

Microgrid Energy Training Using Reinforcement Learning

A Case Study

Moayad Elamin
Fay Elhassan

A thesis presented for the degree of
bachelor of engineering



Electrical and Electronic Engineering
University of Khartoum
Sudan
August 2020

Contents

1	Introduction	4
1.1	Problem Statement	4
1.2	MicroGrids	5
1.3	Reinforcement Learning	5
1.4	Research Objectives	6
1.5	Research Questions	6
2	Literature Review	7
2.1	Microgirds	7
2.1.1	Theoretical Background	7
2.1.2	Technical challenges of Micro-grids	9
2.1.3	Control hierarchy in micro-grids	10
2.2	Reinforcement Learning	13
2.2.1	Machine Learning Introduction	13
2.2.2	Reinforcement Learning Introduction	15
2.2.3	Cross-Entropy Methods	17
2.2.4	Q-Leaning	17
2.2.5	Deep Reinforcement Learning	18
2.2.6	Policy Gradients Methods	19
2.3	Related Work	20

List of Figures

2.1	A Simple Micro-Grid	7
2.2	Hierarchy diagram	11
2.3	A system with two voltage sources	11
2.4	$P - w$ and $Q - V$ droop characteristics	12
2.5	Machine Learning Sections	13
2.6	Simple Neural Network	14
2.7	Reinforcement Learning Parts	15

Chapter 1

Introduction

1.1 Problem Statement

As electricity is the backbone of all development and innovation, Sudan's electricity situation is a challenge that needs to be tackled with excellent efficiency and using innovative solutions. According to governmental sources, the national unified electricity grid covers less than 40% of Sudan; more than 60% of this is residential demand, 75% of power generated goes to the capital state; Khartoum state. Khartoum itself has a 66% connection percentage, while the 3rd and fourth states based on population density; South Darfur and North Kordofan have 2% and 5% connection percentage.

Electricity in Sudan has an average production cost of 18 cents; residential electricity is heavily subsidized and sold at a range of 1.8% to 19.4% of the production cost. Adding this to a usage increase of 10% annually, highest amongst neighboring countries, means that electricity inequality will only continue to grow. The top 1% of earners in Sudan use 49% of electricity produced while the lowest 49% use only 21% of electricity produced. This situation needs a re-evaluation of the system used in Sudan to achieve a fast recovery and adaptation to modern models.

Sudan's system is a conventional system of centralized generation; the government uses a single system to generate electricity at far-away stations or dams; then, it is distributed in the country. This system is dying in the world as the old equipment causes high maintenance cost and reliability problems; it is limited to the addition of new demand areas.

Old engineering, outdated systems of management, and efficiency concerns contribute to Sudan's astronomical loss rate of 25%. Combining this with its environmental concerns and lack of renewable energy utilization make it one to be changed.

The distributed generation where we generate where we will consume and have smaller grids, smaller generation and smaller transportation of electricity is an alternative worth considering. Here we have a flexible system, easy to maintain, efficient, modular, reliable, economically efficient, and above all environmentally responsible because of the renewable energy cornerstone, it stands on.

When looking into Sudan and its high number of villages and lack of major cities, we can see that microgrids present a great model of distributed generation to Sudan's rural areas. Microgrids are a variation of smart grids where we create small scale grids to work on villages, islands, and small residential areas, which can work as a standalone islanded grid or be grid-connected. It uses distributed energy

sources and renewable energy sources to generate electricity locally and then fulfill the local demand and use storage units for night demand and fault cases. As it is a smart grid, then a grid management system is needed.

In the new wave of control, we use machine learning as a way to control different systems robustly through the full cycle from production to generations, and that is the reason we will be using its techniques in this project. Reinforcement learning is our technique of choice due to its usage in sequential control.

1.2 MicroGrids

A microgrid is a collection of energy production units and consumers (load) placed near each other to reduce transport and control costs. These grids use Renewable Energy Resources (RESs) as the production units, those include Solar Energy production units (Photovoltaic units PVs), Wind production units (Wind Turbines), Biomass units as well as Energy storage system. Loads supplied by the Microgrid can be critical loads (industrial factories, hospitals, schools) and non-critical loads (houses).

The supply between this type of loads merely depends on the mode of operation the Microgrid operates in; dual-mode operation on-grid where we connect the Microgrid to the primary utility grid in normal conditions. The other mode is off-grid or known as (islanded) mode, where the microgrid switch to be operating without the back up of the utility grid due to shortage or disturbance in the primary grid.

In some switching mode scenarios, the non-critical load is deactivated from the Microgrid until we reach a stable status. Microgrids components from a control unit, inverters, and batteries help determine the Microgrid's supply-demand profile. We connect The different Microgrids to a primary controller alongside the main grid. This controller receives the supply-demand information of each Microgrid and the mode, on-grid, or off-grid.

1.3 Reinforcement Learning

A reinforcement learning algorithm will handle the high-level control between the microgrids. It works as a black box that trades when needed deciding then with which Microgrid and with what price for our microgrids to achieve equilibrium.

Reinforcement Learning is a field where the problem it solves is an environment affected by an algorithm controlled agent. The agent takes an action that affects the environment, then the agent sees the environment as an observation, and receives a reward on the "goodness" of the action it last took.

We use this to optimize its policy, which decides which action to take based on which state we are in, and a state is a unique set of values that the observation returns that fully or partially describe our environment. We use the Markov decision process concept, where this state is sufficient to predict the future. Through time, the actions taken by the agent traverse the state space until we reach our goal state or goal reward. The optimal policy will do this in the best path possible as we will later introduce concepts that push towards a faster goal-reaching policy.

In the context of our project, we will create an environment that simulates several microgrids with full generation and load profiles. We will then create an agent that

will choose to buy or sell electricity when our generation and stored electricity are insufficient for our load. It will also choose the price of the transaction to achieve the maximum economic gain from trading situations.

1.4 Research Objectives

This thesis proposes working in Microgrid islanded mode and introduces a Reinforcement Learning algorithm to manage an energy trading process between neighboring grids to achieve optimality (equilibrium) in production. The research aims to propose a system for off-grid power management and trading between different microgrids using expected power production (supply), consumption (demand) forecasting, and power storage information. We will apply machine learning techniques in specific reinforcement learning and deep reinforcement learning to solve the trading process. We will test multiple algorithms on our created environment and find the best algorithm to solve it.

1.5 Research Questions

Chapter 2

Literature Review

2.1 Microgrids

2.1.1 Theoretical Background

The environmental and economic conditions, the need to provide a clean environment, and decrease the carbon emissions in the atmosphere and the need to decrease fossil fuels made technological advancements a need. Recent technological developments in micro-generation showed that micro-grids are the future of efficient and fast restoration of the power system.

Micro-grids

They can be defined as "A group of interconnected loads and distributed energy resources(DERs)with set electrical boundaries that act as a single controllable entity concerning the grid that can connect and disconnect itself from the grid based on the mode required."

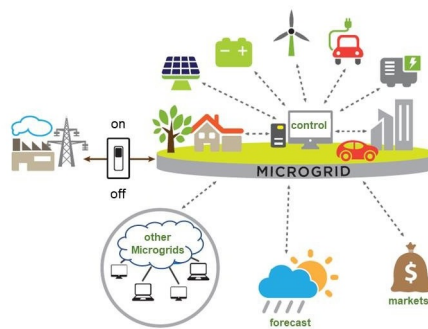


Figure 2.1: A Simple Micro-Grid

The term micro-grid dates back to 1882 when Edison installed 50 DC micro-grid before the operation of the utility grid. With the utilization of utility grid and benefiting from economic and increasing transmission process leading to fade away of micro-grids. Indeed, in the past years, with advancements in power electronics and DER technologies and more engagement with the electricity consumer, the micro-grid concept started seeing the light again.

There would be three different features if we compromised the DER installations could be considered as a micro-grid. There must be a master controller to control the system components as a single controllable entity. The installed generation capacity must exceed the peak critical load; thus, when we can disconnect from the grid and, most importantly, clearly defined electrical boundaries.

The characteristics mentioned earlier present the micro-grid as a small-scale power supply network for a small community; it allows the penetration of distributed generation into the system. One of its major advantages is its ability to work alone during utility grid disturbance or outage; it means that micro-grid can operate in two modes:

1. ON-grid
2. OFF-grid

The on-grid mode is when the microgrid is connected to the primary utility grid and work in synchronization with it. This mode enables bidirectional power flow, and if any disturbance happens to the primary grid, the micro-grid will switch to the off-grid mode or what is known as a standalone grid (islanded). In this mode, the microgrid acts as the primary provider to the specified geographical area, working autonomously with high-quality service by acting as local voltage and frequency regulator [1]. Micro-grid is not a backup generator; a backup generator has been around for quite a while, providing temporary supply to local loads when there is a disturbance in the main utility grid supply. However, micro-grids has a wide range of benefits and noticeably more flexible than a backup generator.

The Micro-grid main components include Loads, DERs, master controller, smart switches, protective devices, communication, control, and automation. The micro-grid load is known to be of two categories; critical and non-critical (fixed and flexible). Critical load (Fixed) must be satisfied at all conditions and is not altered. In contrast, the non-critical load (flexible) can differ and be adjusted based on the economic incentives or the status of the grid (islanded requirements).

DERs consist of distributed generation units(DG) and Energy Storage System (ESS) which can be installed on the utility or consumer premises. The distributed generation units are either dispatchable or non-dispatchable. Dispatchable units can be controlled by the central controller and are subjected to technical constraints depending on the unit type. Non-dispatchable cannot be controlled by the micro-grid controller as its input is changeable, and unrestrained such units are like Solar and wind, mainly renewable sources. The intermittency shows that generation is not always available. Simultaneously, unpredictability reveals that the generation tends to be unstable at different time scales. Those stated characteristics affect our non-dispatchable units negatively and usually increase the forecast error. The right solution is always to reinforce those units with an energy storage system (ESS).

As we know, electricity demand varies based on the time of day and time of year. While in the traditional power system, we are not capable of storing electricity, which leads to a gap between supply and demand. Micro-grid having a mixed power generation will allow as to fill in the mismatch as some generations have significant response times, and others have little flexibility. Some generations can start real quickly to provide more or less depending on demand. Provided the late reasons, the energy storage system is quite beneficial in managing such system .ESS synchronize

with DGs as an assurance to micro-grid generation capability. Its inclusion within the micro-grid system allows the excess energy generated to be stored or in the typical scenario that could be put into the utility grid.

The master controller in the micro-grid performs the scheduling in the microgrid's dual-mode based on economic and security considerations. Usually, the master controller is responsible for interaction with the utility grid, the decision to switch between on-grid and islanded.

With that been said micro-grids benefits are: improving reliability by introducing self-healing at local distribution network, managing local loads due to higher power quality, carbon emission reduction due to diversification usage in renewable energy sources, economically reducing the Transmission and Distribution (T&D) costs [look 2]

2.1.2 Technical challenges of Micro-grids

The integration of DERs units and micro-grid introduces several technical challenges that require addressing the control design and protection system to ensure the level of reliability is not affected. The potential benefits of DG are fully harnessed. Some of these challenges are stability issues arising while at transmission-level, and others are assumptions applied to distribution systems.

The most critical challenges in Protection and control are bidirectional power flow, stability issues, modeling, low inertia, uncertainty. [in 3].

Along with the above, the micro-grid must guarantee the reliable and economical operation of micro-grid while overcoming the challenges above. Henceforth, these are some of the required features in the control system: output control, power balance, DSM, economic dispatch, the transition between mode of operation [see 3].

Furthermore, we can summarize microgrid issues into three points.

1. Islanded mode

This mode represents a future of interconnected grid with a high density of DG. The control strategies of islanding mode are quite essential for the micro-grid to operate in autonomous mode.

We use two kinds of control strategies of islanding to operate the grid. The PQ inverter controls active and reactive power setpoint .furthermore, the VSI control maintains the voltage and frequency feeding the load.

Henceforth the following issues occur within the islanded mode: As beginning as DG supply, the load demand equal sharing is required, but due to various unequal capacities of the DG load sharing tend to be impossible. Along with the harmonics and compensation effort for unbalance and nonlinearity of the load. Secondly, losing a DG in this mode allows the use of load shedding and battery unit to be explored to fulfill the critical load. Finally, guaranteeing Stability in islanded mode is quite challenging with the presence of non-linear load. (An overview on microgrid control strategy).

2. Stability

Stability issues may arise in a micro-grid due to various causes such as islanding the micro-grid and grid reconnection, change in parameters, faults, mismatch

in the generation demand, an immediate connection of DG, or disconnection, and this leads to changes in the voltage and frequency of the system.

Henceforth usage of voltage and frequency controllers or regulators was suggested along with power electronic DGs to give the micro-grid flexibility. Along with ensuring both voltage and frequency are within predefined limit around setpoint values to adjust active and reactive power generated or consumed.

3. Protection

Certain conditions have to be taken into consideration when designing a micro-grid. Its ability to operate under unbalanced conditions such as spacing of overhead transmission and unbalanced impedance from three-phase load any fault within our power system. As the Protection of micro-grid is vital, a new scheme has been introduced that uses ABC-DQ transformation of the system voltage to detect any faults or short circuits. It achieves this by comparing measurements at different locations, thus associating with micro-grid network the faults varieties at different zones.

Unrestrained excess generation results in the voltage profile distortion in an islanded microgrid. Therefore, we should consider the characteristic difference between various DG to develop control strategies to regulate the power output. In cases where active power is not consumed, power oscillations can be used mainly in islanded mode.

2.1.3 Control hierarchy in micro-grids

To understand how the micro-grid is controlled and how it can operate in the two modes, on-grid, and island. Two opposite approaches are identified concerning the architecture of power system control, which is centralized and decentralized.

Centralized control is characterized by having one main central controller responsible for collecting all the required data for decision-making from the various DERs by performing the required calculations and concluding the control actions for each unit at this point.

On the other hand, we have the decentralized control in which we have a local controller for each DERs unit, receiving only local information without being aware of any other system activity.

An interrelated power system is usually characterized by covering large geographical areas. This characteristic means a fully centralized approach is entirely infeasible due to the computation needs and communication needed. Simultaneously, a decentralized approach is not possible either due to its need for a minimum level of coordination and cannot be achieved by using only local variables. Therefore, cooperation between centralized and decentralized control schemes is found in means of a hierarchical control scheme that consists of three control levels: primary, secondary, and tertiary. These control levels vary in their (i) speed of response, (ii) infrastructure requirements. see figure 2.2

1. Primary Control

In local control, it is the first level in our hierarchy featuring the fastest response. It is at the first level; its control is based on local measurements and does not need communication. Given the speed requirements and reliance on

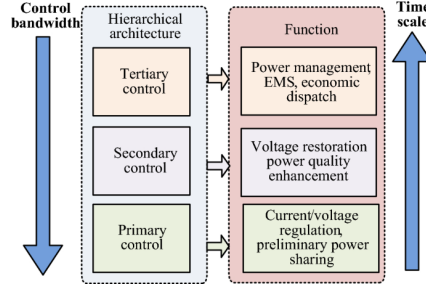


Figure 2.2: Hierarchy diagram

local measurements, islanding detection, inverter output control, and power-sharing balance are all at this level.

- (a) **Inverter Output Control** This control usually contains the outer loop for voltage control and an inner loop for current regulation. Using PI controllers is the typical approach in designing the control loops supported with feed-forward compensation to enhance the current regulator performance; we will look at those control loops further in Chapter 3.
- (b) **Power Sharing Control** A second stage within the primary control level is the power-sharing control concept, which we will cover in two indistinct theories:
 - i. **PQ Control** It is a public control that controls an inverter's voltage output by injecting the active and reactive power in cases the micro-grid cannot give voltage or frequency support .henceforth, and the micro-grid controller is not affected by the unstable voltage and frequency. Usually, when connected to the primary grid, it provides us by the reference frequency, unlike in private mode, it is given by another micro-grid operating on droop control.
 - ii. **Droop Control** The droop method is originally from the power balance of synchronous generators in interrelated power systems. A frequency and voltage deviation occur in our system when there is not inequity between the input mechanical power of the generator and output electrical active power, likely output reactive power. Henceforth in this unit, if we drooped the frequency as a function of active output power, we can then share this power of total load among the various sources. Considering the relationship that dictates power transfer in a two inverter system, droop control applicability is apparent in figure 2.3

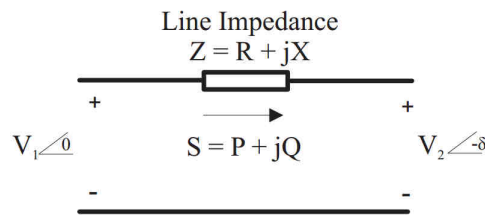


Figure 2.3: A system with two voltage sources

In droop control the relation between real power/frequency and reactive power/voltage can be expressed as:

$$W_0 = W^* - K_p(P_0 - P^*) \quad (2.1)$$

$$V_0 = V^* - Q_p(Q_0 - Q^*) \quad (2.2)$$

Where w^* and V^* are the angular frequency and voltage ,respectively ,and w_0 and V_0 are measured output frequency and voltage of DG system, respectively. The coefficient K_P and K_Q denote the droop coefficients and are determined by the following formulas:

$$K_P = \nabla f / P \quad (2.3)$$

$$K_Q = \nabla V / Q \quad (2.4)$$

$P - w$ droop characteristics are shown in figure 2.4a below while basic $Q - V$ droop characteristics is shown in figure 2.4b

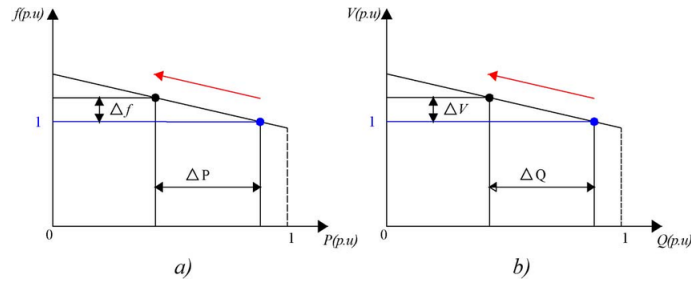


Figure 2.4: $P - w$ and $Q - V$ droop characteristics

Droop control eliminates the need for communication. Its control is based on local measurements, which is outstanding flexibility, in case we are guaranteed a balance between the supply and the demand there is not any need for local controllers. A further illustration will be conducted in Chapter 3

2. Secondary Control

It is known as the Energy Management System (EMS) of the micro-grid, which is in charge of the security and reliability, and economic operation of the micro-grid in its dual-mode. This control level's performance gets more challenging as we switch to isolated mode(islanded) as there are high-variable energy sources, in which the unit dispatch command should be high at a rate enough to keep up with the unexpected changes of load and non-dispatchable DERs.

The EMS works on finding the optimal and unit commitment (UC) and dispatch available DER units; its architecture has two main approaches: centralized and decentralized. With that being said, this level tends to be the highest level of control in the hierarchy for standalone micro-grids.

The centralized approach's architecture contains a central controller that is enriched with the information of every DER and load in the microgrid and network itself as well as forecasting system information. This central controller makes decisions using either online calculation of optimal operation

or databases continuously updated and pre-built with information on proper operation.

Solving energy management related problems while guaranteeing a high level of autonomy for load and DER is one of the decentralized approach benefits. This autonomy is achieved through three levels: Distribution Network Operator (DNO), Microgrid Central Controller(MGCC), and Local Controllers(LC).

DNO controls the communication between the micro-grid and the distribution network and other microgrids, making it part of the tertiary control. MGCC supervise the operation of DERs and load within a micro-grid and in charge of their reliable and economical operation. At the same time, LC control DER units in decentralized architecture, an LC can communicate with MGCC and other LC to share knowledge.

3. Tertiary Control

This control is the highest point in our hierarchical control level, and it works on setting the optimal setpoint based on the power system. It is usually in charge of coordinating multiple micro-grids interacting with one another within the same system and communicating the needs from the primary or host grid.

It works by providing a signal to the secondary level at micro-grid and sub-systems forming the full system. On the contrary, the secondary control coordinates internal primary control leading the primary control to function autonomously and react in predefined ways to identified signals [3].

2.2 Reinforcement Learning

2.2.1 Machine Learning Introduction

Gaining popularity in research and being the centre of technological advancements in all domains, machine learning is a field that focuses on providing data-driven, intelligent answers to research questions. We generally divide it into Supervised Learning, Unsupervised Learning, and Reinforcement Learning.

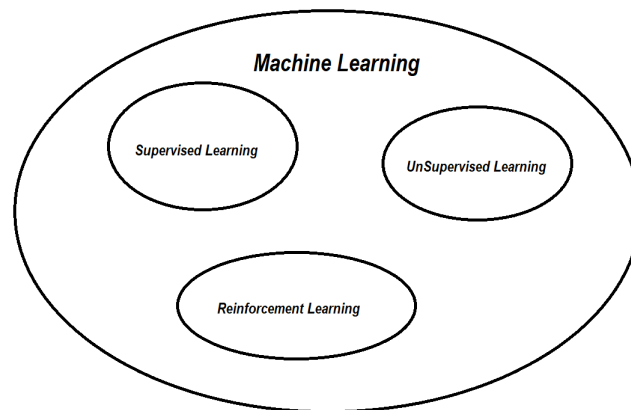


Figure 2.5: Machine Learning Sections

Supervised Learning (SL) is the domain where we try to map inputs to outputs using labeled data. The SL task is either a classification task or a regression task.

In classification, we use a particular object's features to predict its class (i.e., classifying a malignant or benign tumor), meanwhile in regression, where we use features of an object to try to predict another feature of that object, i.e., predicting the price of a house using features of size, location.

Unsupervised Learning (UL) uses unlabeled data to divide it into an unknown set of classes and extract specific structures from the data, ...etc. We divide UL into Parametric UL, where we assume probabilistic distribution of the data based on specific parameters. The mission is to try and get those parameters so that we can predict the future.

The other field of UL is Non-Parametric UL. Here, we make no assumptions about the data, and we only group the data into clusters with resembling features.

Though several problems are solvable using these two methods of learning, the real power and research start when we start talking about the opportunities Artificial Neural Networks ANNs provide. These are trials of modeling the process of thinking that lies inside a human brain.

A neural net is a collection of small processing units called neurons that receive input. It uses its predefined function to calculate a result and then output this result to another neuron or the user. Collecting a group of neurons, we create a neural net that can take multiple inputs and use complex functions to get a single output. We call each level of neurons a layer; we have an input layer, output, hidden layers.

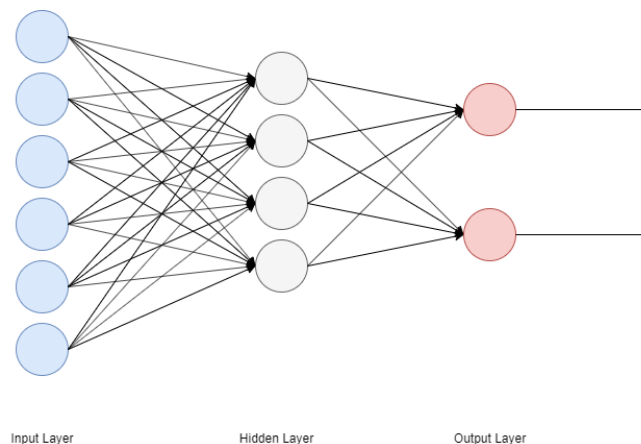


Figure 2.6: Simple Neural Network

Deep Learning DL is an application on ANNs with many levels of hidden layers typically more than three. We use this architecture to solve problems with increasing levels of complexity and individual requirements. The most famous divisions of deep learning are Convolutional Neural Networks CNNs and Recurrent Neural Networks RNNs.

CNNs use filters to tackle problems concerning pictures and images. These problems include object detections, image classification, and much more in-depth domain problems.

The focus of RNNs is time-varying data, i.e., data that require a time-based context to solve. The most important application of this is Natural Language Processing NLP, which needs text context data to understand general speech.

2.2.2 Reinforcement Learning Introduction

Reinforcement Learning RL is the third type of machine learning domain. It is learning what to do so to maximize a digital reward signal [see 4]. The problem is as one where an agent (i.e., a player in a game) is traversing an environment, and he takes actions and collects rewards as he goes.

We can describe the whole problem as an environment *env* that can be described as a state-space S that consists of states s that describe fully the world that can affect or be affected by the agent's decisions. The agent can take action a from an action space where $a \in A$ that will change its state and receive a reward r where $r \in R$.

This Mathematical representation means that an RL problem can be described as a Markov Decision Process MDP. In an MDP, we have the concepts of a reward function $R(s)$, which can map states to rewards achieved when reaching that state.

In an episodic process (one with a clear starting state and a final state), the total reward is the accumulation of each reward received through the journey traversing the environment until reaching the final state. This total reward is the value we are trying to maximize in our RL problem.

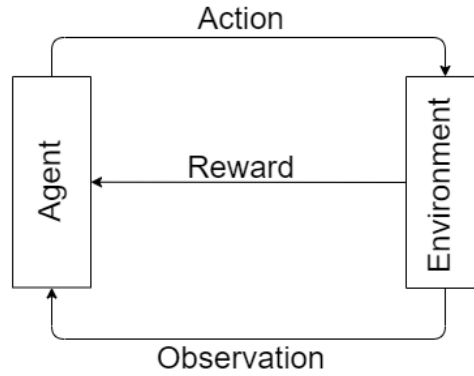


Figure 2.7: Reinforcement Learning Parts

Another concept apparent in MDPs is the concept of a policy π , which defines the path (also known as trajectory τ) that the agent will take during the episode. The policy $\pi(s)$ maps state-action pairs, i.e., what action to take when an agent is in state s . There are two types of policies, either a deterministic policy, one that the agent is told precisely what action to make when arriving at state s where $a = \pi(s)$. The other type of policy is a stochastic policy; here the agent is told in probabilistic values the probability of taking action a when in-state s where

$$\pi(a|s) = P[A_t = a|S_t = s] \quad (2.5)$$

Achieving the maximum reward possible implies that we will take the best possible policy to give the maximum reward at each step. This policy is called the optimal policy π^* , and accordingly, we can say that the main target of RL is to achieve an optimal policy for the environment in which we work. Finding this optimal policy requires a way to measure the optimality or goodness of a certain policy, and this can be done using a Value Function $V(s)$ which is the reward to accumulate over the future starting at current state s [see 4]. This can be expressed as:

$$V_\pi(s) = E_\pi[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s] \quad (2.6)$$

where $V_\pi(s)$ is the set of rewards expected to be accumulated starting at state s and following the policy π . Using the definition of return G_t which is given by:

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{T-1} R_T \quad (2.7)$$

We could re-write $V_\pi(s)$ as:

$$V(s) = E[R_t + \gamma G_t | S_t = s] \quad (2.8)$$

Representing the environment is the model, which is a mapping of what the environment behaviour in response to the agent's action with $P_{ss'}^a$ describing the probability of arriving at state s' when taking action a at state s and R_s^a being the expected reward for arriving at state s taking action a described as follows:

$$P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a] \quad (2.9)$$

and

$$R_s^a = E[R_{t+1} | S_t = s, A_t = a] \quad (2.10)$$

Following, the optimal value function of a state s following a policy π in a number of steps H is given by:

$$V_\pi^* = \max_\pi E\left[\sum_{t=0}^H \gamma^t R_{s_t}^a | S_0 = s\right] \quad (2.11)$$

Moreover, finding the best possible value function of a given policy and then updating the policy to find the best policy is called the policy iteration. It can be given by the algorithm (Value Update or Bellman Update Backup):

Start with $V_0^*(s) = 0$

For $k = 1, 2, \dots, H$:

For all states s in S :

$$V_k^*(s) \leftarrow \max_a \sum_{s'} P_{ss'}^a (R(s, a, s') + \gamma V_{k-1}^*(s')) \quad (2.12)$$

$$\pi_k^*(s) \leftarrow \operatorname{argmax}_a \sum_{s'} P_{ss'}^a (R(s, a, s') + \gamma V_{k-1}^*(s')) \quad (2.13)$$

One alternative to using these dynamic programming methods of value iteration and policy iteration is to introduce a way to learn as the episode is going, using an entire set of visits over states to come up with the value function of each state. These methods are called Monte-Carlo methods and their basic concept is that the value of a state $V(s)$ is calculated using a number of visits to the state in the episode $N(s)$ and total return on that state $S(s)$ by the following algorithm:

To evaluate s

timestep t when s is visited in an episode:

Increment counter $N(s) \leftarrow N(s) + 1$

Increment total return $S(s) \leftarrow S(s) + G_t$

Value is given by $V(s) = S(s)/N(s)$

With $V(s) \rightarrow V_\pi$ as $N(s) \rightarrow \infty$

using this method and generalizing over all episodes:

For all states S_t with return G_t :

$$N(S_t) \leftarrow N(S_t) + 1 \quad (2.14)$$

$$V(S_t) = V(S_t) + (G_t - V(S_t))/N(S_t) \quad (2.15)$$

For non stationary environments:

$$V(S_t) = V(S_t) + \alpha(G_t - V(S_t)) \quad (2.16)$$

2.2.3 Cross-Entropy Methods

One of the most basic solution methods when talking RL problems is cross-entropy methods. Their idea is pretty simple and builds on the intuition gained from looking at the RL main problem.

The idea of gaining as much reward as possible used as we replace all the agent complications with a non-linear trainable NN function, with the input being the observations from the environment s and the output being the policy π . In practice, we represent the policy as a probability distribution over the actions a , making the problem a classification problem. The algorithm describing the method is:

1. Play N number of episodes with initial NN model and environment.
2. Calculate the total reward for each episode and set a reward boundary, usually at the 70th percentile of rewards.
3. Discard all rewards below the boundary.
4. Train the NN model on elite episodes using state s as input and used actions as a target.
5. Repeat 1 until the target mean reward is reached.

2.2.4 Q-Learning

Another value that's important to describe an RL environment is the Q-value $Q(a, s)$ which describes the quality of taking a certain action a when in state s following

policy π and is given by:

$$Q_\pi(s, a) = \sum_{s'} P_{ss'}^a (R_s^a + \gamma Q_\pi(s', a)) \quad (2.17)$$

Finding a solution to this value is inherently a temporal difference problem, which can be defined as a combination between Monte Carlo and Dynamic Programming in which we can learn directly from experience without needing a dynamics model (MC) and we update estimates based in part by other learned without waiting for final estimates. Working from equation 2.16, TD doesn't wait for a whole episode to update $V(S_t)$, it only waits for the next time step $t+1$ to update the value function using both $V(S_{t+1})$ and R_{t+1} updating on the transition to S_{t+1} using the equation:

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)) \quad (2.18)$$

This is called one-step *TD* or *TD(0)*, a special version of the complete *TD*(λ) or *n*-step *TD*. The algorithm for implementing it is:

Loop for each episode:

Initialize S :

Loop for each step of the episode:

$A \leftarrow$ action taken by π for S

Take action A , observe R, S

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

$$S \leftarrow S'$$

Untill S is terminal

Q-learning is an off-policy *TD* algorithm [first introduced 5] which made a major breakthrough in RL in which we use the update 2.19 in the normal *TD*(0) algorithm.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_t, a) - Q(S_t, A_t)) \quad (2.19)$$

2.2.5 Deep Reinforcement Learning

These techniques work on using neural networks to approximate either policy parameters or to get V , Q and A . In value function Deep-RL methods, more specifically *DQNs* [introduced in 6] we can look at the problem as one of regression, using a deep NN to approximate the value of Q function, updating our Q value using the tabular Q learning update:

$$Q(S, A) \leftarrow (1 - \alpha)(Q(S, A) + \alpha(r + \gamma \max_{a' \in A} Q(S', A'))) \quad (2.20)$$

When using this update in practice, we face some problems that limit the usability of the method. The problem of acting randomly or using the Q approximation is solved using an epsilon-greedy method. We start with a completely random action and decaying the randomness with probability ϵ till a small value of 2% randomness.

Another problem is the requirement of SGD on training data, as it requires independent and identically distributed data *i.i.d*; this is not satisfied. We solve this problem by using a large replay buffer, adding new experiments, and pushing old ones out, allowing for independent data that are fresh.

The last problem is the similarity between states after each other, not allowing our NN to distinguish between them. We will solve this by using a target network, keeping a copy of the NN parameters using it for $Q(s', a')$. This network is updated periodically, making the network stable.

Our algorithm for DQN is:

1. Initialize $Q(s, a), Q'(s, a)$ with random weights, $\epsilon = 1.0$, and an empty replay buffer
2. Choose random action a with probability ϵ ;otherwise, $a = \operatorname{argmax}_a Q(s, a)$
3. Execute action a , observe next state s' and reward r
4. Store transition (s, a, r, s') in replay buffer
5. Sample a random mini-batch of transitions from the replay buffer
6. For each entry in the buffer (*Transition*), calculate target $y = r$ if the episode has ended at this step and $y = r + \gamma \max_{a' \in A} \hat{Q}(s', a')$ otherwise
7. Calculate loss $L = (Q(s, a) - r)^2$
8. Update $Q(s, a)$ using SGD
9. Every N steps copy parameters from Q to \hat{Q}
10. Repeat 2 till convergence

2.2.6 Policy Gradients Methods

These can learn a parametrized policy π_θ that can select an action without returning to a value function. It can be used to learn policy parameters but not to select exact actions. Updating the policy parameters can be achieved either using gradient-based or gradient-free methods to maximize expected return. Policy gradients update the policy parameter on each step in the direction of an estimate of the gradient of performance compared to the policy parameter. Given the trajectory τ , which is the set of state action reward sequences, we can define a policy parameter performance:

$$J(\theta) = E[R(\tau)] \quad (2.21)$$

It's obvious that we need to maximize J in order to find the optimal θ^* . Using gradient descent on this problem (the most basic machine learning idea) we find that:

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t) \quad (2.22)$$

The far right of the equation can be approximated to:

$$\nabla E_{\pi_\theta}[R(\tau)] = E_{\pi_\theta}[R(\tau) \nabla \log \pi_\theta(\tau)] \quad (2.23)$$

Using an expansion for $\pi_\theta(\tau)$ we can arrive at:

$$\nabla E_{\pi_\theta}[R(\tau)] = E_{\pi_\theta}[(G_t - b) \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right)] \quad (2.24)$$

The equation tells us that to successfully update the parameters, and we do not need a model, we can only use state-action rewards to make a useful update. We can solve the expectation in the equation by sampling a large number of trajectories while replacing $R(\tau)$ with G_t removes the variance that appears in the standard equation. The baseline b is used to reduce the estimate's bias, b must be independent of the parameters, and a good baseline makes use of the state-value current state. Equation 2.24, without adding a baseline, is called *REINFORCE*, which is the classic policy gradient algorithm.

Actor-critic methods combine policy gradients with model fitting, and we use an actor to model the policy and a critic to model the value function V . By introduce a critic, we reduce the number of samples to collect for each policy update; we do not collect all samples until the end of an episode. We will talk about these specific solution methods and more state-of-the-art solutions in the next chapter to explain the methodology of our work.

2.3 Related Work

We mentioned earlier that the electricity demand is increasing and that it should balance with the constant need to limit global warming, that is why we introduced microgrids. The unidirectional design of national power grids did not allow for surplus microgrid generated power to be redistributed and sold to the primary national grid, so the process of energy trading was devised to accommodate for this problem in low voltage networks between microgrids.

In some cases, the microgrid demand is satisfied through the traditional utility grid when the connection is on-grid. Nevertheless, that means we will be paying an extra cost for them to overcome that shortage. Microgrids stand for economic optimality, which leads us to have efficient backups at reasonable and profitable cost and solve some of the technical problems we mentioned. These problems include better power quality, reduced voltage fluctuations, and a reliable system that is not affected by the utility grid outage through energy trading.

Understanding the problems facing microgrids takes us to the core of our research, the "energy transition between microgrids." The term energy trading is understood as the importing and exporting of energy in a market of retailers, producers, and vendors, including the large industrial consumers in the utility grid. It was

then reinterpreted as the local energy trading that happens amongst users in the microgrid.

Energy trading can be approached from different perspectives. We can refer to it as an optimization problem where we can look at it as we look to the microgrid control either through centralized or decentralized approaches. In centralized, we have one central controller responsible for solving this problem, where it looks into a way to minimize the generation and the transportation cost of the microgrid. In contrast, we have a decentralized approach that looks into studying all the participants and the benefits that fall for each inclination.

In energy trading, what affects one affects all that is any action performed by any participant in the market will affect all. That is why we got introduced to the Game Theory (GT) technique in energy trading that in itself has a different concept to cover or try to solve the energy trading problem. GT acts in competitive situations where it takes into account that one's strategy will affect all the other strategies. We mostly use it on the decentralized structure of our microgrid, which makes it easier to check for each one's behavior. GT games are classified into Direct and Indirect games where one aims to find the ideal policy, and the other is concerned with planning a game that satisfies particular objectives. The former does not have any effect on energy trading so far. We will focus on direct games such as non-cooperative games in which each individual is looking for their benefit.

The work of Pilz [7] covers many studies regarding GT energy trading with its different properties. We find one study taking energy storage with the background of the schedule and reducing the peak to the average ratio where they plan a cost function under certain conditions that lead the system to be balanced. We then assume that each consumer's total load is the sum of the external power delivered from the primary grid. After setting up, he then proposed two different approaches, one is a static non-cooperative game where the utility sets a cost function, and the player plays a scheduled game in order for him to minimize the respective cost. Here we have a user with the advantage of selling energy back to the utility grid, known as a reverse peak. A second game was introduced that takes the utility grid as part of the game (participant) and adjusts the prices and schedules the trade a typical leader-follower structure defined by the "Stackelberg game." which proves that Stackelberg equilibrium is equivalent to minimizing the peak-to-average ratio.

Another study takes a look at situations where the traditional power station could not meet the high demand at some point, so it buys the needed energy from energy consumers (electric vehicles, renewable energy farms, and any participants involved with the central power station as an individual). The researchers proposed a non-cooperative Stackelberg game where we do not deal with each component alone. Instead, they operated a solution that serves the social benefit assuring that each component benefits from participating in energy trading. They introduced a price model where the price can differ for different energy consumers. Henceforth, the authors applied an iterative algorithm to minimize the cost for the central power station and, at the same time, maximize the sum of utility functions of energy consumers.

Another research covered the transition of energy among the MGs. The trade does not happen directly with each other but instead tries to trade surplus energy with the market and request the deficiency as well. This multileader-multifollower Stackelberg game proposed, the sellers act as leaders and the buyers as followers in

which the surplus energy is proposed by the leaders to the followers proportionally to the bids each buyer has placed. This method leads us to know that the best solution for this scenario depends on bids given and the number of players in the game. Because of the expanding rivalry between the purchasers, the worth monotonically diminishes when the quantity of purchasers increases. Simultaneously, the aggregate of the utility qualities for the dealers' increases, since more costumers permit them to sell more.

Another Stackelberg seller-buyer structure among MGs was taken into consideration as the ones before. However, to make the model more expressive, the author encompasses the known structure to the Bayesian game. In this type of game, our knowledge is incomplete, and we do not have full awareness of the game aspects and players' states, meaning that each player is private about their information. In this case, we take the players as normal or abnormal, the emergency state in which the sellers are less profound to sell energy and value the stored energy. From the buyers' point of view, they tend to bid more to ensure the requested energy delivery. We build on the last study by proposing a communication link between respective MGs where a weighting variable is used to express the relation between them. Precisely the conditional probability distribution over the condition of the player is classified as a two-stage technique. In stage one, each Mg estimates the state based on the players' given messages; the second stage updates the estimates based on information gathered from the close neighbors in the structure, looking into increasing the trust within the network showing and the partial trusted information. In the end, a debate is held questioning whether this will increase the power quality but was left for further work.

Unlike the late researches here, the central unit does not only communicate, but it works as a distributor or gatherer for the energy that is traded among the MGs. Also, no scheduling scheme is proposed to pay the sellers. By providing energy to the system, the respective MG collects points that increase the contribution value. If this MG runs into a deficit of energy, its high contribution value will give it a more significant chunk of energy given by the rest. The distributor sets that in order to maximize the social welfare function. Knowing the distribution mechanism, the game here deals with the remark of how much energy to request directly proportional to this and inversely proportional to the contribution value it gets. Furthermore, each buyer is given a stage in the queue in which h should try to be served earlier to minimize what is requested from the utility grid. In this case, not enough surplus energy to serve, we have nash equilibrium property that even if participants deviate from it, the other does not be impacted negatively.

For security reasons, all correspondences are composed of the central unit. Seen from any of the MGs, this prompts a fragmented data game, as no one thinks about the systems and settlements of the others. All the more specifically, the author divides the MGs into merchants and purchasers and plans a two-phase Stackelberg game in which every one of these gatherings attempts to find their best activities by methods or reinforcement learning algorithms. The same classification and giving principle based on proportionality are added; this implies there are two utility capacities, one for each group of buyers and sellers without the knowledge of the other players. It shows that the learning algorithm here converges to the best reply, which is the same as the solution to the sellers' and buyers' optimization, respectively. In comparison, the iteration solutions earlier this take 100 times more to converge to

Nash Equilibrium.

If we looked at the energy exchange proficiencies, we find some focused on selling the energy back to the conventional grid. At the same time, others took into account two types of participants, the sellers and the buyers. In most cases, the energy transition happens in secondary structure as it does not happen between individuals but happens through an operator, a third party that leads us to not fully decentralized scenario.

On the other hand, most of the utility functions taken by the games are focused on the monetary function perspectives from looking into the cost of storing the energy to the cost of transmission of energy between different parties. On the other hand, the utility function did not look into the price function but instead looked into the ratio between allocated energy and requested energy. Other approaches acted upon an auction algorithm where the buyers and the sellers are balanced in the market.

In all scenarios, the customers were referred to as sellers or buyers despite having surplus energy or deficit. In the above models, there was a shortage in models that combine a high-quality demand analysis with the RE generation in energy trading. Most of this researcher proposed what they call blue-sky approaches with "reinforcement learning" and "contribution-based" energy trading. Furthermore, all those authors lacked in the long-term assessable suggestions opposing the merely one-day ahead analyses in energy trading. We use a reinforcement learning algorithm that works on solving the situation without prior information about the microgrid. As achieving Supply-Demand equilibrium is complicated when considering the non-formality of the RES's and many studies were proposed in market-based energy trading among microgrids to utilize DERs across the network fully [8]

The idea of microgrids replacing conventional power grids in rural areas has been the subject of research. B. M. Sivapriya et al. [9] worked with the problem of microgrid design using the center of moment approach to the placement of PV panels on the network providing case studies for their designs on villages in India. Murenzi et al. [10] worked in Africa, introducing Microgrids as a viable method to electrify sub-Saharan Africa. They showed that in a typical Rwandan village, the installation of a microgrid with PV, batteries, and a micro-hydro is a better financial alternative than extending the national power grid transmission to reach the village.

Applications of Reinforcement Learning in smart grids and microgrids vary, A smart building energy management algorithm [11] that uses a Markov decision process to model the smart building. The algorithm controlled included interactions with the utility grid and internal RES. The algorithm used Q-Learning to make decisions on energy dispatch actions achieved better energy costs in the building against multiple pricing policies. Mocanu et al. [12] created a deep belief network that improved the performance of standard reinforcement learning algorithms. They namely worked on SARSA and Q-learning, in the context of predicting energy in a smart building, the algorithm can generalize a learned behavior model into any other building without any specific history of that building. Leo Raju et al.[13] proposed a model-free reinforcement learning algorithm (Q-learning) to solve the optimal dispatch problem, which concerns finding the best combination of available power resources to provide the required load with minimal cost. Their algorithm converged to the optimal solution and provided adaptability in dynamic situations and unforeseen load management.

Fabrice et al.[14] proposed an algorithm to fully control power flow between a multi-storage Microgrid mapping it as a Multi-Agent System (MAS) and using Multi-Agent Reinforcement Learning to solve the problem. They produced results showing that a centralized control; unit for the microgrid is not needed. The algorithm can achieve the minimal cost of drawing power from the primary grid and achieve most grid independence. Finally, Xiao et al.[15] proposed an energy trading game between different microgrids intending to achieve the Nash Equilibrium without knowing the generation and load demand of the other microgrids using a DQN-based energy trading strategy achieving an improvement of 22.3% in the utility of the microgrid.

Bibliography

- [1] Mushtaq N Ahmed et al. “An overview on microgrid control strategies”. In: *International Journal of Engineering and Advanced Technology (IJEAT)* 4.5 (2015), pp. 93–98.
- [2] Sina Parhizi et al. “State of the art in research on microgrids: A review”. In: *Ieee Access* 3 (2015), pp. 890–925.
- [3] Daniel E Olivares et al. “Trends in microgrid control”. In: *IEEE Transactions on smart grid* 5.4 (2014), pp. 1905–1919.
- [4] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [5] Christopher JCH Watkins and Peter Dayan. “Q-learning”. In: *Machine learning* 8.3-4 (1992), pp. 279–292.
- [6] Volodymyr Mnih et al. “Playing atari with deep reinforcement learning”. In: *arXiv preprint arXiv:1312.5602* (2013).
- [7] Matthias Pilz and Luluwah Al-Fagih. “Recent advances in local energy trading in the smart grid based on game-theoretic approaches”. In: *IEEE Transactions on Smart Grid* 10.2 (2017), pp. 1363–1371.
- [8] Sheikh Muhammad Ali. “Electricity trading among microgrids”. In: *Department of Mechanical Engineering, University of Strathclyde* (2009).
- [9] Sivapriya Mothilal Bhagavathy and Gobind Pillai. “PV Microgrid Design for Rural Electrification”. In: *Designs* 2.3 (2018), p. 33.
- [10] Jean Pierre Murenzi and Taha Selim Ustun. “The case for microgrids in electrifying Sub-Saharan Africa”. In: *IREC2015 The Sixth International Renewable Energy Congress*. IEEE. 2015, pp. 1–6.
- [11] Sunyong Kim and Hyuk Lim. “Reinforcement learning based energy management algorithm for smart energy buildings”. In: *Energies* 11.8 (2018), p. 2010.
- [12] Elena Mocanu et al. “Unsupervised energy prediction in a Smart Grid context using reinforcement cross-building transfer learning”. In: *Energy and Buildings* 116 (2016), pp. 646–655.
- [13] Leo Raju et al. “Reinforcement learning in adaptive control of power system generation”. In: *Procedia Computer Science* 46 (2015), pp. 202–209.
- [14] Fabrice Lauri et al. “Managing power flows in microgrids using multi-agent reinforcement learning”. In: *Agent Technologies in Energy Systems (ATES)* (2013).
- [15] Liang Xiao et al. “Reinforcement learning-based energy trading for microgrids”. In: *arXiv preprint arXiv:1801.06285* (2018).