

MASTER THESIS

Thesis submitted in partial fulfillment of the requirements
for the degree of Master of Science in Engineering at the
University of Applied Sciences Technikum Wien
Degree Program Renewable Urban Energysystems

Multi-objective optimization of building control to minimize operational cost with a machine learning approach

Andreas Rippl, BSc
Student Number: 1810578009

Supervisor 1: Simon Schneider, MSc
Supervisor 2: Christoph Gehbauer, MSc

Berkeley, December 04, 2020

Declaration of Authenticity

"As author and creator of this work to hand, I confirm with my signature knowledge of the relevant copyright regulations governed by higher education acts (see Urheberrechtsgesetz/ Austrian copyright law as amended as well as the Statute on Studies Act Provisions / Examination Regulations of the UAS Technikum Wien as amended).

I hereby declare that I completed the present work independently and that any ideas, whether written by others or by myself, have been fully sourced and referenced. I am aware of any consequences I may face on the part of the degree program director if there should be evidence of missing autonomy and independence or evidence of any intent to fraudulently achieve a pass mark for this work (see Statute on Studies Act Provisions / Examination Regulations of the UAS Technikum Wien as amended).

I further declare that up to this date I have not published the work to hand, nor have I presented it to another examination board in the same or similar form. I affirm that the version submitted matches the version in the upload tool."

Place, Date

Signature

Abstract

The transition to more renewable energy sources, will become an ever-greater challenge in the future. Herein, one critical issue, is the natural volatile generation of sun and wind powered power stations, which strain the power grid during peak generation times. Buildings as flexible consumers can help to relieve the stress on the power grids without having to significantly increase their capacity. Current control systems are often realized with proportional-integral-derivative (PID) controllers, which work with a feedback, but do not react to predictions. Therefore, Model-predictive-controllers (MPC) represent an innovation in control, as they can, adapt the control strategy to the needs of the grid by means of weather forecasts and predictions on the occupancy of persons. Nevertheless, their application requires a higher accuracy regarding models of the buildings.

In the thesis, an alternative approach was chosen. A so-called agent with reinforcement learning (RL) is trained to learn the necessary rules to control the room temperature in an office building. RL refers to the fact that the agent improves itself with the help of experience as it takes over the control. The goal of this thesis is to develop an agent that controls the heating and cooling system in a room of a new office building in Berkeley, California and the tint of the electrochromic window used for shading. The room is represented as a resistance and capacitance (RC) model and is controlled by the agent in a way that reduces the total energy costs for heating, cooling, artificial lighting and office equipment compared to PID controller and MPC. Based on past studies of RL-algorithms the Deep Deterministic Policy Gradient (DDPG) is chosen as the algorithm for the agent. The agent interacts with the RC-model during a training process, where the agent learns how to operate with the possible actions to gain the highest reward based on the total energy costs. The goal of the agent is to maximize the reward over all possible timesteps.

The agent uses a forecast including the weather forecast, electricity tariff information and information about occupancy. By combining a DDPG-algorithm in combination with a Multi-Layer Perceptron (MLP) network the agent succeeds in its primary task but is not farsighted enough to lower the maximum demand, what leads to high demand costs. The further improvement with four hours of forecast data as inputs and a reward system based on multiple steps lead to a behavior where the agent precools and preheats the room to lower the demand costs.

The final trained agent uses a DDPG algorithm in combination with a MLP network to achieve the best results and control the room with lower energy costs compared to the PI controller. The "perfect information" model as MPC, developed in this thesis optimizes the room over the entire period and minimizes the energy costs compared to the PI controller by % to the disadvantage of energy consumption. This lies around % higher. The agent

manages to reduce energy costs by % and only consumes % more than the PI controller and thus % less than the MPC.

Kurzfassung

Der Übergang zu mehr erneuerbaren Energiequellen wird in Zukunft zu einer immer größeren Herausforderung werden, da die natürliche volatile Erzeugung von sonnen- und windbetriebenen Kraftwerken das Stromnetz in Spitzenerzeugungszeiten belastet. Gebäude als flexible Verbraucher können dazu beitragen, die Stromnetze soweit zu entlasten, dass die Netzkapazitäten nicht wesentlich erhöht werden müssen. Regelsysteme in Gebäuden werden derzeit häufig mit Proportional-Integral-Derivativ (PID) Reglern realisiert, die zwar mit einer Rückkopplung arbeiten, aber nicht auf Wettervorhersagen reagieren. Modell-Prädiktive-Regler (MPC) stellen eine Regelungsinnovation dar, die z.B. durch Wettervorhersagen und Vorhersagen über die Personenbelegung die Regelstrategie an die Bedürfnisse des Stromnetzes anpassen können, aber stark von genauen Modellen der Gebäude abhängig sind.

In der Masterarbeit wird ein sogenannter agent mit Reinforcement Learning (RL) entwickelt, der die notwendigen Regeln zur Steuerung der Heizung und Kühlung sowie der Beschattung von Gebäuden lernt. RL bezieht sich auf die Tatsache, dass der agent sich mit Hilfe von eigener Erfahrung verbessert, während er die Regelung übernimmt. Ziel dieser Arbeit ist es, einen agenten zu entwickeln, der das Heiz- und Kühlsystem in einem Raum eins neuen Bürogebäudes in Berkeley, Kalifornien und die Tönung des elektrochromen Fensters, das für die Beschattung verwendet wird, regelt. Der Raum wird als Widerstands- und Kapazitätsmodell (RC) dargestellt und wird vom agenten in einer Form geregelt um die Gesamtenergiekosten für Heizung, Kühlung, Kunstlicht und Bürogeräte im Vergleich zu PID Reglern und MPC zu reduzieren. Basierend auf früheren Studien von RL-Algorithmen wird der Deep Deterministic Policy Gradient (DDPG) als Algorithmus für den agenten gewählt. Der agent interagiert mit dem RC-Modell während des Trainingsprozesses, und lernt, wie er mit den möglichen Aktionen die höchste Belohnung auf der Grundlage der Gesamtenergiekosten erhält. Das Ziel des agenten ist es, die Belohnung über alle möglichen Zeitschritte zu maximieren.

Mit der bereitgestellten Vorhersage mit Wettervorhersage, Stromtarifinformationen und Informationen über die Belegung des grundlegenden DDPG-Algorithmus in Kombination mit einem Multi-Layer Perceptron (MLP)-Netzwerk gelingt dem agenten die Aufgabe, ist aber nicht weitsichtig genug, um die maximale Nachfrage zu senken, was zu hohen Nachfragekosten führt. Die weitere Verbesserung mit vier Stunden Vorhersagedaten als Input und ein auf mehreren Schritten basierendes Belohnungssystem führen zu einem Verhalten, bei dem der agent den Raum vorköhlt und vorheizt, um die Nachfragekosten zu senken.

Der fertig ausgebildete agent verwendet einen DDPG-Algorithmus in Kombination mit einem MLP-Netzwerk, um die besten Ergebnisse zu erzielen und den Raum mit geringeren

Energiekosten im Vergleich zum PI-Regler zu steuern. Das in dieser Arbeit entwickelte "Perfect Information"-Modell als MPC optimiert den Raum über den gesamten Zeitraum und minimiert die Energiekosten im Vergleich zum PI-Regler um 70 % zum Nachteil des Energieverbrauchs. Dies liegt um 30 % höher. Dem agenten gelingt es, die Energiekosten um zu senken. % zu senken und verbraucht nur % mehr als der PI-Regler und damit % weniger als der MPC.

Danksagung

Den größten Dank verdient mein Betreuer Christoph Gebauer, der mir die Chance ermöglicht hat, meine Arbeit in Berkeley Kalifornien zu schreiben. Bereits im Herbst 2019 hat mich Christoph dabei unterstützt alle notwendigen Anträge fristgerecht fertig zu stellen. Ohne die Unterstützung bereits vor Beginn meiner Ankunft in den U.S.A wäre mir der Start in das neue Thema nicht so gut gelungen. Auch das Programmieren wurde mir durch die laufende Hilfestellung in den letzten neun Monaten enorm erleichtert.

Danke Christoph!

Mein FH-Betreuer Simon Schneider, hat sich zu Jahresende 2019 sehr spontan und kurzfristig dazu bereit erklärt mein Betreuer zu sein. Besonders in der letzten stressigen Zeit war die Unterstützung bei der Fertigstellung der Masterarbeit sehr wichtig.

Danke Simon!

Da die Pandemie auch Berkeley und das Lawrence Berkley National Laboratory lahmgelegt hat, bevor ich richtig starten konnte war besonders der Beginn sehr anstrengend. Besondere Unterstützung habe ich dabei von Ellen Thomas erhalten, die sich um meine Ergonomie am Heimarbeitsplatz gekümmert hat und auch im Austausch mit Campingtipps immer gute Ratschläge parat hatte. Die Meetings der Windows & Daylighting Group unter der Leitung von Christian Kohler waren mit den „Ice-breakern“ doch sehr erfrischend um dem Home-office zu entkommen.

Danke Ellen, Christian und das gesamte Windows & Daylighting Team!

Meinen beiden Mitbewohnern, Ulirke und Christian in der „Haste Mansion“ danke ich für die Unterstützung als professionelle Korrekturleser und Ablenkung während des gesamten Jahres. Ich habe sehr viel über Chemie lernen dürfen → Das reicht mir jetzt aber auch!

Danke Ulli und Christian

Meinen Eltern danke ich für die Unterstützung während meines gesamten Studiums in Wien und meiner Ausbildung zuvor. Ohne die Zeit mich um meine Bildung zu kümmern hätte ich nicht das alles nicht erreicht.

Danke Mama und Papa!

Ich bin mir sicher, dass auch die Kerzen und Gebete meiner Oma dazu beigetragen haben die Arbeit erfolgreich zu beenden.

Danke Oma!

Acknowledgements

My greatest gratitude goes to my supervisor Christoph Gehbauer, who has given me the opportunity to write my master thesis in Berkeley, California. Already in fall 2019 Christoph supported me in completing all necessary application documents in time. Without his support even before my arrival in the U.S.A. I would not have had such a good start into the new topic. The ongoing support throughout the last nine months especially for programming made the work much easier.

Thanks Christoph!

At the end of 2019, Simon Schneider spontaneously agreed to be my supervisor at short notice. Especially during the last stressful time, the support in completing the master thesis was very important for me.

Thanks Simon!

Since the pandemic also shut down Berkeley and the Lawrence Berkley National Laboratory before I could really start. I received great support from Ellen Thomas, who took care of my ergonomics in my home-office and always had good camping tips at hand. The meetings of the Windows & Daylighting Group under the leadership of Christian Kohler were very refreshing with the "Ice-breakers" to escape the home office.

Thanks Ellen, Christian and the whole Windows & Daylighting Goup!

I special thanks to my two housemates, Ulrike, and Christian in the "Haste Mansion" for their support as professional readers and friends throughout the year. I have been able to learn a lot about chemistry, but that is enough for me now!

Thanks, Ulli and Christian

I would like to thank my parents for their support during my entire studies in Vienna and my studies before. Without the time I got to take care of my education I would not have achieved all this.

Thanks, mom and dad!

I am sure that also the candles and prayers of my grandmother helped to finish the work successfully.

Thanks grandma!

Table of Contents

1	Introduction	11
1.1	Motivation	11
1.2	Aim of the Thesis and Scientific Question.....	12
1.3	Approach	12
2	Methodology	13
2.1	Programming Language – Python 3.8.....	13
2.1.1	Machine Learning Framework.....	14
3	Machine Learning	15
3.1	Applications	17
3.2	Learning Techniques	17
3.3	Reinforcement Learning.....	19
3.4	Reinforcement Learning in Building Technologies	25
3.5	Neural Networks	25
3.5.1	Multi-Layer Perceptron.....	27
3.5.2	Recurrent Neural Network	28
3.5.3	Network features.....	30
4	Results	32
4.1	Deep Deterministic Policy Gradient (DDPG)	33
4.2	Replay Buffer.....	37
4.3	Noise	40
4.4	State of the art Controller.....	41
4.4.1	PID Control	41
4.4.2	Model Predictive Control.....	43
4.5	Room Model	44
4.5.1	Electrochromic Window	46
4.5.2	Solar Position and Radiation.....	48
4.5.3	Electricity Market in California.....	53
4.6	RL-Setup	55
4.6.1	Environment	55
4.6.2	Development	57
5	Discussion and Outlook	74

Bibliography	75
List of Figures	80
List of Tables.....	83
List of Abbreviations.....	84
Appendix A: Setting.....	86
Appendix B: Execution	91
Anhang C: RL-Setup	95

1 Introduction

Building envelopes play a crucial role in the energy performance of buildings, imposing an annual 21.3 quadrillion Btu (6,242.41 TWh) primary energy in the U.S. in 2019, which represents 28 % of the total primary energy consumption (U.S. Energy Information Administration 2020). The initiative of the U.S. Department of Energy launched an initiative for Grid-interactive Buildings whose aim is to optimize the interplay of energy efficiency, demand response, behind-the-meter generation and energy storage to enable more demand-side management possibilities. State-of-the-art control systems, such as Proportional Integral Derivative (PID) controls use conventional feedback and are rule-based. Specifically, they are reactionary (cannot consider future operation) and largely univariate (only consider a single variable) and often fail to deliver sustained performance over the time of the installation (Wang and Hong 2020). These controllers cannot consider future climatic conditions like predicted hot outside air temperatures and only react to the outside conditions, which leads to high peak loads for heating and cooling.

Model-predictive controls (MPC) on the other hand can take the future outside conditions into account and have proven the potential to save energy in simulations, as well as in real life buildings. The disadvantage of the MPC is the fact that as the name implicates a detailed model of the building must be programmed. Therefore, the development and calibration are cost intensive as every building is unique. That is the main reason for the limited application of, predictive control in real buildings.

Current investigations at the Lawrence Berkeley National Laboratory (LBNL) in California induce the application of machine learning (ML) in a building life cycle and show that ML is applicable in many stages of this life cycle (Hong et al. 2020). These studies already demonstrate the potential of ML to benefit the performance of the buildings. ML and the field of reinforcement learning (RL) is especially suited for the desired control strategies and can help to eliminate the developing and calibrating of detailed building models as known from MPC.

1.1 Motivation

The rising requirements for energy management, occupant interactions, on-site renewable generation, on-site storage, electric grid interfacing, etc., demand innovative control methods to integrate multiple subsystems. Furthermore, it becomes necessary to address the number of high-performance objectives, such as minimizing the use of energy, energy cost, increasing the demand response capacity, while satisfying the occupant comfort. Therefore, the control methods need to be responsive to real-time and forecasted conditions, consider the interaction of multiple subsystems, require minimal to no set-up and commissioning, and have to be adaptable over the life of the installation.

First studies already indicate the high potential for energy savings, the current challenge hereby is the implementation in buildings to enable more electric loads and distributed Energy systems without reinforcement of the power grid.

Here, the latest study of the LBNL focusing on model predictive controllers (MPC) showed that a total energy saving of 28% is possible compared to state-of-the-art heuristic controllers (Gehbauer et al. 2020). The complexity of the building model necessary for the development of the controllers must be decreased to enable more buildings to have advanced building controls in order to path the way for renewable energy systems.

1.2 Aim of the Thesis and Scientific Question

The LBNL investigates the potential of MPC in an environment, where the shading system, and the heating ventilating air conditioning (HVAC) system is controlled. Herein, the constraints in form of occupancy comfort (e.g. indoor temperature control) and cost savings have to be considered. Therefore, the aim of this thesis is the implementation of an agent that aims to minimize the total energy costs and the peak electricity load, while ensuring the comfort parameters for the occupants.

Within this framework, the following question needs to be answered to improve existing control strategies:

- Which Reinforcement Learning (RL) methodology is best suited for the control of building technology to further reduce total energy costs compared to state-of-the-art controllers and MPC controllers?

1.3 Approach

At the beginning of the work, a fundamental understanding of state-of-the-art ML approaches, and of RL in particular, needs to be gained. RL is a powerful deep learning (DL) technique in the field of artificial intelligence (AI). The most renowned successes of DL were achieved in the video game and board game sector. For example, an agent trained with RL defeated the world champion in the game “Go” which was considered to be impossible due to the complexity of the game (DeepMind 2016). Based on the gained knowledge the RL agent shall be developed in open-source based programming language Python which supports modules and packages that make it suitable for ML applications.

The development of the agent will be performed entirely in Python with the ML framework TensorFlow. The RL agent has to regulate the heating and cooling of the building, as well as the control of the dynamic façade as a shading device. The costs for electricity will be compared with a MPC system programmed in python with the module pyomo. The

comparison will be performed in the context of a California using the electricity tariff for medium office buildings of Pacific Gas and Electricity (E-19) as a time-of-use (TOU) tariff.

2 Methodology

A plethora of research and development has already been conducted in the field of ML. In the following chapter, the relevant methods employed to answer the scientific question posed in this thesis are summarized and explained.

Some of the tools used to develop the agent are prescribed by the LBNL to enable the communication with existing programs and environments. All used tools and programs used are freeware to ensure reproducibility.

2.1 Programming Language – Python 3.8

Python is an open-source programming language which is administrated by the non-profit corporation Python Software Foundation (Python Software Foundation 2020). The language use is widely spread amongst industry due to its flexibility. It is used for web-development, scientific and numeric or software development. Python is an interpreted, object-oriented high-level programming language with a clear syntax.

Figure 1 shows the popularity of programming languages based on raw data based on Google Trends. The numbers show the share of how often a programming tutorial for the corresponding language has been searched in 2020. By a large margin, python is the most popular programming language with a share of more than 30 % of searches on Google.

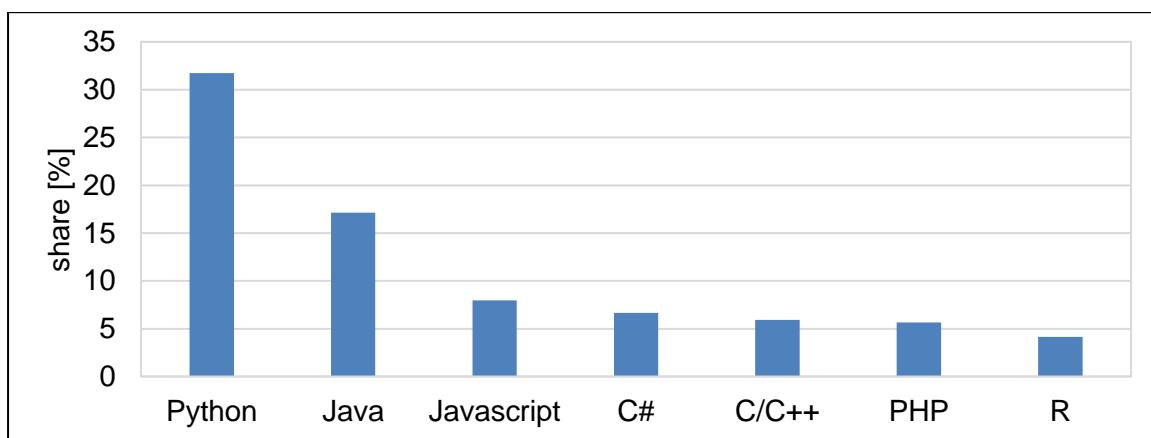


Figure 1: Worldwide PYPL Popularity of Programming Language in 2020 (modified according to (Pierre 2020))

2.1.1 Machine Learning Framework

The variety of ML frameworks was studied by Jeff Hale who described the popularity of different ML frameworks with a power ranking based on online Job Listings, Google Search Volume, Medium Articles, ArXiv Articles, GitHub Activity and others (Figure 2) (Hale 2019). With applied weights, Tensorflow is the most popular framework for machine learning followed by Keras and Pytorch.

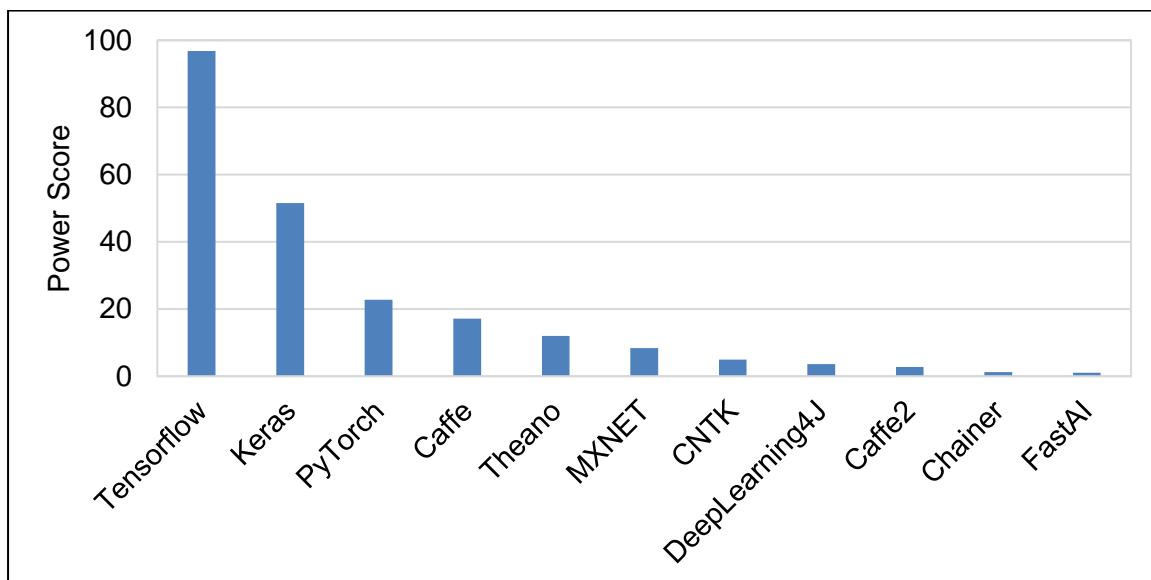


Figure 2: ML Framework Power Scores 2018 (modified according to (Hale 2019))

The further development of deep learning frameworks lead to a new survey by Hale where the growth of the leading frameworks in 2019 was observed as presented in Figure 3 (Hale 2020). The leading frameworks currently are Tensorflow with Keras as the high-level application programming interface (API) and Pytorch with fast.ai. According to these results Tensorflow is both most in demand framework, as well as fastest growing.

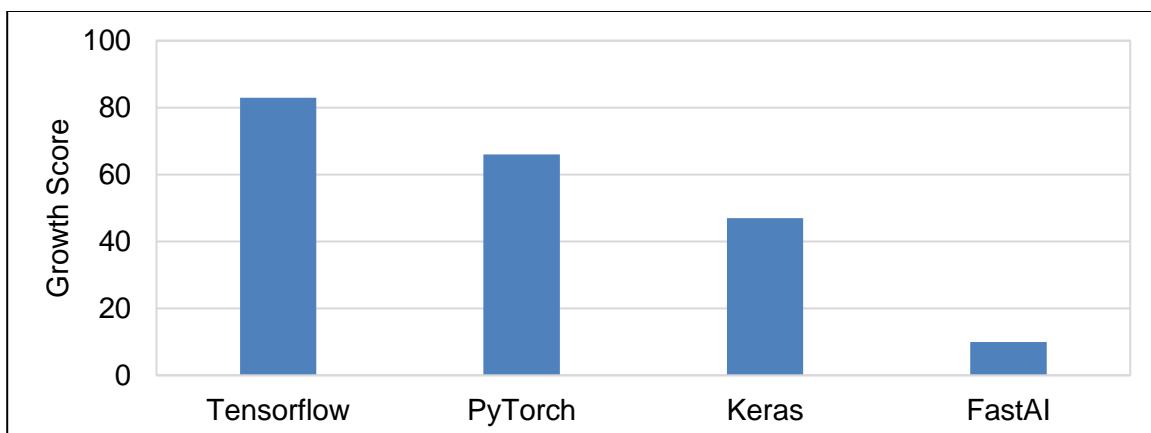


Figure 3: DL Framework Six-Month Growth Scores 2019 (modified according to (Hale 2020))

The open-source library **Tensorflow (version 2.3.0)** was developed by Google and was intended for the spam filter of Gmail before it was available to the public in 2015 (Open Data Science 2019). As shown before, TensorFlow is currently one of the most popular ML frameworks and is widely employed for DL. TensorFlow can run on various platforms, such as Linux, macOS, Windows and on the mobile platforms iOS, Android or on Raspberry Pi. For performance reasons, the library is written in C++, but the API is also available in Python, and others. TensorFlow can be executed on the central processing unit (CPU) or on the graphics processing unit (GPU) with enabled multiprocessing to boost the performance. One of the biggest advantages of TensorFlow is the possibility to work with low-level, as well as with high-level API.

Keras (version 2.4.0) as a high-level API was launched in 2015 and became the framework for developers due a clean API and the possibility to use it with different DL libraries as the backend such as TensorFlow, Theano or CNTK (Google Inc. 2019). In 2019 with TensorFlow 2.0, Keras was integrated and now is the standard interface, when developing DL environments.

The python package **pyomo (version 5.7)** is used for developing the MPC controller and is an is the open-source package which provides a variety of different optimization models (Sandia National Laboratories 2019). The high-level programming language has the advantage of usability over other algebraic modelling languages like AMPL, AIMMS or GAMS.

3 Machine Learning

In 1942, the idea of AI was born in the USA when it was mentioned in the science fiction short story called “Runaround” by Isaac Asimov (Haenlein and Kaplan 2019). At the same time, a machine called “The Bombe” for deciphering Enigma, an encryption device used for secure communications by the German military in the second world war was developed by the English mathematician Alan Turing. The ability to decipher Enigma led to Turing’s seminal paper “Computing Machinery and Intelligence” in 1950 which stated, that, for a machine to be intelligent, it needs to respond in a manner that it is not differentiable from a human being (Turing 1950). These criteria are seen as a benchmark for the intelligence of machines considered to be AI-systems. The first machine that matched this criterion was called ELIZA, it was able to simulate a conversation with a human and was developed between 1964 and 1966 at MIT. The System used for ELIZA was a so-called “Expert System” in which rules are programmed assuming that human intelligence can be formalized with a top-down “if-then” approach. The same system was used in IBM’s Deep Blue in 1997 which was able to beat the reigning chess world champion Gary Kasparov.

A more technical definition for ML was stated by Tom M. Mitchell in 1997: “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.” (Mitchell 1997). Ethem Alpaydin describes the task of ML as a optimization of a performance criterion using example data and experience (Alpaydin 2010). These definitions are still valid today and based on them different algorithms and approaches have evolved.

The next big milestone for Artificial Intelligence was made in 2015 by Google with the program “Alpha-Go” which can play the board game Go and was able to beat Lee Sedol the reigning world champion (Haenlein and Kaplan 2019). Figure 4 shows two children playing Go on a board with black and white stones which are placed anywhere on the grid and cannot be moved afterwards (DeepMind 2016). The goal of Go is to capture as much free space and surround as many of the opponent’s stones until no more move is possible. This leads to 10^{17} possible board configurations.

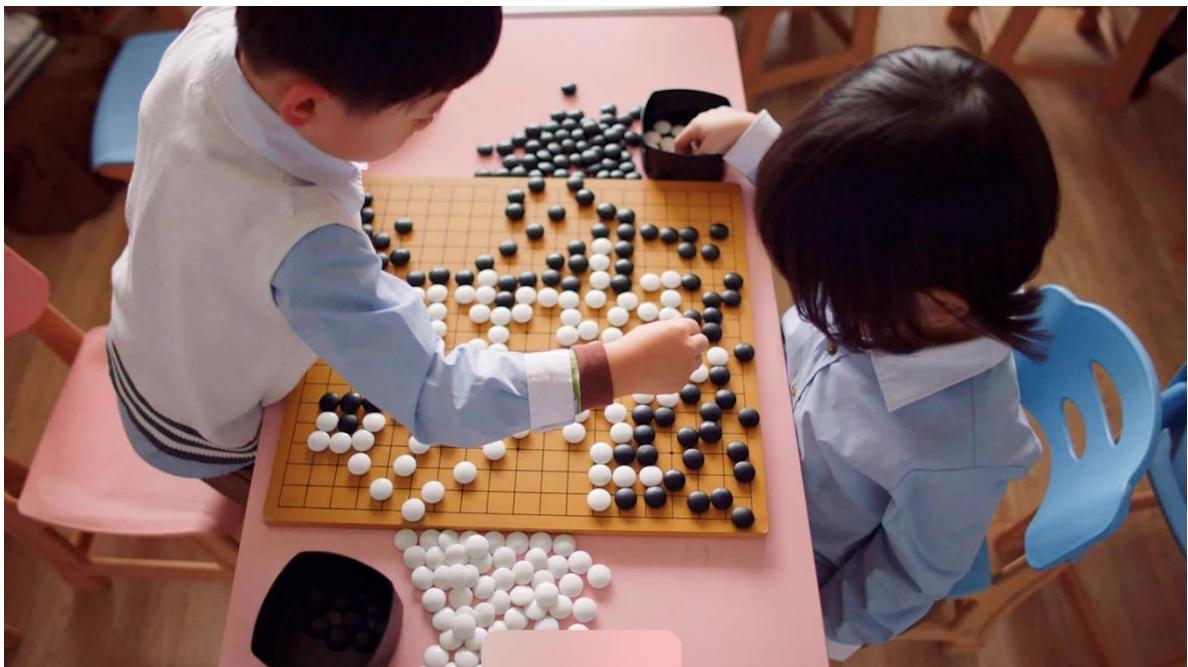


Figure 4: Children playing Go on a regular Go board (DeepMind 2016)

The possibility of 361 first moves in Go makes the game more complex than chess with only 20 possible starting moves. The brute-force of analyzing all possible moves as used in IBM’s Deep Blue chess gets infeasible for that number of possible actions with an exponentially increasing cost of calculation. Therefore, the team of DeepMind chose an approach of in form of a deep neural network. By playing against amateur players and against itself AlphaGo developed an understanding of how humans play and ultimately outplayed them.

3.1 Applications

The use of AI in industry is strongly driven by information technology companies like Google, Microsoft, Apple and Intel (Pan 2016). Google as an example uses Deep Learning to improve their picture search or develop their unmanned ground vehicle. The research in AI is shifting from academia-related research to research which addresses social demands like intelligent cities, medicine, transportation, logistics, manufacturing, as well as driverless automobiles. AI nowadays supports us constantly in everyday life. Google uses AI to sort the emails into different categories and, most importantly, to filter spam emails. Moreover, it recommends search queries based on the first words typed into the search field and then tries to find the best matches for the question (Bradley 2018). The business-focused social media platform LinkedIn uses AI to find best matches of employees to employers, by observing the behavior of applicants and the outcome of hiring processes. Facebook is helping to prevent suicide and to save lives by detecting suicidal thinking patterns and sending resources to help.

In the specific field of building technologies the research in RL started as early as 1997 and gained more interest since 2015 (Wang and Hong 2020). Wang and Hong found in their study that the main focus in building technologies was Heating-Ventilating-Air Conditioning (HVAC) with a margin of 35 % papers released in this topic in 2015. In 2015 Barrett and Lindner introduced a learning thermostat where the desired room temperature is set by the user and the learning thermostat controls the heating or cooling signal with on or off signals by learning the time schedule of the occupants (Barrett and Linder 2015). In comparison, Wei et. al. introduced a system which controls the air flow of the HVAC system with an agent (Wei et al. 2017). An RL-algorithm for the combination of HVAC control and window control was developed by Chen et. al. in 2018 (Chen et al. 2018). The similarity of these approaches is the cost saving potential in comparison to a heuristic control system.

These examples show that ML can be applied to a variety of integral tasks and different learning techniques are necessary to solve these problems.

3.2 Learning Techniques

The different ML techniques can be classified in the four categories of supervised, unsupervised, semi-supervised and reinforcement learning, depending on the required data (Figure 5) (Mohammed et al. 2017). The designation of the data with classified data refers to whether the data have a specific label, e.g. the picture of the dog has the name dog. With unclassified data the name of the picture has a name that is not referred to the content.

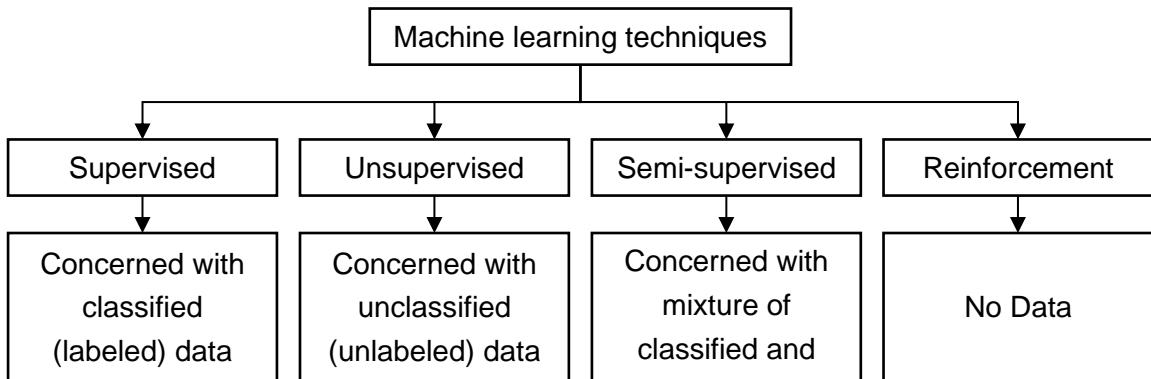


Figure 5: Different machine learning techniques and their required data (modified according to (Mohammed et al. 2017)

Supervised learning

The goal of supervised learning in its basic form is to find a correlation between input data and output data (Brownlee 2019). The two main types of supervised learning are classification and regression. A classification problem could be e.g. a dataset of handwritten digits with pixel data for which the learner should recognize the digits representing numbers from 0 to 9. The regression problem deals with numerical numbers as output, for example house prices could be calculated by given variables that describe the house itself and the neighborhood.

Unsupervised learning

Problems are solved without labelled input data as a reference for learning. In contrast to supervised learning, the model tries to describe or extract relationships in the data. The two main problems it is being used for are clustering and density estimation which are performed to find patterns in data. Another method where unsupervised learning is visualization for graphing or data plotting, as well as projection for reducing the complexity of multidimensional data.

Semi-supervised learning

This technique is a hybrid of supervised- and unsupervised learning where the training dataset contains more unlabeled than labeled data. This method is common for real-world supervised learning problems as in computer vision, natural language processing and automatic speech recognition, due to the lack of training data.

Reinforcement learning

The reinforcement learning technique does not have an initial works dataset available at the start of the training process. In this case, an agent operates in an environment and learns how to operate using feedback. Google's AlphaGo is an example for the most recognized example of reinforcement learning problem to date. Reinforcement learning is the technique used in this work and will be discussed in greater detail in the next chapter.

3.3 Reinforcement Learning

The goal of this thesis to control the temperature and illuminance in a room for minimal cost. For the problem at hand, no initial dataset exists prior to the training, making this an obvious candidate for RL.

RL is based on the process by which humans naturally learn (Sutton and Barto 2018). Gaining experience by interacting with our environment is one of the major sources for our knowledge. Figure 6 shows the basic agent-environment setup for RL. The agent operating with the environment selects the actions a_t to take in the current state s_t to reach the next state s_{t+1} and get the reward r_{t+1} as a feedback.

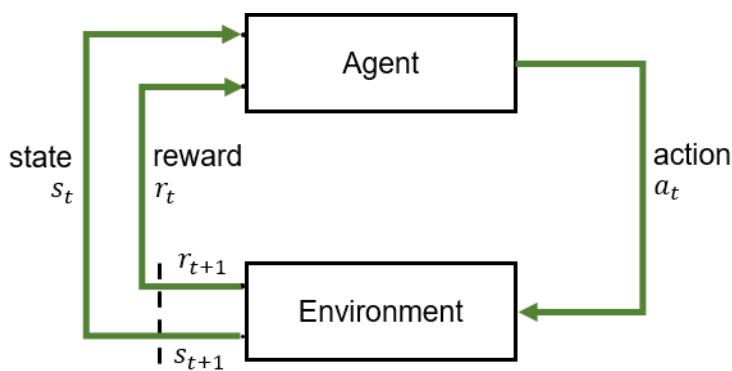


Figure 6: The agent–environment interaction in a Markov decision process. (modified according to (Sutton and Barto 2018, p.48))

The main elements of the RL-system are the policy, the reward function and the value function which are built into the agent and the Environment.

Environment

The environment can be a variety of problems, such as a car or boardgames like chess. In this thesis, the given environment is a thermal room model. The room temperature is the state of the environment and acts as the output. It should be ensured by the agent. The possible actions for the agent are the energy input by heating or cooling, as well as the control of the shading system. The room reacts to actions taken by the agent and creates an output in form of the next state and the immediate reward.

The reward signal is a single number calculated with a **reward function** in the environment. Its design is crucial for the learning success of the agent, as discussed by Sutton and Barto (Sutton and Barto 2018). The reward in the context of the given room model in this thesis contain the cost of energy and every exceedance of any comfort parameter. How the agent can learn from this reward function is described chapter 4.1 in greater detail.

Agent

The agent is responsible for selecting actions according to the current state of the environment following a **policy** as the main element of the process. This policy can be a look-up table, a function, or a search policy. Recent algorithms make use of parameterized policy by introducing a neural network. The actions can be selected with a stochastic function, with probabilities for each action, or deterministic with the output of the policy being the real value of the action. The optimization goal in of the agent in this thesis is to save energy costs while maintaining to meet the needs of the occupants. To achieve an optimal control strategy, the agent tries to maximize a cumulated reward (return) over all viewed timesteps. In its simplest case, the return can be the sum of the rewards equation 1.

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T \quad (1)$$

G_t return, cumulated reward

R_{t+i} ... reward of timestep

R_T reward of the terminal timestep, last, timestep in the viewed timeperiod

Heating or cooling a building is a continuous task without a terminal state which would lead to an infinite return with the formulation in equation 1. Adding a discount rate γ to future rewards prevents this behavior with equation 2. The initially received reward is worth more than the reward received after the next step.

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad (2)$$

G_t ... return, cumulated reward

R_{t+i} ... reward of timestep

γ discounting factor

Value-function

In RL, the two value functions used are called **state-value** and the **action-value** function. The value functions are used to estimate the return, because the rewards for each timestep are not known prior to the state visitation. The value functions are used to train the agent to achieve the optimization goal of maximizing the return.

The **state-value** $v_\pi(s)$ is defined as the total expected reward achievable in the future starting from this state. The state-value noted as $v_\pi(s)$ indicates what the best option for the long run is and takes into account the next states which are most likely to follow. That means that the reward in a specific state can be low, however the value of this state can still be high if the following states can gain a high reward. In the simple maze depicted in Figure 7, the goal is to move from the start in the top left corner to the goal in the bottom right corner. Following the orange line with the highest reward in every single box and summing up the rewards (green numbers), the total return is 120 whereas following the red line by taking the future rewards into account results in the total reward of 155, making it the better option.

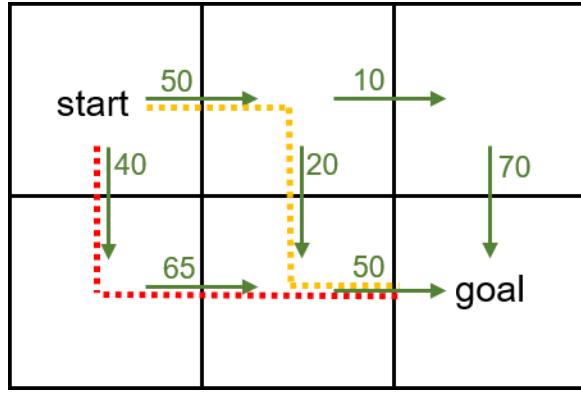


Figure 7: Simple Maze displaying the difference of reward and state-value

This simple example shows that the state-value is crucial for the performance of RL algorithms but is not as straightforward as the reward estimation which is a direct feedback from the environment. The state-value of each state is dependent on the possible actions and the probability these actions are taken following the current policy. Figure 8 shows the state-value and is the visualization of equation 3. Starting from a specific state s the agent can perform any action a . The transition from state s to the next state s' and the immediate reward is expressed as the probability $p(s', r|s, a)$. The reward r is added to the discounted state-value of the next state $v_\pi(s')$. The sum of all possible state-values by taking different actions is then averaged over all possible actions by multiplying the probabilities $\pi(a|s)$ of taking each action a in state s .

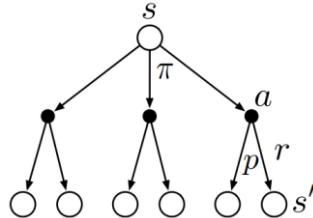


Figure 8: backup diagram for $v_\pi(s)$ (Sutton and Barto 2018, p.59)

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) [r + \gamma v_\pi(s')] \quad (3)$$

v_π ... state-value

π policy

s state

s' next state

a action

r reward

Like the state-value, the **action-value** $q_\pi(s, a)$ is the estimated return with respect to the state s and the action a . The action value indicates how good it is to take a specific action in a specific state Figure 9.

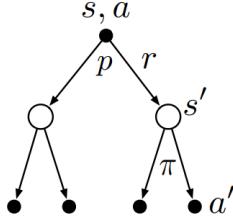


Figure 9: q_π backup diagram (Sutton and Barto 2018, p.61)

RL-algorithms with a model are called model-based algorithms which are planning the future situations without ever experiencing them. Model-free methods are simple trial and error learners learning from the experience they obtain during operation.

Figure 10 displays an overview of modern RL-Algorithms divisible into model-based- and model-free algorithms (OpenAI 2018). Model-based algorithms do know the model dynamics or learn the dynamics of the environment model with the advantage for planning ahead and seeing what will happen when choosing certain actions. While Google's AlphaZero is model-based with the agent being provided with the ruleset of the game, whereas algorithms like Stochastic Value Gradient learn the model dynamics as part of the learning process. This has the disadvantage that biased models are possibly learnt during the training with the result of sub-optimal performance in the real environment.

Model-free algorithms are separated into the two main approaches of policy optimization and Q-learning. The RL-system either learns policies, action-values (Q-functions) or value functions. As can be seen in Figure 10, the DDPG, TD3 or SAC algorithm are a combination of Policy-Optimization and Q-learning. These policy optimization algorithms are so-called **actor-critic** algorithms and are characterized by using a critic and an actor for training and selecting an action. The state-value of the current step is the estimated state-value of the next step added to the actual reward given by the environment. The reward with the estimated state-value of the second step is then called the one-step return $G_{t:t+1}$ which is used to assess the action (Sutton and Barto 2018). The use of the state-value function in this way is called a critic and the function which takes actions is called the actor. Actor-critic algorithms take the one-step return to update and improve the policy.

Other well-known algorithms like PPO or TRPO are pure policy gradient algorithms. Another group of algorithms are Q-learning algorithms, because they learn the Q-value which is in fact the action. These algorithms are not feasible for large action spaces because of the discrete actions they use and the necessity to calculate the action-value of all possible actions in the specific state. Policy optimizations can perform in continuous action spaces and are therefore the preferred algorithms for the problem discussed in this work and are further described below.

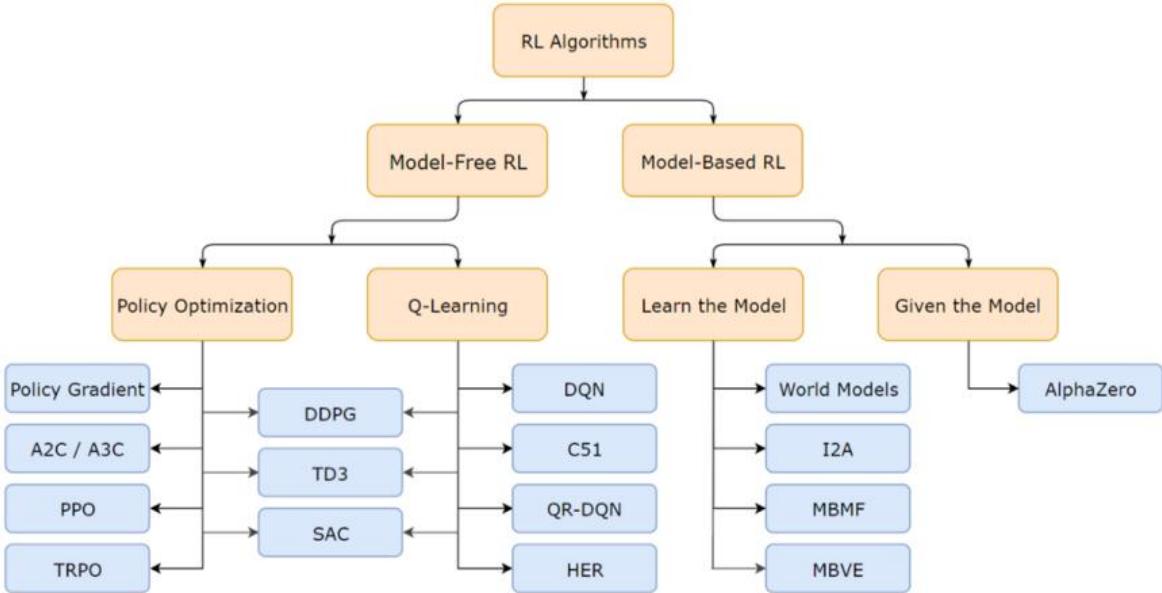


Figure 10: A non-exhaustive, but useful taxonomy of algorithms in modern RL (OpenAI 2018)

The key differences of the Policy-optimization algorithms are described in terms of the action space and policy, the performance measure, and the question if the algorithm is an on- or off-policy algorithm. Off-policy for example means, that the experience used for training the policy and value functions is not produced by the current policy and can be used multiple times which makes these algorithms more sample efficient. On-policy methods use the experience from the last episode or steps and compare the new policy with the old one to find out if this episode is better. The performance measure for the comparison is called the advantage of the policy. The key points of the policy optimization algorithms are:

DDPG: Deep Deterministic policy gradient(Lillicrap et al. 2015) (Lillicrap et al. 2019)

- Continuous action spaces with a deterministic policy
- Learns policy and action-value function
- Off-policy

TD3: Twin Delayed Deep Deterministic policy gradient (Fujimoto et al. 2018)

- Continuous action spaces with a deterministic policy
- Learns policy and stabilized action-value function
- Off-policy

SAC: Soft Actor-Critic (Haarnoja et al. 2018)

- Continuous or discrete action spaces with a stochastic policy
- Learns policy and stabilized action-value function
- Off-policy

A2C/A3C: Asynchronous Advantage Critic (Mnih et al. 2016)

- Continuous or discrete action spaces with a stochastic policy

- Advantage function
- On-policy

PPO: Proximal Policy Optimization (Schulman, Wolski, et al. 2017)

- Continuous action spaces with a stochastic policy
- clipped, advantage function
- On-policy

TRPO: Trust Region Policy Optimization (Schulman, Levine, et al. 2017)

- Continuous or discrete action spaces with a stochastic policy
- KL-divergence advantage function
- On-policy

The **policy gradient** algorithms select the actions based on a parameterized policy (e.g. neural network) that uses a performance measure (e.g. value function) to update the parameters and improve the performance (Sutton and Barto 2018). The policy is learned based on the gradient of an accumulated reward, as the performance measure with respect to the policy parameters referred to as $J(\theta)$ which can be written as equation 4 with μ as the distribution of the states and π as the policy corresponding to the parameter vector θ .

$$\nabla J(\theta) = \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a|s, \theta) \quad (4)$$

$\nabla J(\theta)$... gradient of the performance measure

s state

a action

θ parameter vector

$\mu(s)$ state distribution

q_π action-value

π policy

The update of policy gradient methods is based on gradient ascent with respect to the current policy parameters θ_t in equation 5.

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t) \quad (5)$$

θ_{t+1} ... current parameter vector at timestep t+1

α learning rate (step size of the gradient)

θ_t parameter vector at timestep t

∇J gradient of performance measure

The policy π selects actions a based on the current state s with the current parameters θ and can be noted as $\pi(a|s, \theta)$.

3.4 Reinforcement Learning in Building Technologies

Wang and Hong have very recently published a review giving a detailed analysis of publications since 1997 in the field of RL in building technologies (Wang and Hong 2020). The algorithms which have been used so far are to 76.6 % based on the number of publications value-based algorithms (Q-learning) which were already excluded for this thesis because of their disability to work in continuous action spaces. Actor-critic algorithms got more popular in recent years with a total share of 15.1 % of all publications. The popularity of actor-critic algorithms is due to the possibility for transfer-learning which means, that a trained behavior from one building can be generalized to other buildings as well. The policy function is suitable for transfer learning, because the task of ensuring the room temperature is the same in every building, whereas the mapping from states to actions is not transferable due to different control goals and structures in building technologies.

The methods used to represent the policy and value function shift more and more to neural network estimators which were used in all publications in 2019 listed by Wang and Hong. The study concludes that the majority of utilized RL controllers adopted supervisory control which they describe as setpoint control where conventional controllers are still needed to track this setpoint.

Given the analysis by Wang and Hong this thesis will focus on an actor-critic algorithm for developing a RL-controller applicable for transfer learning. Sutton and Barto state, that the advantage of an approximation policy is that it can approach a deterministic policy. Together with the advantages of the policy gradient algorithm the Deep Deterministic Policy Gradient fulfills the approach of a deterministic policy which is described in detail in the section 4.1.

3.5 Neural Networks

The chosen DDPG-algorithm uses neural network for the actor to select the actions and the critic to estimate the actor-value. The idea of a neural network is based on the functionality of a brain (Ertel 2016). The big step towards an AI of neural networks was taken by McCulloch and Pitts in 1943 with the mathematical model of the neuron as a basic switching element for brains. This formulation laid the foundation for the construction of artificial neural networks.

The neuron of a brain is comparable with a conductor which get charged by incoming impulses and sends a signal if the voltage exceeds a certain threshold to all connected neurons where the same process is repeated. A neuron can have multiple inputs and outputs and are connected to other neurons (Figure 11).

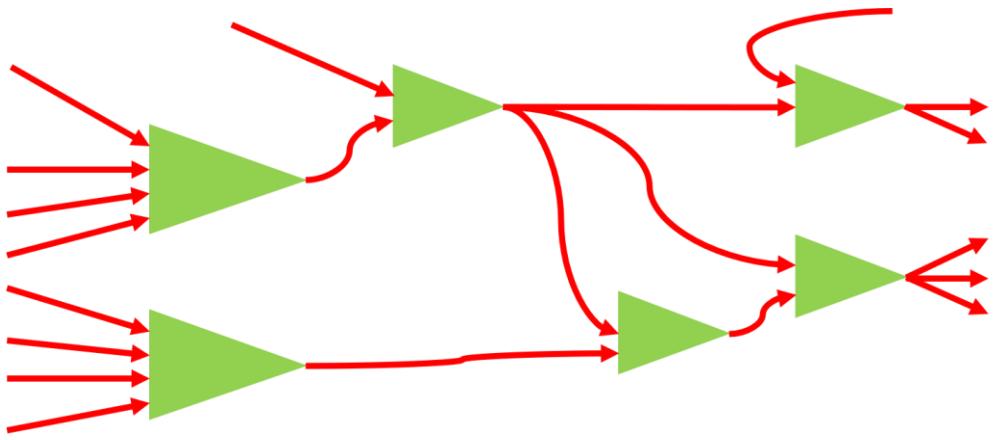


Figure 11: Formal model with neurons and directed connections between them (modified according to (Ertel 2016, p.267))

The mathematical formulation for this process replaces the continuous process of the brain with a discrete time scale and the charging of the activation potential is the sum of the weighted output values with weight ω_{ij} of all input values x_j with an applied activation function f (equation 6 and Figure 12). There are several options for the activation function which are explained in the section 3.5.3.

$$x_i = f \left(\sum_{j=1}^n \omega_{ij} x_j \right) \quad (6)$$

x_i output of neuron

f activation function

ω_{ij} ... weights of the connections

x_j inputs pf neuron

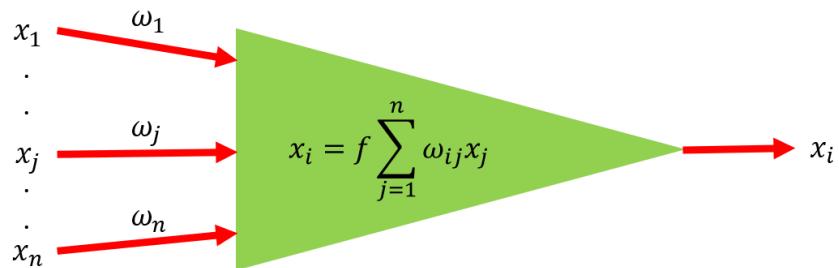


Figure 12: The structure of a formal neuron that applies the activation function f to the weighted sum of all inputs (modified according to (Ertel 2016, p.269))

The most used neural network model is the backpropagation algorithm because of its universal applicability for any approximation task. Figure 13 shows a backpropagation network with an input layer, a hidden layer, and an output layer. The values x_j^p of the output layer are compared with the values of the targets t_j^p . In the tables in Figure 13 the values of

the inputs outputs and the target values of the neural network used in this thesis are shown. The state input is the room temperature, and the forecast input is the weather forecast with the outside air temperature, solar radiation, cost of energy and the occupancy of the room. The actor output Q is the heating-or cooling energy input and T_v is the value for the shading system.

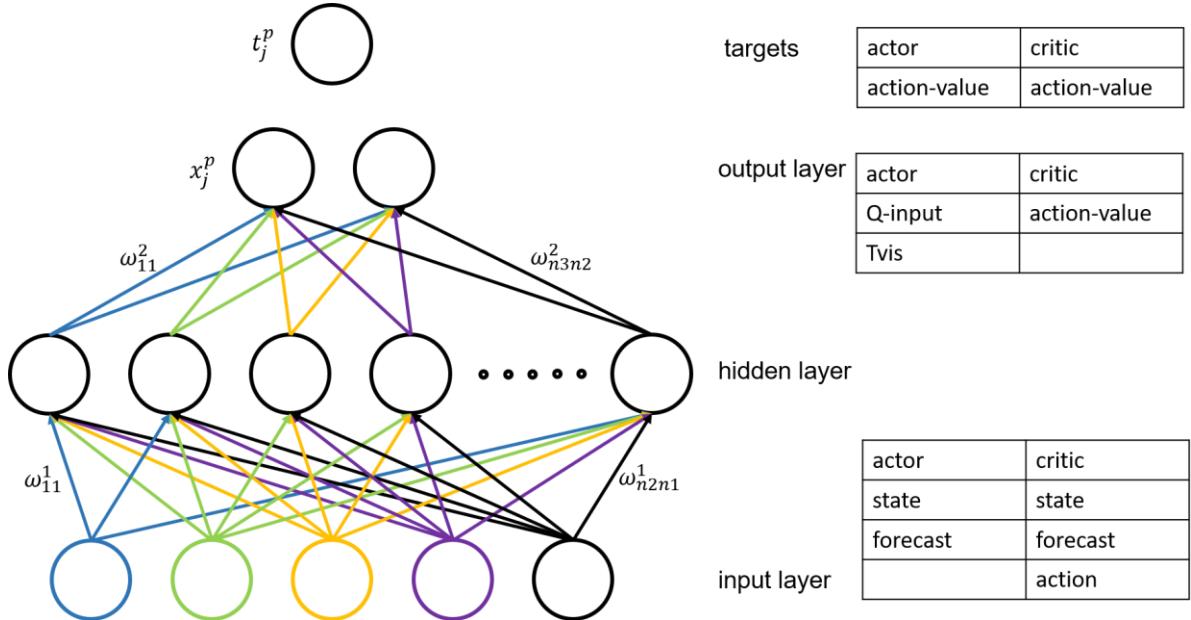


Figure 13: A three-layer backpropagation network with n_1 neurons in the first layer, n_2 neurons in the second and n_3 neurons in the third layer (modified according to (Ertel 2016, p.291))

The target value for the critic network is compared with the value of the output layer and the error is calculated with the preferred function. This error is then used to calculate the negative gradient of the weights and further tune the weights to minimize the error and make accurate estimations of the action-value. The actor network with the actions as an output is not trained to minimize an error and get accurate predictions but trained to minimize the action-value function.

The following section shows two neural network architectures based on the backpropagation model for RL which have already proven their usefulness in a wide range of problems. The structure of Multi-layer perceptron models and Recurrent neural network models is described in the following section.

3.5.1 Multi-Layer Perceptron

The Multi-Layer Perceptron network (MLP) is viewed as the classical neural network (Brownlee 2016). The basic structure of this network class is an input layer followed by one or multiple hidden layers and an output layer. Figure 14 shows this structure and displays that the layer size of the layers can vary.

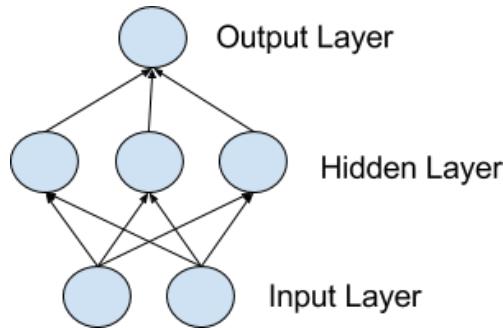


Figure 14: Model of a Simple Network (Brownlee 2016)

The input layer is not constructed with neurons and passes the input to the first hidden layer in the network. The network can have multiple hidden layer which is referred to as Deep Learning. The properties of the output layer as the final layer depends on the problem the neural network is used for. The output layer in this thesis has one output neuron for the critic network estimating the action-value function and the actor has two outputs for two actions. The properties and what range of values this neuron can output is depending on the activation function described in 3.5.3.

3.5.2 Recurrent Neural Network

MLP networks are not able to learn time related dependencies, because they have no knowledge of what happened in the timestep before (Olah 2015). Recurrent Neural Network (RNN) address this issue with loops in the neurons of the RNN layers. A RNN neuron look like the left-hand side of Figure 15 with a loop that allows to use the past information to be used in the current step. The right-hand side shows the unrolled neuron where the output h_0 of timestep zero is passed to the next timestep and is the input together with X_0 .

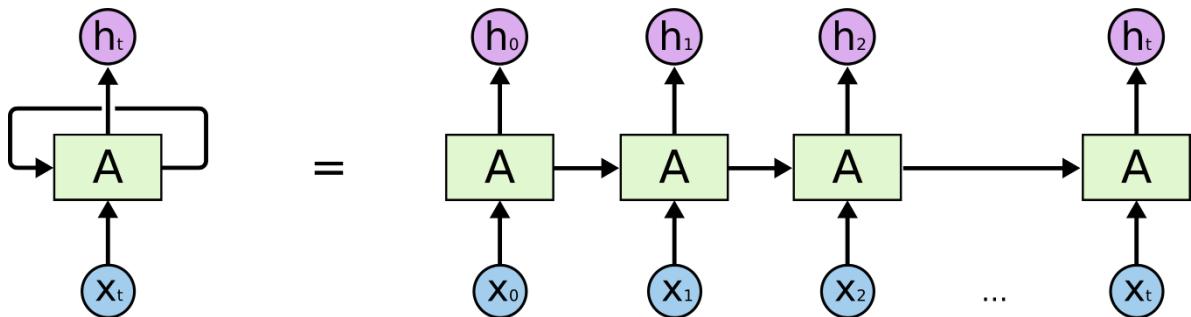


Figure 15: An unrolled recurrent neural network (Olah 2015)

This additional knowledge lead to success in speech recognition, language modelling, image captioning or timeseries forecasting. In 2015 Heess et al. used an RNN approach in the DDPG algorithm to conquer problems with partial observable environments like a way sign in a navigation task which is only temporary available (Heess et al. 2015). In the task of room conditioning the interesting value to remember is the past actions that were taken.

The idea to use an RNN in such a task is to connect previous information (way sign) to the present task (navigation). Unfortunately, basic RNNs have a problem with long term dependencies where not only the information of the last time-step is needed but also the information of a few timesteps back (Olah 2015). Following example by Olah makes this issue clear: I grew up in France I speak fluent "?". For a human it is clear, that the missing word is French. The bigger this gap grows it gets more likely for the RNN to fail.

This problem is solved with Long-Short-Term Memory (LSTM) networks which are designed to learn these long-term dependencies. LSTMs were introduced by Hochreiter and Schmidhuber in 1997 (Hochreiter and Schmidhuber 1997). The difference between the RNN and the LSTM is how the information of past timesteps is passed to the next timestep. In RNN the repeating modules responsible for the forward pass of past information is a simple structure with an activation function (Figure 16).

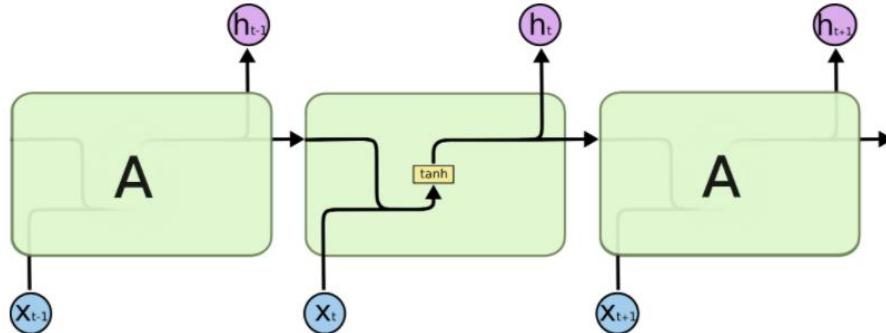


Figure 16: The repeating module in a standard RNN contains a single layer (Olah 2015).

The improved LSTM network layers repeating module is built with four interacting network layers shown in the middle of Figure 17.

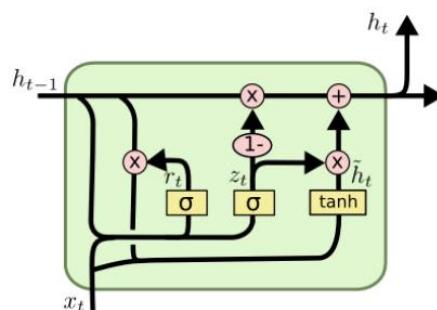


Figure 17: The repeating module in an LSTM contains four interacting layers (Olah 2015).

The architecture of LSTM decides and learns what information to keep of the past information, what information to store as the state of the layer and what information to pass as the output. This is done with the sigmoid (sig) or hyperbolic tangent (tanh) layers called gate layers. The first layer is the forget gate layer which decides what information is thrown away and what to keep followed from the input gate layer which decides which values are updated. Together with the tanh layer the state is updated. The output of the LSTM is a

filtered version of the state which is put through a tanh layer to push the values between -1 and 1 multiplied by a sigmoid gate.

3.5.3 Network features

For both presented network architectures the network features like the number of layers, number of neurons of each layer, activation function of the layer and what loss-function should be used to train the network have to be set.

Activation functions

The most common activation functions in neural networks are the sig, tanh and variants of rectified linear units (relu) (Ding et al. 2018). Ding et al. analysed the different activation functions based on their characteristics in neural networks.

The **sig function** is the most used activation function because the calculation is easy. The problem with the sigmoid function is that while backpropagating the derivative will reduce to zero around saturation, as shown in Figure 18 and that leads to a vanishing gradient. The gradient vanishes, when more layers with the same activation function are added to a neural network (Wang 2019). The weights are not updated effectively which can lead to an inaccurate NN. The output of the sigmoid function is between 0 and 1.

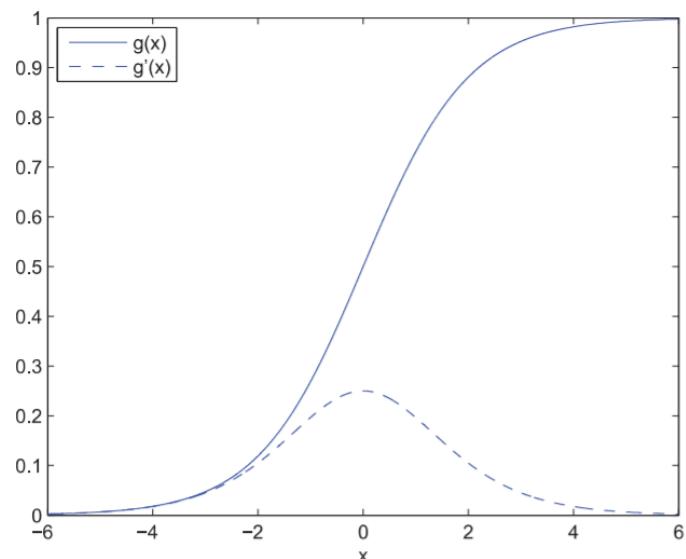


Figure 18: The graphic depiction of Sigmoid function and its derivative (Ding et al. 2018, p.1837).

Similar to the sigmoid function is the **tanh** with output values between -1 and 1 (Figure 19). The symmetric nature of the function makes it more likely to be used than sigmoid because the average of the layer is close to zero and the neural network converges faster. The problem with the vanishing gradient also exists with the tanh activation and is more complicated to calculate what makes the computing of the gradient and the update of the weights more time consuming.

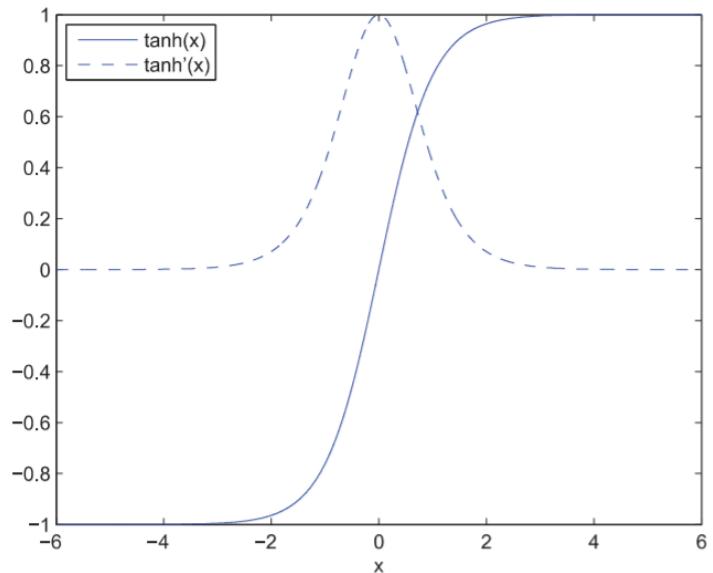


Figure 19: The graphic depiction of hyperbolic tangent function and its derivative (Ding et al. 2018, p.1838)

The **relu** activation and its improvements are currently the most used activation functions in neural networks. Values smaller than zero which are passed to the activation are always zero and values bigger than zero are activated with a linear function (Figure 20). The relu function has advantage of being less computational demanding. With a derivative of 1 the neural network converges faster and avoids local optimizations and a vanishing gradient. The disadvantage of relu function is the dying neuron problem. The output of negative values as zero lead to so called dead neurons which will never be activated.

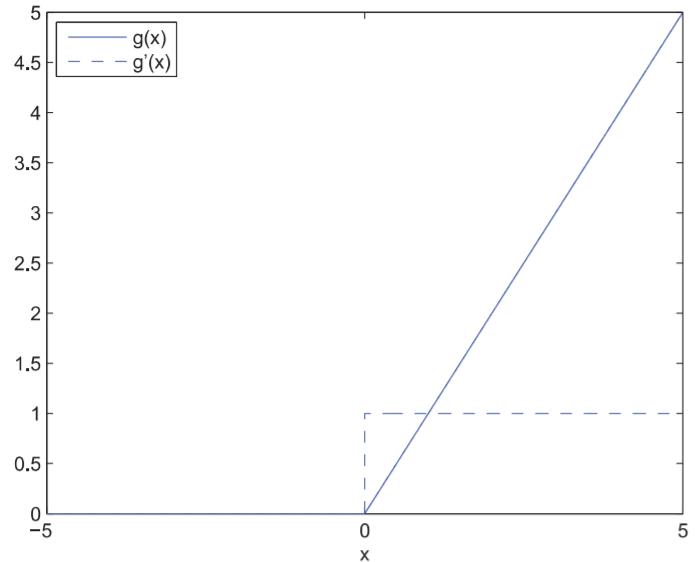


Figure 20: The graphic depiction of ReLU function and its derivatives (Ding et al. 2018, p.1838)

This dying neuron issue can be solved with the **leaky relu (lrelu)** activation function where the negative values of the neuron are not zero and are calculated with a fixed scale for the negative slope, shown in Figure 21. For other activation functions like the **prelu** and the **rrelu**, the negative slope is not fixed but trainable or selected randomly.

Ding et. al. tested these activation functions with a classification problem where the neural network with the **ReLU** function performed the best.

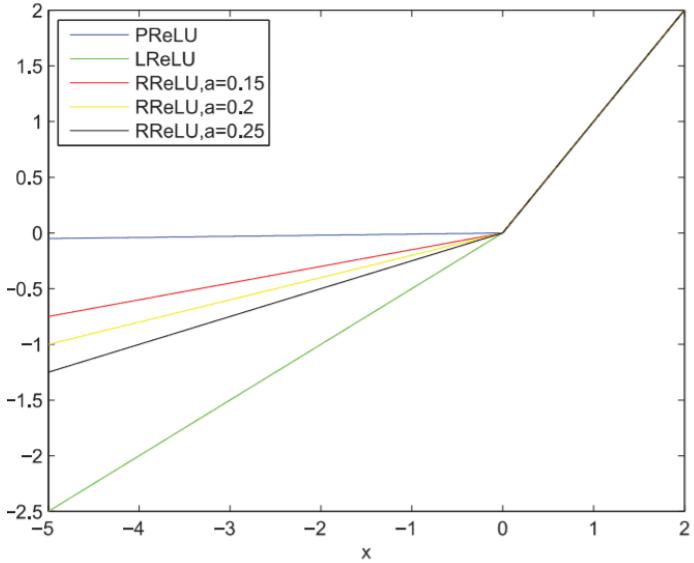


Figure 21: The graphic depiction of LReLU, PReLU, and RReLU function (Ding et al. 2018, p.1839)

Loss function

The loss function is the measure of how accurate the model of the neural network predicts the target values (Seif 2019). The Mean Squared Error (MSE) loss used as a default in the DDPG (equation 7) is the right choice when the aim for a neural network is to be accurately in the majority of situation.

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - q_\pi(s, a))^2 \quad (7)$$

N number of samples

y_i action-value as the target value

$q_\pi(s, a)$... estimated action-value

4 Results

The basis of the algorithm used to solve the control problem for heating, cooling and controlling of the shading device in an office room is the DDPG, an improvement of the initial Deterministic Policy Gradient by Silver et al which implements a Deep Neural Network, was proposed in 2015 by Lillicrap et al. (Lillicrap et al. 2015) (Lillicrap et al. 2019). Numerous improvements have been made to this particular RL-algorithm since it was introduced. Improvements considered in this thesis are optimized replay buffer approaches and ways to manage the exploration of the agent using different noise processes.

For a better understanding of the nomenclature in the following chapter and for linking the algorithm to the use case, the state properties and the action space is defined as follows.

state	room temperature
forecast	forecast data for air temperature, solar irradiation, cost of energy and a Boolean variable if the room is occupied or not.
observation	the state and forecast
actions	thermal heating/cooling power and the shading factor
action space	heating/cooling power is bound between -1 and 1 with a scaling factor depending on the room properties shading factor is set between 0.01 and 0.6.

4.1 Deep Deterministic Policy Gradient (DDPG)

The DDPG is a model-free, off-policy, actor-critic algorithm which can solve problems with high dimensional, continuous action spaces (Lillicrap et al. 2015) (Lillicrap et al. 2019). Lillicrap et al. showed, that DPG is unstable for challenging problems and therefore combined the DPG algorithm with a Deep Q Network algorithm. The advantage of the Deep Q Network algorithm is given by the replay buffer which is replayed in an off-policy way to reduce the correlation between the samples and the use of target networks to reduce the variance of targets while calculating the temporal difference errors. The implementation of DDPG follows a straight-forward actor-critic architecture and is therefore easy to implement and to scale to different tasks and network sizes.

The main elements, visualized in Figure 22 of this algorithm are the replay buffer, the environment, the actor network initialized as $\mu(o|\theta^\mu)$ and the critic network as $Q(o,a|\theta^Q)$. The weights θ^Q, θ^μ of both networks are used to initialize the target networks μ', Q' as copies of the actor and critic with the respective weights $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$ which are introduced to stabilize training.

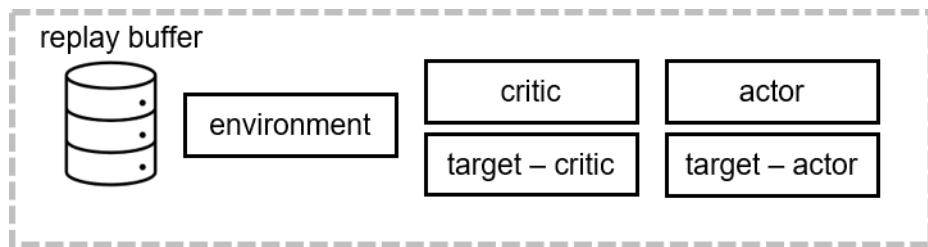


Figure 22: Elements of the DDPG algorithm

The off-policy algorithm explores the action space by selecting an action following the current policy μ in the current observation o_t with an action noise N_t added to the selected action at

equation 8. In the DDPG the action noise for exploring the action space can be handled independently of the learning algorithm.

$$a_t = \mu(o_t | \theta^\mu) + N_t \quad (8)$$

a_t ... selected actions (Q-input, Tvis)

o_t ... observation (state and forecast values)

μ deterministic policy (actor)

θ^μ ... parameters of the actor

N_t ... action noise

The trajectory following the execution of the action is stored with the transition from one state s_t with a forecast f_t and action a_t to the next state s_{t+1} with the next forecast f_{t+1} and the reward r_t for the current timestep as the trajectory $(s_t, f_t, a_t, s_{t+1}, f_{t+1}, r_t)$. Figure 23 shows the process starting from selecting the action until storing the trajectory.

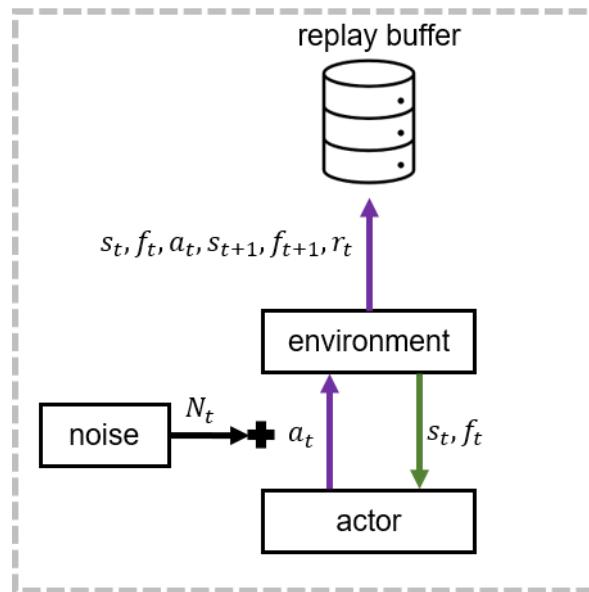


Figure 23: DDPG – agent-environment interaction

The training of the actor-critic networks is executed after each timestep with a minibatch of trajectories, which are sampled randomly from the replay buffer. The training process starts with calculating the action-value with the target networks. With the observation of the timestep t+1 from the sampled minibatch the target actor selects an action and passes it to the target critic to calculate the action-value by adding it to the reward from timestep t. In Figure 24 the green arrows show input data from the replay buffer and the purple arrows are outputs of neural networks. Equation 9 depicts the mathematical formulation of the process.

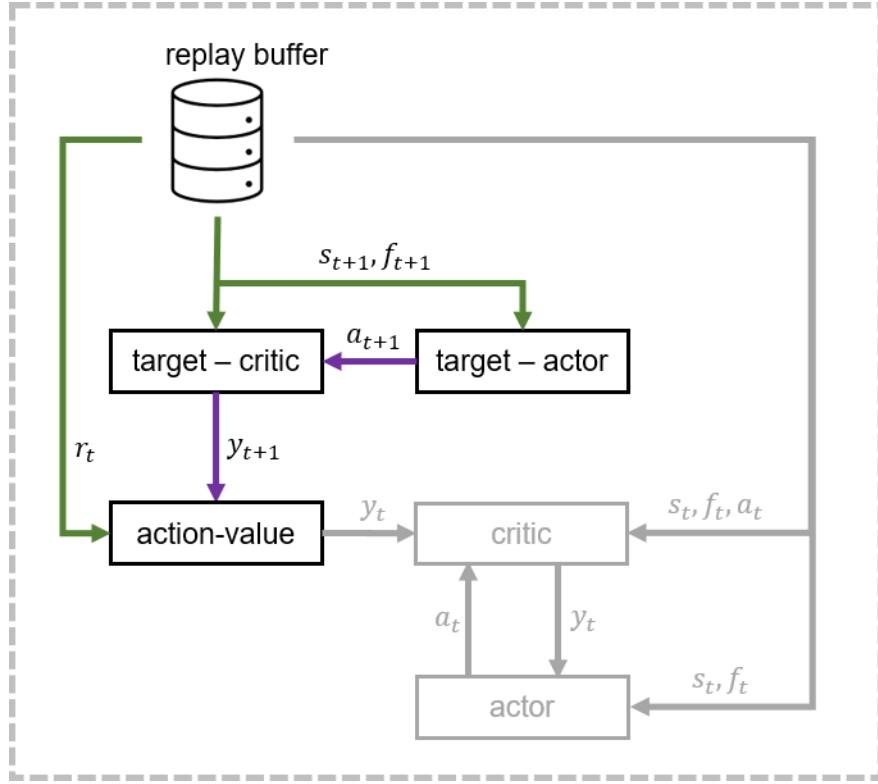


Figure 24: DDPG – calculating the action-values

$$y_t = r_t + \gamma * Q'(o_{t+1}, \mu'(o_{t+1} | \theta^{\mu'}) | \theta^{Q'}) \quad (9)$$

y_t action-value as target

o_{t+1} ... observation of timestep t+1 (state and forecast values)

r_t reward of timestep t

γ discount factor

Q' target critic

$\theta^{Q'}$ parameters of target critic

μ' target actor

$\theta^{\mu'}$ parameters of target actor

The loss L of the critic-network by estimating the action-value is minimized during training of the critic with the mean squared error between the action-value y_t and the approximation of the critic (equation 10).

$$L = \frac{1}{N} \sum_t (y_t - Q(o_t, a_t | \theta^Q))^2 \quad (10)$$

L critic loss

N number of samples

y_t action-value as target

o_t observation of timestep t (state and forecast values)

a_t selected actions (Q-input, Tvis)

Q critic

θ^Q ... parameter of critic

The training process presented in Figure 25 illustrates the off-policy training of DDPG. Green arrows are the inputs from the replay buffer, purple arrows are the outputs from the neural network, and the blue arrows are the values used for backpropagation through the network. The critic is trained with actions selected by an old policy and the action value calculated before. The training of the actor starts with selecting actions with the sampled inputs according to the new policy. The new observation and action inputs are feed into the critic. The objective for optimizing the actor policy is the sampled policy gradient following the updated critic network. The mean of the estimated actor-value is used to calculate the gradients which are applied to the actor policy. Equation 11 depicts the mathematical formulation of the process.

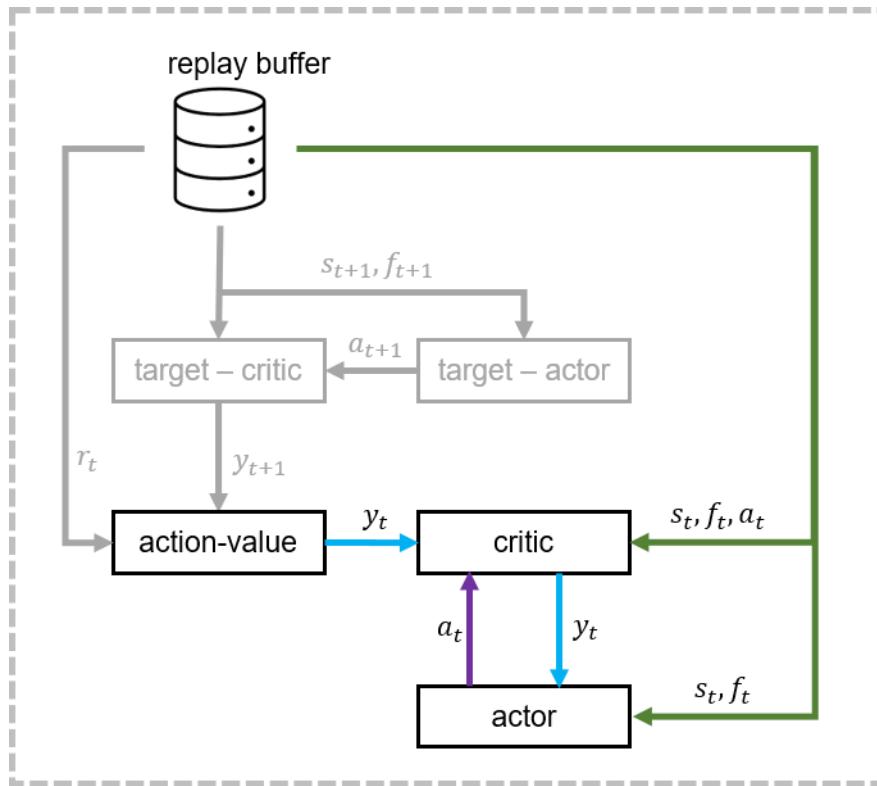


Figure 25: DDPG – training of the critic and actor network

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_t \nabla_\alpha Q(o_t, \theta^Q)|_{o=o_t, a=\mu(o_t)} \nabla_{\theta^\mu} \mu(o_t | \theta^\mu)|_{o_t} \quad (11)$$

$\nabla_{\theta^\mu} J$... gradient of the performance measure

N Number of samples

o_t observation of timestep t (state and forecast values)

Q critic
 θ^Q parameter of critic
 μ deterministic policy (actor)
 θ^μ parameters of the actor
 α learning rate (stepsize of the gradient)

The training of the neural network is stabilized with a soft update, which means that the parameter of the actor- and critic network are decreased with the factor τ before copying the parameters to the target networks, calculated with equation 12 for the critic network and with equation 13 for the actor network. The disadvantage of this soft update is the slower propagation of the action-value estimation of the critic.

$$\theta^{Q'} \leftarrow \tau\theta^Q + (1 - \tau)\theta^{Q'} \quad (12)$$

$$\theta^{\mu'} \leftarrow \tau\theta^\mu + (1 - \tau)\theta^{\mu'} \quad (13)$$

$\theta^{Q'}$... parameters of the target critic
 θ^Q parameters of the critic
 $\theta^{\mu'}$... parameters of the target actor
 θ^μ parameters of the actor
 τ soft constraint

4.2 Replay Buffer

Experience from the interaction of the agent with the environment is stored in the replay buffer with state, forecast, action, next state, next forecast, and the reward. In the original DDPG-algorithm by Lillicrap et al a subset of experiences is randomly sampled from the replay buffer to train the networks (Lillicrap et al. 2019).

Schaul et. al proved that the learning process can be improved by sampling the experiences according to a priority for each experience. **Prioritized Experience Replay (PER)** uses the absolute value of the temporal difference (TD) error (equation 14) of the estimation of the action-value by the critic network (Schaul et al. 2016). The priority of each sample is updated after these experiences are used for training the neural network for future training steps. Since new experience do not have a priority, and thus, would never be selected for training the priority is set to the clipped maximum priority set by the user.

$$\delta_i = |r_t + \gamma * Q'(o_{t+1}, \mu'(o_{t+1} | \theta^{\mu'}) | \theta^{Q'}) - Q(o_{t-1}, a_{t-1})| \quad (14)$$

δ_i TD error
 o_{t+1} ... observation of timestep t+1 (state and forecast values)
 o_{t-1} ... observation of timestep t-1 (state and forecast values)

a_{t-1} ... actions of timestep t-1
 r_t reward of timestep t
 γ discount factor
 Q' target critic
 $\theta^{Q'}$ parameters of target critic
 μ' target actor
 $\theta^{\mu'}$ parameters of target actor

In PER, the TD error shrinks slowly, which leads to a frequent replay of experiences with an initial high TD error. This lack of variety in the training data for the neural network can lead to over-fitting, meaning that the agent is able to solve the problem in specific states with specific forecasts only. To overcome this issue Schaul et al. introduces a stochastic sampling method, which interpolates between a pure greedy-sampling and random sampling of the experiences with equation 15. The exponent α sets how much prioritization is used with $\alpha = 1$ as the prioritized case with no randomness.

$$P(i) = \frac{\delta_i^\alpha}{\sum_i \delta_i^\alpha} \quad (15)$$

$P(i)$... priority of sample i
 δ_i^α scaled TD error of sample i
 α prioritization of randomness

Prioritized sampling introduces a bias in the network because experiences with high priorities are used more often for training. Importance sampling (IS) weight (equation 16) is a way to correct the bias. An unbiased sampling is especially important at the end of training, therefore the exponent β sets the amount of correction and increases over time to one. Another benefit of IS weights are the lower magnitudes of the gradients of samples with a high TD error, which enables the use of a higher global step size of the optimizer.

$$\omega_j = \frac{(N * P(i))^{-\beta}}{\max_i \omega_i} \quad (16)$$

$P(i)$... priority of sample
 ω_j weight of sample
 N batch size
 β amount of importance correction

The weight change with IS weights is set according to equation 17.

$$\Delta \leftarrow \Delta + \omega_i * \delta_i * \nabla_\theta Q(o_{i-1}, a_{i-1}) \quad (17)$$

ω_i ... importance sampling weight
 δ_i ... TD error

o_{i-1} ... observation of timestep t-1
 a_{i-1} ... actions of timestep t-1
 θ parameter of critic network

Another priority sampling algorithm introduced by Cao et al. in 2019 called the **High-Value Prioritized Experience Replay (HPER)** builds on PER but combines the action-value and the TD-error for each sample (Cao et al. 2019). The high TD-errors in the first episodes of training do not improve the agent because the optimal policy will not reach these states. The IS weight, as well as the TD error are calculated the same way as in equation 16 and equation 17, respectively. The priority calculation is extended by the variable u_i , which is updated with $u_i = u_0 * \mu$ every time the experience is used for training.

The priority value for the action-value and the TD-error are often not in the same range. Therefore, these values must be normalized. Cao et al. used the sigmoid function (equation 18) to do so and updates the priorities the action-value and TD error-priority with equation 19 and equation 20.

$$y = \frac{1}{(1 + e^{-x})} \quad (18)$$

$$p_{q_\pi}(i) = \text{sig}(q_\pi(o_i, a_i)) \quad (19)$$

$$p_{TD}(i) = \text{sig}(|\delta_i|) * 2 - 1 \quad (20)$$

$p_Q(i)$ priority of action-value

$p_{TD}(i)$... priority of TD-error

q_π action-value

o_i observation of sample

a_i actions of sample

δ_i TD error

The full calculation of the priority is presented in equation 21. The variable λ shifts the weight of the priorities from the start with a higher weight for the Q-priority until the end with a higher weight of the TD error to speed up the convergence of the neural network. The value of u_i declines every time this experience is used, which leads to a smaller priority.

$$p(i) = \lambda * p_{q_\pi}(i) + (1 - \lambda) * p_{TD}(i) * u_i \quad (21)$$

$p(i)$ priority of sample

$p_{q_\pi}(i)$ priority of action-value

$p_{TD}(i)$ priority of TD error

λ Prioritization of action-value/TD-error

u_i scale of priority according to the number of using this sample

The sampling of the experience is a combination of random sampling and priority sampling to reduce the time overhead for updating every priority in the replay buffer with a capacity of up to 10^6 samples. The first step is to randomly select samples with a size of $k * n$ and then select samples via HPER sampling with a size of n .

These three different approaches for the replay buffer are investigated within the research environment.

4.3 Noise

The noise in RL-algorithms prevents the algorithm to converge to a local optimum and can be applied as an action noise or as a parameter noise (Plappert et al. 2018). In the original DDPG algorithm an Ornstein-Uhlenbeck noise-process is initialized at the start of each episode and added to the selected action (Lillicrap et al. 2019). The Ornstein-Uhlenbeck noise is a temporally correlated noise visualized in Figure 26 by the blue line compared to a Gaussian noise. As discovered by Barth-Maron et al. the correlated noise has no impact on the performance of the algorithm compared to a fixed Gaussian noise. (Barth-Maron et al. 2018).

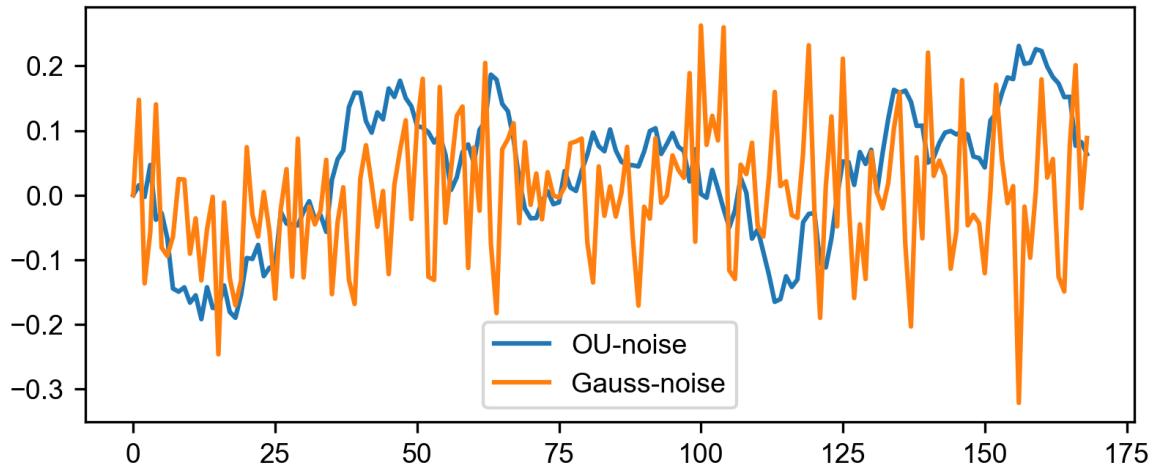


Figure 26: Action-noise process Ohrnstein-Uhlenbeck and Gaussian

An alternative to action noise is to perturbate the network parameters of the actor (Plappert et al. 2018). Gaussian noise is applied to the parameter vector of the policy network at the beginning of every episode. The action obtained by the policy with action space noise is different with a fixed observation as the input because the noise is independent of the observation. With parameter space noise the obtained action will always be same when passing a fixed observation.

Especially in environments with a sparse reward, means not providing a reward at every timestep, the algorithm with parameter space noise succeeded in the task, whereas the algorithm with action noise failed completely. Scaling the Gaussian noise for the perturbed actor is not as intuitive as scaling the actor noise. Plappert et al. introduced an adaptive noise scaling suitable for all RL-algorithm where the scale over time changes over time with a measure depending on the distance between the actor and the perturbed actor.

4.4 State of the art Controller

PID -controller and MPC can be considered as state of the art controller with MPC (Wang et al. 2017). The simplicity and reliability of PID controllers makes them still widely used, even though MPC has proven to perform better for energy savings and cost savings as Gehbauer et al. demonstrated in their study (Gehbauer et al. 2020).

4.4.1 PID Control

PID controller are a simple form of feedback-controller seen in the control loop displayed in Figure 27 shows the PID controller with the three main elements of P, I, and D (Heinrich et al. 2020).

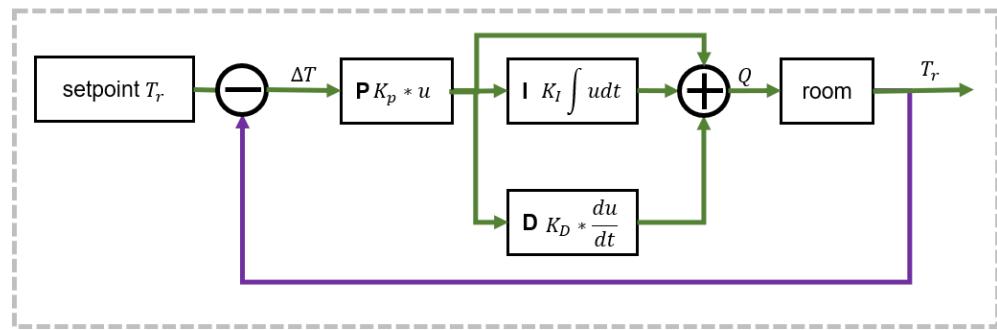


Figure 27: Functional diagram of a PID-controller (modified according to(Heinrich et al. 2020, p.163))

The following icons show the step function after a step for the target value. The step of the target value could be a change of the temperature setpoint, due a time schedule or occupancy sensor. In the control terminology the elements are described with the unit-step response $G(s)$.

Proportional

The P element is a multiplication by a proportional constant with the error between setpoint and the target value (room temperature). This element follows the error without delay (equation 22).

$$G(s) = K_P \quad (22)$$

Integral

With the I element the controller gets more accurate, due to the nature of integration, the control value is not zero if the error is not zero. The target value is reached accurately but the minimization of the error takes longer than with the P element. The unit-step response is given in equation 23.

$$G(s) = \frac{K_I}{s} \quad (23)$$

Derivative

The unit-step response calculated with equation 24 gives the step function of the D element which is an impulse function with a value of zero except at timestep t=0. In combination with a P-element as a PD controller the performance is fast, but the controller is inaccurate, produces high frequent malfunctions.

$$G(s) = K_D * s \quad (24)$$

The combination of P- and I-element or of P-, I-, and D-element is a classic combination for controller as PI-controller or PID-controller. The unit-step response of the PID-controller is specified in equation 25. For a PID controller the equation remains the same, but the derivative constant is set to zero.

$$\begin{aligned} G(s) &= K_P \left(1 + \frac{K_I}{K_P s} + \frac{K_D}{K_P} s \right) \\ G(s) &= K_P \left(1 + \frac{1}{T_I s} + T_D s \right) \end{aligned} \quad (25)$$

$G(s)$... unit-step function

K_P proportional constant

K_I integration constant

K_D derivative constant

s operator for the derivative by time d/dt

T_I reset time

T_D rate time

Setting

The configuration of the PID controller parameters can be done empirically by analyzing the step response and apply the equations 26 of Ziegler and Nichols with the tuning parameters given by the step response in Figure 28.

$$\begin{aligned} K_P &= 0.9 * \frac{T_b}{K_S T_e} \\ T_I &= 3.3 * T_e \\ T_D &= 0.5 * T_e \end{aligned} \quad (26)$$

T_b time constant

T_e delay time

K_S gain

T_I reset time

T_D rate time

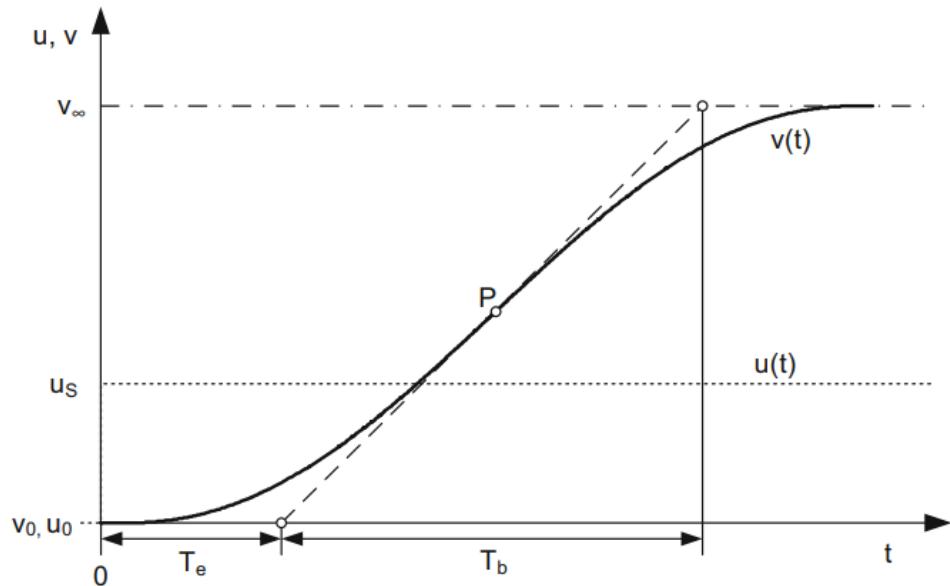


Figure 28: Step response with aperiodic course (Heinrich et al. 2020, p.174)

In this thesis a PI-controller is used to compare it with the developed agent. The parameters of the PI controller in this thesis are:

$$K_P = 10,000$$

$$T_I = 30$$

$$T_D = 0$$

4.4.2 Model Predictive Control

In a perfect world, the predictive control model has the knowledge of all relevant information and optimizes its strategy based on this knowledge. The model built in this thesis is a so-called perfect information model and is used to evaluate the agent in the development process. The perfect information model is built with numerical functions and is a twin of the RC-model built in python as the environment of the RL-setup. In Figure 29 the information flow in the model and the constraints and penalties are shown.

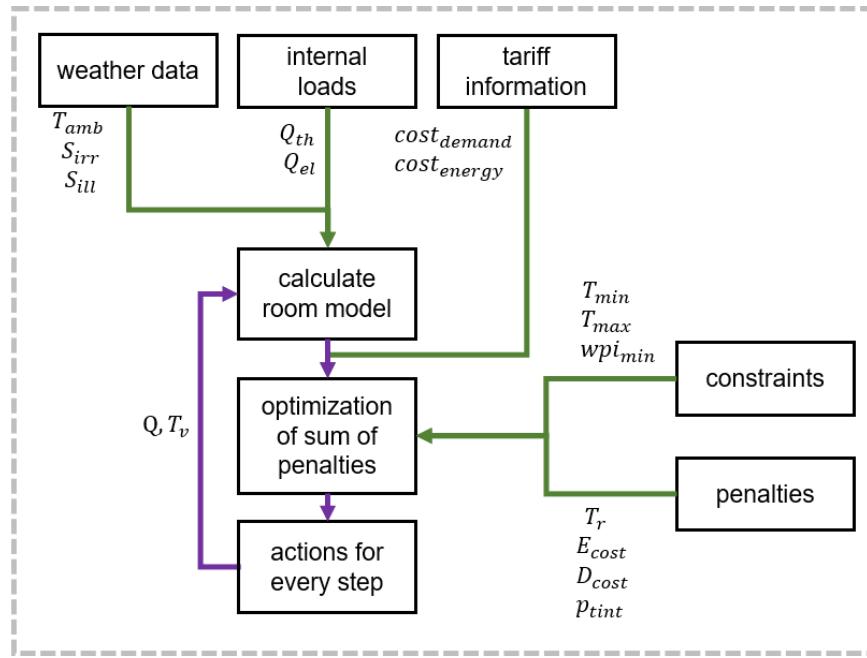


Figure 29: Information flow in the perfect knowledge MPC model

- T_{amb} outside air temperature
- S_{irr} solar irradiation on the tilted window
- S_{ill} global horizontal illuminance
- Q_{th} thermal internal load
- Q_{el} electrical internal load
- $cost_{energy}$ tariff information energy costs
- $cost_{demand}$... tariff information demand costs
- T_r room temperature
- E_{cost} energy costs
- D_{cost} demand costs
- p_{tint} penalty for tinting the window
- T_{min} minimum room temperature
- T_{max} maximum room temperature
- wpi_{min} minimum workplace illuminance
- Q energy input
- T_v visibility through EC-window

4.5 Room Model

For this thesis a medium office building, based on a study conducted by the National Renewable Energy Laboratory is the basis of the building properties used for developing the agent (Deru et al. 2011). The reference building has the form parameters of a medium office building which corresponds to a mass or steel construction. These parameters are summarized in Table 1.

Table 1: Reference Building Form Assignments (Deru et al. 2011, p.19)

Floor Area		Aspect Ratio	No. of Floors	Floor-to-Floor Height		Floor-to-ceiling Height		Glazing Fraction
ft ²	m ²			ft	m	ft	m	
53,628	4,982	1.5	3	13	3.96	9	2.74	0.33

The energy relevant specifications of medium office buildings are shown in Table 2.

Table 2: U-Value by Reference Building Vintage - Standard 90.1-2004 (Deru et al. 2011, p.26)

	Btu/h*ft ^{2*} °F	W/m ^{2*K}
Roof	0.034	0.1936
Wall	0.580	3.294
Window	1.22	6.927

The single office room controlled in this thesis (Figure 30 in green) has an area of 14 m² and a window with a size of 5.2 m² which corresponds to a typical window to wall ratio according to a study conducted by the U.S. Department of Energy of 33 % (Deru et al. 2011).

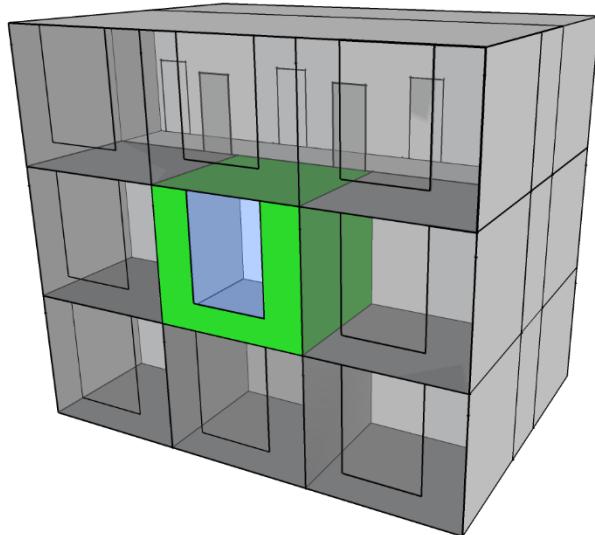


Figure 30: Room model (green)

The resistance value (R-value) is calculated applying the U-values and the respective wall- and window area. The room has no heat loss through ceiling, floor or inside walls. The total capacity of the room is calculated by taking the air properties at 20 °C and by calculating the effective thermal mass of the walls, floor and ceiling following the standard EN-ISO 13786 with the calculation tool developed by HTflux (Rüdisser 2018). The specifications of the room regarding the building envelope is stated in Table 3.

Table 3: specification of the room model

area	14 m ² (150 ft ²)
height	3.95 m (13.12 ft)
window area	5.2 m ²
exterior wall area	10.6 m ²
U-value wall	3.294 W/m ² K (1.22 Btu/h·ft ² °F)
U-value window	6.923 W/m ² K (1.22 Btu/h·ft ² °F)
R-value room	0.014 K/W
C Room	2205 kJ/K

The HVAC system is modelled with a fixed coefficient of performance with 3.5 for cooling and 1 for heating.

4.5.1 Electrochromic Window

The shading device, controlled by the agent is integrated in the glazing of the window as an Electrochromic Window (EC-window). EC-windows are coated with a switchable nanometer-thick (1×10^{-9} m) thin-film which tint can be reversibly changed by applying a small direct current voltage (Lee et al. 2006). The thin film is formed with the following layers:

1. transparent conductor
2. active electrochromic
3. counter-electrode
4. ion-conducting electrolyte

When a bipolar potential is applied to the outside layer (transparent conductor) where lithium ions migrate across the ion-conducting layer from the counter electrode layer to the electrochromic layer. The EC-window is tinted to a Prussian Blue and can be reversed to a clear state by reversing the potential. The window only needs power while changing its tint state and remains unchanged until a voltage is applied. In Figure 31 the principle of an EC-window is shown for the clear and colored state.

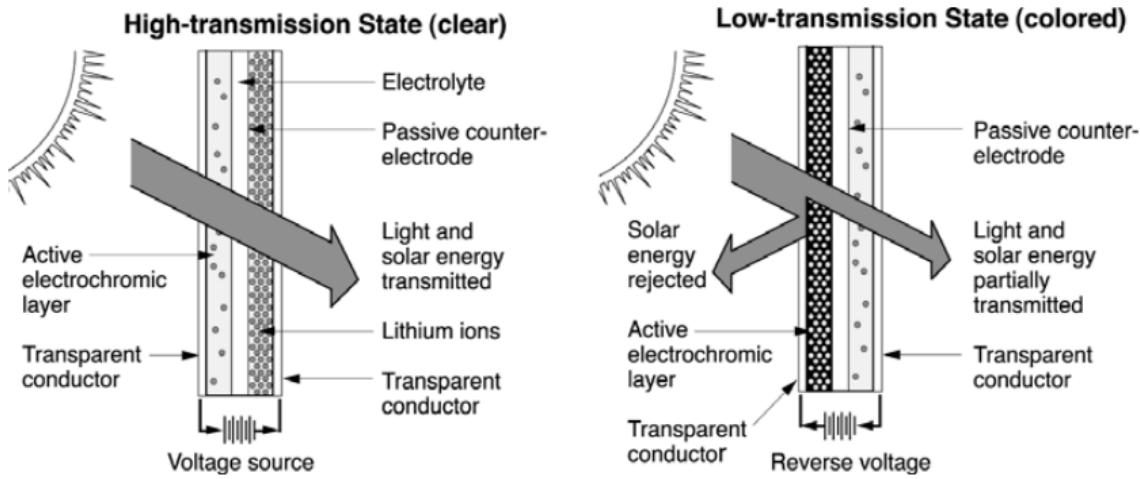


Figure 31: Diagram of a typical tungsten-oxide electrochromic coating (Lee et al. 2006, p.6)

The window can be controlled by changing the visibility transmittance (T_v) in a range of $T_v = 0.6 - 0.01$. Consequently, the solar heat gain coefficient (SHGC) changes accordingly ranging from $SHGC = 0.48 - 0.09$. EC-windows are considered to have the potential for real time optimization in buildings regarding the total energy-and demand costs, the stress on the power grid and occupant comfort due to an undistorted view to the sky.

In Figure 32 the EC-window is shown installed in an office building in Sacramento, CA (Fernandes et al. 2018).



Figure 32: Each window pane had three sub-zones that could be independently controlled (Fernandes et al. 2018, p.14)

The three independent subpanels of the glass enable a better glare control. The Subpanels can be tinted in four discreet states with the glazing properties for the EC-windows used in this study shown in Table 4.

Table 4: Name and visible transmittance of the four tint levels. (Fernandes et al. 2018, p.15)

Tint name	Visible transmittance [%]	Solar transmittance [%]	SHGC [-]	U-value	
				[W/m ² K]	[BTU/ft ² F]
Clear	60	33	0.42	1.816	0.32
Light tint	18	7	0.16		
Medium tint	6	2	0.12		
Full tint	1	0.4	0.1		

The dependency of T_v to SHGC is shown in Figure 33 as a linear and a quadratic function. The SHGC is calculated after taking the action T_v to calculate the solar heat gain. The linear function is chosen to calculate SHGC because the quadratic function would slow the simulation down and has no further advantage over the linear function. The action taken by the agent is continuous and can be any number between 0.6 and 0.01.

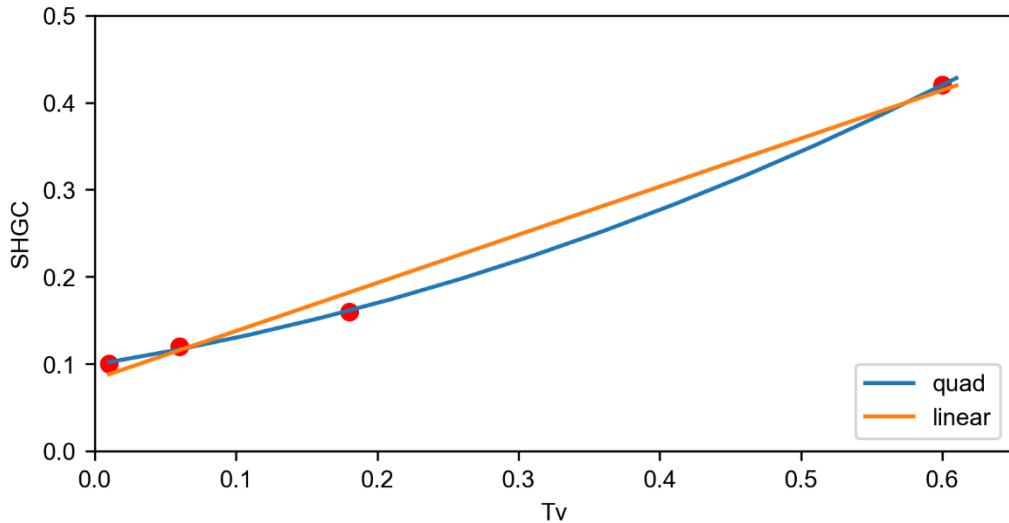


Figure 33: EC-window properties

4.5.2 Solar Position and Radiation

The determination of the solar position and thus the calculation of the incident radiation on the window is necessary for the calculation of the room model. The global horizontal irradiance (GHI), the diffuse horizontal irradiance (DHI) and the direct normal irradiance (DNI) together with the geographical position and the time zone are needed as inputs for the calculation. Starting with the calculation of the real location time t_{WOZ} (equation 29) and the hour-angle ω (equation 30) (Duffie and Beckman 2013). Equation 27 describes the time the

earth traveled on the orbit so far this year in degrees and is used in the equation of time (equation 28) which describes the variable length of the days in the year.

$$B = \frac{360}{365} * (N - 1) \quad (27)$$

B ... travelled distance in degree

N ... day of year

$$E = 229,2 * (0,000075 + 0,001868 * \cos B - 0,032077 * \sin B - 0,014615 * \cos 2B - 0,04089 * \sin 2B) \quad (28)$$

E ... equation of time

B ... distance in degree of the earth on the earth orbit

The real location time is referenced to the standard meridian of the timezone and the latitude of the location. With E the elliptic orbit of the earth is also included in the equation 29.

$$t_{WOZ} = t_{LZ} - DST + \frac{\phi_{BZ} - \phi}{15} + E * \frac{1h}{60min} \quad (29)$$

t_{WOZ} ... real location time

t_{LZ} local time

DST daylight saving time

ϕ_{BZ} standard meridian

ϕ latitude

E equation of time

The hour angle is referenced to the real location time and is negative before noon and positive in the afternoon.

$$\omega = (t_{WOZ} - 12) * 15 \quad (30)$$

ω hour angle

t_{LZ} ... real local time

The orbit of the sun and thus also the position of the sun can be described over several angles, some of them are shown in Figure 34 and described further on.

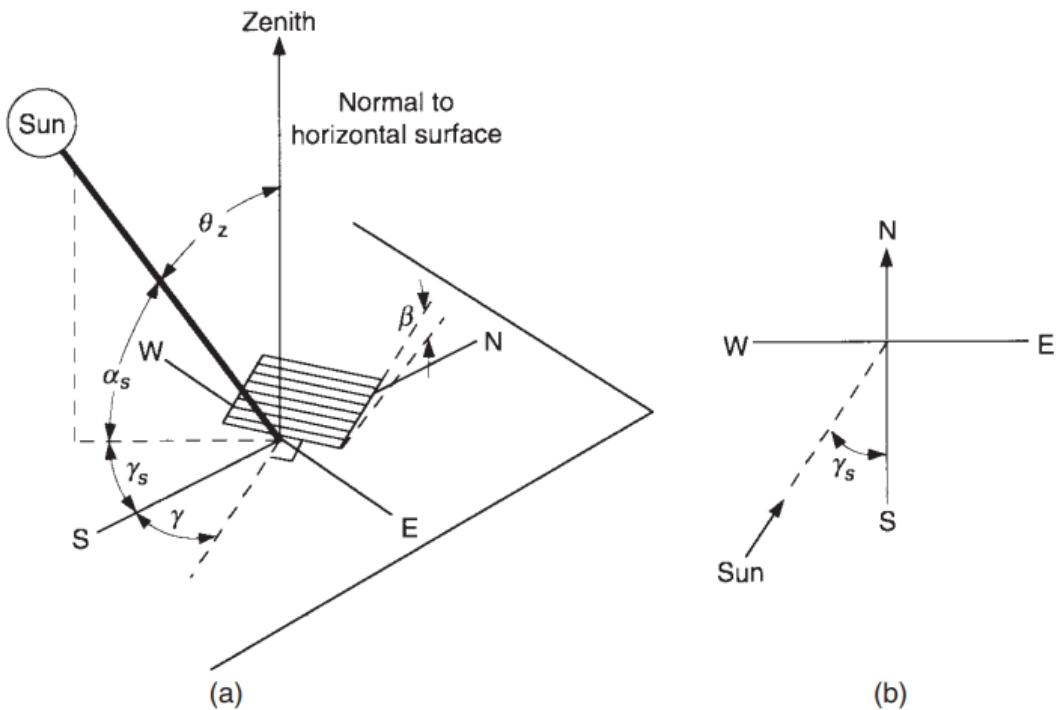


Figure 34: (a) Zenith angle, slope, surface azimuth angle, and solar azimuth angle for a tilted surface. (b) Plan view showing solar azimuth angle (Duffie and Beckman 2013, p.13)

Another angle, not shown in Figure 34 is the declination of the earth which varies between -23° and 23° as seen in Figure 35 and can be described with the approximation by Cooper in equation 31 (Duffie and Beckman 2013). The declination is the angle between the sun at solar noon and a plane on the equator.

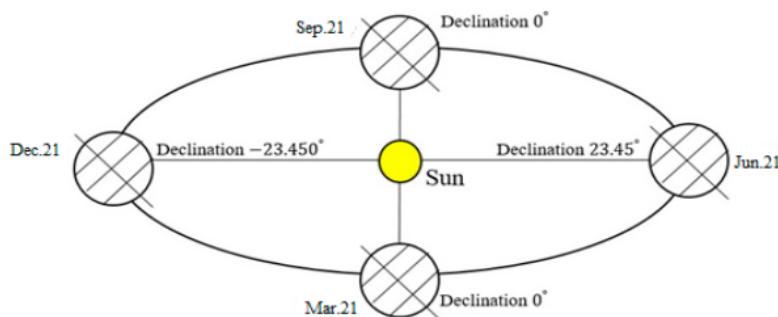


Figure 35: Maximum and minimum value of declination angle (Mousavi Maleki et al. 2017, p.2)

$$\delta = 23.45 \times \sin \left(360 \times \frac{284 + N}{365} \right) \quad (31)$$

δ ... declination

N ... day of year

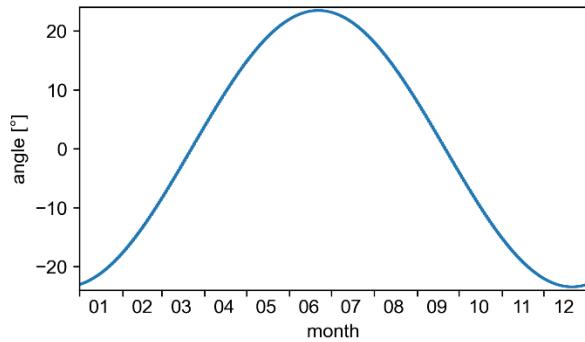


Figure 36: Declination angle in Oakland, CA

The zenith – angle θ_z shown in Figure 34 is a function of the declination δ , latitude ϕ as well as, the hour angle ω .

$$\cos(\theta_z) = \cos(\phi) * \cos(\delta) * \cos(\omega) + \sin(\phi) * \sin(\delta) \quad (32)$$

θ_z ... zenith angle

ϕ latitude

δ declination

ω hour angle

The azimuth angle γ_s is related to south and varies between -180° and 180° which represents before noon and after noon.

$$\gamma_s = \text{sign}(\omega) \left| \arccos \left(\frac{\cos(\theta_z) * \sin(\phi) - \sin(\delta)}{\sin(\theta_z) * \cos(\phi)} \right) \right| \quad (33)$$

γ_s ... azimuth angle

θ_z ... zenith angle

ϕ latitude

δ declination

ω hour angle

With the calculated angles the angle of incidence θ_{Di} can be calculated according to equation 34.

$$\cos \theta_{Di} = \cos(\theta_z) * \cos(\beta) + \sin(\theta_z) * \sin(\beta) * \cos(\gamma_s - \gamma) \quad (34)$$

θ_{Di} ... angle of incidence

γ_s ... azimuth angle

θ_z ... zenith angle

ϕ latitude

δ declination

ω hour angle

γ surface azimuth angle

β slope of the surface (window 90 °)

The GHI is a product of DHI and the DNI dependent on the zenith angle.

$$GHI = DHI + DNI * \cos(\theta_Z) \quad (35)$$

GHI ... global horizontal irradiation

DHI ... diffuse horizontal irradiation

DNI ... direct normal irradiation

θ_Z azimuth angle

The product of equation 36 is the DNI on the tilted surface, calculated with the angle of incidence.

$$DNI_T = DNI * \cos \theta_{Di} \quad (36)$$

DNI_T ... direct normal irradiation on the surface (window)

DNI direct normal irradiation

θ_{Di} ... angle of incidence

The total irradiation on the tilted surface, calculated with equation 37 is the sum of the DNI on the tilted surface, the DHI depending on the angle of the surface in respect to the sky and the GHI depending on the angle of the surface in respect to the ground and the value for ground reflection.

$$I_T = DNI_T + DHI * \left(\frac{1 + \cos(\beta)}{2} \right) + \rho_B * GHI * \left(\frac{1 - \cos(\beta)}{2} \right) \quad (37)$$

I_T total irradiation on the tilted surface

DNI_T ... direct normal irradiance on the tilted surface

DHI ... diffuse horizontal irradiation

GHI ... global horizontal irradiation

β slope of the surface (window 90 °)

ρ_B reflectance of the ground (albedo)

Figure 37 summarizes the calculated solar angles and displays the total solar irradiation on a window oriented to the south for Oakland, CA with a longitude of -122.22 and latitude 37.72 for a window with an orientation with 0° off south and the slope with 90° of the window. The figure shows the calculated values for January 1st and August 1st with the weather data from 2019.

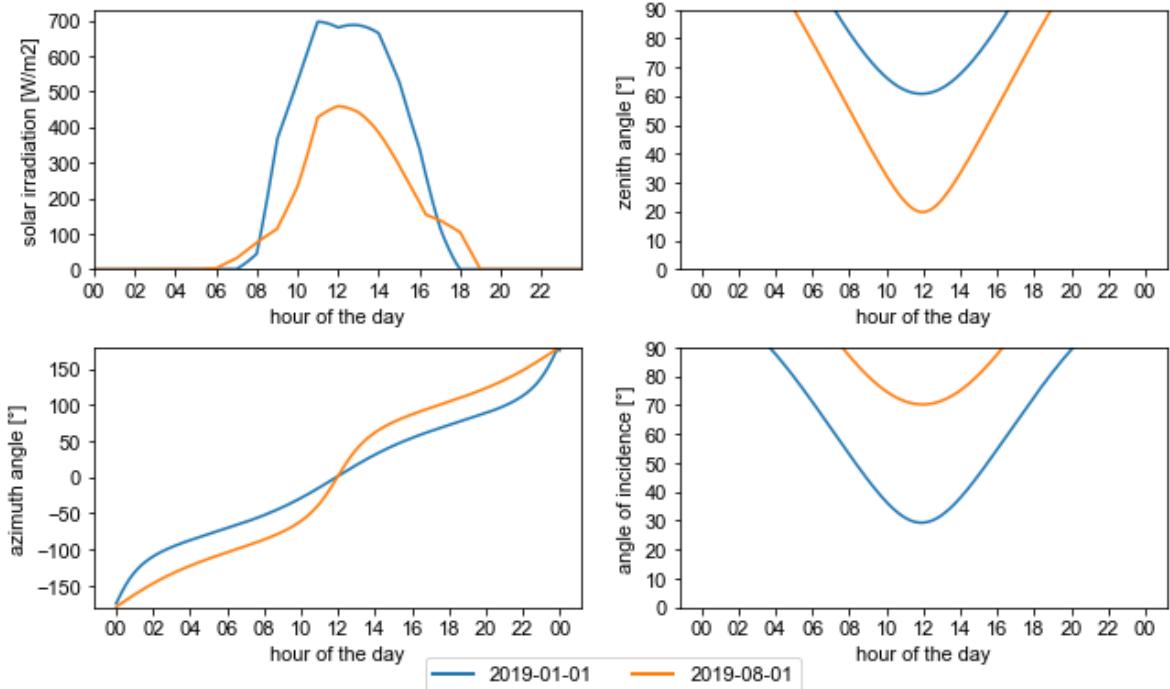


Figure 37: Solar angles and solar irradiation on tilted surface

4.5.3 Electricity Market in California

The master thesis focuses on cost savings for electricity demand. Therefore, the tariff structure of the electricity market in California for economic criteria are analyzed. The focus in this thesis is on the electricity market in Berkeley, Alameda County. This area of California is in the electric utility area of Pacific Gas and Electricity (PG&E) (CEC 2020). The electric power industry is deregulated since 1992 when the U.S. Congress passed the Energy Policy Act and opened the transmission networks to independent energy producers and dissolved the natural monopole of electric utilities (State of California 2018). Due to an energy crisis in 2001 the customer choice has a limited availability. Customers can enter a lottery system if they intend to choose their energy service provider and opt out of from PG&E as the default energy provider in the city of Berkeley.

For commercial customers PG&E offers two rate options with time-of-use (TOU) or peak day pricing (PG&E 2020b). With the PDP rate plans the customer gets discounted electricity rates in the summer in exchange of higher priced peak periods during peak events from 2-6 p.m., which occur during the summer months on the hottest days of the year. PG&E proposes the TOU rates with “Maximize your savings with time-of-use rates”. Since the thesis focuses on reducing the electricity bill the electricity rate is chosen from the TOU plans portfolio. The representative electricity rate is the PG&E E-19 tariff with a winter and summer period with different time schedules and energy prices (PG&E 2020a).

Table 5: E-19 definition of time periods, energy- and demand-costs

SUMMER	May 1 st	October 31 st	Energy cost [\$/kWh]	Demand cost [\$/kW]
Peak	12:00 p.m. - 06:00 p.m.	workdays	0.16225	19.63
Partial-peak	08:30 a.m. - 12:00 p.m. 6:00 p.m. to 09:30 p.m.	weekdays weekdays	0.11734	5.37
Off-peak	09:30 p.m. - 08:30 a.m. 24 hours	weekdays weekends and holidays	0.08846	0.00
WINTER	November 1 st	April 30 th		
Partial-peak	08:30 a.m. - 09:30 p.m.	workdays	0.11127	0.18
Off-peak	09:30 p.m.- 08:30 a.m. 24 hours	weekdays weekends and holidays	0.09559	0.00
Base rate	All year			17.63

The maximum demand is averaged over 15-minute intervals and is calculated and charged monthly. For the demand calculation PG&E uses the maximum demand for each period multiplied with the corresponding costs. The base rate is multiplied with the maximum demand in the month. The bill for the demand costs of one month in the summer could look like Table 6.

Table 6: Example for the demand cost calculation

	Demand [kW]	Demand cost [\\$]
Peak	0.75	14.7225
Partial-peak	1.12	6.0144
Off-peak	0.64	0.00
Base rate	1.12	19.7456
Total demand cost		40.4825

The energy costs are calculated according to the energy consumption every hour corresponding to the TOU-tariff. Figure 38 shows the four different cases occurring in a year.

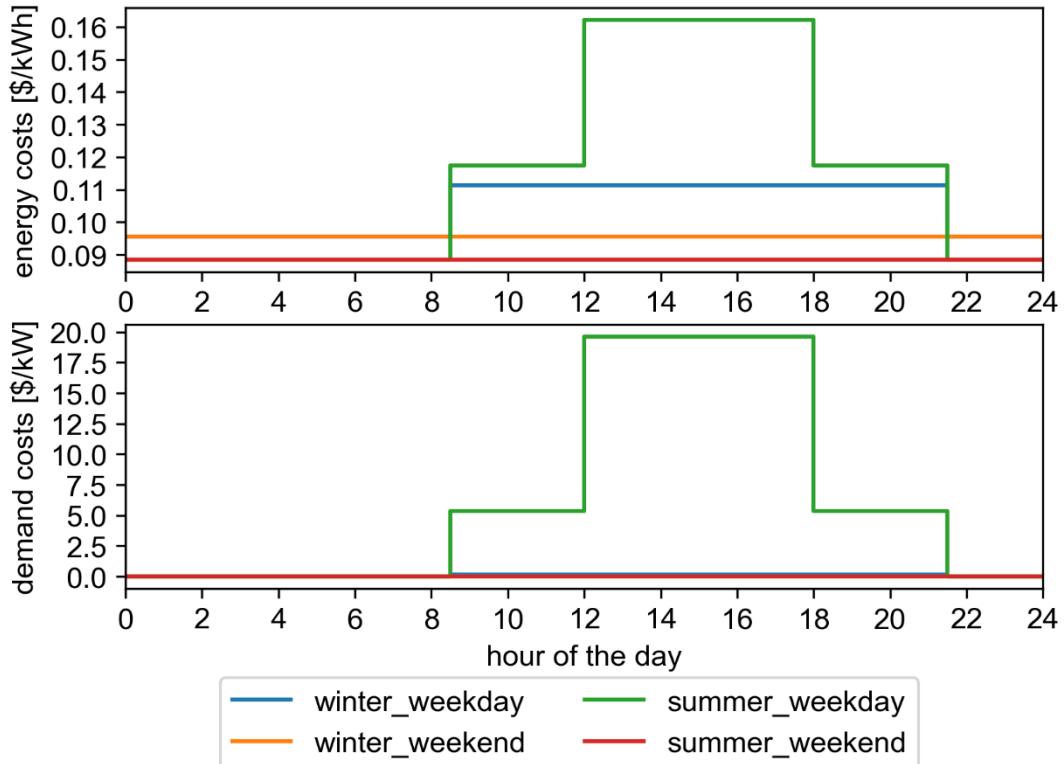


Figure 38: E-19 tariff with time dependent energy- and demand costs (PG&E 2020a)

4.6 RL-Setup

The task of the agent is defined as follows:

Ensure the room temperature within the boundaries of 21 – 24 °C while the room is occupied and 15,5 – 26,5 °C while the room is empty. The workplace illuminance (WPI) should be at least 350 lx for office work. The goal hereby is, to lower the total costs for energy and demand while the constraints are met. The agent can control the shading system by setting the visibility with a linear dependency to the applied current and the heating – and cooling system by controlling the thermal power distributed to the room.

4.6.1 Environment

In the environment in the RL-setup the thermal model and the reward function are defined and calculated. The environment for the development of the agent is a simplified resistance and capacitance (RC) model.

Room- Model

The equation for the RC-model is given with equation 38 and considers the outside air temperature, the room temperature of the previous timestep and the current room

temperature. The solar heat gain on the tilted surface I_T is calculated with equation 37 as described in chapter 4.5.2.

$$Q = \frac{(T_{r(t)} - T_{amb})}{R_r} + C_r * (T_{r(t)} - T_{r(t-1)}) + I_T * SHGC - Q_{int} \quad (38)$$

Q heating- or cooling energy (action of agent)

I_T solar irradiation on the window glazing

$SHGC$ solar heat gain coefficient

Q_{int} internal loads (people, power consumers, artificial lights)

$T_{r(t)}$ current room temperature

$T_{r(t-1)}$... room temperature of last timestep

T_{amb} outside air temperature

C_r capacitance of the room

R_r thermal resistance of the room

For the second task of the agent to ensure the WPI the illuminance in the room has to be calculated. A detailed calculation of the WPI using raytracing is a computer intensive work. In this thesis the goal is to develop an agent and a raytracing calculation exceeds the scope. The Building Research Establishment on behalf of the Department for Communities and Local Government of the United Kingdom developed an analyzing tool for building's energy consumption (BRE 2015). Building Research Establishment calculates the average daylight factor (DF) with total window area and the area of all surfaces in the room (equation 39).

$$DF = 45 * \frac{A_w * T_v}{0.76 * A_{surf}} \quad (39)$$

DF average daylight factor

A_w window area

A_{surf} ... area of all room surfaces (ceiling, floor, walls, and windows)

T_v visibility (action of the agent)

The DF per definition is the ratio between global horizontal illuminance and the average illuminance in the room. Therefore, by calculating the DF with equation 39 the available illuminance in the room is calculated by multiplying the global horizontal illuminance with DF.

The internal loads are the sum of artificial light, power consumers, and the people in the room. The artificial light ensures the minimum level of WPI, therefore the only signal the agent gets for the tint status is the energy consumption of the artificial light. The power consumers per workplace are assumed to be 10.78 W/m^2 with 10 % of standby energy consumption (Deru et al. 2011). The thermal internal load per workplace, equivalent to one

person is 100 W The time schedule for the power consumers and people's presence on weekdays is 07:00 am to 06:00 pm and no occupancy on the weekends is assumed.

Reward

The reward (equation 40) is calculated out of the total costs for energy and demand for all energy consumers as the optimization goal. Furthermore, a penalty for exceeding the room temperature boundaries and a penalty for tinting the EC-window with no solar radiation are included in the calculation. The demand is charged on a monthly basis, therefore the costs per month are scaled to represent the ratio between one hour of energy costs and the monthly demand costs. The penalty for the room temperature is limited to a maximum value of two to prevent the reward deviate too much from the optimal policy especially at the beginning of the training process. The tint penalty is one when the visibility is set to a lower level than 99 % of the maximum visibility level which means no tinting.

$$r = -\frac{|E_{cost}|}{maxE_{cost}} - \frac{|D_{cost}|}{maxD_{cost}} * scale_D - max(|T_{r(t)} - T_{const}|, 2) - p_{tint} \quad (40)$$

r reward

E_{cost} energy costs

E_{max_cost} ... maximum energy costs

D_{cost} demand costs

D_{max_cost} ... maximum demand costs

$scale_D$ scale factor for the demand costs

$T_{r(t)}$ current room temperature

T_{const} temperature boundary (min, max)

p_{tint} penalty for tinting the window

4.6.2 Development

The development of the agent includes the selection of the algorithm, the network architecture with its input values, and the replay buffer to improve the agent. The action space for the energy input Q is set with a maximum specific heat- and cooling power of 100 W/m² and the visibility T_v , with the properties shown in 4.5.1 with action boundaries of 0.01 to 0.6. The observation of the agent contains the state of the room and a forecast including the outside air temperature, solar radiation, costs of energy, and the occupancy of the room.

The training process of the agent runs within episodes with a length of one day and a total of 3,000 episodes. For each episode, the start day is selected randomly from the selected weather dataset and a random start state (room temperature) within the temperature boundaries. The agent is trained with the weather data set of Oakland Intl AP 724930, distributed by EnergyPlus (EnergyPlus 2019).

The basic setup of the agent is based on the experimental setup of the DDPG with a MLP Network with 2 hidden layers with a layer size of 400 and 300 respectively (Lillicrap et al. 2019). The structure of both neural networks is the same, with the difference that the action is added to the critic network after the first hidden layer. The activation function for the hidden layers is the relu function and for the output layer of the critic a linear function is applied. For the actor, the activation function for the output of Q is tanh with values between -1 and 1 and the output for T_v , with a relu function with a maximum value of 1 is utilized. The learning rate was chosen to be 10^{-4} for the actor and 10^{-3} regarding the critic to ensure, that the critic converges faster than the actor. The soft update for the target networks is set to 0.001.

Furthermore, the magnitude of the input values is a critical parameter for the neural network. This in regards, all input values are normalized between -1 and 1.

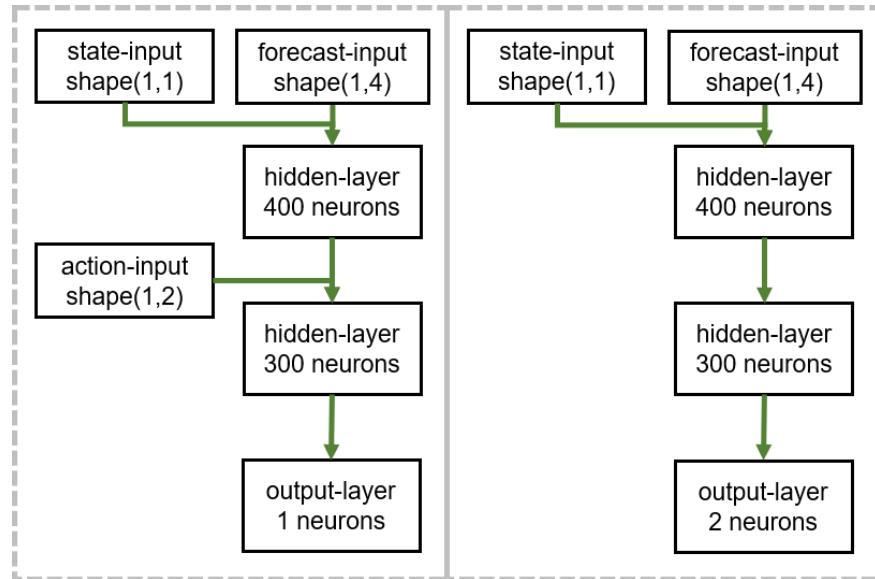


Figure 39: Neural network architecture; critic left and actor on the right with the shape of the input vectors

The default action noise in the DDPG algorithm is the Ornstein-Uhlenbeck process, which is initialized for both actions with a different scale, due to the varying action spaces with 0.15 for Q and 0.1 for T_v . 64 samples for each training step are selected randomly from the replay buffer.

The following figures represent the agent after the training process for one test week starting from August 1st or January 1st. The figures are structured as follows:

- The weather data is represented in the first graph including:
 - The solar radiation with (GHI, DHI, DNI)
 - The outside air temperature (T-out) on the right y-axis

- The second graph shows the thermal power of the HVAC system whether its heating or cooling.
- The third graph is the tint state of the EC-window
- The fourth graph shows the room temperature and the temperature boundaries with the setpoints for the time the room is occupied and not.

In the first training run the agent is limited to one action with a fixed T_v to 0.6 to proof if it can succeed. After the training run the agent with the basic settings of the DDPG algorithm is able to maintain the room temperature most of the time, but has problems in the morning hours when the minimal room temperature increases from 15.5 °C to 21 °C (Figure 40).

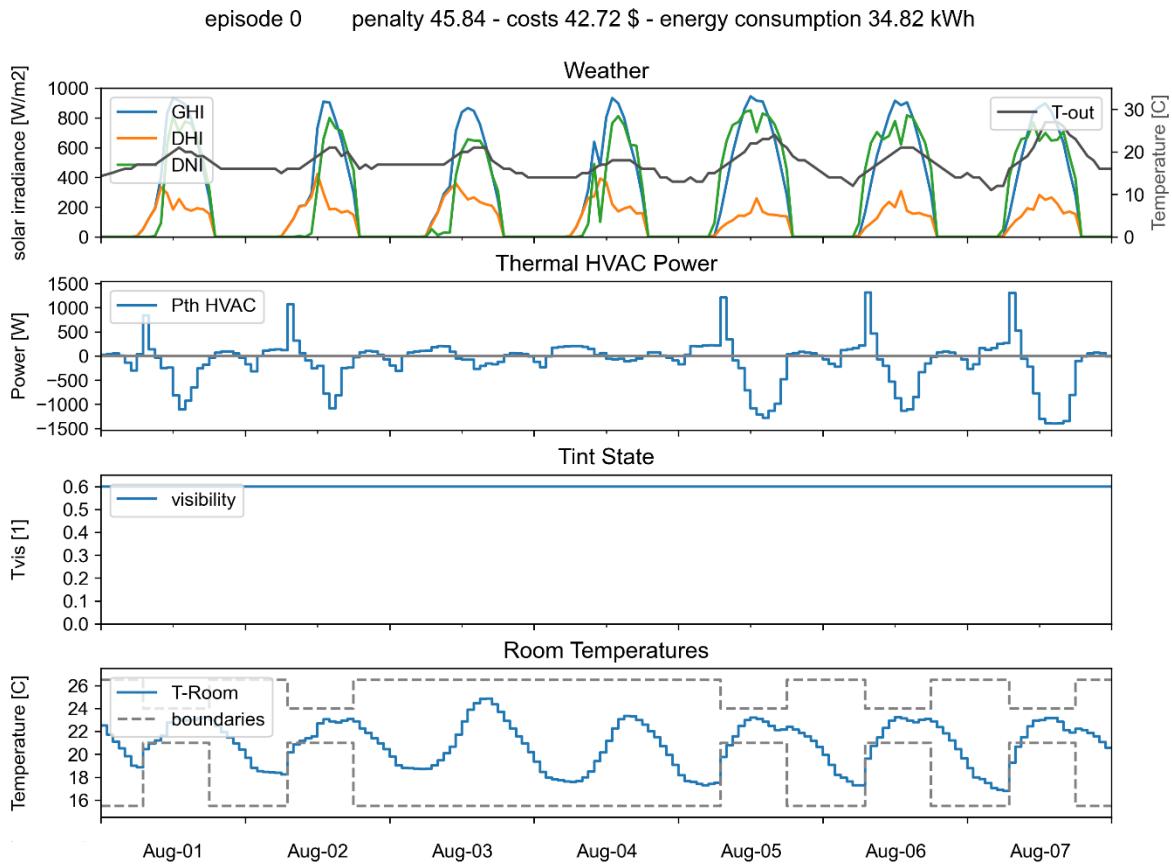


Figure 40: First training result of a week starting on August 1st with HVAC control and a fixed Tint state

Moving on with the development the agent must control both possible actions. With the same setup as before the agent does not succeed in its task (Figure 41). The agent is not eager to heat the building, despite the fact that the temperature constraints are not met, and only does that on the weekends.

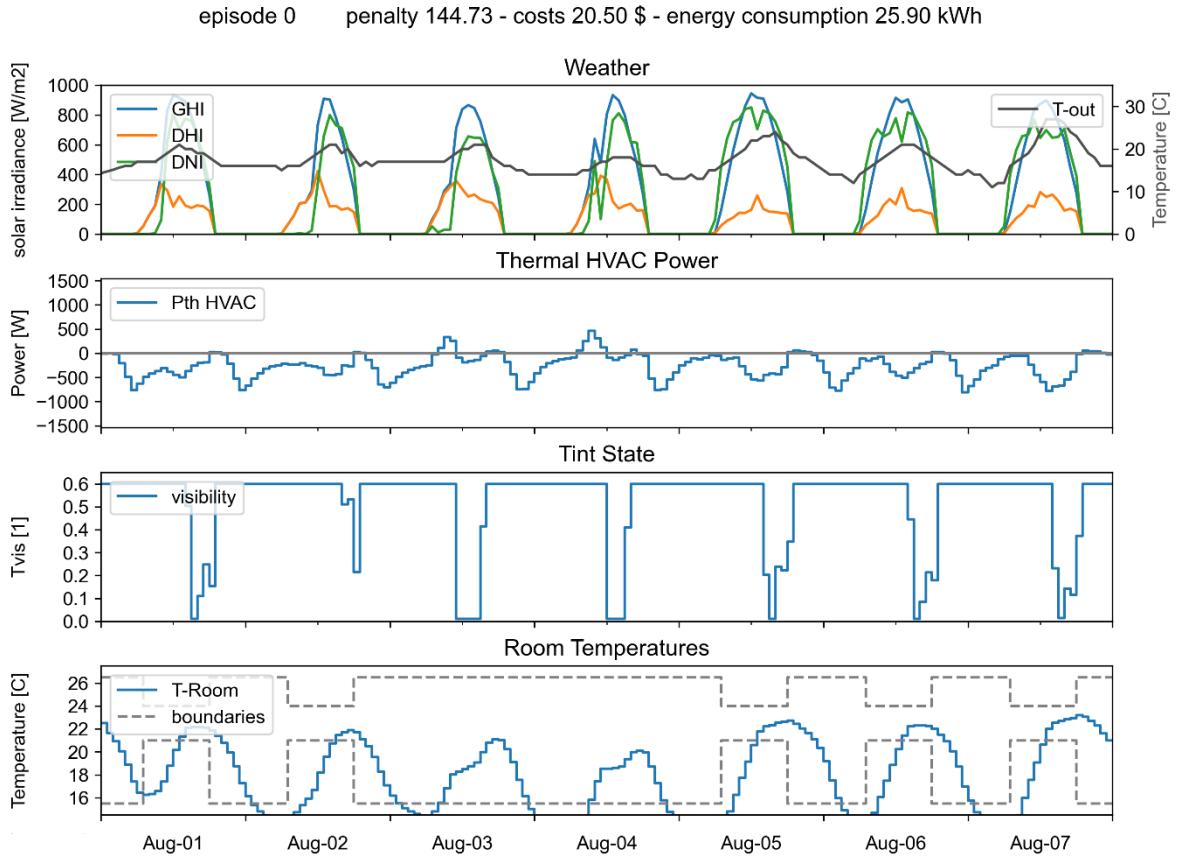


Figure 41: Training result of a week starting on August 1st with HVAC and Tint state control

Investigating the reason, why the agent failed, the loss as an accuracy measurement of the critic was analyzed. Figure 42 indicates that the critic loss increases over time but stabilizes at the end of the run.

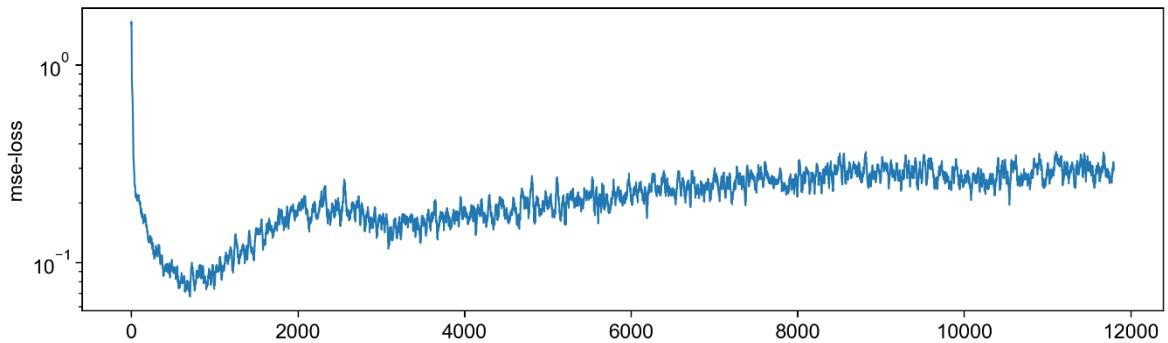


Figure 42: Critic loss of a week starting on August 1st with HVAC and Tint state control

To get a more accurate critic which is leading to a more successful agent, the neural network architecture of Fujimoto et. al. proposed in the TD3 algorithm in 2018 is implemented with the critic architecture including the action in the same layer as the state-and forecast input (Figure 43) while the rest of the network remains unchanged (Fujimoto et al. 2018).

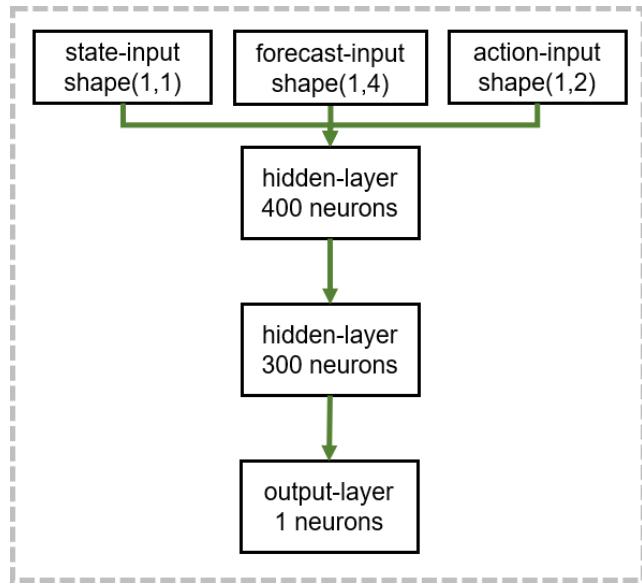


Figure 43: New critic network based on the TD3 algorithm with the shape of the input vectors

After another training run with the new critic network, the critic loss is lower (Figure 44), which indicates that the critic is more accurate and stabilizes after 3,800 training steps.

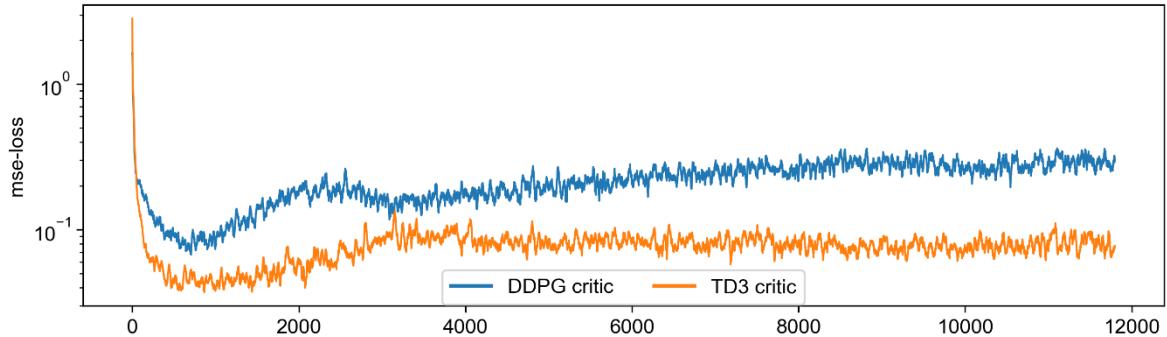


Figure 44: Comparison of critic loss between DDPG critic and TD3 critic architecture

The agent is successful, regarding the temperature constraints by controlling the energy input but fails to control the tint state Figure 45. Following the reward function, the agent should select the brightest tint state during the night to avoid getting penalized for tinting the window. The agent does not find the correct way to handle the reward function. The behavior of the agent, regarding the HVAC system is not energy saving by any means. The agent heats the room starting in the night until the room temperature gets close to the upper boundary and then starts cooling the building Figure 45. The positive thing out of this analysis is, that the agent recognizes the constraints.

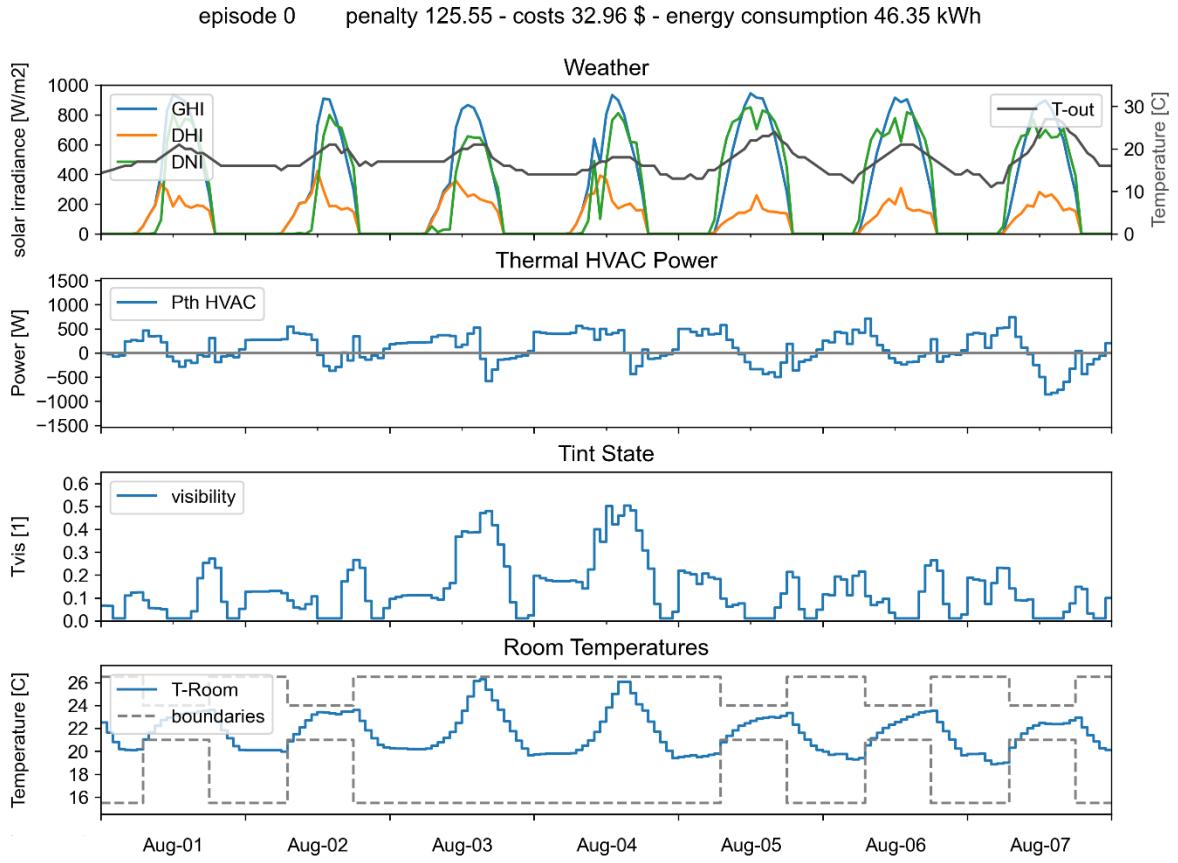


Figure 45: Training result with the new critic of a week starting on August 1st with HVAC and Tint state control

A reason for this behavior may be the lack of forecast data. With forecast data, the agent should be able to change its behavior based on future data. The training process should guide the agent, which timestep is the most important to learn a control strategy that optimizes the energy costs and keep the room temperature within the boundaries. Due to the low capacitance of the room the long-term dependency is low, and thus, leads to the decision of four hours of forecast. The network architecture remains the same, but the forecast input now has 16 values for the four-hour forecast (Figure 46).

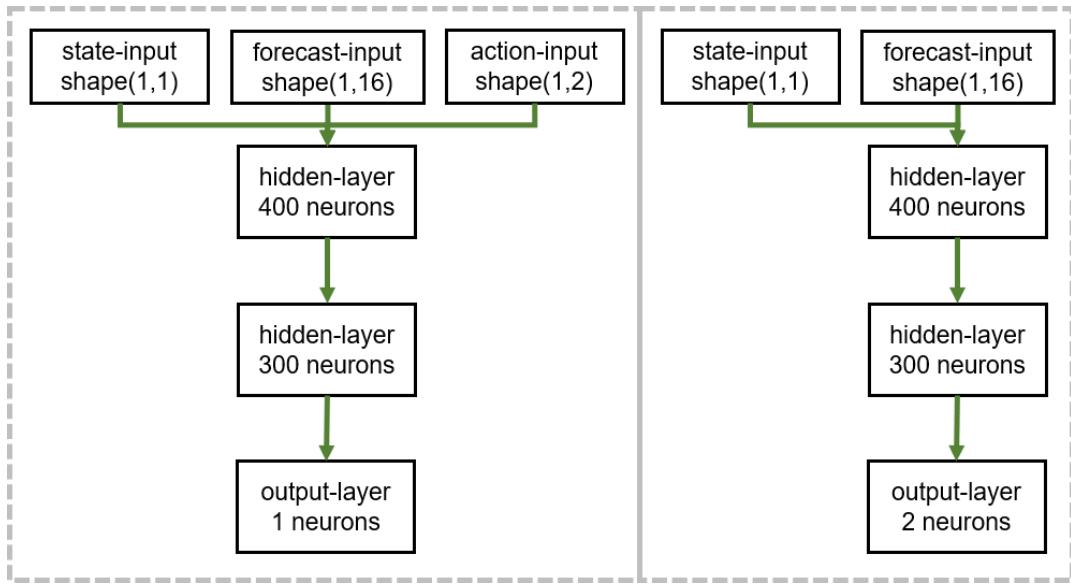


Figure 46: Neural network setup with 4 forecast hours with the shape of the input vectors

The critic loss, of the first training run with multiple hours of forecast is lower, than with only current values of the forecast inputs.

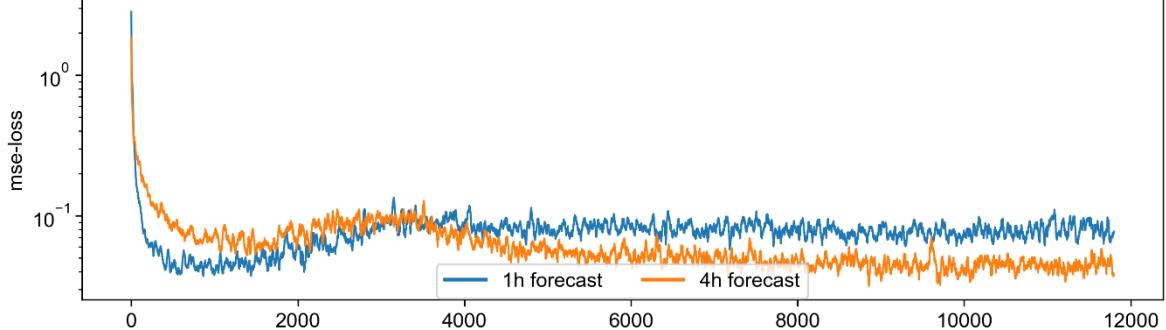


Figure 47: Comparison of critic loss between 1h and 4h forecast

It is not clear why the agent fails in the control task since the loss of the four hour forecast is lower (Figure 47). The agent failed to maintain the room temperature before noon but performs better in terms of penalty and the taken actions do not fluctuate as much as with one hour of forecast, as can be seen in Figure 48.

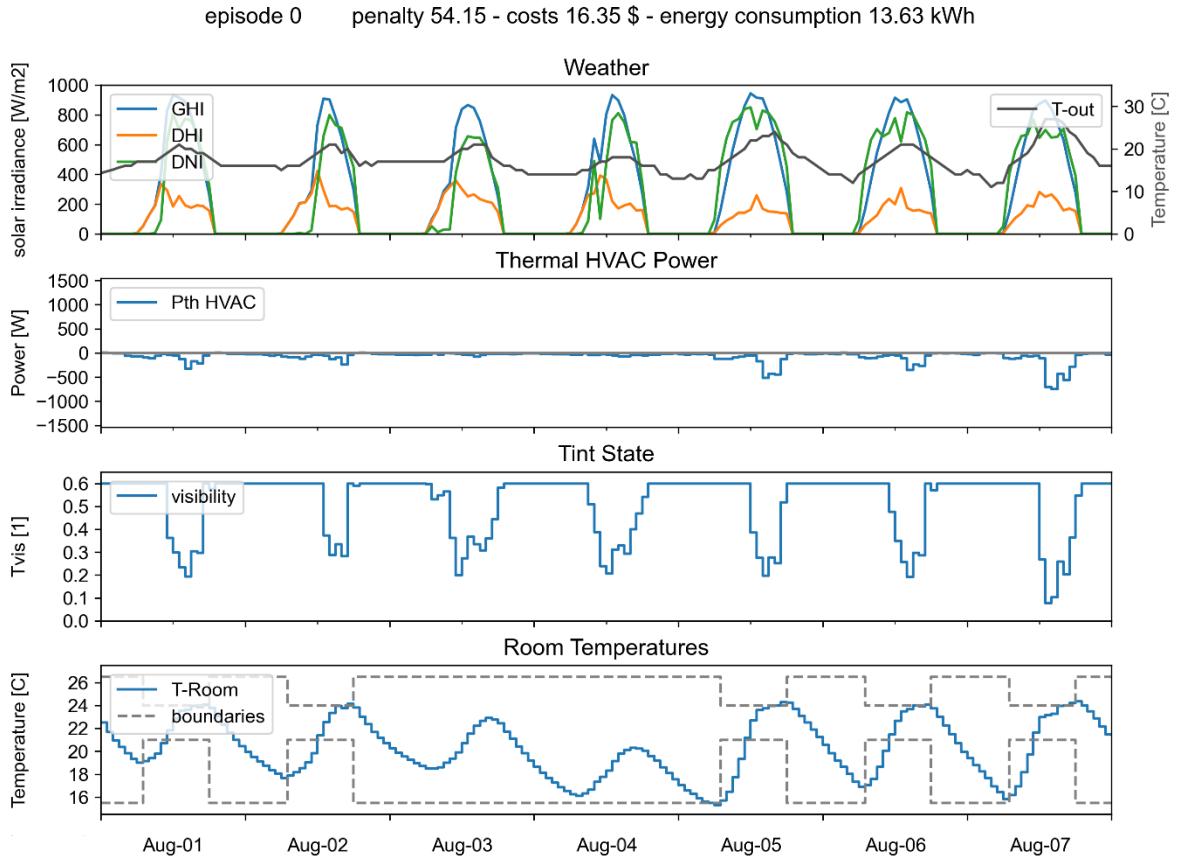


Figure 48: Training run with four forecast hours of a week starting on August 1st with HVAC and Tint state control

The solution in this case is not obvious, therefore a Gridsearch is performed, where the network configuration in terms of neurons per layer and number of hidden layers is tested in all possible combinations of one or two hidden layers and a layer size of 300 to 600 neurons with a step size of 100. The best results of the Gridsearch in Table 7 compared with the reference of the best run so far, show that the critic with two hidden layers tends to be more accurate and the penalty for the test run with a larger first layer than the second layer is lower.

Table 7: Gridsearch results of best neural network configurations

ID	hidden layer	layer size 1	layer size 2	critic loss	test Aug 1 st	test Jan 1 st
Ref	1	400	300	0.026	54.15	164.58
30	2	500	300	0.101	32.81	71.78
07	1	600	400	0.491	35.19	102.41
27	2	400	300	0.245	35.42	81.78
06	1	600	300	0.232	37.83	141.14
25	1	600	400	0.277	39.42	110.48

The results of the Gridsearch in Figure 49 indicate, that long term dependencies were not taken into account by the agent and thus, tend to react too slow on changes of outside air temperatures. Herein, all agents manage to keep the room temperature within the boundaries with a similar behavior. The best EC-window control was achieved by the agent with jobID 07 which keeps a brighter tint state of the window from August 3rd to August 4th to keep the room temperature higher compared to the other agents.

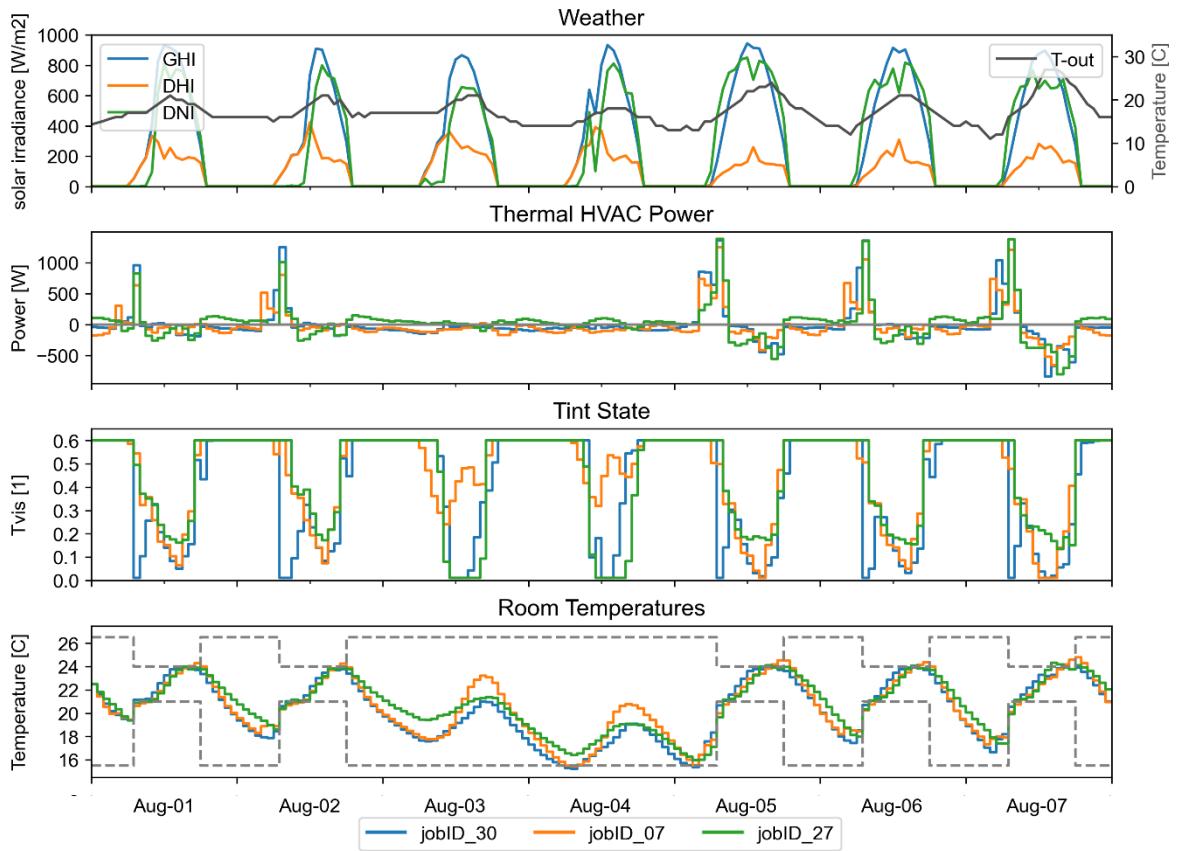


Figure 49: Comparison of best Gridsearch results of a week starting on August 1st with HVAC and Tint state control

The results of the performance in terms of the costs, penalty and the maximum peak load, of the best results of the Gridsearch, are summarized in Figure 50. The agent with jobID 07 has the lowest peak load with 1.41 kW with a small difference to the other two runs which have a peak load of 1.54 kW and 1.58 kW. Based on these performances and general policy of the agent during the test week, the next improvement step is conducted with all three agents.

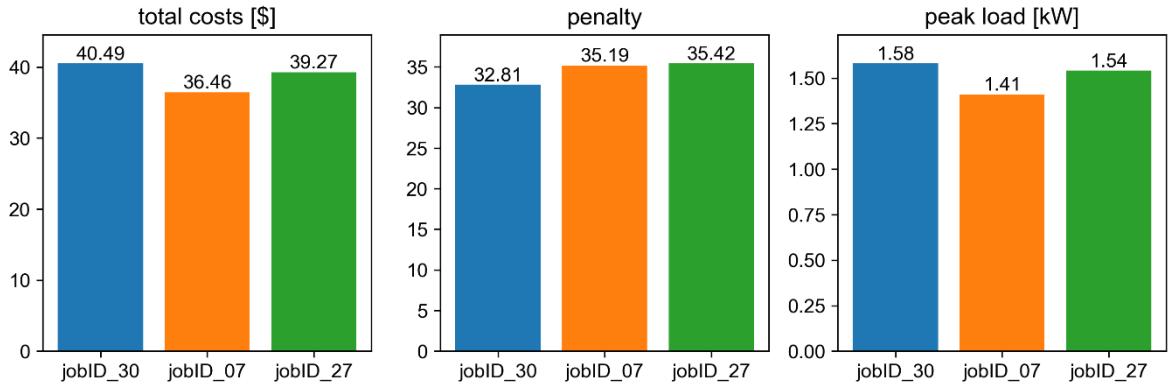


Figure 50: Performance measure of the best Gridsearch results

As already described the four-hour forecast does not lead the agent to a successful strategy and increasing the forecast to eight hours does not contribute to an improvement either. A possible way to overcome the issue of a too short dependence, is to use the N-step reward for calculating the action-values with the critic. A variation of the DDPG algorithm was proposed in 2018 called the Distributed Distributional Deterministic Policy Gradient which outperforms the DDPG algorithm (Barth-Maron et al. 2018). The use of the N-step reward had the greatest influence on their performance and was most successful with a length of 5 steps. The N-step reward is the sum of discounted rewards of a fixed length and is calculated with equation 41.

$$Y_t = \left(\sum_{n=0}^{N-1} \gamma^n r_{t+n} \right) + \gamma^N * Q'(o_{N+1}, \mu'(o_{N+1} | \theta^{\mu'}) | \theta^{Q'}) \quad (41)$$

Y_t action-value as target

N length of N-step reward

n step in N-step reward

γ^n discount factor

o_{N+1} ... observation of timestep t+1 (state and forecast values)

r_{t+n} reward of timestep t

Q' target critic

$\theta^{Q'}$ parameters of target critic

μ' target actor

$\theta^{\mu'}$ parameters of target actor

The calculation of a three N-step reward looks like following example with the transition from the start step with time t to the time step t+2 as the last step.

$$Y_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2}$$

The trajectories with the N-step reward are stored with the observation of the current timestep (s_t, f_t) with the timestep after the N-step reward (s_{t+N+1}, f_{t+N+1}) as $(s_t, f_t, a_t, s_{t+N+1}, f_{t+N+1}, Y_t)$. The trajectories are stored for every timestep, to gather as many trajectories as possible for the training process and not have any gaps in the stored data.

The Gridsearch for the best fitting N-step reward will be run with a possible N-steps of 2, 3, 4, 5. The best results indicate, that the optimal length of the N-step for this problem is four steps of the run jobID 27, by taking both, the test week starting on August 1st and the week starting on January 1st into account. However, the critic loss for the run is higher, than of the run with 2 N-steps because a longer N-step reward makes it hard for the critic to estimate the action-value. The critic has no further information, of the next states and which actions are taken to reach the current state.

Table 8: Gridsearch results of best N-step reward

ID	hidden layer	layer size 1	layer size 2	N-step	critic loss	test Aug 1 st	test Jan 1 st
Ref	2	500	300	1	0.101	32.81	71.78
30	2	500	300	2	0.108	33.15	106.79
27	2	400	300	4	0.431	39.04	86.60
07	1	600	400	2	0.283	41.92.30	129.78

With the N-step reward system, the agent uses the forecast data to his advantage and precools or preheats the room displayed in Figure 51. The maximum peak demand can be lowered by all agents compared to the 1 step reward used in the basic DDPG algorithm. The agent with jobID 07 has promising behavior for tinting the EC-window and the oscillation of the thermal HVAC power on the weekend is the lowest but fails to keep the room temperature in the boundaries. This agent reacts always a bit slower than the two others. A non-optimal behavior, regarding the tinting of the EC-window is seen by the agent with jobID 30. This agent tints the window on the weekend what leads to lower room temperatures and a higher demand to heat up the building. The N-step length of four of the agent with jobID 27 leads to a farsighted behavior and a similar good tinting behavior as the agent with jobID 07.

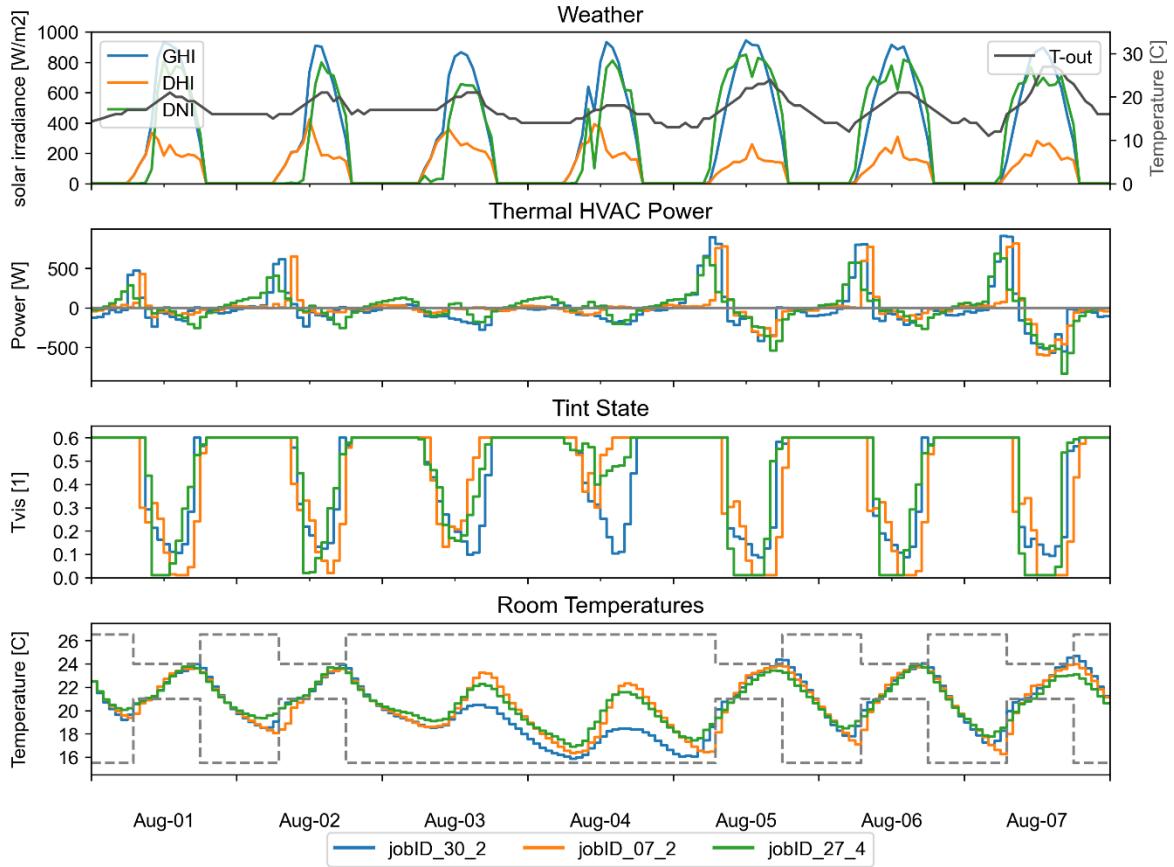


Figure 51: Comparison of Gridsearch results for different N-step rewards starting on August 1st

The performance measures are compared in Figure 52 and show that the agent with jobID 27 with a N-step length of four has the highest cost saving potential and has the lowest impact on the power grid.

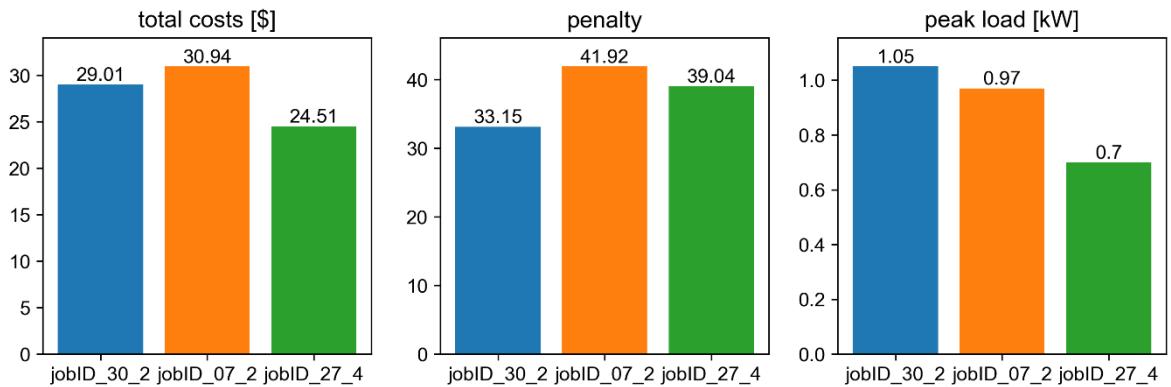


Figure 52: Performance measure of the best Gridsearch results for different N-step rewards

The improvement introduced in chapter 4.1 with the replay buffer, noise process and the activation function of the hidden layers are applied in the final run with the agent with jobID 27 as the most promising. Following options are possible for the improvements:

- Activation function
 - Rectified Linear Unit – relu
 - Leaky Rectified Linear Unit – lrelu
- Replay Buffer
 - Uniform
 - Prioritized Experience Replay – PER
 - High-Value Prioritized Experience Replay – HVPER
- Noise process
 - Ornstein Uhlenbeck noise – OU
 - Gaussian noise – Gauss
 - Parameter noise – Param

Table 9: Gridsearch results of best improvements to the agent

ID	activation function	replay buffer	noise process	critic loss	test Aug 1 st	test Jan 1 st
Ref	relu	Uniform	OU	0.026	39.04	86.60

The N-step reward introduced for the MLP network manages to take long term dependencies into account but misses knowledge of the steps taken after the initial step. An algorithm developed for time series dependent problems published by Google DeepMind is the Recurrent Deterministic Policy Gradient with a LSTM network for both neural networks (Heess et al. 2015). For this algorithm the entire history $(o_1, a_1, o_2, a_2, \dots, a_{t-1}, o_t)$ is used for selecting actions with the deterministic policy μ . The critic network therefore is initialized as $Q(h, a|\theta^Q)$ and the actor as $\mu(h|\theta^\mu)$. Same as in the DDPG algorithm noise is added to the selected actions to explore the continuous action space. The value function introduced in chapter 4.1 DDPG stays unchanged and is calculated for every step.

The neural network architecture is based on the experimental setup of Song et al. published in 2019. In their paper, the inputs for the neural networks were based on a combination of pictures and numbers. Since that is not the case for this thesis, the layers dedicated to the pictures are not used. The adapted architecture is shown in Figure 53 with the critic and actor network. The inputs for the critic are the observation-and action history and for the

actor only the observation history is passed. The forecast in the observation is passed with the current values.

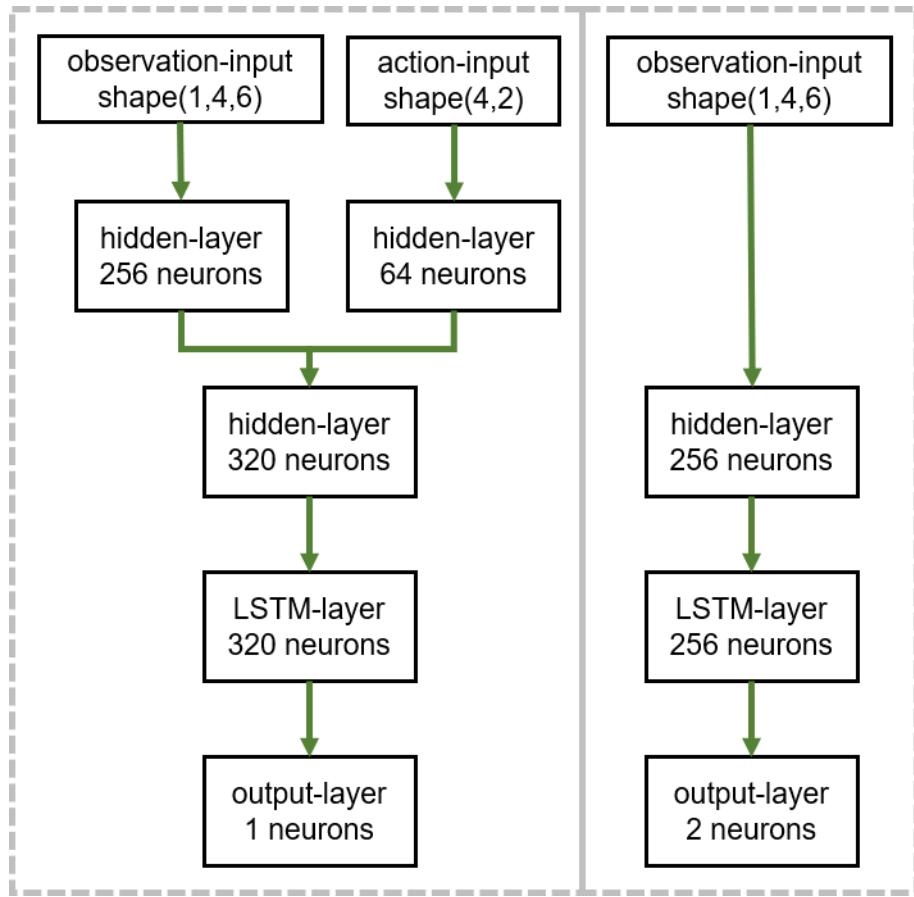
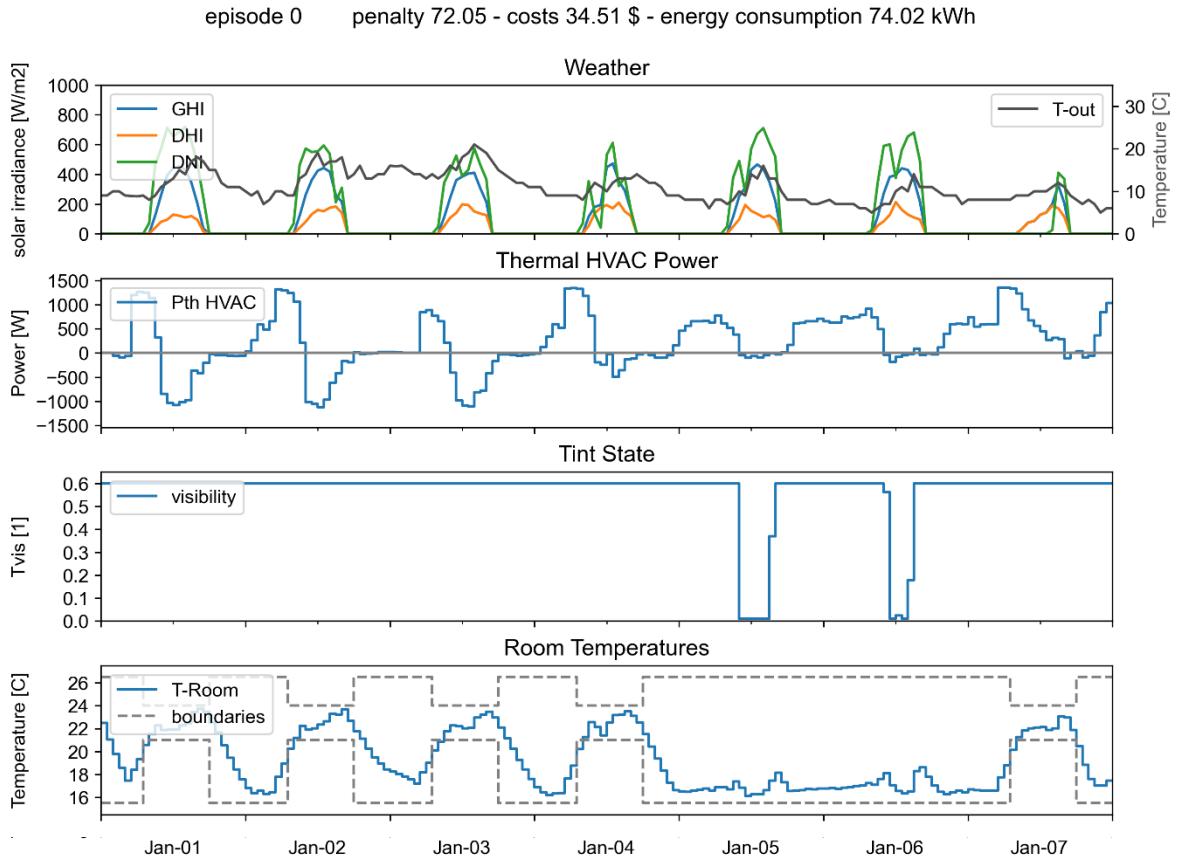


Figure 53: Neural network setup for the RDPG with the shape of the input vectors

With the described setup and the beforehand selected improvements for the activation function, the replay buffer, and the noise process, the agent with the RDPG algorithm is successful in keeping the room temperature between the boundaries. The taken actions for the EC-windows, however, are not beneficial for cost saving. The lack of forecast information also leads to high peak loads for heating and cooling.



The same approach as with the DDPG algorithm of four forecast hours as inputs does not lead to any improvements but leads to a failure of the agent.

Therefore, the latest DDPG agent is the best performing agent and is compared with the PI-controller and the MPC in Figure 54. It gets clear, that the agent has not the same foresight, as the MPC but can decrease the maximum peak compared to the PI-controller.

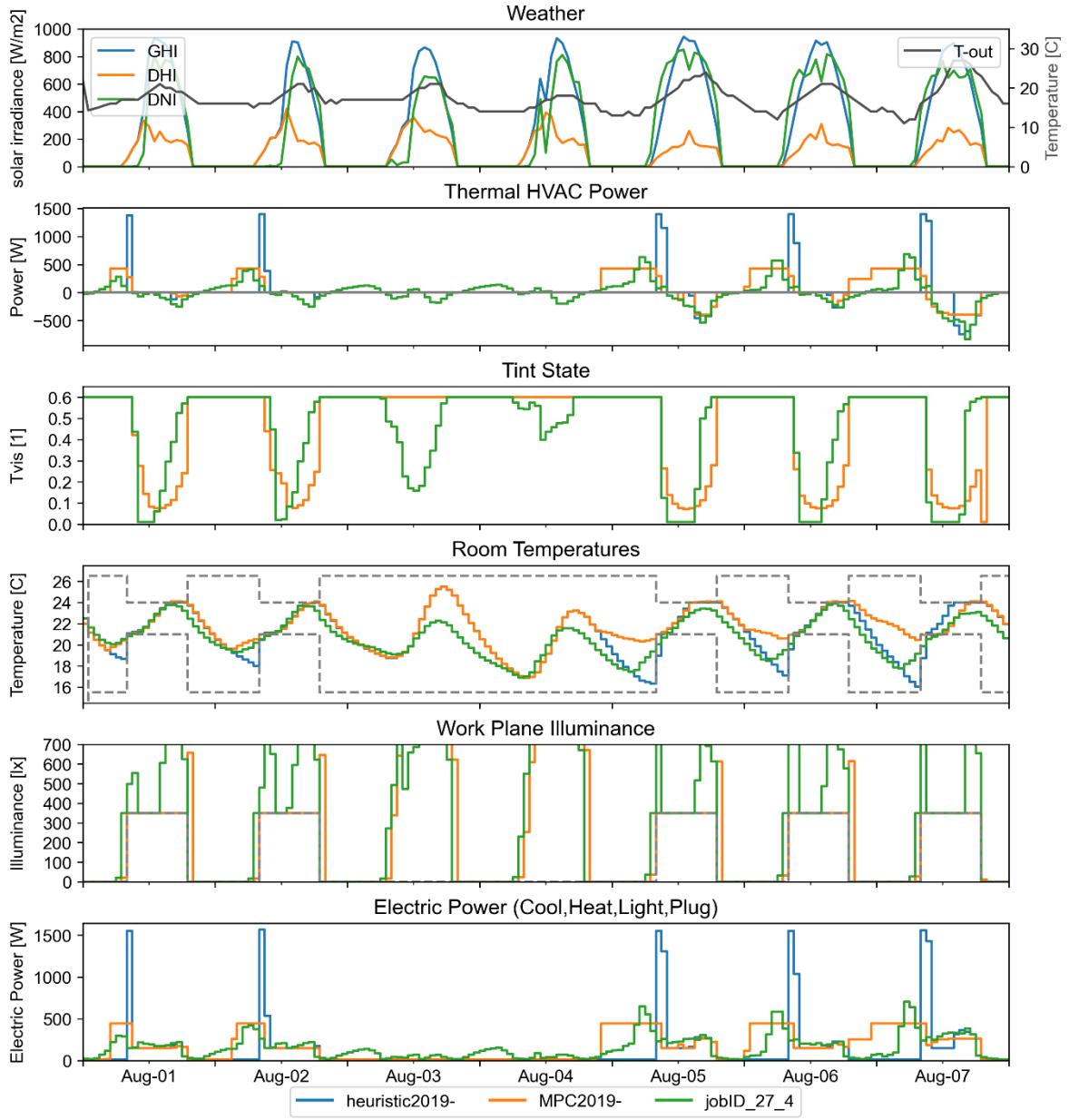


Figure 54: Results of the final RL agent (due to some errors this is an outdated version from 11.15.2020)

The MPC with its optimizing nature precools or preheats the room, which leads to a 75 % lower peak load for heating (Figure 55). With the strategy of precooling the MPC consumes 28.65 kWh of electrical energy with a total cost for energy and demand in this week of 17.49 \$. Whereas the PI-controller consumes less energy with 22.28 kWh but with total costs of 38.04 \$. The week controlled by the agent costs 24.51 \$. The Clearly recognizable is the impact of the TOU-tariff.

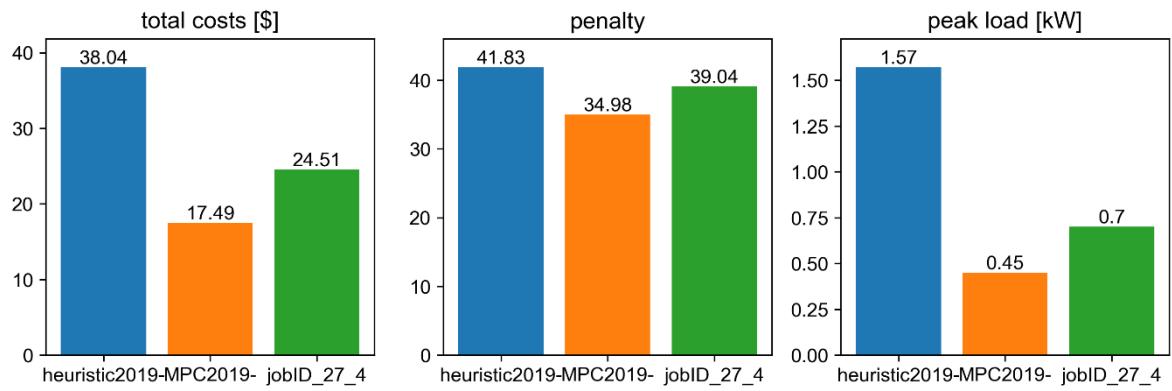


Figure 55: Performance measure compared between PI-controller, MPC and best agent (due to some errors this is an outdated verion from 11.15.2020)

5 Discussion and Outlook

The aim of this thesis was the implementation of a machine learning agent which strives to minimize the total energy costs of a room, while ensuring the comfort parameters for the occupants. One of the main tasks was the question which Reinforcement Learning (RL) methodology would be best suited for the control of building technology to further reduce total energy costs compared to state-of-the-art controllers and MPC controllers.

In this thesis, an agent was developed for the heating and cooling control of an office building, as well as its window shading system with input values for the weather-forecast, occupancy level and TOU-tariff. The latter is the most crucial factor for a cost-effective control system. The TOU-tariff as a main input value for the agent enables the power grid operator to actively manage the energy load of the building by changing energy costs for a short period of time. This possibility for the power grid operator will help to increase the share of renewable energy systems, without the necessity to reinforce the power grid. The advantage of an agent for building owners are the significantly lower total energy costs compared to state-of-the-art PI-controller. That is due to the possibility of minimizing the maximum peak load and energy consumption during high priced periods since the agent is learning a control strategy without rules and is continuously trained during operation. With a controller that takes the total energy costs into account, including for the HVAC-system, as well as for the artificial light and all equipment, the illuminance level of the room remains unknown to the agent and is not required for training or operating. The final comparison with the MPC controller shows that the agent misses farsightedness and cannot decrease the maximum demand as much as the MPC does.

Even if the agent can reduce the energy costs, while maintaining the room temperature, it does not work perfectly. The actions taken are never zero, but rather oscillate on days on which the room temperature would stay within the set boundaries even if no actions were required from the agent. A solution for this problem could be a hierarchical agent setup. An agent would set the goal, e.g. in form of the room temperature and the illuminance level, and the underlying agent would try to take actions to reach these goals while an additional threshold would prevent the agent from performing unnecessary actions. By including the TOU-tariff in the neural network inputs, the controller is affected by changes of the tariff or the tariff structure of the utilities. It is important to recalibrate the normalization of the TOU tariff to ensure a successful behavior of the mode. The learning setup during operation in a real building is not included yet but could be implemented the same way as the learning in the simulation was. The training of an agent would have needed to be done prior to deploying it to a building, since the training time with 3000 episodes or days, respectively, is too long for learning from the real building.

Bibliography

- [1] Alpaydin, E., 2010, Introduction to machine learning. 2nd ed. London, England, The MIT Press
- [2] Barrett, E., Linder, S., 2015, Autonomous HVAC Control, A Reinforcement Learning Approach. in: A. Bifet, M. May, B. Zadrozny, R. Gavalda, D. Pedreschi, F. Bonchi, J. Cardoso, & M. Spiliopoulou (eds.), Machine Learning and Knowledge Discovery in Databases. Lecture Notes in Computer Science. Cham, Springer International Publishing; 3–19
- [3] Barth-Maron, G., Hoffman, M.W., Budden, D., Dabney, W., Horgan, D., TB, D., Muldal, A., Heess, N., Lillicrap, T., 2018, Distributed Distributional Deterministic Policy Gradients. arXiv:1804.08617 [cs, stat]; 16
- [4] Bradley, R., 2018, 16 Examples of Artificial Intelligence (AI) in Your Everyday Life | The Manifest.; <https://themanifest.com/development/16-examples-artificial-intelligence-ai-your-everyday-life/>; 19.11.2020
- [5] BRE, 2015, A Technical Manual for SBEM.; https://www.uk-ncm.org.uk/filelibrary/SBEM-Technical-Manual_v5.2.g_20Nov15.pdf; 26.11.2020
- [6] Brownlee, J., 2019, 14 Different Types of Learning in Machine Learning. Machine Learning Mastery; <https://machinelearningmastery.com/types-of-learning-in-machine-learning/>; 18.8.2020
- [7] Brownlee, J., 2016, Crash Course On Multi-Layer Perceptron Neural Networks. Machine Learning Mastery; <https://machinelearningmastery.com/neural-networks-crash-course/>; 23.11.2020
- [8] Cao, X., Wan, H., Lin, Y., Han, S., 2019, High-Value Prioritized Experience Replay for Off-Policy Reinforcement Learning. in: 2019 IEEE 31st International Conference on Tools with Artificial Intelligence (ICTAI). Portland, OR, USA; 1510–1514
- [9] CEC, 2020, Electric Utility Service Area.; <https://cecgis-caenergy.opendata.arcgis.com/pages/pdf-maps>; 26.11.2020
- [10] Chen, Y., Norford, L.K., Samuelson, H.W., Malkawi, A., 2018, Optimal control of HVAC and window systems for natural ventilation through reinforcement learning. Energy and Buildings, Volume 169, Number; 195–205
- [11] DeepMind, 2016, AlphaGo: The story so far. Deepmind; </research/case-studies/alphago-the-story-so-far>; 13.7.2020
- [12] Deru, M., Field, K., Studer, D., Benne, K., Griffith, B., Torcellini, P., Liu, B., Halverson, M., Winiarski, D., Rosenberg, M., Yazdanian, M., Huang, J., Crawley, D., 2011, U.S. Department of Energy Commercial Reference Building Models of the National Building Stock.
- [13] Ding, B., Qian, H., Zhou, J., 2018, Activation functions and their characteristics in deep neural networks. in: 2018 Chinese Control And Decision Conference (CCDC). Shenyang; 1836–1841

- [14] Duffie, J.A., Beckman, W.A., 2013, Solar Engineering of Thermal Processes. fourth Edition. Madison, John Wiley & Sons
- [15] EnergyPlus, 2019, Weather Data by Location. EnergyPlus; https://energyplus.net/weather-location/north_and_central_america_wmo_region_4/USA/CA/USA_CA_Oakland.Intl.AP.724930_TMY3; 28.11.2020
- [16] Ertel, W., 2016, Neuronale Netze. in: Grundkurs Künstliche Intelligenz. Wiesbaden, Springer Fachmedien Wiesbaden; 265–311
- [17] Fernandes, L.L., Lee, E.S., Dickerhoff, D., Thanachareonkit, A., Wang, T., Gebauer, C., 2018, Electrochromic Window Demonstration at the John E. Moss Federal Building, 650 Capitol Mall, Sacramento, California.
- [18] Fujimoto, S., van Hoof, H., Meger, D., 2018, Addressing Function Approximation Error in Actor-Critic Methods. arXiv:1802.09477 [cs, stat]; <http://arxiv.org/abs/1802.09477>; 22.9.2020
- [19] Gebauer, C., Blum, D.H., Wang, T., Lee, E.S., 2020, An assessment of the load modifying potential of model predictive controlled dynamic facades within the California context. Energy and Buildings, Volume 210, Number; 109762
- [20] Google Inc., 2019, Keras: the Python deep learning API.; <https://keras.io/>; 13.11.2020
- [21] Haarnoja, T., Zhou, A., Abbeel, P., Levine, S., 2018, Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. arXiv:1801.01290 [cs, stat]; 14
- [22] Haenlein, M., Kaplan, A., 2019, A Brief History of Artificial Intelligence: On the Past, Present, and Future of Artificial Intelligence. California Management Review, Volume 61, Number 4; 5–14
- [23] Hale, J., 2019, Deep Learning Framework Power Scores 2018. Medium; <https://towardsdatascience.com/deep-learning-framework-power-scores-2018-23607ddf297a>; 23.7.2020
- [24] Hale, J., 2020, Which Deep Learning Framework is Growing Fastest? Medium; <https://towardsdatascience.com/which-deep-learning-framework-is-growing-fastest-3f77f14aa318>; 23.7.2020
- [25] Heess, N., Hunt, J.J., Lillicrap, T.P., Silver, D., 2015, Memory-based control with recurrent neural networks. arXiv:1512.04455 [cs]; 11
- [26] Heinrich, B., Linke, P., Glöckler, M., 2020, Grundlagen Automatisierung: Erfassen - Steuern - Regeln. 3rd ed. Springer Vieweg
- [27] Hochreiter, S., Schmidhuber, J., 1997, Long Short-Term Memory.; <http://www.bioinf.jku.at/publications/older/2604.pdf>

- [28] Hong, T., Wang, Z., Luo, X., Zhang, W., 2020, State-of-the-art on research and applications of machine learning in the building life cycle. *Energy and Buildings*, Volume 212, Number; 109831
- [29] Lee, E.S., Selkowitz, S., Clear, R., DiBartolomeo, D., Klems, J., Fernandes, L.L., Ward, G., Inkarojrit, V., Yazdanian, M., 2006, A Design Guide for Early-Market Electrochromic Windows. California Energy Commission, PIER. 500-01-023. LBNL-59950
- [30] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D., 2015, Continuous control with deep reinforcement learning. arXiv:1509.02971 [cs, stat]; 10
- [31] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D., 2019, Continuous control with deep reinforcement learning. arXiv:1509.02971 [cs, stat]; <http://arxiv.org/abs/1509.02971>; 20.5.2020
- [32] Mitchell, T.M., 1997, Machine Learning. New York, McGraw-Hill
- [33] Mnih, V., Badia, A.P., Mirza, M., Graves, A., Harley, T., Lillicrap, T.P., Silver, D., Kavukcuoglu, K., 2016, Asynchronous Methods for Deep Reinforcement Learning. arXiv:1602.01783v2 [cs.LG]; 10
- [34] Mohammed, M., Khan, M.B., Bashier, E.B.M., 2017, Machine Learning: Algorithms and Applications. 0 ed. Boca Raton : CRC Press, 2017., CRC Press
- [35] Mousavi Maleki, S., Hizam, H., Gomes, C., 2017, Estimation of Hourly, Daily and Monthly Global Solar Radiation on Inclined Surfaces: Models Re-Visited. *Energies*, Volume 10, Number 1; 134
- [36] Olah, C., 2015, Understanding LSTM Networks -- colah's blog.; <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>; 22.11.2020
- [37] Open Data Science, 2019, What is TensorFlow? Medium; <https://medium.com/@ODSC/what-is-tensorflow-13200525e852>; 10.11.2020
- [38] OpenAI, 2018, Part 2: Kinds of RL Algorithms — Spinning Up documentation.; https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html#a-taxonomy-of-rl-algorithms; 14.7.2020
- [39] Pan, Y., 2016, Heading toward Artificial Intelligence 2.0. *Engineering*, Volume 2, Number 4; 409–413
- [40] PG&E, 2020a, Electric Schedule E-19 Medium General Demand-Metered TOU Service, effective April 19,2020.; https://www.pge.com/tariffs/assets/pdf/tariffbook/ELEC_SCHEDS_E-19.pdf; 2.10.2020
- [41] PG&E, 2020b, PG&E, Pacific Gas and Electric - Gas and power company for California.; <https://www.pge.com/>; 26.11.2020
- [42] Pierre, C., 2020, PYPL PopularitY of Programming Language index.; <http://pypl.github.io/PYPL.html>; 23.7.2020

- [43] Plappert, M., Houthooft, R., Dhariwal, P., Sidor, S., Chen, R.Y., Chen, X., Asfour, T., Abbeel, P., Andrychowicz, M., 2018, Parameter Space Noise for Exploration. arXiv:1706.01905 [cs, stat]; 18
- [44] Python Software Foundation, 2020, Python. Python.org; <https://www.python.org/>; 12.5.2020
- [45] Rüdisser, D., 2018, A brief guide and free tool for the calculation of the thermal mass of building components (according to ISO 13786).; <http://rgdoi.net/10.13140/RG.2.2.18312.72967>; 25.11.2020
- [46] Sandia National Laboratories, 2019, Pyomo. Pyomo; <http://www.pyomo.org>; 13.11.2020
- [47] Schaul, T., Quan, J., Antonoglou, I., Silver, D., 2016, Prioritized Experience Replay. arXiv:1511.05952 [cs]; <http://arxiv.org/abs/1511.05952>; 13.5.2020
- [48] Schulman, J., Levine, S., Moritz, P., Jordan, M.I., Abbeel, P., 2017, Trust Region Policy Optimization. arXiv:1502.05477 [cs]; 16
- [49] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O., 2017, Proximal Policy Optimization Algorithms. arXiv:1707.06347 [cs]; <http://arxiv.org/abs/1707.06347>; 14.7.2020
- [50] Seif, G., 2019, Understanding the 3 most common loss functions for Machine Learning Regression. Medium; <https://towardsdatascience.com/understanding-the-3-most-common-loss-functions-for-machine-learning-regression-23e0ef3e14d3>; 23.11.2020
- [51] State of California, 2018, California Costumer Choice.; https://www.cpuc.ca.gov/uploadedFiles/CPUC_Public_Website/Content/Utilities_and_Industries/Energy_-_Electricity_and_Natural_Gas/Cal%20Customer%20Choice%20Report%208-7-18%20rm.pdf; 22.9.2020
- [52] Sutton, R.S., Barto, A.G., 2018, Reinforcement learning: an introduction. Second edition. Cambridge, Massachusetts, The MIT Press
- [53] Turing, A.M., 1950, I.—COMPUTING MACHINERY AND INTELLIGENCE. Mind, Volume LIX, Number 236; 433–460
- [54] U.S. Energy Information Administration, 2020, U.S. energy facts explained - consumption and production - U.S. Energy Information Administration (EIA).; <https://www.eia.gov/energyexplained/us-energy-facts/>; 9.11.2020
- [55] Wang, C.-F., 2019, The Vanishing Gradient Problem. Medium; <https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484>; 23.11.2020
- [56] Wang, Y., Kuckelkorn, J., Liu, Y., 2017, A state of art review on methodologies for control strategies in low energy buildings in the period from 2006 to 2016. Energy and Buildings, Volume 147, Number; 27–40
- [57] Wang, Z., Hong, T., 2020, Reinforcement learning for building controls: The opportunities and challenges. Applied Energy, Volume 269, Number; 115036

- [58] Wei, T., Wang, Y., Zhu, Q., 2017, Deep Reinforcement Learning for Building HVAC Control. in: Proceedings of the 54th Annual Design Automation Conference 2017. DAC '17: The 54th Annual Design Automation Conference 2017. Austin TX USA, ACM; 1–6

List of Figures

Figure 1: Worldwide PYPL PopularitY of Programming Language in 2020 (modified according to (Pierre 2020)).....	13
Figure 2: ML Framework Power Scores 2018 (modified according to (Hale 2019))	14
Figure 3: DL Framework Six-Month Growth Scores 2019 (modified according to (Hale 2020))	14
Figure 4: Children playing Go on a regular Go board (DeepMind 2016)	16
Figure 5: Different machine learning techniques and their required data (modified according to (Mohammed et al. 2017)	18
Figure 6: The agent–environment interaction in a Markov decision process. (modified according to (Sutton and Barto 2018, p.48))	19
Figure 7: Simple Maze displaying the difference of reward and state-value.....	21
Figure 8: backup diagram for $v\pi(s)$ (Sutton and Barto 2018, p.59)	21
Figure 9: $q\pi$ backup diagram (Sutton and Barto 2018, p.61)	22
Figure 10: A non-exhaustive, but useful taxonomy of algorithms in modern RL (OpenAI 2018)	23
Figure 11: Formal model with neurons and directed connections between them (modified according to (Ertel 2016, p.267))	26
Figure 12: The structure of a formal neuron that applies the activation function f to the weighted sum of all inputs (modified according to (Ertel 2016, p.269))	26
Figure 13: A three-layer backpropagation network with n_1 neurons in the first layer, n_2 neurons in the second and n_3 neurons in the third layer (modified according to (Ertel 2016, p.291))	27
Figure 14: Model of a Simple Network (Brownlee 2016).....	28
Figure 15: An unrolled recurrent neural network (Olah 2015)	28
Figure 16: The repeating module in a standard RNN contains a single layer (Olah 2015). 29	29
Figure 17: The repeating module in an LSTM contains four interacting layers (Olah 2015).	29
Figure 18: The graphic depiction of Sigmoid function and its derivative (Ding et al. 2018, p.1837).....	30
Figure 19: The graphic depiction of hyperbolic tangent function and its derivative (Ding et al. 2018, p.1838).....	31
Figure 20: The graphic depiction of ReLU function and its derivatives (Ding et al. 2018, p.1838).....	31
Figure 21: The graphic depiction of LReLU, PReLU, and RReLU function (Ding et al. 2018, p.1839).....	32
Figure 22: Elements of the DDPG algorithm.....	33
Figure 23: DDPG – agent-environment interaction	34
Figure 24: DDPG – calculating the action-values	35
Figure 25: DDPG – training of the critic and actor network	36

Figure 26: Action-noise process Ohrnstein-Uhlenbeck and Gaussian	40
Figure 27: Functional diagram of a PID-controller (modified according to(Heinrich et al. 2020, p.163))	41
Figure 28: Step response with aperiodic course (Heinrich et al. 2020, p.174)	43
Figure 29: Information flow in the perfect knowledge MPC model	44
Figure 30: Room model (green)	45
Figure 31: Diagram of a typical tungsten-oxide electrochromic coating (Lee et al. 2006, p.6)	47
Figure 32: Each window pane had three sub-zones that could be independently controlled (Fernandes et al. 2018, p.14)	47
Figure 33: EC-window properties	48
Figure 34: (a) Zenith angle, slope, surface azimuth angle, and solar azimuth angle for a tilted surface. (b) Plan view showing solar azimuth angle (Duffie and Beckman 2013, p.13)	50
Figure 35: Maximum and minimum value of declination angle (Mousavi Maleki et al. 2017, p.2)	50
Figure 36: Declination angle in Oakland, CA	51
Figure 37: Solar angles and solar irradiation on tilted surface	53
Figure 38: E-19 tariff with time dependent energy- and demand costs (PG&E 2020a)	55
Figure 39: Neural network architecture; critic left and actor on the right with the shape of the input vectors.....	58
Figure 40: First training result of a week starting on August 1 st with HVAC control and a fixed Tint state	59
Figure 41: Training result of a week starting on August 1 st with HVAC and Tint state control	60
Figure 42: Critic loss of a week starting on August 1 st with HVAC and Tint state control ...	60
Figure 43: New critic network based on the TD3 algorithm with the shape of the input vectors	61
Figure 44: Comparison of critic loss between DDPG critic and TD3 critic architecture	61
Figure 45: Training result with the new critic of a week starting on August 1 st with HVAC and Tint state control	62
Figure 46: Neural network setup with 4 forecast hours with the shape of the input vectors	63
Figure 47: Comparison of critic loss between 1h and 4h forecast.....	63
Figure 48: Training run with four forecast hours of a week starting on August 1 st with HVAC and Tint state control.....	64
Figure 49: Comparison of best Gridsearch results of a week starting on August 1st with HVAC and Tint state control	65
Figure 50: Performance measure of the best Gridsearch results	66
Figure 51:Comparison of Gridsearch results for different N-step rewards starting on August 1st.....	68
Figure 52: Performance measure of the best Gridsearch results for different N-step rewards	68

Figure 53: Neural network setup for the RDPG with the shape of the input vectors.....	70
Figure 54: Results of the final RL agent (due to some errors this is an outdated verion from 11.15.2020).....	72
Figure 55: Performance measure compared between PI-controller, MPC and best agent (due to some errors this is an outdated verion from 11.15.2020).....	73

List of Tables

Table 1: Reference Building Form Assignments (Deru et al. 2011, p.19)	45
Table 2: U-Value by Reference Building Vintage - Standard 90.1-2004 (Deru et al. 2011, p.26).....	45
Table 3: specification of the room model	46
Table 4: Name and visible transmittance of the four tint levels. (Fernandes et al. 2018, p.15)	48
Table 5: E-19 definition of time periods, energy- and demand-costs	54
Table 6: Example for the demand cost calculation	54
Table 7: Gridsearch results of best neural network configurations.....	64
Table 8: Gridsearch results of best N-step reward.....	67
Table 9: Gridsearch results of best improvements to the agent.....	69

List of Abbreviations

AI	Artificial Intelligence
DDPG	Deep Deterministic Policy Gradient
DF	Daylight Factor
DHI	Diffuse Horizontal Irradiance
DL	Deep Learning
DNI	Direct Normal Irradiance
GHI	Global Horizontal Irradiance
HVAC	Heating Ventilating Air Conditioning
HVPER	High-Value Prioritized Experience Replay
ID	Importance Sampling
LSTM	Long-Short Term Memory
MAE	Mean Absolute Error
ML	Machine Learning
MLP	Multi-Layer Perceptron
MPC	Model Predictive Control
MSE	Mean Squared Error
NN	Neural Network
OU	Ohrstein-Uhlenbeck
PER	Prioritized Experience Replay
PG&E	Pacific Gas and Electricity
PI	Proportional-Integra
PID	Proportional-Integral-Derivative
RC	Resistance and Capacitance
relu	Rectified Linear Unit
RL	Reinforcement Learning
SHGC	Solar Heat Gain Coefficient

sig	Sigmoid
tanh	Hyperbolic tangent
TD	Temporal Difference
TOU	Time-Of-Use
T-out	Outside Air Temperature
Tv	Visibility Transmittance
WPI	Workplace Illuminance

Appendix A: Setting

Parameter setting

```
1 import numpy as np
2 import os
3
4 def get_parameter(l_wall, U_wall, U_window, A_room, h_room, A_window, step_size, model, \
5                   hvac_control, radiance, max_power):
6     rho_air = 1.1894 # Density, in kg/m3
7     cp_air = 1.0086 # Air Specific heat capacity, in kJ/kgK
8     R_wall = (U_wall * (l_wall * h_room - A_window))
9     R_window = (U_window * A_window)
10    C_int_wall = ((l_wall + A_room/l_wall)*2*h_room - l_wall*h_room)*8.184 # Drywall
11    Capacitance in kJ/K
12    C_int_floor_ceiling = 2*A_room * 43.7545#[kJ/K] light concrete floor and ceiling
13    C_Air = A_room * h_room * rho_air * cp_air #[kJ/K]
14    C_ext_wall = l_wall*h_room * 35.284 #[kJ/K]
15    C_room = (C_int_wall + C_Air + C_ext_wall + C_int_floor_ceiling)
16    operatingsys = 'windows' if os.name == 'nt' else 'linux'
17
18    weather_columns = ['weaCelHei','weaCloTim','weaHDifHor','weaHDirNor',
19                      'weaHGloHor','weaHHorIR','weaNOpA','weaNTot',
20                      'weaPAtm','weaRelHum','weaSoltTim','weaSolZen',
21                      'weaTBlaSky','weaTDewPoi','weaTDryBul',
22                      'weaTWetBul','weaWinDir','weaWinSpe']
23
24    # Inputs
25    parameter['inputs'] = {}
26    parameter['inputs']['labels'] = ['Q_hvac', 'Tvis', 'start_time', 'T_out', 'S_irr',
27                                    'S_ill', 'Q_int_th', 'Q_int_el',
28                                    'Occ', 'C_energy', 'C_demand', 't_min', 't_max',
29                                    'wpi_min'] + weather_columns
30
31    # input data setting
32    parameter['input_data'] = {}
33    parameter['input_data']['month'] = 8
34    parameter['input_data']['day'] = 1
35    parameter['input_data']['step_size'] = step_size # hour
36
37    # Window
38    tints = np.array([[0.42,0.6],[0.16,0.18],[0.12,0.06],[0.1,0.01]]) # [shgc, Tvis]
39    coeff = np.poly1d(np.polyfit(tints[:,1], tints[:,0], 1))
40    parameter['window'] = {}
41    parameter['window']['area'] = A_window # Window area, in m2
42    parameter['window']['coeff_a'] = coeff[0] # Window tint fit funciton
43    parameter['window']['coeff_b'] = coeff[1] # Window tint fit funciton
44
45    # Room configuraiton
46    parameter['zone'] = {}
47    parameter['zone']['area'] = A_room # Room area, in m2
48    parameter['zone']['height'] = h_room # Room height, in m
49    parameter['zone']['length'] = l_wall # Room length, in m
50    parameter['zone']['surface_area'] = ((parameter['zone']['length'] +
51                                         parameter['zone']['area']) *
52                                         parameter['zone']['length']) *
53                                         parameter['zone']['height'] \
54                                         + parameter['zone']['area']) * 2
55    parameter['zone']['t_init'] = None # Initial room temperature, in K (None = random)
56    parameter['zone']['t_init_min'] = 21 + 273.15 # Minimal initial room temperate;
57    Minimum room temperature when occupied, in K
58    parameter['zone']['t_init_max'] = 24 + 273.15 # Maximal initial room temperate;
59    Maximum room temperature when occupied, in K
60    parameter['zone']['eff_lights'] = 5 * A_room / 500 # 5 W/m2 => LPD of 0.5
61    W/ft2
62    parameter['zone']['eff_heat'] = 1 # Efficieny of heating
63    parameter['zone']['eff_cool'] = 1/3.5 # Efficiency of cooling
```

```

57     parameter['zone']['int_th_load'] = 100 * np.round(parameter['zone']['area']/18.58)
58     # NREL --> 1 person on 18.58 m2
59     parameter['zone']['int_el_load'] = 10.76*parameter['zone']['area'] # --> 10.76 W/m2
60     parameter['zone']['office_hours'] = [[7,18]]
61     if model == 'RC':
62         parameter['zone']['control_hvac'] = False
63     else:
64         parameter['zone']['control_hvac'] = hvac_control # Flag to control HVAC system
65
66     # Constraints
67     parameter['constraints'] = {}
68     parameter['constraints']['max_t_penalty'] = 2 # Maximal reward penalty
69     parameter['constraints']['night_tint_penalty'] = 1 # penalty for tinting the
70     windows during the night
71     parameter['constraints']['cool_max'] = max_power * A_room # Maximal cooling power,
72     in W
73     parameter['constraints']['heat_max'] = max_power * A_room # Maximal heating power,
74     in W
75     parameter['constraints']['wpi_min'] = 350 # Minimum work place illuminance (wpi)
76     when occupied, lux
77     parameter['constraints']['t_min'] = 15.5 + 273.15 # Minimum temperature when not
78     occupied, K
79     parameter['constraints']['t_max'] = 26.5 + 273.15 # Maximum temperature when not
80     occupied, K
81
82     # demand charge properties
83     parameter['demand_charge'] = {}
84     parameter['demand_charge']['period'] = 24 * 30.436875 # demand is charged every
85     month, in h
86     parameter['demand_charge']['n_step'] = 4 # length of demand charge period for
87     penalty, in h
88     parameter['demand_charge']['base_charge'] = 17.63 # demand is charged every month,
89     in h
90
91     # Solar-optical model selection
92     path_emulator = os.path.join(os.getcwd(), '..', 'DynamicFacades', 'emulator')
93
94     parameter['somodel'] = {}
95     parameter['somodel']['use_radiance'] = radiance # Flag to use Radiance or simplified
96     parameter['somodel']['path_radiance_handler'] = os.path.join(path_emulator,
97     'emulator')
98     path_resources_emulator = os.path.join(path_emulator, 'resources')
99     parameter['somodel']['config_file'] = os.path.join(path_resources_emulator,
100    'radiance', 'room0.4WWR_ec.cfg')
101    parameter['somodel']['regenerate_matrices'] = False # Flag to regenerate Radiance
102    matrices
103    parameter['somodel']['orientation'] = 0 # Radiance orientation S-0, W-90, E-270,
N-180
104    parameter['somodel']['location'] = {}
105    parameter['somodel']['location']['latitude'] = 37.72 # Radiance latitude
106    parameter['somodel']['location']['longitude'] = -122.22 * -1 # Radiance longitude
107    parameter['somodel']['location']['timezone'] = -8 * 15 * -1 # Radiance timezone
108    parameter['somodel']['location']['orient'] = parameter['somodel']['orientation']
109    parameter['somodel']['filestruct'] = {}
110    parameter['somodel']['filestruct']['resources'] =
111    os.path.join(path_resources_emulator, 'radiance', 'BSDFs')
112    parameter['somodel']['filestruct']['matrices'] =
113    os.path.join(path_resources_emulator, 'radiance',
114                                'matrices', 'ec',
115                                '0.4')
116    parameter['somodel']['tvvis_to_state'] = {0.6:3,0.18:2,0.06:1,0.01:0}
117
118    # Model selection (Modelica RC model)
119    parameter['model'] = {}

```

```

104 if model == 'RC':
105     parameter['model']['use_fmu'] = False # Flag to use FMU
106 else:
107     parameter['model']['use_fmu'] = True # Flag to use FMU
108 parameter['model']['fmu_loglevel'] = 0 # Loglevel, 0-none, 5-debug
109 #'''
110 # Modelica RC model
111 if model != 'Room_modelica':
112     parameter['model']['fmu_path'] =
113         r'resources\fmus\{}\RCmodelwithHvac.fmu'.format(operatingsys)
114     parameter['model']['include_solar_gains'] = True
115     parameter['model']['param'] = {}
116     parameter['model']['param']['R1'] = 1 / (R_wall + R_window)
117     parameter['model']['param']['C1'] = C_room * 1e3
118     parameter['model']['param']['timestep'] = parameter['input_data']['step_size']
119     *60*60 # Timestep in seconds
120     parameter['model']['param']['gaiHea.k'] = parameter['constraints']['heat_max']
121     # Limit HVAC thermal power, in W
122     parameter['model']['inputs_Q_thermal'] = 'Q1'
123     parameter['model']['inputs'] = {}
124     parameter['model']['inputs'][T_ctrl] = 1 if parameter['zone']['control_hvac']
125     else 0
126     parameter['model']['inputs_map'] = {}
127     parameter['model']['inputs_map'][T_out] = 'T_out'
128     parameter['model']['inputs_map'][T_set_heat] = 't_min'
129     parameter['model']['inputs_map'][T_set_cool] = 't_max'
130     parameter['model']['outputs_T_in'] = 'T_in'
131     parameter['model']['outputs_Q_hvac'] = 'Q_hvac'
132 elif model == 'Room_modelica':
133 # Modelica Room model
134     parameter['model']['fmu_path'] =
135         r'resources\fmus\{}\AWTBModelica.fmu'.format(operatingsys)
136     parameter['model']['include_solar_gains'] = False
137     parameter['model']['param'] = {}
138     parameter['model']['param']['timestep'] =
139         parameter['input_data']['step_size']*60*60 # Timestep in seconds
140     parameter['model']['param']['gaiHea.k'] = parameter['constraints']['heat_max']
141     # Limit HVAC thermal power, in W
142     parameter['model']['inputs_Q_thermal'] = 'Q_con'
143     parameter['model']['inputs'] = {}
144     parameter['model']['inputs'][T_ctrl] = 1 if parameter['zone']['control_hvac']
145     else 0
146     parameter['tariff'] = {}
147     parameter['tariff']['periods'] = {}
148     # month
149     parameter['tariff']['periods'][summer] = {}
150     parameter['tariff']['periods'][winter] = {}
151     parameter['tariff']['periods'][summer]['month'] = [5,6,7,8,9,10]
152     parameter['tariff']['periods'][winter]['month'] = [1,2,3,4,11,12]
153     # day of week
154     parameter['tariff']['dayofweek'] = {}
155     parameter['tariff']['dayofweek'][weekday] = [0,1,2,3,4]
156     parameter['tariff']['dayofweek'][weekend] = [5,6]
157     # hour of day
158     parameter['tariff']['periods'][summer]['s_peak'] = [[12,6+12]]
159     parameter['tariff']['periods'][summer]['s_part_peak'] = [[8.5,12],[6+12,9.5+12]]
160     parameter['tariff']['periods'][summer]['s_off_peak'] = [[0,8.5],[9.5+12,12+12]]

```

```

160 parameter['tariff']['periods']['winter']['w_part_peak'] = [[8.5,9.5+12]]
161 parameter['tariff']['periods']['winter']['w_off_peak'] = [[0,8.5],[9.5+12,12+12]]
162 # energy rate
163 parameter['tariff']['C_energy'] = {}
164 parameter['tariff']['C_energy']['s_peak'] = 0.16225
165 parameter['tariff']['C_energy']['s_part_peak'] = 0.11734
166 parameter['tariff']['C_energy']['s_off_peak'] = 0.08846
167 parameter['tariff']['C_energy']['w_part_peak'] = 0.11127
168 parameter['tariff']['C_energy']['w_off_peak'] = 0.09559
169 # demand rate
170 parameter['tariff']['C_demand'] = {}
171 parameter['tariff']['C_demand']['s_peak'] = 19.63
172 parameter['tariff']['C_demand']['s_part_peak'] = 5.37
173 parameter['tariff']['C_demand']['s_off_peak'] = 0.0
174 parameter['tariff']['C_demand']['w_part_peak'] = 0.18
175 parameter['tariff']['C_demand']['w_off_peak'] = 0.0
176 parameter['tariff']['C_demand']['base_rate'] = 17.63
177
178 # parameters for RL Agent
179 parameter['agent']={}
180 parameter['agent']['stepsize'] = parameter['input_data']['step_size'] # hour
181 parameter['agent']['Agent'] = 'DDPG' # DDPG, RDPG
182 parameter['agent']['dtype'] = 'float64'
183 parameter['agent']['actions'] = 0 # 0 = Q_hvac+Tvis; 1 = Q_hvac; 2 = Tvis
184
185 # set the hyperparameters for the neural network
186 parameter['agent']['NN'] = {}
187 parameter['agent']['NN']['network_architecture'] = 'MLP' # MLP, CNN, LSTM
188 parameter['agent']['NN']['hidden_layers'] = 3
189 parameter['agent']['NN']['layer_size_1'] = 300
190 parameter['agent']['NN']['layer_size_2'] = 832
191 parameter['agent']['NN']['activation'] = 'relu'
192 parameter['agent']['NN']['tow'] = 0.001
193 parameter['agent']['NN']['discount_factor'] = 0.99
194 parameter['agent']['NN']['demand_charge_scale'] = 1
195 parameter['agent']['NN']['act_learning_rate'] = 0.0001
196 parameter['agent']['NN']['crit_learning_rate'] = 0.001
197
198 parameter['agent']['setting'] = {}
199 parameter['agent']['setting']['n_step'] = 4
200 parameter['agent']['setting']['forecast_hours'] = 4
201 parameter['agent']['setting']['training_days'] = 7
202 parameter['agent']['setting']['test_days'] = 7
203 parameter['agent']['setting']['noise_process'] = 'OU_noise' # method for OU_noise,
Gauss_noise or param_noise
204 parameter['agent']['setting']['forecast_col'] = ['weaTDryBul', 'S_irr',
'TOU_tariff', 'occupancy']
205 parameter['agent']['setting']['episodes'] = 600
206 parameter['agent']['setting']['episodes_with_noise'] = 0.5
207 parameter['agent']['setting']['exploration_episodes'] = 50
208
209 parameter['agent']['replay_buffer'] = {}
210 parameter['agent']['replay_buffer']['Buffer'] = 'Uniform' # Uniform, PER, HVPER
211 parameter['agent']['replay_buffer']['max_buffer_size'] = int(1e9)
212 parameter['agent']['replay_buffer']['batch_size'] = int(1e3)
213
214 parameter['agent']['scaling'] = {}
215 parameter['agent']['scaling']['forecast_norm_data'] =
np.array([[0.+273.15,0.,-0.16225,0.],[34.+273.15,1000.,0.16225,2.1]])
216 parameter['agent']['scaling']['state_norm_data'] =
np.array([[parameter['constraints']['t_min']],[parameter['constraints']['t_max']]])
217 parameter['agent']['scaling']['actions_norm_data'] =
np.array([[-parameter['constraints']['heat_max'],
-tints[:,1].max(),[parameter['constraints']['heat_max']] ,tints[:,1].max()]])

```

```

218     parameter['agent']['scaling']['noise_scale'] =
219         parameter['constraints']['heat_max']*0.15, 0.1 #if action_noise 2 dim, if param
220         noise 1 dim
221     parameter['agent']['scaling']['reward_0'] = 24 * 30.436875 /
222     parameter['agent']['setting']['n_step']
223
224     # set flags
225     parameter['agent']['flags'] = {}
226     parameter['agent']['flags']['valuefunction'] = True
227     parameter['agent']['flags']['plot'] = True
228     parameter['agent']['flags']['print'] = True
229     parameter['agent']['flags']['hp_tune'] = False
230     parameter['agent']['flags']['load_model'] = False
231     parameter['agent']['flags']['save_models'] = False
232
233     return parameter

```

Appendix B: Execution

Training

```
1 import numpy as np
2 import os
3 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
4 from tf_agents.environments import tf_py_environment
5
6 from environment import Room
7 from input_data import get_weather_files
8 from input_data_handler import RL_results_handler, files_handler
9 from parameter_handler import get_parameter
10 from RL_Agent import agent
11
12 l_wall = 4 # Length of wall, in m
13 U_wall = 3.294 # U-value wall, in W/m2K
14 U_window = 6.923 # U-value window, in W/m2K
15 A_room = 14 # Room area, in m2
16 h_room = 3.95 # Room height, in m
17 A_window = l_wall * h_room * 0.33 # Window area, in m2
18 step_size = 60/60 # hour
19 model = 'RC_modelica' #RC = simple function; RC_modelica; Room_modelica
20 hvac_control = False # if True modelica controls the heating and cooling
21 radiance = False
22 max_power = 100 # specific power input W/m2
23 parameter = get_parameter(l_wall, U_wall, U_window, A_room, h_room, \
24                             A_window, step_size, model, hvac_control, radiance, max_power)
25
26 # parameters for RL_Agent
27 parameter['agent']['Agent'] = 'DDPG' # DDPG, RDPG
28 parameter['agent']['actions'] = 0 # 0 = Q_hvac+Tvis; 1 = Q_hvac; 2 = Tvis
29
30 # set the hyperparameters for the neural network
31 parameter['agent']['NN']['network_architecture'] = 'MLP' # MLP, CNN, LSTM
32 parameter['agent']['NN']['hidden_layers'] = 2
33 parameter['agent']['NN']['layer_size_1'] = 400
34 parameter['agent']['NN']['layer_size_2'] = 300
35 parameter['agent']['NN']['activation'] = 'relu'
36 parameter['agent']['NN']['demand_charge_scale'] = 1
37
38 parameter['agent']['setting']['n_step'] = 1
39 parameter['agent']['setting']['forecast_hours'] = 1
40 parameter['agent']['setting']['training_days'] = 1
41 parameter['agent']['setting']['test_days'] = 7
42 parameter['agent']['setting']['noise_process'] = 'OU_noise' # method for OU_noise,
Gauss_noise or param_noise
43 parameter['agent']['setting']['forecast_col'] = ['weaTDryBul', 'S_irr', 'TOU_tariff',
'occupancy']
44 parameter['agent']['setting']['episodes'] = 3000
45 parameter['agent']['setting']['episodes_with_noise'] = 0.5
46 parameter['agent']['setting']['exploration_episodes'] = 50
47
48 parameter['agent']['replay_buffer']['Buffer'] = 'Uniform' # Uniform, PER, HPER
49 parameter['agent']['replay_buffer']['max_buffer_size'] = int(1e9)
50 parameter['agent']['replay_buffer']['batch_size'] = 64
51
52 parameter['agent']['scaling']['forecast_norm_data'] =
np.array([[0.+273.15,0.,-0.16225,0.],[34.+273.15,1000.,0.16225,2.1]])
53
54 # set flags
55 parameter['agent']['flags']['plot'] = True
56 parameter['agent']['flags']['print'] = True
57 parameter['agent']['flags']['hp_tune'] = False
58 parameter['agent']['flags']['load_model'] = False
59 parameter['agent']['flags']['save_models'] = True
60
61 weather_config = {}
```

```

62 weather_config['start_time'] = '2019-01-01 00:00:00'
63 weather_config['stepsize'] = parameter['input_data']['step_size'] * 60*60 #set hourly
64 timestep in parameter
65 weather_config['weather_dir'] = r'resources\weather\\'
66 weather_config['location'] = ['Oakland'] # enter the name of the city as a list for the
67 weather file, as string (None = random)
68 weather_config['weather_path'] = get_weather_files(weather_config)
69 while not weather_config['weather_path']:
70     print('location not available in this path')
71     weather_config['location'] = [input('enter new location: ')]
72     weather_config['weather_path'] = get_weather_files(weather_config)
73     print(weather_config['weather_path'])
74 weather_config['weather_columns'] = ['waeCelHei','waeCloTim','waeHDifHor','waeHDirNor',
75                                     'waeHGloHor','waeHHorIR','waeNOpa','waeNTot',
76                                     'waePAtm','waeRelHum','waeSolTim','waeSolZen',
77                                     'waeTBLaSky','waeTDewPoi','waeTDryBul',
78                                     'waeTWetBul','waeWinDir','waeWinSpe']
79
80 job_id = 0
81 mode = input("enter train or test: ")
82 input_files = files_handler(weather_config, parameter)
83 if mode == 'train':
84     train_results = RL_results_handler()
85     test_results = RL_results_handler()
86     env = tf_py_environment.TFPyEnvironment(Room(parameter))
87     Agent = agent(env, job_id, parameter, input_files, train_results, test_results)
88     train_results, test_results = Agent.train()
89 elif mode == 'test':
90     parameter['agent']['flags']['load_model'] = True
91     parameter['agent']['flags']['save_models'] = False
92     parameter['zone']['t_init'] = 22.5+273.15
93     test_results = RL_results_handler()
94     env = tf_py_environment.TFPyEnvironment(Room(parameter))
95     Agent = agent(env, job_id, parameter, input_files, None, test_results)
96     testing = 'select' #input('"select" for new location or "rand" for random
97     initialized location: ') #rand to choose randomly from weather data set and set
98     random date, select when setting new location and date
99     if testing == 'select':
100         weather_config['location'] = ['Oakland']#[input('Enter City for the location
101         file: ')]# select city for weather file, as string (None = random)
102         weather_config['weather_path'] = get_weather_files(weather_config)
103         while not weather_config['weather_path']:
104             print('location not available in this path')
105             weather_config['location'] = [input('enter new location: ')]
106             weather_config['weather_path'] = get_weather_files(weather_config)
107             print('What date (month,day) should be the first day of the episode?')
108             date = [8,1]#[int(input('month: ')), int(input('day: '))]
109             else: date = [None,None]
110
111 data = files_handler(weather_config, parameter).input_files['file0']
112 training_days = 7 #int(input('how many days should be tested? '))
113 test = Agent.test(testing ,data, training_days, date[0], date[1])

```

Gridsearch

```
1 import copy
2 import multiprocessing as mp
3 import os
4 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
5 import pandas as pd
6 from sklearn.model_selection import ParameterGrid
7 import tensorflow as tf
8 from tf_agents.environments import tf_py_environment
9 from parameter_handler import get_parameter
10
11
12 def ddpg_worker(job):
13     from environment import Room
14     from input_data import get_weather_files
15     from input_data_handler import files_handler
16     from RL_Agent import agent
17     from parameter_handler import get_parameter
18
19     job_id = job[0]
20     params = job[1]
21
22     weather_config = {}
23     weather_config['start_time'] = '2019-01-01 00:00:00'
24     weather_config['stepsize'] = params['input_data']['step_size'] * 60*60 #set hourly
25     timestep in parameter
26     weather_config['weather_dir'] = r'resources\weather\\'
27     weather_config['location'] = ['Oakland'] # enter the name of the city as a list for
28     the weather file, as string (None = random)
29     weather_config['weather_path'] = get_weather_files(weather_config)
30     weather_config['weather_columns'] =
31         ['weaCelHei','weaCloTim','weaHDifHor','weaHDirNor',
32             'weaHGloHor','weaHHorIR','weaNOpa','weaNTot',
33             'weaPAtm','weaRelHum','weaSolTim','weaSolZen',
34             'weaTBlasSky','weaTDewPoi','weaTDryBul',
35             'weaTWetBul','weaWinDir','weaWinSpe']
36
37     input_files = files_handler(weather_config, params)
38     ## load the environment and wrap it to tf_environment
39     envpy = Room(params)
40     env = tf_py_environment.TFPyEnvironment(envpy)
41
42     Agent = agent(env, job_id, params, input_files, None, None)
43
44     return Agent.train()
45
46 if __name__ == '__main__':
47     from parameter_handler import get_parameter
48     l_wall = 4 # Length of wall, in m
49     U_wall = 3.293 # U-value wall, in W/m2K
50     U_window = 6.923 # U-value window, in W/m2K
51     A_room = 14 # Room area, in m2
52     h_room = 3.95 # Room height, in m
53     A_window = l_wall * h_room * 0.33 # Window area, in m2
54     step_size = 60/60 # hour
55     model = 'RC_modelica' #RC = simple function; RC_modelica; Room_modelica
56     hvac_control = False # if True modelica controls the heating and cooling
57     radiance = False
58     max_power = 100 # specific power input W/m2
59     parameter = get_parameter(l_wall, U_wall, U_window, A_room, h_room,
60                               A_window, step_size, model, hvac_control, radiance, max_power)
61
62     # Make grid
63     agent_algorithm = list(['DDPG', 'RDPG'])
64     agent_architecture = list(['MLP', 'LSTM'])
```

```

62     forecast_hours = list([1,4])
63     noise_process = list(['OU_noise', 'Gauss_noise', 'param_noise'])
64     buffer = list(['Uniform', 'PER', 'HUPER'])
65     hidden_layer = list([1,2])
66     layer1 = list([400,500,600])
67     layer2 = list([300,400,500])
68
69     parameterset = dict(agent_algorithm = agent_algorithm, agent_architecture =
70                           agent_architecture,\n                               forecast_hours = forecast_hours, noise_process = noise_process,\n                               buffer = buffer)
71
72
73     params = list(ParameterGrid([parameterset]))
74     jobs = copy.deepcopy(params)
75     removed = 0
76     for dict_p in params:
77         if dict_p['agent_algorithm'] == 'RDPG' and dict_p['agent_architecture'] !=\n             'LSTM':
78             jobs.remove(dict_p)
79             removed += 1
80             continue
81         if dict_p['agent_algorithm'] == 'DDPG' and dict_p['agent_architecture'] ==\n             'LSTM':
82             jobs.remove(dict_p)
83             removed += 1
84             continue
85         if dict_p['agent_algorithm'] == 'DDPG' and dict_p['forecast_hours'] == 1:
86             jobs.remove(dict_p)
87             removed += 1
88             continue
89     parameter_list = []
90     for params in jobs:
91         job_params = {}
92         job_params['inputs'] = copy.deepcopy(parameter['inputs'])
93         job_params['input_data'] = copy.deepcopy(parameter['input_data'])
94         job_params['window'] = copy.deepcopy(parameter['window'])
95         job_params['zone'] = copy.deepcopy(parameter['zone'])
96         job_params['constraints'] = copy.deepcopy(parameter['constraints'])
97         job_params['somodel'] = copy.deepcopy(parameter['somodel'])
98         job_params['model'] = copy.deepcopy(parameter['model'])
99         job_params['tariff'] = copy.deepcopy(parameter['tariff'])
100        job_params['agent'] = copy.deepcopy(parameter['agent'])
101        job_params['agent']['flags']['plot'] = False
102        job_params['agent']['flags']['print'] = False
103        job_params['agent']['flags']['hp_tune'] = True
104        job_params['agent']['flags']['load_model'] = False
105        job_params['agent']['flags']['save_models'] = True
106        job_params['agent']['Agent'] = params['agent_algorithm']
107        job_params['agent']['NN']['network_architecture'] = params['agent_architecture']
108        job_params['agent']['NN']['hidden_layers'] = params['hidden_layer']
109        job_params['agent']['NN']['layer_size_1'] = params['layer1']
110        job_params['agent']['NN']['layer_size_2'] = params['layer2']
111        job_params['agent']['setting']['forecast_hours'] = params['forecast_hours']
112        job_params['agent']['setting']['noise_process'] = params['noise_process']
113        job_params['agent']['replay_buffer']['Buffer'] = params['buffer']
114        parameter_list.append(job_params)
115
116    # Run in parallel
117    pool = mp.Pool(mp.cpu_count()-1)
118    physical_devices = tf.config.list_physical_devices()
119    print(physical_devices)
120    data = pool.map(ddpg_worker, jobs)
121    pool.close()
122
123    # Convert to pandas
124    data = pd.DataFrame(data)
125    data.to_csv('logs/job_results.csv')

```

Anhang C: RL-Setup

Agent

```
1 import calendar
2 import copy
3 from datetime import datetime, timedelta
4 import json
5 import numpy as np
6 import random
7 import tensorflow as tf
8 from tensorflow.keras.layers.experimental.preprocessing import Normalization
9
10 import tf_agents
11 import time
12 print('tensorflow', tf.__version__)
13 print('tf_agents', tf_agents.__version__)
14
15 from AC_NN import AC_network
16 from input_data_handler import RL_results_handler
17 from ReplayBuffer import Uniform, PER, HVPER
18 from plot import episodeplot, runplot
19
20 class agent():
21     def __init__(self, env, job_id, parameter, input_files, train_results,
22                  test_results):
23
24         ''' agent in the RL framework
25             Inputs: env ..... initialized environment (tf_Py_environment)
26                     parameter ..... parameters set in parameter_handler.py and Main.py
27                     input_files ..... all selected weather files in a dict as
28                         input_handlers
29                     train_results .... Results handler for training data
30                     test_results ..... Results handler for test data
31
32         ...
33
34         self.parameter = parameter
35         self.env = env
36
37         self.input_files = input_files
38
39         if parameter['agent']['flags']['hp_tune']:
40             # tf.config.threading.set_inter_op_parallelism_threads(1)
41             # tf.config.threading.set_intra_op_parallelism_threads(1)
42             self.job_id = job_id
43             self.train_results = RL_results_handler()
44             self.test_results = RL_results_handler()
45         else:
46             self.train_results = train_results
47             self.test_results = test_results
48
49         self.T_train = int(parameter['agent']['setting']['training_days']* 24 / \
50                            parameter['agent']['stepsize']) # length of train episode
51         self.T_test = int(parameter['agent']['setting']['test_days']*24 / \
52                            parameter['agent']['stepsize']) # length of train episode
53         self.update_freq = int(6) # update frequency of neural networks
54         self.n_Step = parameter['agent']['setting']['n_step']
55         self.dtype = parameter['agent']['dtype']
56
57         # decrease noise scale to 0 after episodes with noise
58         self.noise_decrease = 1-(1/(parameter['agent']['setting']['episodes'] * \
59                               parameter['agent']['setting']['episodes_with_noise']))
59
59
60         # select noise scale for action noise shape(action_dim) or param noise shape(1)
61         if parameter['agent']['setting']['noise_process'] == 'param_noise':
62             self.noise_scale = np.array(0.6)
63             #np.array(parameter['agent']['scaling']['noise_scale'][1])
64             self.init_noise_scale = np.array(0.6)
```

```

60         #np.array(parameter['agent']['scaling']['noise_scale'][1])
61     else:
62         self.noise_scale = np.array(parameter['agent']['scaling']['noise_scale'])
63         self.init_noise_scale =
64         np.array(parameter['agent']['scaling']['noise_scale'])
65
66     if parameter['agent']['flags']['print']:
67         print('update frequency ',self.update_freq, ' steps')
68
69     # NN Paramters
70     self.Agent = parameter['agent']['Agent']
71     if self.Agent == 'RDPG':
72         self.parameter['agent']['NN']['network_architecture'] = 'LSTM'
73
74     # set dimension of NN output
75     self.action_dim = env._action_spec.shape[0]
76
77     # get max action set in environment
78     self.action_bound_range = env.action_spec().maximum
79     self.state_dim = env.reset()[3].shape[0]
80
81     # shape for forecast data
82     if self.parameter['agent']['NN']['network_architecture'] == 'CNN':
83         self.forecast_dim = (parameter['agent']['setting']['forecast_hours'],
84                             len(parameter['agent']['setting']['forecast_col']))
85     else:
86         self.forecast_dim = (parameter['agent']['setting']['forecast_hours']\*
87                             * len(parameter['agent']['setting']['forecast_col']))
88
89     # normalization layers prior to NN input
90     self.norm_state_layer = Normalization()
91     self.norm_state_layer.adapt(parameter['agent']['scaling']['state_norm_data'])
92     self.norm_forecast_layer = Normalization()
93
94     self.norm_forecast_layer.adapt(parameter['agent']['scaling']['forecast_norm_data'])
95
96     self.norm_action_layer = Normalization()
97     if parameter['agent']['actions'] == 0:
98
99         self.norm_action_layer.adapt(parameter['agent']['scaling']['actions_norm_data'])
100
101    elif parameter['agent']['actions'] == 1:
102
103        self.norm_action_layer.adapt(self.parameter['agent']['scaling']['actions_norm_data']\
104                                     [:,0].reshape(2,self.action_dim))
105
106    else:
107
108        self.norm_action_layer.adapt(self.parameter['agent']['scaling']['actions_norm_data']\
109                                     [:,1].reshape(2,self.action_dim))
110
111    # initialize replay buffer according to replay
112    if parameter['agent']['replay_buffer']['Buffer'] == 'PER':
113        self.buffer = PER(parameter['agent']['replay_buffer']['max_buffer_size'],
114                           parameter['agent']['dtype'], self.forecast_dim)
115    elif parameter['agent']['replay_buffer']['Buffer'] == 'HUPER':
116        self.buffer = HUPER(parameter['agent']['replay_buffer']['max_buffer_size'],
117                           parameter['agent']['dtype'], self.forecast_dim)
118
119    else:
120        self.buffer =
121        Uniform(parameter['agent']['replay_buffer']['max_buffer_size'],
122                           parameter['agent']['dtype'], self.forecast_dim)

```

```

110     # importance sampling for prioritized Replay Buffer (PER and HPER)
111     if isinstance(self.buffer, Uniform):
112         self.importance_sampling = False
113     else:
114         self.importance_sampling = True
115         self.beta = 0.5 # beta corrects the prioritized sampling probability and
116         changes linear over time to 1
117         self.beta_increase = 1 + (1/(parameter['agent']['setting']['episodes'] * 
118             self.T_train) * self.update_freq / 2)
119
120     # initialize actor and critic network
121     self.agent_network = AC_network(self.state_dim, self.forecast_dim,
122         self.action_dim, self.action_bound_range,\ 
123             self.norm_forecast_layer, self.norm_state_layer,
124             self.norm_action_layer, parameter)
125
126     #load saved network model
127     if parameter['agent']['flags']['load_model'] == True:
128         self.agent_network.load_model()
129     else:
130
131         self.agent_network.build_actor(self.parameter['agent']['NN']['network_archite
132         cture'])
133
134         self.agent_network.build_critic(self.parameter['agent']['NN']['network_archit
135         ecture'])
136
137     if parameter['agent']['flags']['print']:
138         self.agent_network.actor.summary()
139         self.agent_network.critic.summary()
140
141     def test(self, testing, data, episode_length, month, day):
142         ''' testing of the Agent with different weather conditions as in training '''
143         if testing == 'rand':
144             if self.parameter['agent']['flags']['hp_tune']:
145                 test_episodes = 1
146             else:
147                 test_episodes = 10
148         else:
149             test_episodes = 1
150             self.T_test = episode_length * 24
151
152         for episode in range(test_episodes):
153             if testing == 'rand':
154                 if len(self.input_files.input_files) > 1 or episode == 0:
155                     self.inputs_data, self.parameter['somodel'], year = \
156                         self.input_files.select_file(self.parameter['somodel'])
157
158                     # randomly select the start date for this episode
159                     if self.parameter['agent']['flags']['hp_tune']:
160                         month = 8
161                         day = 1
162                     else:
163                         month = random.choice(np.arange(1,13))
164                         day =
165                             random.choice(np.arange(1,calendar.monthrange(year,month)[1]+1))
166                     if month >= 12:
167                         month = min(month,12)
168                         day = min(day,
169                             calendar.monthrange(year,month)[1]-self.parameter['agent']['setting']
170                             ['test_days']+1)
171                     start_time = datetime(year, month, day,0,0)
172                     save_figure = False
173
174                 else:
175                     self.inputs_data, self.parameter['somodel'], year = \
176                         self.input_files.choose_file(data,self.parameter['somodel'])

```

```

163     if month >= 12:
164         month = min(month,12)
165         day = min(day, calendar.monthrange(year,month)[1]-episode_length-1)
166         start_time = datetime(year, month, day,0,0)
167         save_figure = True
168         current_time = start_time
169
170     # Append mapping to fmu inputs (for Buildings Library only)
171     if not 'RCmodel' in self.parameter['model']['fmu_path']:
172         weather_offset = len(self.inputs_data.cols_inputs) +
173             len(self.inputs_data.cols_data)
174         for i, c in enumerate(self.inputs_data.cols_weather):
175             #parameter['model']['inputs_map'][c] = weather_offset + i
176             self.parameter['model']['inputs_map'][c] = c
177
178     add_noise = False
179     # reset the environment to start setting
180     state_t = np.array(self.env.reset()[:3]).reshape(1,1)
181     # get the init forecast
182     forecast_t = self.inputs_data.get_forecast(current_time)
183     n_step_forecast = np.zeros(shape=(forecast_t.shape))
184     n_step_state = np.zeros(shape=(1,1))
185     for t in range(self.T_test):
186         if self.Agent == 'RDPG':
187             if n_step_state.shape[0] == 1:
188                 action_t =
189                     self.agent_network.take_action_lstm(state_t[0],forecast_t,
190                     add_noise)
191             elif 1 < n_step_state.shape[0] < 4 or t < 4:
192                 action_t =
193                     self.agent_network.take_action_lstm(n_step_state[1:],n_step_forecast[1:],add_noise)
194             else:
195                 action_t =
196                     self.agent_network.take_action_lstm(n_step_state,n_step_forecast,
197                     add_noise)
198             else:
199                 action_t = self.agent_network.take_action(state_t, forecast_t,
200                     add_noise)
201
202             if self.parameter['agent']['actions'] == 1:
203                 Q = action_t[0]
204                 Tvis = np.array([0.6])
205             elif self.parameter['agent']['actions'] == 2:
206                 Q = np.array([0.1])
207                 Tvis = action_t[0]
208             else:
209                 Q = action_t[0][0]
210                 Tvis = action_t[0][1]
211             env_input = self.inputs_data.get_inputs(Q, Tvis, current_time,
212             start_time)
213
214             _, rwd_t, _, state_t_pls_n = self.env.step(env_input)
215             self.test_results.append_res(self.env, 0, episode, current_time)
216             current_time += timedelta(hours=self.parameter['input_data']['step_size']))
217
218             # set the next forecast
219             forecast_t_pls_n = self.inputs_data.get_forecast(current_time)
220
221             n_step_state = np.append(n_step_state,state_t, axis=0)
222             n_step_forecast = np.append(n_step_forecast,forecast_t, axis=0)
223
224             if (t+1) >= self.n_Step:

```

```

217         # store the experience into the buffer for updating the critic
218         network
219         n_step_state = np.delete(n_step_state, 0, axis=0)
220         n_step_forecast = np.delete(n_step_forecast, 0, axis=0)
221
222         state_t = state_t_pls_n
223         forecast_t = forecast_t_pls_n
224
225         if (t+1) % self.T_test == 0:
226             rr = self.test_results.data['reward_0'].iloc[-self.T_test: ].sum() + \
227                 self.test_results.data['reward_1'].iloc[-self.T_test: ].sum()
228             if self.parameter['agent']['flags']['print']:
229                 print('Episode %d : Total Penalty = %f' % (episode+1, rr))
230             if self.parameter['agent']['flags']['plot']:
231                 episodeplot(self.parameter,
232                             self.test_results.data.iloc[-int(self.T_test):], save =
233                             save_figure, fig_name =
234                             str(self.agent_network.AC_name)+str(start_time).split(
235                             ')')[0])
236
237         return self.test_results
238
239     def train(self):
240         ''' main training function of the Agent w
241             maximum action taken in the episode/day/n_step gets added to every step
242             delayed storage of the experience'''
243         time_start = time.time()
244         experience_cnt = 0
245         logs = 0
246         rand = True
247         add_noise = True
248         distance = 0.6
249
250         for episode in range(self.parameter['agent']['setting']['episodes']):
251             # reset the environment to start setting
252             if len(self.input_files.input_files) > 1 or episode == 0:
253                 self.inputs_data, self.parameter['somodel'], year = \
254                     self.input_files.select_file(self.parameter['somodel'])
255
256             # randomly select the start date for this episode
257             month = random.choice(np.arange(1,13))
258             day = random.choice(np.arange(1,calendar.monthrange(year,month)[1]+1))
259             if month >= 12:
260                 month = min(month,12)
261                 day = min(day,
262                           calendar.monthrange(year,month)[1]-self.parameter['agent']['setting']['trainning_days']-1)
263             start_time = datetime(year, month, day,0,0)
264             current_time = start_time
265
266             # Append mapping to fmu inputs (for Buildings Library only)
267             if not 'RCmodel' in self.parameter['model']['fmu_path']:
268                 weather_offset = len(self.inputs_data.cols_inputs) +
269                 len(self.inputs_data.cols_data)
270                 for i, c in enumerate(self.inputs_data.cols_weather):
271                     #parameter['model']['inputs_map'][c] = weather_offset + i
272                     self.parameter['model']['inputs_map'][c] = c
273
274             # reset the environment to the startvalues
275             state_t = np.array(self.env.reset()[3]).reshape(1,1)
276
277             # get the init forecast
278             forecast_t = self.inputs_data.get_forecast(current_time)
279             done_t = False

```

```

273     # init add OUA-noise or Guassian noise process at the start of every
274     # episode or
275     # add param noise to the actor network
276     if add_noise:
277         if self.parameter['agent']['setting']['noise_process'] == 'param_noise':
278             self.agent_network.param_noise_process.calc_scale(distance)
279             self.agent_network.parameter_noise_handling()
280         else:
281             self.noise_scale = self.noise_scale * self.noise_decrease
282             self.agent_network.action_noise_handler(self.noise_scale)
283
284     # initialize n_step buffer
285     n_step_state = np.zeros(shape=(1,1))
286     n_step_forecast = np.zeros(shape=(forecast_t.shape))
287     n_step_forecast_t_pls_n = np.zeros(shape=(forecast_t.shape))
288     n_step_actions = np.zeros(shape=(1, self.action_dim))
289     n_step_rwrd_t = np.zeros(shape=(1,1))
290     n_step_demand = np.zeros(shape=(1,1))
291     n_step_state_t_pls_n = np.zeros(shape=(1,1))
292
293     # start of the episode
294     for t in range(self.T_train):
295         if self.Agent == 'RDPG':
296             if n_step_state.shape[0] == 1:
297                 action_t =
298                     self.agent_network.take_action_lstm(state_t[0], forecast_t,
299                     add_noise)
300             elif 1 < n_step_state.shape[0] < 4 or t < 4:
301                 action_t =
302                     self.agent_network.take_action_lstm(n_step_state[1:], n_step_forec
303                     ast[1:], add_noise)
304             else:
305                 action_t =
306                     self.agent_network.take_action_lstm(n_step_state, n_step_forecast,
307                     add noise)
308         else:
309             action_t = self.agent_network.take_action(state_t, forecast_t,
310             add_noise)
311
312             if self.parameter['agent']['actions'] == 1:
313                 Q = action_t[0]
314                 Tvis = np.array([0.6])
315             elif self.parameter['agent']['actions'] == 2:
316                 Q = np.array([0.])
317                 Tvis = action_t[0]
318             else:
319                 Q = action_t[0][0]
320                 Tvis = action_t[0][1]
321             env_input = self.inputs_data.get_inputs(Q, Tvis, current_time,
322             start_time)
323             _, rwrd_t, _, state_t_pls_n = self.env.step(env_input)
324             self.train_results.append_res(self.env, logs, episode, current_time)
325             logs=0
326             current_time +=
327             timedelta(hours=self.parameter['input_data']['step_size']))
328
329     # set the next forecast
330     forecast_t_pls_n = self.inputs_data.get_forecast(current_time)
331
332     if (t+1) == self.T_train:
333         done_t = True
334
335     # store the step in the n_step_memory for calculation of the
336     n_step_reward

```

```

326     n_step_state = np.append(n_step_state,state_t, axis=0)
327     n_step_forecast = np.append(n_step_forecast,forecast_t, axis=0)
328     n_step_forecast_t_pls_n =
329         np.append(n_step_forecast_t_pls_n,forecast_t_pls_n, axis=0)
330     n_step_actions = np.append(n_step_actions, action_t, axis=0)
331     n_step_rwrd_t = np.append(n_step_rwrd_t,np.array([rwrd_t[0][0]]))
332     n_step_demand = np.append(n_step_demand,np.array([rwrd_t[0][1]]))
333     n_step_state_t_pls_n = np.append(n_step_state_t_pls_n,state_t_pls_n,
334                                     axis=0)
335
336     # Calculation of the n_step_reward
337     if (t+1) >= self.n_Step:
338         # store the experience into the buffer for updating the critic
339         # network
340         n_step_state = np.delete(n_step_state,0, axis=0)
341         n_step_forecast = np.delete(n_step_forecast,0, axis=0)
342         n_step_actions = np.delete(n_step_actions,0, axis=0)
343         if self.Agent == 'RDPG':
344             n_step_action = n_step_actions
345         else:
346             n_step_action = np.array([n_step_actions[0]])
347             n_step_rwrd_t = np.delete(n_step_rwrd_t,0, axis=0)
348             n_step_demand = np.delete(n_step_demand,0, axis=0)
349             n_step_state_t_pls_n = np.delete(n_step_state_t_pls_n,0, axis=0)
350             n_step_forecast_t_pls_n = np.delete(n_step_forecast_t_pls_n,0,
351                                               axis=0)
352             reward = 0.
353
354         if self.Agent != 'RDPG':
355             for j, (rwrd_j) in enumerate(n_step_rwrd_t):
356                 reward += rwrd_j *
357                     self.parameter['agent']['NN']['discount_factor']**(j+1)
358             demand_charge = n_step_demand[-1]
359             forecast = n_step_forecast[0]
360             next_forecast = n_step_forecast_t_pls_n[-1]
361
362         else:
363             reward = n_step_rwrd_t + n_step_demand[-1]
364             reward = reward.reshape(self.n_Step,1)
365             demand_charge = n_step_demand
366             forecast = n_step_forecast
367             next_forecast = n_step_forecast_t_pls_n
368
369             self.buffer.add_experience(n_step_state, forecast, n_step_action,
370                                       reward, demand_charge,
371                                         n_step_state_t_pls_n, next_forecast,
372                                         done_t)
373
374     # TRAINING AND UPDATING THE NETWORKS
375     if not rand and experience_cnt % self.update_freq == 0:
376         if isinstance(self.buffer,Uniform):
377             states_batch, forecast_batch, actions_batch, rewards_batch,
378             demand_batch, next_states_batch, next_forecast_batch,
379             done_batch =
380                 self.buffer.sample_batch(self.parameter['agent']['replay_buffer']
381                                         ['batch_size'])
382             indices = None
383             importance_weight = np.array([1])
384         else:
385             self.beta = min(self.beta * self.beta_increase,1)
386             states_batch, forecast_batch, actions_batch, rewards_batch,
387             demand_batch, next_states_batch, next_forecast_batch,
388             done_batch, indices, importance_weight =
389                 self.buffer.sample_batch(self.parameter['agent']['replay_buffer']
390                                         ['batch_size'], self.beta)

```

```

375
376     num_samples =
377     min(len(states_batch), self.parameter['agent']['replay_buffer']['batch
378         _size'])
379
380     states_batch = self.norm_state_layer(states_batch)
381     forecast_batch = self.norm_forecast_layer(forecast_batch)
382     actions_batch = self.norm_action_layer(actions_batch)
383     next_states_batch = self.norm_state_layer(next_states_batch)
384     next_forecast_batch = self.norm_forecast_layer(next_forecast_batch)
385
386     if self.parameter['agent']['NN']['network_architecture'] != 'CNN'
387         and self.Agent != 'RDPG':
388         forecast_batch = tf.reshape(forecast_batch, (num_samples,
389             self.forecast_dim))
390         next_forecast_batch =
391             tf.reshape(next_forecast_batch, (num_samples, self.forecast_dim))
392         actions_batch =
393             tf.reshape(actions_batch, (num_samples, self.action_dim))
394
395
396     # states_batch ..... shape(num_samples, n_Step, 1)
397     # forecast_batch ..... shape(num_samples, forecast_dim)
398     # actions_batch ..... shape(num_samples, actions_dim)
399     # rewards_batch ..... shape(num_samples, 1)
400     # demand_batch ..... shape(num_samples, 1)
401     # next_states_batch .... shape(num_samples, n_Step, 1)
402     # next_forecast_batch ... shape(num_samples, forecast_dim)
403     # done_batch ..... shape(num_samples, 1)
404     if self.Agent != 'RDPG':
405         logs, td_error =
406             self.agent_network.train_critic_network(num_samples, states_batch,
407                 forecast_batch, actions_batch,
408                     rewards_batch, demand_batch,
409                     next_states_batch, next_forecast_batch,
410                     done_batch, indices, importance_weight)
411             self.agent_network.train_actor_network(num_samples,
412                 states_batch, forecast_batch)
413             if self.parameter['agent']['setting']['noise_process'] ==
414                 'param_noise' and add_noise == True:
415                 actions =
416                     self.agent_network.actor([states_batch[:, 0], forecast_batch])
417                     actions[0] = actions[0]/self.action_bound_range
418                     p_actions =
419                         self.agent_network.perturbed_actor([states_batch[:, 0], forecas
420                             t_batch])
421                     p_actions[0] = p_actions[0]/self.action_bound_range
422                     actions =
423                         np.dstack(actions).reshape(num_samples, self.action_dim)
424                     p_actions =
425                         np.dstack(p_actions).reshape(num_samples, self.action_dim)
426                     distance = tf.math.sqrt(1/self.action_dim *
427                         tf.reduce_mean(tf.math.square(actions-p_actions)))
428
429             else:
430                 forecast_batch = tf.reshape(forecast_batch, (num_samples,
431                     self.n_Step, self.forecast_dim))
432                 next_forecast_batch =
433                     tf.reshape(next_forecast_batch, (num_samples, self.n_Step,
434                         self.forecast_dim))
435                     hist_t = np.concatenate((states_batch, forecast_batch), axis=2)
436                     hist_t_1 =
437                         np.concatenate((next_states_batch, next_forecast_batch), axis=2)
438                     logs, td_error =
439                         self.agent_network.train_critic_lstm(num_samples, hist_t,
440                             actions_batch, rewards_batch,

```

```

415                                         hist_t_1, done_batch, indices,
416                                         importance_weight)
417                                         self.agent_network.train_actor_lstm(num_samples, hist_t,
418                                         actions_batch)
419                                         if self.parameter['agent']['setting']['noise_process'] ==
420                                         'param_noise' and add_noise == True:
421                                         actions = self.agent_network.actor([hist_t])
422                                         actions[0] = actions[0]/self.action_bound_range
423                                         p_actions = self.agent_network.perturbed_actor([hist_t])
424                                         p_actions[0] = p_actions[0]/self.action_bound_range
425                                         actions =
426                                         np.dstack(actions).reshape(num_samples,self.action_dim)
427                                         p_actions =
428                                         np.dstack(p_actions).reshape(num_samples,self.action_dim)
429                                         distance = tf.math.sqrt(1/self.action_dim *
430                                         tf.reduce_mean(tf.math.square(actions-p_actions)))
431                                         self.buffer.batch_update(indices, td_error)
432                                         experience_cnt += 1
433                                         if episode >=
434                                         self.parameter['agent']['setting']['exploration_episodes']:
435                                         rand = False
436                                         if episode == self.parameter['agent']['setting']['episodes_with_noise']
437                                         *\
438                                         self.parameter['agent']['setting']['episodes']:
439                                         add_noise = False
440                                         state_t = state_t_pls_n
441                                         forecast_t = forecast_t_pls_n
442                                         state_t = state_t_pls_n
443                                         forecast_t = forecast_t_pls_n
444                                         #### Plot of 1 episode
445                                         if (t+1) % self.T_train == 0:
446                                         rr = self.train_results.data['reward_0'].iloc[-self.T_train:].
447                                         sum() + \
448                                         self.train_results.data['reward_1'].iloc[-self.T_train:].
449                                         sum()
450                                         if self.parameter["agent"]['flags']['print']:
451                                         print('Episode %d : Total Penalty = %f' % (episode, rr))
452                                         if self.parameter['agent']['flags']['plot'] and episode >
453                                         self.parameter['agent']['setting']['episodes'] - 15:
454                                         episodeplot(self.parameter,
455                                         self.train_results.data.iloc[-int(self.T_train):])
456                                         if experience_cnt % (10*self.T_train) == 0 and
457                                         self.parameter['agent']['flags']['save_models']:
458                                         self.agent_network.save_model()
459                                         # Plot of 1 run
460                                         if self.parameter['agent']['flags']['plot']:
461                                         runplot(self.train_results.data, save = True, fig_name =
462                                         str(self.agent_network.AC_name))
463                                         self.test('rand',0,0,0,0)
464                                         time_end = time.time()
465                                         if self.parameter['agent']['flags']['hp_tune']:
466                                         filename = str('logs/' + str(self.job_id) + '_'+
467                                         datetime.now().strftime("%Y_%m_%d-%I_%M_%S") + '.json')
468                                         run_data = {'train_data': self.train_results.data.to_json(), 'test_data':
469                                         self.test_results.data.to_json()}
470                                         mean_critic_loss = self.train_results.data['critic_loss']
471                                         mean_critic_loss = mean_critic_loss.drop(mean_critic_loss[mean_critic_loss
472                                         == 0].index).values
473                                         episode_reward =

```

```

465             self.train_results.data[['episode', 'reward_0']].groupby(['episode']).sum().va
466             lues +\
467             self.train_results.data[['episode', 'reward_1']].groupby(['ep
468             isode']).sum().values
469             res = {}
470             res['job_id'] = self.job_id
471             res['critic_loss'] = mean_critic_loss[-1]
472             res['mean_loss'] = mean_critic_loss[-100:].mean()
473             res['reward'] = episode_reward[-1].item()
474             res['mean_reward'] = episode_reward[-30:].mean()
475             res['test_reward'] =
476                 self.test_results.data['reward_0'].iloc[-self.T_test: ].sum() + \
477                     self.test_results.data['reward_1'].iloc[-self.T_test: ].s
478                     um()
479             res['energy_use [kWh]'] =
480                 abs(self.train_results.data['grid_import'].iloc[-self.T_train:]).sum() / 1e3 * \
481                     self.parameter['input_data']['step_size']
482             res['episodes'] = self.parameter['agent']['setting']['episodes']
483             res['n_step'] = self.parameter['agent']['setting']['n_step']
484             res['forecast_hours'] = self.parameter['agent']['setting']['forecast_hours']
485             res['batch_size'] = self.parameter['agent']['replay_buffer']['batch_size']
486             res['episodes_with_noise'] =
487                 self.parameter['agent']['setting']['episodes_with_noise']
488             res['exploration_episodes'] =
489                 self.parameter['agent']['setting']['exploration_episodes']
490             res['Agent'] = self.parameter['agent']['Agent']
491             res['network'] = self.parameter['agent']['NN']['network_architecture']
492             res['replay'] = self.parameter['agent']['replay_buffer']['Buffer']
493             res['num_layer_1'] = self.parameter['agent']['NN']['layer_size_1']
494             res['num_layer_2'] = self.parameter['agent']['NN']['layer_size_2']
495             res['hidden_layers'] = self.parameter['agent']['NN']['hidden_layers']
496             res['activation'] = self.parameter['agent']['NN']['activation']
497             res['act_learn'] = self.parameter['agent']['NN']['act_learning_rate']
498             res['crit_learn'] = self.parameter['agent']['NN']['crit_learning_rate']
499             res['noise_scale'] = self.parameter['agent']['scaling']['noise_scale']
500             res['discount_factor'] = self.parameter['agent']['NN']['discount_factor']
501             res['demand_charge_scale'] =
502                 self.parameter['agent']['NN']['demand_charge_scale']
503             res['value_function'] = self.parameter['agent']['flags']['valuefunction']
504             res['duration'] = str(timedelta(seconds=time_end - time_start))
505             res['resname'] = filename
506
507             with open(filename, 'a') as json_file:
508                 json_file.write(json.dumps(run_data))
509
510             return res
511
512     return self.train_results, self.test_results

```

Environment:

```

1 import pandas as pd
2 import sys
3 import random
4 import numpy as np
5 from tf_agents.environments import py_environment
6 from tf_agents.specs import array_spec
7 from tf_agents.trajectories import time_step as ts
8
9 class Room(py_environment.PyEnvironment):
10     """
11         Training environment of a thermal zone (room) for controls development and
12         evalauiton.
13     """
14     def __init__(self, parameter):
15
16         self.parameter = parameter
17         self.tvis_to_shgc = np.poly1d([[self.parameter['window']['coeff_b'],
18                                         self.parameter['window']['coeff_a']]])
19         # initiate the array for the calculation of the demand charge with
20         # [grid_import,C_demand]
21         self.demand_charge_period = np.zeros(shape=(1,2))
22         self.max_energy_cost = self.parameter['constraints']['heat_max']/1e3 *
23             max(parameter['tariff']['C_energy'].values())
24         self.max_demand_cost = self.parameter['constraints']['heat_max']/1e3 *
25             max(parameter['tariff']['C_demand'].values())
26
27         self.reward = 0
28         self.demand_costs_calc = np.zeros(shape=(1,2))
29         if self.parameter['agent']['actions'] == 0:
30             action_dim = 2
31         else:
32             action_dim = 1
33         self.action_spec = array_spec.BoundedArraySpec(shape=(action_dim,),,
34             dtype=np.float64, minimum=-self.parameter['constraints']['cool_max'],
35             maximum=self.parameter['constraints']['heat_max'], name='action')
36         self.observation_spec = array_spec.BoundedArraySpec(shape=(1,),,
37             dtype=np.float64, name='observation')
38
39         # Radiance for solar-optical model
40         if self.parameter['somodel']['use_radiance']:
41             sys.path.append(parameter['somodel']['path_radiance_handler'])
42             import radiance.radiance_emulator as radiance_handler
43             self.radiance_handler = radiance_handler
44             self.radiance = None
45
46         # FMI interface (load FMU)
47         if self.parameter['model']['use_fmu']:
48             from pyfmi import load_fmu
49             self.load_fmu = load_fmu
50             self.fmu_loaded = False
51
52     def setup_fmu(self, T_init, start_time):
53         """
54             Load and setup the FMU.
55
56             Inputs
57             T_init (float): Initial temperature for model, in K.
58             start_time (float): Start time of the model, in sceonds.
59         """
60         # Load FMU
61         self.fmu = self.load_fmu(self.parameter['model']['fmu_path'],
62                             log_level=self.parameter['model']['fmu_loglevel'])
63         # Initizlaise FMU
64         self.fmu.setup_experiment(start_time=start_time,

```

```

63                         stop_time=start_time+1,
64                         stop_time_defined=False)
65
66     # Parameterize FMU
67     param = self.parameter['model']['param']
68     param['T_init'] = T_init
69     self.fmu.set(list(param.keys()), list(param.values()))
70
71     self.fmu.initialize()
72     self.fmu_loaded = True
73
74     def action_spec(self):
75         return self._action_spec
76
77     def observation_spec(self):
78         return self._observation_spec
79
80     def get_info(self):
81         '''function returns variables calculated in the environment'''
82         return self.data.index, self.data.values
83
84     def heat_balance(self, inputs, T_room):
85         ''' heat balance model '''
86
87         Q_thermal = inputs['Q_int_th'] + inputs['Q_int_el'] + inputs['P_lights']
88         if not self.parameter['zone']['control_hvac']:
89             Q_thermal += inputs['Q_hvac']
90         if self.parameter['model']['include_solar_gains']:
91             Q_thermal += inputs['Q_solar']
92
93         if self.parameter['model']['use_fmu']:
94             if not self.fmu_loaded:
95                 # Careful, always initialized with 0 as start_time
96                 self.setup_fmu(T_room, 0)
97
98             # Set Inputs
99             fmu_inputs = {}
100            for k, v in self.parameter['model']['inputs_map'].items():
101                fmu_inputs[k] = inputs.loc[v]
102
103            # Add Q_thermal as input
104            if self.parameter['model']['inputs_Q_thermal']:
105                fmu_inputs[self.parameter['model']['inputs_Q_thermal']] = Q_thermal # W
106
107            # Add inputs from parameters
108            for k, v in self.parameter['model']['inputs'].items():
109                fmu_inputs[k] = v
110
111            # Set inputs
112            self.fmu.set(list(fmu_inputs.keys()), list(fmu_inputs.values()))
113
114            # Compute FMU
115            start_time = self.fmu.time
116            self.fmu.do_step(current_t=start_time,
117                             step_size=inputs['start_time']-start_time)
118
119            # Get Outputs
120            T_in = self.fmu.get(self.parameter['model']['outputs_T_in'])[-1]
121            if self.parameter['model']['outputs_Q_hvac']:
122                Q_hvac = self.fmu.get(self.parameter['model']['outputs_Q_hvac'])[-1]
123            else:
124                Q_hvac = 0
125
126        else:
127            T_in = ((Q_thermal * self.parameter['input_data']['step_size'])+ \
128                  inputs['T_out'] * 1/self.parameter['model']['param']['R1']+\
129                  T_room * self.parameter['model']['param']['C1'] / 3600)/\
130                  (1/self.parameter['model']['param']['R1'] + \
131                   self.parameter['model']['param']['C1'] / 3600)
132            Q_hvac = 0
133
134        return T_in, Q_hvac

```

```

127
128     def calc_illuminance(self, solar_Illumination, Tvis):
129         # The calculation is using the daylight factor from
130         # https://www.uk-ncm.org.uk/filelibrary/SBEM-Technical-Manual_v5.2.g_20Nov15.pdf
131         return (solar_Illumination * Tvis * (45*self.parameter['zone']['area']) / \
132                 (self.parameter['zone']['surface_area']*0.76))/100
133
134     def calculate_solar_power(self, inputs):
135         outputs = {}
136         outputs['shgc'] = self.tvis_to_shgc(inputs['Tvis'])
137         outputs['Q_solar'] = inputs['S_irr'] * self.parameter['window']['area'] * \
138             outputs['shgc']
139         tvis_to_state = pd.Series(self.parameter['somodel']['tvis_to_state'])
140         outputs['tint'] = tvis_to_state.iloc[tvis_to_state.index.get_loc(inputs['Tvis']),
141                                             method='nearest')]
142         outputs['uWin'] = 1 - outputs['tint'] / tvis_to_state.max()
143         if self.parameter['somodel']['use_radiance']:
144             if not self.radiance:
145                 self.radiance = \
146                     self.radiance_handler.Radiance(self.parameter['somodel']['config_file'],
147                                         regenerate=self.parameter['somodel']['regenerate_matrices'],
148                                         orient=self.parameter['somodel']['orientation'],
149                                         location=self.parameter['somodel']['location'],
150                                         filestruct=self.parameter['somodel']['filestruct'],
151                                         use_gendaymtx=False)
152         weather = pd.DataFrame({k:[v] for k,v in inputs.items() if
153                                k.startswith('wea') or k == 'start_time'})
154         weather.index = [pd.to_datetime('2020-01-01') + pd.DateOffset(seconds=ix)
155                           for ix in weather['start_time']]
156         outputs['uShade'] = [outputs['tint']] * 3 # homogenous control of ECs
157         radiance_outputs = self.radiance.compute(weather, outputs['uShade'],
158                                                 weather.index[0])
159         outputs['daylight'] = radiance_outputs[0][0]
160         outputs['glare'] = radiance_outputs[1][0]
161         # outputs['radiance'] = radiance_outputs
162         # Rough approximation of solar heat gain
163         outputs['Q_solar_radiance'] = sum(radiance_outputs[2:7]) +
164             radiance_outputs[8] + 0 * radiance_outputs[7]
165     else:
166         outputs['daylight'] = self.calc_illuminance(inputs['S_ill'],
167                                         inputs['Tvis']) # lux
168
169     return outputs
170
171     def calculate_lighting_power(self, inputs):
172         if inputs['daylight'] > inputs['wpi_min']:
173             P_light = 0
174             Ill_light = 0
175         else:
176             P_light = (inputs['wpi_min'] - inputs['daylight']) *
177                         self.parameter['zone']['eff_lights']
178             Ill_light = inputs['wpi_min'] - inputs['daylight']
179
180     return P_light, Ill_light
181
182     def calculate_demand_charge(self,inputs):
183         # for the reward calculation the demand charge is calculated for the n_step
184         # period
185         demand_cost = 0
186         self.demand_charge_period = np.append(self.demand_charge_period,\n187             (np.array([[inputs['grid_import']/1e3,inputs['C_demand']]])),axis=0)
188         if self.demand_charge_period.shape[0] > \

```

```

180     self.parameter['agent']['setting']['n_step']:
181         # delete the first row in the array to have the latest steps with length
182         # (n_step)
183         self.demand_charge_period = np.delete(self.demand_charge_period ,0, axis=0)
184         # take the maximum grid_import of each unique period and calculate the
185         # resulting demand costs
186         periods_power = np.split(self.demand_charge_period[:,0],
187             np.unique(self.demand_charge_period[:,1], return_index =
188             True)[1]))[1:]
189         periods_cost = np.unique(self.demand_charge_period[:,1])
190         for i in range(periods_cost.shape[0]):
191             demand_cost += np.max(periods_power[i]) * periods_cost[i]
192         demand_cost /= periods_cost.shape[0]
193         # add the base demand charge with the maximum grid_import of the latest
194         # steps with length (n_step)
195         demand_cost += np.max(self.demand_charge_period[:,0]) *
196         self.parameter['tariff']['C_demand']['base_rate']
197     return demand_cost
198
199     def _step(self, inputs):
200         # Parse inputs
201         data = pd.Series(inputs, index=self.parameter['inputs']['labels'])
202
203         # Calculate solar gains and daylighting in room
204         data = data.append(pd.Series(self.calculate_solar_power(data)))
205
206         # Calculate lighting requirement in room
207         data['P_lights'],data['Ill_light'] = self.calculate_lighting_power(data)
208
209         # Calculate heat balance
210         data['T_now'], data['Q_hvac_env'] = self.heat_balance(data, self.state[0])
211         self.state[0] = data['T_now']
212
213         # Electricity balance
214         if self.parameter['zone']['control_hvac']:
215             data['Q_hvac'] = data['Q_hvac_env']
216             data['P_hvac'] = (data['Q_hvac'] * self.parameter['zone']['eff_heat']) if
217             data['Q_hvac'] > 0 else \
218                 (abs(data['Q_hvac']) * self.parameter['zone']['eff_cool'])
219             data['grid_import'] = data['P_hvac'] + data['P_lights'] + data['Q_int_el']
220
221         # Cost calculation
222         data['energy_cost'] = data['grid_import'] *
223         self.parameter['input_data']['step_size'] * data['C_energy'] / 1e3
224         data['demand_cost'] = self.calculate_demand_charge(data)
225
226         # Temperature constraint
227         if data['T_now'] < data['t_min']:
228             data['Penalty_T_room'] = min(abs(data['t_min'] - data['T_now']),
229                                         self.parameter['constraints']['max_t_penalty'])
230         elif data['T_now'] > data['t_max']:
231             data['Penalty_T_room'] = min(abs(data['T_now'] - data['t_max']),
232                                         self.parameter['constraints']['max_t_penalty'])
233         else:
234             data['Penalty_T_room'] = 0
235
236         # penalty for tint status in the night
237         if data['S_ill'] == 0 and data['Tvis'] < 0.59:
238             data['Penalty_tint'] = self.parameter['constraints']['night_tint_penalty']
239         else:
240             data['Penalty_tint'] = 0
241
242         data['reward_0'] = data['energy_cost']/self.max_energy_cost +
243             (data['Penalty_T_room'] + data['Penalty_tint'])

```

```

235     data['reward_1'] =
236     data['demand_cost']/self.max_energy_cost/self.parameter['agent']['scaling']['rewa
237     rd_0']
238
239     reward = np.array([data['reward_0']**-1, data['reward_1']**-1])
240
241     self.data = data
242     return ts.transition(self.state, reward = reward , discount = 0.0)
243
244     def _reset(self):
245         '''reset is called at the start of the episode'''
246         if self.parameter['zone']['t_init']:
247             self.state = np.array([self.parameter['zone']['t_init']])
248         else:
249             self.state = np.array([random.uniform(self.parameter['zone']['t_init_min'],
250                                         self.parameter['zone']['t_init_max'])])
250
251         self.fmu_loaded = False
252         self.demand_costs_calc = np.zeros(shape=(1,2))
253
254         return ts.restart(self.state)

```

Actor Critic Network

```
1  from datetime import datetime
2  import numpy as np
3  import os
4  import pandas as pd
5  import sys
6  import tensorflow as tf
7  import tensorflow.keras
8  from tensorflow.keras.optimizers import Adam
9  from tensorflow.keras.layers import Dense, ReLU, LeakyReLU, ELU, Input, concatenate,
BatchNormalization, Dropout
10 from tensorflow.keras.layers import Conv1D, Flatten, Multiply, Add
11 from tensorflow.keras.layers import LSTM
12 from tensorflow.keras.losses import MeanSquaredError, Huber
13 from tensorflow.keras import Model
14 from tensorflow.keras.models import load_model
15 tf.keras.backend.set_floatx('float32')
16
17 from Noise import OU_Noise, Gauss_Noise, Param_Noise, add_OU_noise, add_Gauss_noise
18
19 class _actor_MLP():
20     def __init__(self, state_dim, forecast_dim, action_bound_range, parameter):
21         self.state_dim = state_dim
22         self.action_bound_range = tf.constant(action_bound_range, shape=(1,), dtype='float64')
23         self.forecast_dim = forecast_dim
24         self.num_ly1 = parameter['NN']['layer_size_1']
25         self.num_ly2 = parameter['NN']['layer_size_2']
26         self.hidden_layers = parameter['NN']['hidden_layers']
27         self.normalize_a = tf.constant(0.6-0.01, shape=(1,), dtype='float64')
28         self.normalize_b = tf.constant(0.01, shape=(1,), dtype='float64')
29         if parameter['NN']['activation'] == 'relu':
30             self.activation = ReLU
31         elif parameter['NN']['activation'] == 'leakyrelu':
32             self.activation = LeakyReLU
33         elif parameter['NN']['activation'] == 'elu':
34             self.activation == ELU
35
36         self.action = parameter['actions']
37
38     def model(self):
39
40         state = Input(shape=(self.state_dim,), name='state_input')
41         forecast = Input(shape=(self.forecast_dim,), name='forecast_input')
42
43         sf = concatenate([state, forecast])
44
45         sf = Dense(self.num_ly1, bias_initializer = 'zeros', name='param_noise')(sf)
46         sf = self.activation()(sf)
47         sf = BatchNormalization()(sf)
48         for i in range(self.hidden_layers):
49             sf = Dense(self.num_ly2, bias_initializer = 'zeros', name =
50             'hidden'+str(i+1))(sf)
51             sf = self.activation()(sf)
52             sf = BatchNormalization()(sf)
53
54         if self.action != 0:
55             if self.action == 1:
56                 action = Dense(1, activation='tanh')(sf)
57                 action = Multiply(name = 'Q')([action, self.action_bound_range])
58             elif self.action == 2:
59                 action = Dense(1)(sf)
60                 action = ReLU(max_value=1)(action)
61                 action = Multiply()( [action, self.normalize_a])
62                 action = Add(name = 'Tvis')([action, self.normalize_b])
```

```

62         return Model(inputs=[state, forecast], outputs=[action], name='actor')
63     else:
64         action1 = Dense(1, activation='tanh')(sf)
65         action1 = Multiply(name = 'Q')([action1, self.action_bound_range])
66
67         action2 = Dense(1)(sf)
68         action2 = ReLU()(action2)
69         action2 = Multiply()([action2, self.normalize_a])
70         action2 = Add(name = 'Tvis')([action2,self.normalize_b])
71
72     return Model(inputs=[state, forecast], outputs=[action1, action2],
73                  name='actor')
73
74 class _critic_MLP():
75     def __init__(self, state_dim, forecast_dim, action_dim, n_Step, parameter):
76         self.state_dim = state_dim
77         self.action_dim = action_dim
78         self.forecast_dim = forecast_dim
79         self.num_ly1 = parameter['NN']['layer_size_1']
80         self.num_ly2 = parameter['NN']['layer_size_2']
81         self.hidden_layers = parameter['NN']['hidden_layers']
82         if parameter['NN']['activation'] == 'relu':
83             self.activation = ReLU
84         elif parameter['NN']['activation'] == 'leakyrelu':
85             self.activation = LeakyReLU
86         elif parameter['NN']['activation'] == 'elu':
87             self.activation == ELU
88
89     def model(self):
90         state = Input(shape=(self.state_dim,), name='state_input')
91         forecast = Input(shape=(self.forecast_dim,), name='forecast_input')
92         action = Input(shape=(self.action_dim,), name='action_input')
93
94         sfa = concatenate([state, forecast, action])
95
96         sfa = Dense(self.num_ly1, bias_initializer = 'zeros')(sfa)
97         sfa = self.activation()(sfa)
98         sfa = BatchNormalization()(sfa)
99
100        for i in range(self.hidden_layers):
101            sfa = Dense(self.num_ly2, bias_initializer = 'zeros')(sfa)
102            sfa = self.activation()(sfa)
103            sfa = BatchNormalization()(sfa)
104
105        value = Dense(1, activation='linear', name = 'value')(sfa)
106        return Model(inputs=[state, forecast, action], outputs=value, name='critic')
107
108 class _actor_LSTM():
109     def __init__(self, state_dim, forecast_dim, action_bound_range, parameter):
110         self.observation_dim = (None,forecast_dim+state_dim)
111         self.action_bound_range = tf.constant(action_bound_range, shape=(1,), dtype='float64')
112         self.num_ly1 = parameter['NN']['layer_size_1']
113         self.num_ly2 = parameter['NN']['layer_size_2']
114         self.hidden_layers = parameter['NN']['hidden_layers']
115         self.normalize_a = tf.constant(0.6-0.01, shape=(1,), dtype='float64')
116         self.normalize_b = tf.constant(0.01, shape=(1,), dtype='float64')
117         if parameter['NN']['activation'] == 'relu':
118             self.activation = ReLU
119         elif parameter['NN']['activation'] == 'leakyrelu':
120             self.activation = LeakyReLU
121         elif parameter['NN']['activation'] == 'elu':
122             self.activation == ELU
123

```

```

124         self.action = parameter['actions']
125
126     def model(self):
127         observation = Input(shape=(self.observation_dim), name='state_input')
128
129         x = Dense(256, bias_initializer = 'zeros', name='param_noise')(observation)
130         x = self.activation()(x)
131         x = BatchNormalization()(x)
132
133         x = LSTM(256)(x)
134
135         if self.action != 0:
136             if self.action == 1:
137                 action = Dense(1, activation='tanh')(x)
138                 action = Multiply(name = 'Q')([action, self.action_bound_range])
139             elif self.action == 2:
140                 action = Dense(1)(x)
141                 action = ReLU(max_value=1)(action)
142                 action = Multiply()( [action, self.normalize_a])
143                 action = Add(name = 'Tvis')([action, self.normalize_b])
144
145             return Model(inputs=observation, outputs=[action], name='actor')
146
147         else:
148             action1 = Dense(1, activation='tanh')(x)
149             action1 = Multiply(name = 'Q')([action1, self.action_bound_range])
150
151             action2 = Dense(1)(x)
152             action2 = ReLU(max_value=1)(action2)
153             action2 = Multiply()( [action2, self.normalize_a])
154             action2 = Add(name = 'Tvis')([action2, self.normalize_b])
155
156         return Model(inputs=observation, outputs=[action1, action2], name='actor')
157
158     class _critic_LSTM():
159         def __init__(self, state_dim, forecast_dim, action_dim, n_Step, parameter):
160             self.observation_dim = (n_Step, forecast_dim+state_dim)
161             self.action_dim = (n_Step, action_dim)
162             self.num_ly1 = parameter['NN'][['layer_size_1']]
163             self.num_ly2 = parameter['NN'][['layer_size_2']]
164             self.hidden_layers = parameter['NN'][['hidden_layers']]
165             if parameter['NN'][['activation']] == 'relu':
166                 self.activation = ReLU
167             elif parameter['NN'][['activation']] == 'leakyrelu':
168                 self.activation = LeakyReLU
169             elif parameter['NN'][['activation']] == 'elu':
170                 self.activation = ELU
171
172         def model(self):
173             observation = Input(shape=(self.observation_dim), name='obs_input')
174             action = Input(shape=(self.action_dim), name='action_input')
175
176             observation_i = Dense(256, activation='relu')(observation)
177             observation_i = BatchNormalization()(observation_i)
178             observation_i = Dense(256, activation='relu')(observation_i)
179             observation_i = BatchNormalization()(observation_i)
180
181             action_i = Dense(64, activation='relu')(action)
182             action_i = BatchNormalization()(action_i)
183
184             x = concatenate([observation_i, action_i], name='concat')
185
186             x = LSTM(320, return_sequences = True)(x)
187
188             value = Dense(1, activation='linear', name = 'value')(x)

```

```

188
189     return Model(inputs=[observation,action], outputs=value, name='critic')
190
191 class AC_network():
192     def __init__(self, state_dim, forecast_dim, action_dim, action_bound_range, \
193                  forecast_norm_layer, state_norm_layer, action_norm_layer, parameter):
194         self.parameter = parameter
195         self.state_dim = state_dim
196         self.forecast_dim = forecast_dim
197         self.action_dim = action_dim
198         self.n_Step = self.parameter['agent']['setting']['n_step']
199         self.action_bound_range = action_bound_range
200
201         self.norm_forecast_layer = forecast_norm_layer
202         self.norm_state_layer = state_norm_layer
203         self.norm_action_layer = action_norm_layer
204
205         self.actor_network_types = {'MLP':_actor_MLP, 'CNN':_actor_CNN, 'LSTM': \
206                                     _actor_LSTM}
207         self.critic_network_types = {'MLP':_critic_MLP, 'CNN':_critic_CNN, 'LSTM': \
208                                     _critic_LSTM}
209
210     # set optimizer for network updates
211     self.actor_opt = Adam(self.parameter['agent']['NN']['act_learning_rate'])
212     self.critic_opt = Adam(self.parameter['agent']['NN']['crit_learning_rate'])
213     self.loss_function = MeanSquaredError()
214
215     self.lamb = 1. # initial weight for Q-value priority which changes over time to 0
216     self.lamb_decrease = 1-(1/(self.parameter['agent']['setting']['episodes']) * \
217                           int(parameter['agent']['setting']['training_days']* 24 / \
218                               parameter['agent']['stepsize'])) * 6)
219
220     self.AC_name = datetime.now().strftime("%Y-%m-%d-%I-%M-%S") + '-' + \
221     str(self.parameter['agent']['Agent'])+'_\'' \
222     + str(self.parameter['agent']['NN']['network architecture'])+'\'' \
223     + 'action' + str(self.parameter['agent']['actions'])+'\'' \
224     + 'layer' + \
225     str(self.parameter['agent']['NN']['layer_size_1'])+'\'' \
226     + str(self.parameter['agent']['NN']['layer_size_2'])+'\'' \
227     + str(self.parameter['agent']['NN']['hidden_layers'])+'\'' \
228     + str(self.parameter['agent']['NN']['activation'])+'\'' \
229     + str(self.parameter['agent']['setting']['n_step'])+'\'' \
230     + 'forecast' + \
231     str(self.parameter['agent']['setting']['forecast_hours'])+'\'' \
232     + str(self.parameter['agent']['replay_buffer']['Buffer'])+'\'' \
233     + 'batch' + \
234     str(self.parameter['agent']['replay_buffer']['batch_size'])+'\'' \
235     + str(self.parameter['agent']['setting']['noise_process'])+'\'' \
236     + 'valuefct' + \
237     str(self.parameter['agent']['flags']['valuefunction'])+'\'' \
238     + 'demand_scale' + \
239     str(self.parameter['agent']['NN']['demand_charge_scale'])

240     if os.path.exists('actor/' + self.AC_name + '.h5') == True and \
241     self.parameter['agent']['flags']['save_models'] == True:
242         job = input('Do you want to overwrite the existing network? (y/n) ')
243         if job == 'y':
244             pass
245         elif job == 'n':
246             sys.exit()
247
248     def build_actor(self, architecture):
249         self.actor = self.actor_network_types[architecture]\ \
250                     (self.state_dim, self.forecast_dim, self.action_bound_range,

```

```

242         self.parameter['agent']).model()
243     self.actor_target = self.actor_network_types[architecture]\ 
244         (self.state_dim, self.forecast_dim, self.action_bound_range,
245         self.parameter['agent']).model()
246     self.actor_target.set_weights(self.actor.get_weights())
247 
248     if self.parameter['agent']['setting']['noise_process'] == 'param_noise':
249         self.perturbed_actor = self.actor_network_types[architecture]\ 
250             (self.state_dim, self.forecast_dim, self.action_bound_range,
251             self.parameter['agent']).model()
252         self.perturbed_actor.set_weights(self.actor.get_weights())
253         actor_weights = self.actor.get_layer(name='param_noise').get_weights()
254         n_weights = actor_weights[0].size
255         layer_shape = actor_weights[0].shape
256         self.param_noise_process = Param_Noise(n_weights, layer_shape,
257             self.action_dim)
258         self.actor.compile(optimizer = self.actor_opt)
259 
260     self.actor.save('actor/' + self.AC_name + '.h5')
261     del self.actor
262     self.actor = load_model('actor/' + self.AC_name + '.h5')
263 
264     def build_critic(self, architecture):
265         self.critic = self.critic_network_types[architecture]\ 
266             (self.state_dim, self.forecast_dim, self.action_dim, self.n_Step,
267             self.parameter['agent']).model()
268         self.critic_target = self.critic_network_types[architecture]\ 
269             (self.state_dim, self.forecast_dim, self.action_dim, self.n_Step,
270             self.parameter['agent']).model()
271         self.critic_target.set_weights(self.critic.get_weights())
272         self.critic.compile(optimizer = self.critic_opt)
273 
274         self.critic.save('critic/' + self.AC_name + '.h5')
275     del self.critic
276     self.critic = load_model('critic/' + self.AC_name + '.h5')
277 
278     def parameter_noise_handling(self):
279         """ takes actor NN initialized in this class and
280         changes the weights of the NN with gaussian noise """
281         actor_weights = self.actor.get_layer(name='param_noise').get_weights()
282         noisy_weights = actor_weights.copy()
283         noise = self.param_noise_process.perturb_actor()
284         noisy_weights[0] = noisy_weights[0] + noise
285         self.perturbed_actor.get_layer(name='param_noise').set_weights(noisy_weights)
286 
287     def action_noise_handler(self,noise_scale):
288         if self.parameter['agent']['setting']['noise_process'] == 'OU_noise':
289             self.noise_Q = OU_Noise(mu = np.zeros(1), sigma = noise_scale[0], theta =
290                 0.2, dt=1e-1)
291             self.noise_Tvis = OU_Noise(mu = np.zeros(1), sigma = noise_scale[1], theta
292                 = 0.2, dt=1e-2)
293         else:
294             self.noise_Q = Gauss_Noise(noise_scale[0])
295             self.noise_Tvis = Gauss_Noise(noise_scale[1])
296 
297     def save_model(self):
298         self.actor.save('actor/' + self.AC_name + '.h5')
299         self.critic.save('critic/' + self.AC_name + '.h5')
300 
301     def load_model(self):
302         self.AC_name = input('filename:')
303         try:
304             self.actor = load_model('actor/' + self.AC_name + '.h5')
305             self.actor_target = load_model('actor/' + self.AC_name + '.h5')

```

```

298     self.critic = load_model('critic/' + self.AC_name + '.h5')
299     self.critic_target = load_model('critic/' + self.AC_name + '.h5')
300     print('actor-critic model successfully loaded')
301 except:
302     print('network not available')
303
304     job = input('Network not available! continue with new network (y/n)? ')
305     if job == 'y':
306         pass
307     else:
308         sys.exit()
309
310 def take_action(self, state_t, forecast_t, noise):
311     ''' Input for selecting an action
312         state_t..... current state (Room Temperature, cost of energy)
313             (np.array(shape=1,2))
314         forecast_t... forecast (np.array(shape=forecast_dim))
315         noise..... flag indicating if noise should be added during
316             exploration (boolean)
317         output
318             actn..... selected action (Energy Input, mode of facade/tvis)
319                 (np.array(shape=1,2))'''
320
321     # normalization for NN input
322     state = self.norm_state_layer(state_t)
323     forecast = self.norm_forecast_layer(forecast_t)
324
325     if self.parameter['agent']['NN']['network_architecture'] == 'MLP':
326         forecast = tf.reshape(forecast, (1, self.forecast_dim))
327     if self.parameter['agent']['setting']['noise_process'] == 'param_noise' and
328         noise:
329         action = self.perturbed_actor([state, forecast])
330     else:
331         action = self.actor([state, forecast])
332
333     if self.parameter['agent']['actions'] == 0:
334         Q = action[0][0]
335         Tvis = action[1][0]
336     elif self.parameter['agent']['actions'] == 1:
337         Q = action[0]
338         Tvis = 0.6
339     else:
340         Q = 0
341         Tvis = action[0]
342
343     if noise:
344         if self.parameter['agent']['setting']['noise_process'] == 'OU_noise':
345             Q, Tvis = add_OU_noise(self.noise_Q, Q, self.noise_Tvis, Tvis)
346         elif self.parameter['agent']['setting']['noise_process'] == 'Gauss_noise':
347             Q, Tvis = add_Gauss_noise(self.noise_Q, Q, self.noise_Tvis, Tvis)
348
349     # clip the values after adding the action noise to min and max values
350     if self.parameter['agent']['actions'] == 0:
351         Q = tf.clip_by_value(action[0], -self.action_bound_range,
352             self.action_bound_range)[0].numpy()
353         Tvis = tf.clip_by_value(action[1], 0.01, 0.6)[0].numpy()
354         actn = np.array([Q, Tvis])
355     elif self.parameter['agent']['actions'] == 1:
356         Q = tf.clip_by_value(action[0], -self.action_bound_range,
357             self.action_bound_range)[0].numpy()
358         actn = np.array([Q])
359     elif self.parameter['agent']['actions'] == 2:
360         Tvis = tf.clip_by_value(action[0], 0.01, 0.6)[0].numpy()
361         actn = np.array([Tvis])
362         actn = actn.reshape(1, self.action_dim)
363
364 return actn

```

```

356
357     def take_action_lstm(self, state_history_t, forecast_history_t, noise):
358         ''' Input for selecting an action
359             state_t..... current state (Room Temperature, cost of energy)
360             (np.array(shape=1,2))
361             forecast_t... forecast (np.array(shape=forecast_dim))
362             noise..... flag indicating if noise should be added during
363             exploration (boolean)
364             output
365             actn..... selected action (Energy Input, mode of facade/tvis)
366             (np.array(shape=1,2))'''
367         # normalization for NN input
368         state_history = self.norm_state_layer(state_history_t)
369         forecast_history = self.norm_forecast_layer(forecast_history_t)
370
371         history =
372         np.concatenate((state_history, tf.reshape(forecast_history, (state_history.shape[0]
373             ,self.forecast_dim))),axis=1)\n
374             .reshape(1,state_history.shape[0],self.forecast_dim+self.state_dim)
375
376         if self.parameter['agent']['setting']['noise_process'] == 'param_noise' and
377             noise:
378             action = self.perturbed_actor(history)
379         else:
380             action = self.actor(history)
381
382         if self.parameter['agent']['actions'] == 0:
383             Q = action[0][0]
384             Tvis = action[1][0]
385         elif self.parameter['agent']['actions'] == 1:
386             Q = action[0]
387             Tvis = 0.6
388         else:
389             Q = 0
390             Tvis = action[0]
391
392         if noise:
393             if self.parameter['agent']['setting']['noise_process'] == 'OU_noise':
394                 Q, Tvis = add_OU_noise(self.noise_Q, Q, self.noise_Tvis, Tvis)
395             elif self.parameter['agent']['setting']['noise_process'] == 'Gauss_noise':
396                 Q, Tvis = add_Gauss_noise(self.noise_Q, Q, self.noise_Tvis, Tvis)
397             # clip the values after adding the action noise to min and max values
398             if self.parameter['agent']['actions'] == 0:
399                 Q = tf.clip_by_value(action[0], -self.action_bound_range,
400                     self.action_bound_range)[0].numpy()
401                 Tvis = tf.clip_by_value(action[1], 0.01, 0.6)[0].numpy()
402                 actn = np.array([Q,Tvis])
403             elif self.parameter['agent']['actions'] == 1:
404                 Q = tf.clip_by_value(action[0], -self.action_bound_range,
405                     self.action_bound_range)[0].numpy()
406                 actn = np.array([Q])
407             elif self.parameter['agent']['actions'] == 2:
408                 Tvis = tf.clip_by_value(action[0], 0.01, 0.6)[0].numpy()
409                 actn = np.array([Tvis])
410             actn = actn.reshape(1,self.action_dim)
411         return actn
412
413     def train_critic_network(self,num_samples, states_batch, forecast_batch,
414         actions_batch, rewards_batch, demand_batch, next_states_batch, next_forecast_batch,
415         done_batch, indices, importance_weight):
416         ''' off-policy training of the critic network with stored experience
417             DDPG by (Lillicrap et.al.)
418
419             after training the priority of the sampled expirience based on the

```

```

410         td_error is changed (for Prioritized_experience_replay or sumTree)
411
412         the loss function for the actor is the negative value function
413         (critic network) as we want to maximize this value ''
414
415     if self.parameter['agent']['flags']['valuefunction']:
416         target_actions =
417             self.actor_target([next_states_batch[:,self.parameter['agent']['setting']['n_'
418                 step']-1], next_forecast_batch])
419         target_actions =
420             np.dstack(target_actions).reshape(num_samples,self.action_dim)
421         target_actions = self.norm_action_layer(target_actions)
422         target_critic_value =
423             self.critic_target([next_states_batch[:,self.parameter['agent']['setting']['n_'
424                 step']-1], next_forecast_batch, target_actions])
425
426         y_i = rewards_batch + demand_batch
427         y_i = np.reshape(y_i,(num_samples,1))
428         for i in range(num_samples):
429             y_i[i] = y_i[i] +
430                 self.parameter['agent']['NN']['discount_factor']**self.parameter['agent']
431                     ['setting']['n_step'] * target_critic_value[i]
432         with tf.GradientTape(watch_accessed_variables=False) as tape:
433             tape.watch(self.critic.trainable_variables)
434             critic_value = self.critic([states_batch[:,0], forecast_batch,
435                 actions_batch])
436             critic_loss =
437                 self.loss_function(y_i,critic_value,sample_weight=importance_weight)
438             critic_grad = tape.gradient(critic_loss,
439                 self.critic.trainable_variables)
440             self.critic_opt.apply_gradients(zip(critic_grad,
441                 self.critic.trainable_variables))
442
443     else:
444         y_i = rewards_batch + demand_batch
445         y_i = tf.reshape(y_i,(num_samples,1))
446         with tf.GradientTape(watch_accessed_variables=False) as tape:
447             tape.watch(self.critic.trainable_variables)
448             critic_value = self.critic([states_batch[:,0], forecast_batch,
449                 actions_batch])
450             critic_loss =
451                 self.loss_function(y_i,critic_value,sample_weight=importance_weight)
452             critic_grad = tape.gradient(critic_loss,
453                 self.critic.trainable_variables)
454             self.critic_opt.apply_gradients(zip(critic_grad,
455                 self.critic.trainable_variables))
456
457     if self.parameter['agent']['replay_buffer']['Buffer'] == 'PER':
458         priority = np.abs(y_i - self.critic([states_batch[:,0], forecast_batch,
459             actions_batch])[0])
460
461     elif self.parameter['agent']['replay_buffer']['Buffer'] == 'HVPER':
462         td_priority = tf.math.sigmoid(np.abs(y_i - self.critic([states_batch[:,0],
463             forecast_batch, actions_batch])[0])) * 2 -1
464         q_priority = tf.math.sigmoid(critic_value)
465         self.lamb = max(self.lamb * self.lamb_decrease,0)
466         priority = self.lamb * q_priority + (1-self.lamb) * td_priority
467
468     else:
469         priority = 0
470
471     if self.parameter['agent']['flags']['print'] :
472         print('critic_loss:\t', critic_loss.numpy())
473
474     return pd.Series(critic_loss.numpy()), priority
475
476 def train_actor_network(self,num_samples, states_batch, forecast_batch):

```

```

457     with tf.GradientTape(watch_accessed_variables=False) as tape:
458         tape.watch(self.actor.trainable_variables)
459         actions = self.actor([states_batch[:,0], forecast_batch])
460         actions = tf.reshape(tf.concat(actions, -1), (num_samples, self.action_dim))
461         actions = self.norm_action_layer(actions)
462         critic_value = self.critic([states_batch[:,0], forecast_batch, actions])
463         actor_loss = -tf.math.reduce_mean(critic_value)
464         actor_grad = tape.gradient(actor_loss, self.actor.trainable_variables)
465         self.actor_opt.apply_gradients(zip(actor_grad, self.actor.trainable_variables))
466
467     if self.parameter['agent']['flags']['print']:
468         print('actor_loss:\t', actor_loss.numpy())
469
470     self.update_target(self.critic_target, self.critic,
471     self.parameter['agent']['NN']['tow'])
472     self.update_target(self.actor_target, self.actor,
473     self.parameter['agent']['NN']['tow'])
474
475 def train_critic_lstm(self, num_samples, hist_t, actions_batch, rewards_batch,
476 hist_t_1, done_batch, indices, importance_weight):
477     ''' off-policy training of the critic network with stored experience
478     RDPG by (Hess et.al.)
479
480     after training the priority of the sampled expirience based on the
481     td-error is changed (for Prioritized_experience_replay or sumTree)
482
483     the loss function for the actor is the negative value function
484     (critic network) as we want to maximize this value '''
485
486     target_actions = self.actor_target([hist_t_1])
487     target_actions =
488     tf.reshape(tf.stack(target_actions, axis=self.action_dim), (num_samples, 1, self.action_dim))
489     target_actions = self.norm_action_layer(target_actions)
490     target_actions = tf.concat((actions_batch[:,1:,:], target_actions), axis=1)
491     target_critic_value = self.critic_target([hist_t_1, target_actions])
492
493     y_i = rewards_batch
494     for i in range(num_samples):
495         y_i[i] = y_i[i] + self.parameter['agent']['NN']['discount_factor'] *
496             target_critic_value[i] #* (1 - done_batch[i])
497
498     with tf.GradientTape(watch_accessed_variables=False) as tape:
499         tape.watch(self.critic.trainable_variables)
500         critic_value = self.critic([hist_t, actions_batch])
501         critic_loss =
502             self.loss_function(y_i, critic_value, sample_weight=importance_weight)
503         critic_grad = tape.gradient(critic_loss,
504         self.critic.trainable_variables)
505         self.critic_opt.apply_gradients(zip(critic_grad,
506         self.critic.trainable_variables))
507
508     if self.parameter['agent']['replay_buffer']['Buffer'] == 'PER':
509         priority = np.abs(y_i - self.critic([hist_t, actions_batch])[0])
510     elif self.parameter['agent']['replay_buffer']['Buffer'] == 'HPER':
511         td_priority = tf.math.sigmoid(np.abs(y_i - self.critic([hist_t,
512         actions_batch])[0])) * 2 - 1
513         q_priority = tf.math.sigmoid(critic_value)
514         self.lamb = max(self.lamb * self.lamb_decrease, 0)
515         priority = self.lamb * q_priority + (1-self.lamb) * td_priority
516     else:
517         priority = 0
518
519     if self.parameter['agent']['flags']['print']:

```

```

511         print('critic_loss:\t', critic_loss.numpy())
512
513     return pd.Series(critic_loss.numpy()), priority
514
515     def train_actor_lstm(self, num_samples, hist_t, actions_batch):
516         with tf.GradientTape(watch_accessed_variables=False) as tape:
517             tape.watch(self.actor.trainable_variables)
518             actions =
519                 tf.reshape(tf.stack(self.actor(hist_t), axis=self.action_dim), (num_samples, 1, self.action_dim))
520             actions = self.norm_action_layer(actions)
521             actions = tf.concat((actions_batch[:, 0:self.n_Step-1], actions), axis=1)
522             critic_value = self.critic([hist_t, actions])
523             actor_loss = -tf.math.reduce_mean(critic_value)
524             actor_grad = tape.gradient(actor_loss, self.actor.trainable_variables)
525             self.actor_opt.apply_gradients(zip(actor_grad, self.actor.trainable_variables))
526
527             if self.parameter['agent']['flags']['print']:
528                 print('actor_loss:\t', actor_loss.numpy())
529
530             self.update_target(self.critic_target, self.critic,
531                               self.parameter['agent']['NN']['tow'])
532             self.update_target(self.actor_target, self.actor,
533                               self.parameter['agent']['NN']['tow'])
534
535     def update_target(self, target, online, tow):
536         ''' soft update of target networks'''
537         init_weights = online.get_weights()
538         update_weights = target.get_weights()
539         weights = []
540         for i in tf.range(len(init_weights)):
541             weights.append(tow * init_weights[i] + (1 - tow) * update_weights[i])
542         target.set_weights(weights)
543
544     return target

```

Replay Buffer

```
1  from collections import deque
2  import numpy as np
3  import random
4
5  class Uniform:
6      def __init__(self, max_buffer_size, dflt_dtype, forecast_dim):
7          self.buffer = deque(maxlen=max_buffer_size)
8          self.dflt_dtype = dflt_dtype
9          self.forecast_dim = forecast_dim
10
11     def add_experience(self, state, forecast, action, reward, demand, next_state,
12                         next_forecast, done):
13         self.buffer.append([state, forecast, action, reward, demand, next_state,
14                            next_forecast, done])
15
16     def batch_update(self, indices, priorities):
17         pass
18
19     def sample_batch(self, batch_size):
20         num_samples = min(len(self.buffer), batch_size)
21         replay_buffer = np.array(random.sample(self.buffer, num_samples))
22         arr = np.array(replay_buffer)
23         states_batch = np.stack(arr[:, 0])
24         forecast_batch = np.stack(arr[:, 1])
25         actions_batch = np.stack(arr[:, 2])
26         rewards_batch = np.stack(arr[:, 3])
27         demand_batch = np.stack(arr[:, 4])
28         next_states_batch = np.stack(arr[:, 5])
29         next_forecast_batch = np.stack(arr[:, 6])
30         done_batch = np.vstack(arr[:, 7]).astype(bool)
31
32     return states_batch, forecast_batch, actions_batch, rewards_batch,
33           demand_batch, next_states_batch, next_forecast_batch, done_batch
34
35     class HVPER:
36         def __init__(self, max_buffer_size, dflt_dtype, forecast_dim):
37             self.max_buffer_size = max_buffer_size
38             self.buffer = deque(maxlen=max_buffer_size)
39             self.priorities = deque(maxlen=max_buffer_size)
40             self.usage = deque(maxlen=max_buffer_size)
41             self.indexes = deque(maxlen=max_buffer_size)
42             self.dflt_dtype = dflt_dtype
43             self.forecast_dim = forecast_dim
44             self.n_k = 5
45
46         def add_experience(self, state, forecast, action, reward, demand, next_state,
47                           next_forecast, done):
48             self.buffer.append([state, forecast, action, reward, demand, next_state,
49                               next_forecast, done])
50             self.priorities.append(1)
51             self.usage.append(1)
52             ln = len(self.buffer)
53             if ln < self.max_buffer_size : self.indexes.append(ln)
54
55         def batch_update(self, indices, priorities):
56             for indx, priority in zip(indices, priorities):
57                 self.usage[indx-1] += 1
58                 self.priorities[indx-1] = priority * (1*pow(0.95, self.usage[indx-1]))
59
59         def sample_batch(self, batch_size, beta):
60             num_samples = min(len(self.buffer), batch_size*self.n_k)
61             n_k = min(self.n_k, int(num_samples/batch_size))
62             if n_k == 0:
63                 n_k = 1
```

```

60     num_samples = min(len(self.buffer),batch_size)
61     nk_indices = random.sample(self.indexes,k = num_samples * n_k)
62     nk_priorities = np.array([self.priorities[indx-1] for indx in nk_indices])
63     indices = random.choices(nk_indices,weights=nk_priorities, k = num_samples)
64
65     importance_weight = np.array([(self.priorities[indx-1] * num_samples)**-beta
66     for indx in indices])
67     importance_weight = importance_weight / max(importance_weight)
68
69     replay_buffer = [self.buffer[indx-1] for indx in indices]
70     arr = np.array(replay_buffer)
71     states_batch = np.stack(arr[:, 0])
72     forecast_batch = np.stack(arr[:, 1])
73     actions_batch = np.stack(arr[:, 2])
74     rewards_batch = np.stack(arr[:, 3])
75     demand_batch = np.stack(arr[:, 4])
76     next_states_batch = np.stack(arr[:, 5])
77     next_forecast_batch = np.stack(arr[:, 6])
78     done_batch = np.vstack(arr[:, 7]).astype(bool)
79
80     return states_batch, forecast_batch, actions_batch, rewards_batch,
81     demand_batch, next_states_batch, next_forecast_batch, done_batch, indices,
82     importance_weight
83
84 class PER:
85     def __init__(self, max_buffer_size, dflt_dtype, forecast_dim):
86         self.max_buffer_size = max_buffer_size
87         self.buffer = deque(maxlen=max_buffer_size)
88         self.priorites = deque(maxlen=max_buffer_size)
89         self.indexes = deque(maxlen=max_buffer_size)
90         self.dflt_dtype = dflt_dtype
91         self.forecast_dim = forecast_dim
92
93     def add_experience(self, day):
94         self.buffer.append(day)
95         self.priorites.append(1)
96         ln = len(self.buffer)
97         if ln < self.max_buffer_size : self.indexes.append(ln)
98
99     def update_priorities(self,indices,priorites):
100        for indx,priority in zip(indices,priorites):
101            self.priorites[indx-1] = priority
102
103    def sample_batch(self, batch_size, beta):
104        num_samples = min(len(self.buffer),batch_size)
105        indices = random.sample(self.indexes,k = num_samples)
106        priorities = np.array([self.priorites[indx-1] for indx in indices])
107        indices = random.choices(indices,weights=priorities, k = num_samples)
108
109        importance_weight = np.array([(self.priorites[indx-1] * num_samples)**-beta
110        for indx in indices])
111        importance_weight = importance_weight / max(importance_weight)
112
113        replay_buffer = [self.buffer[indx-1] for indx in indices]
114        arr = np.array(replay_buffer)
115        states_batch = np.stack(arr[:, 0])
116        forecast_batch = np.stack(arr[:, 1])
117        actions_batch = np.stack(arr[:, 2])
118        rewards_batch = np.stack(arr[:, 3])
119        demand_batch = np.stack(arr[:, 4])
120        next_states_batch = np.stack(arr[:, 5])
121        next_forecast_batch = np.stack(arr[:, 6])
122        done_batch = np.vstack(arr[:, 7]).astype(bool)
123
124        return states_batch, forecast_batch, actions_batch, rewards_batch,
125        demand_batch, next_states_batch, next_forecast_batch, done_batch, indices,
126        importance_weight

```

Noise

```
1 import numpy as np
2 import tensorflow as tf
3 import tensorflow_probability as tfp
4
5 class OU_Noise(object):
6     '''OU noise process'''
7     def __init__(self, mu, sigma = 0.15, theta = 0.2, dt=1e-1, x0 = None):
8         self.theta = theta
9         self.mu = mu
10        self.dt = dt
11        self.sigma = sigma
12        self.x0 = x0
13        self.reset()
14
15    def __call__(self):
16        x = self.x_prev + self.theta*(self.mu-self.x_prev)*self.dt +
17            self.sigma*np.sqrt(self.dt)*np.random.normal(size=self.mu.shape)
18        self.x_prev = x
19        return x
20
21    def reset(self):
22        self.x_prev = self.x0 if self.x0 is not None else np.zeros_like(self.mu)
23
24    class Param_Noise(): #Gaussian Noise  $\theta_e = \theta + N(0, \sigma^2 I)$  from
25        https://arxiv.org/pdf/1706.01905.pdf
26        ''' Parameter noise as gaussian normal distribution
27            Input is the shape of the layer weights of the NN
28            Output is a gaussian distribution in the shape of the weights'''
29        def __init__(self, size, shape, action_dim):
30            self.shape = shape
31            self.scale = 0.6
32            self.size = size
33            self.action_dim = action_dim
34            self.alpha = 1.01
35
36        def calc_scale(self, distance):
37            if distance < self.scale:
38                self.scale = self.alpha * self.scale
39            else:
40                self.scale = 1/self.alpha * self.scale
41
42        def perturb_actor(self):
43            print('noise_scale ', self.scale)
44            paramnoise = np.random.normal(loc = 0, scale = self.scale, size = self.size)
45            paramnoise = paramnoise.reshape(self.shape)
46            return paramnoise
47
48    def Gauss_Noise(scale): #Gaussian Noise  $\theta_e = \theta + N(0, \sigma^2 I)$  from
49        https://arxiv.org/pdf/1706.01905.pdf
50        ''' action noise as gaussian normal distribution
51            Output is a gaussian distribution in the shape of the weights'''
52        gaussnoise = tfp.distributions.Normal(loc = 0, scale = scale)
53        return gaussnoise
54
55    def add_OU_noise(noise_Q, Q, noise_Tvis, Tvis):
56        ''' Input: Q ..... Energy Input selected by actor
57                  Tvis ... facade/tvis selected by actor
58                  noise .. current initiated noise
59                  Output: Q ..... Energy Input + noise
60                  Tvis ... facade/tvis + noise'''
61        noise_Q = noise_Q()
62        noise_Tvis = noise_Tvis()
63        Q = Q + noise_Q
64        Tvis = Tvis + noise_Tvis
```

```
62     return Q, Tvis
63
64 def add_Gauss_noise(noise_Q, Q, noise_Tvis, Tvis):
65     ''' Input:  Q ..... Energy Input selected by actor
66             Tvis ... facade/tvis selected by actor
67             noise .. current initiated noise
68         Output: Q ..... Energy Input + noise
69             Tvis ... facade/tvis + noise'''
70     noise_Q = tf.cast(noise_Q.sample(), dtype = 'float32')
71     noise_Tvis = tf.cast(noise_Tvis.sample(), dtype='float32')
72     Q = Q + noise_Q
73     Tvis = Tvis + noise_Tvis
74
75     return Q, Tvis
```