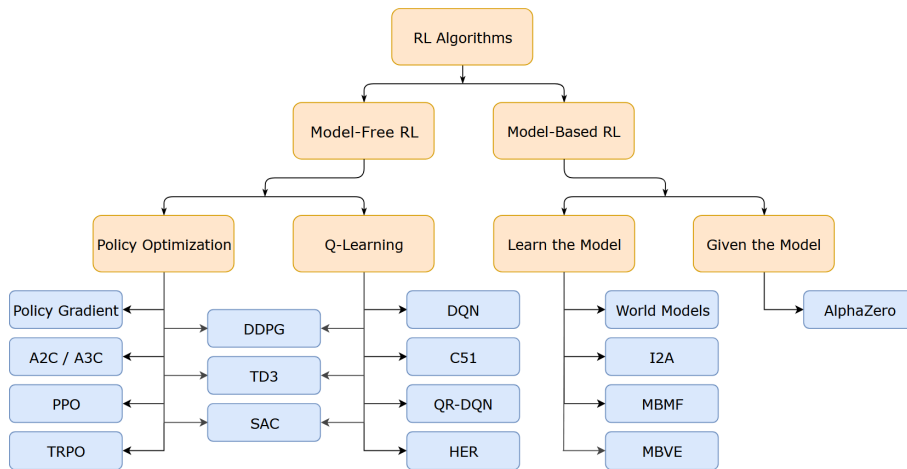


Figure 3.4: Partial taxonomy of algorithms in modern RL.

Source: [OpenAI, 2018]

The primary advantage of having a model is that it enables the agent to plan by predicting future states and rewards. However, the major drawback is that an accurate model of the environment is often unavailable to the agent. Algorithms which use a model are known as **model-based** methods, while those that do not are referred to as **model-free** methods. Our focus will be on the latter.

Model-free methods can be categorized into two main approaches for representing and training agents:

- **Policy Optimization:** explicitly represents a policy as  $\pi_\theta(a \mid s)$  and optimizes the parameters  $\theta$  either directly via gradient ascent on the performance objective  $J(\pi_\theta)$  or indirectly by maximizing local approximations of  $J(\pi_\theta)$ . Typically, this optimization is performed **on-policy**, meaning that each update uses data collected while the agent is acting according to the most recent version of the policy.
- **Q-Learning:** involves learning an approximator  $Q_\theta(s, a)$  for the optimal action-value function  $Q^*(s, a)$ . Unlike policy optimization, Q-learning is generally performed **off-policy**, allowing updates to use data collected at any time during training, irrespective of the agent's policy at the time of data collection. The corresponding policy is derived from the relationship between  $Q^*$  and  $\pi^*$ , where the actions taken by the Q-learning agent are determined by

$$a(s) = \arg \max_a Q_\theta(s, a) \quad (3.9)$$

An example of a Q-learning algorithm is Deep Q-Network (**DQN**), which uses a deep neural network to approximate the optimal action-value function in environments with large state spaces, and its variants such as Double DQN (**DDQN**) and Dueling Double DQN (**D3QN**).

There are also hybrid approaches that combine elements of both policy optimization and Q-learning. In this spectrum of algorithms we can find **actor-critic methods**, which consist of two components: an **actor** that makes actions and a **critic** that evaluates them. An example of an actor-critic algorithm is Twin-Delayed DDPG (**TD3**).

If samples are collected during training the reinforcement learning is considered **online**, otherwise, if the training set is fixed, it is **offline**.

### 3.3 Generative Adversarial Networks

**Generative Adversarial Networks** (GANs) is a framework for estimating generative models via an adversarial process, initially introduced by [Goodfellow et al., 2014]. This framework involves training two models: a **generative** model  $G$ , which captures the data distribution, and a **discriminative** model  $D$ , which predicts samples as either originating from the true data distribution or the generative model. The

training objective for  $G$  is to maximize the likelihood of  $D$  making incorrect classifications. This adversarial process can be conceptualized as a minimax two-player game. In the theoretical space of arbitrary functions  $G$  and  $D$ , a unique solution exists where  $G$  accurately recovers the training data distribution and  $D$  outputs  $\frac{1}{2}$  for all inputs (in other words, it reaches maximum uncertainty). When  $G$  and  $D$  are defined by multilayer perceptrons, the entire system can be trained using backpropagation.

To learn the generator's distribution  $p_g$  over data  $x$ , we define a prior  $p_z(z)$  on input noise variables  $z \in \Omega_z$ , called **latent distribution**, and represent a mapping to the data space as  $G(z; \theta_g)$ , where  $G : \Omega_z \rightarrow \Omega_x$  is a differentiable function parameterized by  $\theta_g$ . Additionally, we define a second network  $D(x; \theta_d)$  that outputs a single scalar, representing the probability that  $x$  originated from the true data distribution  $p_x$  rather than from  $p_g$ . The discriminative model  $D$  is trained to maximize the probability of correctly labeling both the training examples and the samples generated by  $G$ . Simultaneously, the generative model  $G$  is trained to minimize  $\log(1 - D(G(z)))$ .

Formally,  $D$  and  $G$  engage in the following two-player minimax game with the value function  $V(G, D)$ :

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]. \quad (3.10)$$

## 3.4 Causal Inference

[Neal, 2020]

## Chapter 4

# Methodology



## Chapter 5

# Experiments





Chapter 6

Conclusion



# Bibliography

- [Bishop, 2008] Bishop, C. M. (2008). *Pattern Recognition and Machine Learning*. Springer-Verlag, New York.
- [Goodfellow et al., 2016] Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. Adaptive computation and machine learning. MIT Press.
- [Goodfellow et al., 2014] Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial networks.
- [Hornik et al., 1989] Hornik, K., Stinchcombe, M., and White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366.
- [Neal, 2020] Neal, B. (2020). Introduction to causal inference.
- [OpenAI, 2018] OpenAI (2018). Part 2: Kinds of RL Algorithms - Spinning Up documentation — spinningup.openai.com. [https://spinningup.openai.com/en/latest/spinningup/rl\\_intro2.html](https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html). [Accessed 30-06-2024].
- [Sutton and Barto, 1998] Sutton, R. S. and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. The MIT Press.
- [Zhao, 2024] Zhao, S. (2024). *Mathematical Foundations of Reinforcement Learning*. Springer Nature Press.