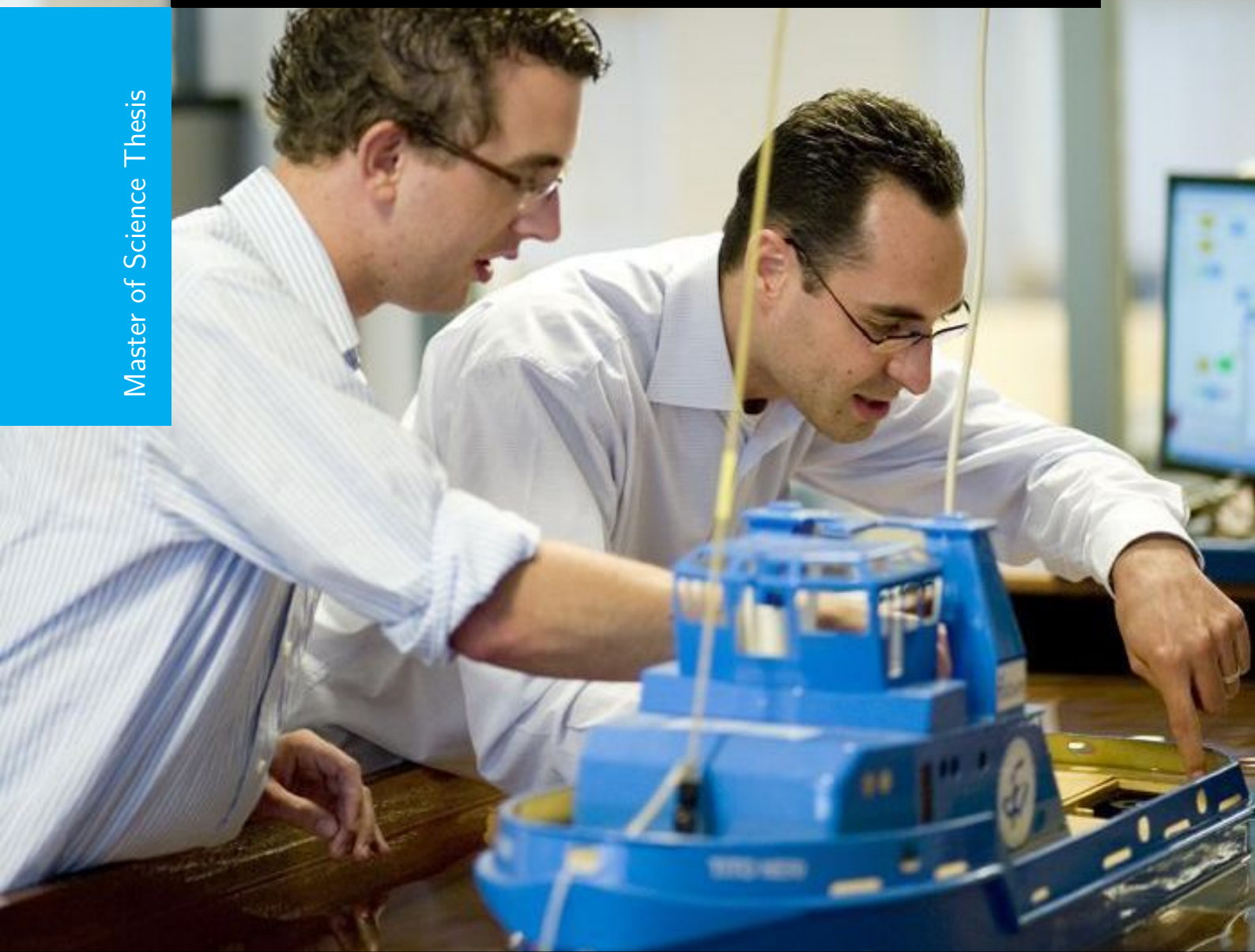


# Learning to Control Multi-Dimensional Autonomous Agents using Hebbian Learning

A Global Reward Approach

Ajdin Husić

Master of Science Thesis





# **Learning to Control Multi-Dimensional Autonomous Agents using Hebbian Learning**

## **A Global Reward Approach**

MASTER OF SCIENCE THESIS

For the degree of Master of Science in Systems and Control at Delft  
University of Technology

Ajdin Husić

December 18, 2018

Faculty of Mechanical, Maritime and Materials Engineering (3mE) · Delft University of  
Technology



Copyright © Delft Center for Systems and Control (DCSC)  
All rights reserved.



---

# Abstract

The novelty-raahn algorithm has been shown to effectively learn a desired behavior from raw inputs by connecting an autoencoder with a Hebbian network. Hebbian learning is compelling for its biological plausibility and simplicity. It changes the weight of a connection based only on the activations of neurons it connects, and can effectively reinforce good behaviors when combined with neuromodulation. These low-level synaptic weight changes make for a better merge of the three learning tasks of perception, prediction and action. However, the state-of-the-art algorithm requires the design of a highly detailed modulation scheme designed for a specific system, which is disconnected from the overall objective it optimizes. In this thesis, we will propose that similar learning behavior can be achieved, by making the autonomous agent react to longer-term rewards, and thus implicitly introducing prediction capabilities. In doing so, the required modulation scheme becomes connected to the global optimization objective.



---

# Table of Contents

<b>Acknowledgements</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Background</b>	<b>3</b>
2-1 The Hebbian Learning Process . . . . .	3
2-2 The state-of-the-art implementation . . . . .	5
2-2-1 The environment . . . . .	6
2-2-2 The autoencoder . . . . .	7
2-2-3 The Hebbian layer . . . . .	9
2-2-4 Experiments and results . . . . .	10
2-2-5 Current drawbacks . . . . .	11
2-2-6 General proposal . . . . .	12
<b>3 Proposal of the Thesis</b>	<b>13</b>
3-1 Defining the general global reward . . . . .	13
3-2 Relating the global reward to modulated Hebbian learning . . . . .	14
3-3 Chapter conclusion: proposal of the thesis . . . . .	15
3-3-1 application and use-case . . . . .	15
<b>4 Implementation</b>	<b>17</b>
4-1 Implementing the environment for the autonomous car . . . . .	17
4-2 Implementing the agent for the autonomous car . . . . .	19
4-3 Implementing the cart-pole environment . . . . .	19
4-4 Implementing the cart-pole agent . . . . .	21

<b>5 Experiments</b>	<b>23</b>
5-1 Experiments with the autonomous car environment . . . . .	23
5-1-1 Convergence test . . . . .	24
5-1-2 Performance test . . . . .	24
5-2 Experiments with the cart-pole . . . . .	26
<b>6 Discussion and Future Work</b>	<b>31</b>
<b>7 Conclusion</b>	<b>33</b>
<b>Bibliography</b>	<b>35</b>
<b>Glossary</b>	<b>37</b>
List of Acronyms . . . . .	37
List of Symbols . . . . .	37
<b>Index</b>	<b>39</b>
<b>Source Code</b>	<b>41</b>
-1 Python code for the autonomous car agent . . . . .	41
-1-1 The environment . . . . .	41
-1-2 Xmap.xml . . . . .	52
-1-3 Agent . . . . .	66
-2 Python Code for the cartpole implementation . . . . .	93
-2-1 Cartpole.py . . . . .	93



---

# List of Figures

2-1	Block diagram of a general reinforcement learning framework using Hebbian learning.	5
2-2	Graphical representation of the simulated environment as designed in Python: an autonomous agent traverses the lap shaped map. The autonomous car is required to make as many laps as possible in an allotted amount of time. The non-uniformly shaped map ensures that the agent does not simply repeat one behavior, but rather is required to generalize its policy to any input it experiences.	6
2-3	Depiction of the car's rangefinder sensors. A sensor's activation denotes a measure of wall intersection.	7
2-4	Artificial neural network as used in [1]. The network receives an activation $y_k$ of input neurons (top part), from which the autoencoder (first layer) computes the high-level encoded features $\hat{x}_k$ in the hidden neurons (middle part). Finally, these features are used by the Hebbian network (second layer) to compute the control value $u_k$ in the final output neuron (bottom part).	8
4-1	Example showcasing AABB collision in 2D. Collision between two rectangular objects only occurs if they overlap in all dimension axes. Adopted from [2].	18
5-1	Graph of the average global reward occurring at each episode.	25
5-2	Average number of crashes for each algorithm during training. From this figure we can conclude firstly that the global reward method learns the optimal behavior of not crashing along with the other algorithms shown. Secondly, global reward has indeed learned a similar performance, because its final performance is in the range of optimal performances as derived from the other algorithms when they do not crash (not shown in this figure).	26
5-3	Depiction of the amount of laps completed over a certain amount of experiments. The figure shows that the global reward method achieves similar performances as the state-of-the-art algorithms after sufficient training, since the satisfactory threshold of 9.64 laps is achieved on average.	27
5-4	Graphical representation of the cart-pole environment, with four known states. Objective is to stabilize the pole in its upper equilibrium (unstable point).	28



---

# List of Tables

4-1	Summary of the observation and action spaces of the cart-pole environment. The real-valued fully observable state is shown in the left part, and the discrete-valued one-dimensional action space is shown in the right part. . . . .	20
-----	---	----



---

# Acknowledgements

I would like to thank my supervisor Martijn Wisse for his assistance during my master's research. I would further like to thank Wouter Wolfslag for his useful ideas and feedback during my literature research.

Delft, University of Technology  
December 18, 2018

Ajdin Husić



---

# Chapter 1

---

## Introduction

Robotic tasks are becoming more complex ever since the progressive developments in computational resources permitted the additional freedom of building more autonomy into the robotic agents. But unlike in industry, the environments in which these agents are employed are not static, and so the complexity grows further in that these agents require not only better hardware, but intelligence to maximally utilize its hardware potential. Ideally, agent intelligence allows not only the generalization of its behavior to unperceived situations, but also the optimal adaptation to it. Since such environments can be very complex (many sensor inputs, many actuator outputs and varying dynamics), the main challenge lies in the design of *scalable learning* algorithms.

An agent is hypothesized to exhibit intelligent behavior when they effectively minimize their prediction errors [3, 4]. It is straightforward that in order to optimize its **predictions** with respect to its **sensation** of the world, which is caused by its **actions**, it needs to accurately *learn to*:

- represent the state (sensory perception),
- predict the state given an action (prediction),
- control to achieve the state (action).

Interestingly, within these three tasks, we recognize three major types of learning. Namely unsupervised, supervised and reinforcement learning respectively. Of course, implementing an algorithm that combines these task must be **meaningful**, and so we must consider a few assumptions:

- A.1** The agent features unknown (and uncertain) dynamics. Otherwise, there is no need for prediction, and the point of minimizing an error becomes futile.
- A.2** The agent must sufficiently explore the unknown environment, in order to effectively learn the three tasks. If an autonomous agent does not excite sufficiently, it will most likely get stuck at a local equilibrium.

**A.3** The algorithm must be applicable to a wide range of system classes, up to higher dimensional state and input spaces, from which follows that it must be not only computationally simple, but also scalable. Indeed, if we have high order complex and unknown system dynamics, learning must be efficient to be applied in either real-time learning, or episodic optimization tasks.

Hebbian learning is a learning type, compelling for its simplicity and biological plausibility. It straightforwardly changes the weight of a synapse based only on the activations of the neurons it connects and can effectively reinforce good behaviors and repel bad behaviors when combined with neuromodulation. Hebbian learning in general is scalable and computationally efficient, since the agent behavior is shaped through low-level neuromodulation as opposed to high-level value-function approximation as in other reinforcement learning types like Q-learning and SARSA. The biologically plausible and local learning rule, makes for a better merge of the three tasks.

The major obstacle to useful applications with Hebbian learning — which is that the performance of Hebbian plasticity is highly sensitive to the choice of inputs — has recently been solved by [1]. Their Novelty-raahn algorithm enables a neural network to **represent** accurately the features of the input domain by inserting a certain type of autoencoder that extrapolates higher level features. It is a real-time learning algorithm that makes the agent effectively learn a good **state representation**, combined with good **feedback control**. A *local* modulation term makes the agent highly adaptive to its environment. These attributes makes that their algorithm is currently the best in the eyes of the aforementioned assumptions and requirements.

The drawbacks of such method however, is that this algorithm still requires a considerable **design effort** to appropriately construct the modulation scheme, which demands detailed logic in order for the agent to learn correctly. The modulation term is also highly local for the current time tick in which the agent fares, and therefore the agent is never really learning to anticipate, i.e. the prediction part of the three learning tasks is missing.

**We will propose in this thesis**, that if we add a general global reward function from which the modulation scheme is derived, that it will enable the agent to regulate any type of control variable, and also regulate multiple control variables simultaneously. Moreover, with this novel method of deriving the modulation from global rewards, we will show that the agent can react to more distal rewards, and as such implicitly develop prediction capabilities. The utilization of a self-defined global reward function then gives more freedom to the engineers to decide how important each control aspect is. This generalized method poses the research question whether it can achieve similar results as the state-of-the-art algorithm.

Firstly, we will give an introduction to the Hebbian learning process in chapter 2 to familiarize the reader with its mathematical functionality. In chapter 3 we will present the proposal of how to use global rewards with Hebbian learning, and thereby yield a generalized version of novelty-Raahn. In chapter 5 we demonstrate that indeed the new algorithm can compete with the original novelty-Raahn algorithm while using a more complex set of dynamics, and prove its generality by applying the new algorithm on a cart-pole system, where we attempt to control an unstable equilibrium. Finally in chapter 7, we will conclude the thesis with our main findings in this research. Relevant code can further be found in appendix 7.



---

## Chapter 2

---

# Background

In this chapter we will explain briefly the idea of Hebbian learning, and how that idea translates into a mathematical formulation. We will then proceed to demonstrate the state-of-the-art literature and the use case for this algorithm as currently implemented. We will also show some experimental results done in literature and by myself, and finally argue about its drawbacks and propose how to improve this algorithm by eliminating these drawbacks.

### 2-1 The Hebbian Learning Process

Hebbian, or associative learning, is a theory about learning, hypothesized by psychologist Donald Hebb [5], that explains how biological neurons shape the learning process. The following statement is often quoted among neuroscientists to explain how Hebbian learning works:

“When an axon of Cell A is near enough to excite a Cell B and repeatedly and persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells, such that A’s efficiency, as one of the cells firing B, is increased.” [5]

The idea is that any two cells that are repeatedly active at the same time, tend to become associated, so that activity in one facilitates activity in the other. This is often informally summarized as neurons that fire together, wire together [6]. Haykin [7] generalized this statement and rephrased it as a two-part rule:

1. If two neurons on either side of a synapse (connection) are activated simultaneously (i.e., synchronously), then the strength of that synapse is selectively increased.
2. If two neurons on either side of a synapse are activated asynchronously, then that synapse is selectively weakened or eliminated.

This is called a Hebbian synapse. In [7], Haykin explained this concept of synaptic weakening, where he included the processes of *synaptic enhancement* and *synaptic depression*. In this way, a Hebbian synapse is recognized to produce synaptic strengthening (enhancement) by positively correlated activity, or synaptic weakening (depression) by negatively or uncorrelated activity. Opposite from the Hebbian synapse is the anti-Hebbian synapse which weakens the synapse with positively correlating activity, whereas it strengthens the synapse with negatively correlating activity.

Hence, these statements spark an inspiration for the development of biologically plausible learning algorithms in optimization and artificial intelligence.

Although Hebb did not provide a precise mathematical formulation of his postulate, a general form can be considered from the previous statements in a feedforward artificial neural network, when the weight between an input and output neuron is strengthened, when the input neuron causes a strong activation in the output neuron [8]. In that sense, the simplest form of Hebbian plasticity is described by **Hebb's hypothesis**, where the *plasticity* (change in weight) is given as

$$\Delta w_{ij} = \eta \cdot y_i x_j, \quad (2-1)$$

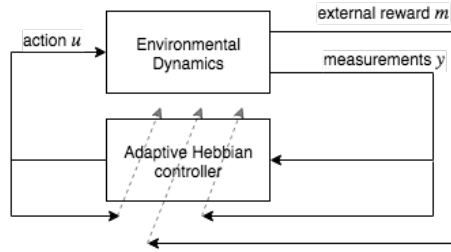
where  $\eta$  denotes a step size (or learning rate), and  $x_j$  denotes the activation of input neuron  $j$ , and  $y_i$  denotes the activation of output neuron  $i$ . This equation represents the standard mathematical form of a Hebbian synapse.

Of course this type of plasticity does not say anything about the neuronal dynamics themselves. If we consider for example that neurons can have activations in the range of real numbers, and their connections is a linear one, then Hebbian plasticity would suffer from exponential growth. That is, if presynaptic neuron fires strongly with the postsynaptic neuron, then their weight will grow strongly. This leads to an even higher value for the rate of weight change, driving the weight into saturation, after which selectivity is lost, i.e. no meaningful information will be stored in the synapse any longer. This is the sensitivity problem of Hebbian learning as mentioned before. I.e., if we feed a synapse too much of a same input, its weight will simply saturate, and this synaptic connection will hold only useless information in its memory.

To counteract this problem, a modulatory signal can be added to the learning rule, as is done by various methods in literature, see e.g. [9, 10] but also by [1] of which the algorithm is the basis for this thesis. This modulatory signal acts as a reinforcement-type feedback, and acts as a reward or punishment for an experienced input-output relationship, in order to develop the desired behavior. Hebbian plasticity then takes the following form:

$$\Delta w_{ij} = \eta \cdot m y_i x_j, \quad (2-2)$$

where  $m$  denotes the modulation in this scenario. Note that  $m$  can not only vary to change the amount of plasticity, but it can switch to the negative range, at which point anti-Hebbian (or unlearning) is turned on. Now, we can recognize that modulated Hebbian plasticity is in line with selectively reinforcing good behaviors as is the case with animals, where the modulation corresponds to dopamine in the brain. In we consider modulated Hebbian plasticity in a reinforcement learning framework, then the presynaptic activity is regarded as the agent's sensory perception  $y$  or state  $x$ , and the postsynaptic activity is equivalent to the agent's action  $u$ . This means that the Hebbian synapse is simply an adaptive feedback controller.



**Figure 2-1:** Block diagram of a general reinforcement learning framework using Hebbian learning.

For the purpose of consistency in our conventions, let us define the terminology and notations that goes unchanged for this entire thesis. The adaptive Hebbian feedback controller is a nonlinear combination of measurements  $y$  (perception) fed back as a signal  $u$  to its actuators (action). Generally, we have that

$$u_k = f(w_k y_k). \quad (2-3)$$

Here we will denote the subscript  $k$  as a discrete time-index. Then  $y_k \in \mathbb{R}^n$  denotes the  $n$ -dimensional sensor inputs,  $w_k \in \mathbb{R}^{m \times n}$  a matrix contain the plastic Hebbian weights and  $u_k \in \mathbb{R}^m$  denotes the  $m$ -dimensional control outputs. The function  $f(\cdot)$  is a nonlinear activation function. Hebbian adaptation is then generally of the form

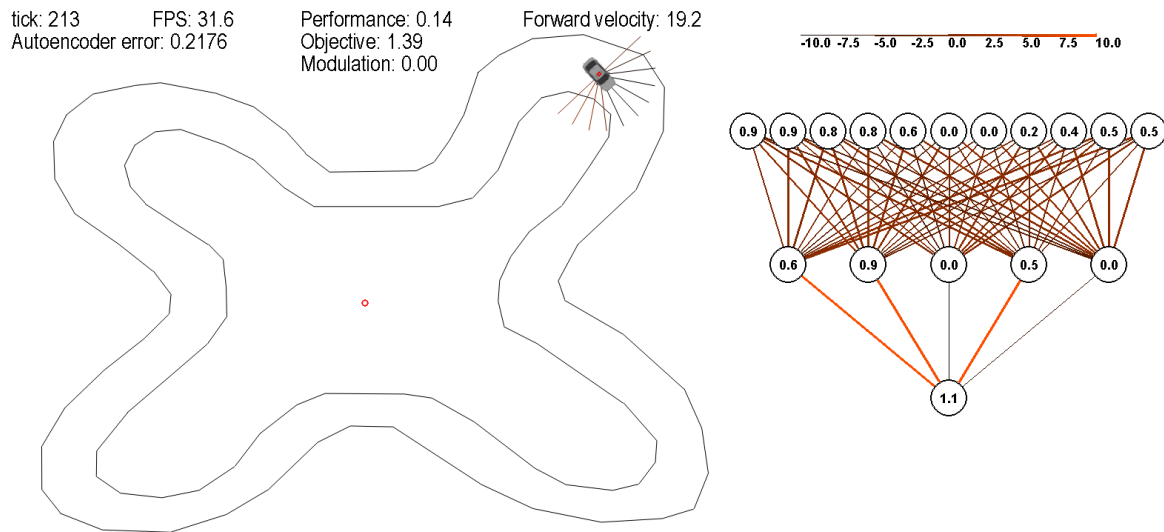
$$w_{k+1} = w_k + \eta \cdot m_k u_k y_k^T, \quad (2-4)$$

where  $m_k$  is the scalar modulation signal, which is the same for each Hebbian connection. This modulation value assumes a measure of 'correctness' that the input-output sample  $(y_k, u_k)$  turned out to cause based on later observations ( $k + 1$  in this thesis). A schematic overview of the closed loop structure is depicted in Figure 2-1.

The modulation scheme in this algorithm is the most important aspect for proper learning, as it directly changes the importance and sign of each training sample. But apart from the modulation term, good exploitation of the input space is still important. If, for example, the network gets excited with only *bad* experiences where the modulation is negative, the agent will proceed to unlearn everything it sees, and no viable progress may still occur. So to avoid feeding the agent with redundant and repetitive experiences, [1] includes an autoencoder, where the input space is represented by the autoencoder's higher-level features. In this way, the input as fed to the Hebbian synapses are more carefully selected, so that the entire input space of a given environment is better exploited.

## 2-2 The state-of-the-art implementation

In [1], the Novelty-Raahn algorithm is proposed, to solve the main obstacle of plain Hebbian learning. The novelty-Raahn algorithm enables an artificial neural network to **represent** accurately the features of the input domain (sensory perceptions) by extrapolating higher level features, based on the collection of novel agent experiences. These features of the autoencoder are fully connected with a second layer, which is considered the Hebbian **control** component of the agent. In the following subsections we will inform about the simulation environment in which the algorithm has been tested, and the functioning of the algorithm itself. Finally,



**Figure 2-2:** Graphical representation of the simulated environment as designed in Python: an autonomous agent traverses the lap shaped map. The autonomous car is required to make as many laps as possible in an allotted amount of time. The non-uniformly shaped map ensures that the agent does not simply repeat one behavior, but rather is required to generalize its policy to any input it experiences.

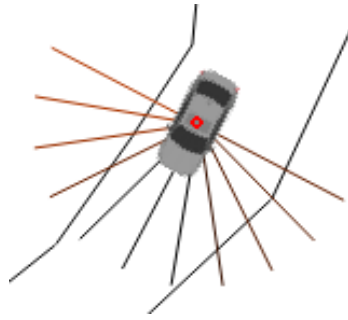
we will summarize some drawbacks still present, and propose what aspects need to change for improvement.

### 2-2-1 The environment

The *system* comprises both the *environment* and the *agent*. The environment entails the dynamics which govern the world's state, and the agent entails that which acts upon the world, to change its state, based on its perception of the world. The system in this description corresponds to a closed loop in feedback control.

Specifically, the simulated environment in [1] concerns an autonomous car agent, driving around through a static map in laps, see Figure 2-2. The entire track from the leftmost wall to the rightmost wall has a width of approximately 4120 units, whereas the lowermost wall up to the uppermost wall has a height of approximately 3074 units. The car has no dimensions, and is considered a point mass, which always drives at a constant speed of approximately 15 units<sup>1</sup> per time tick, unless it gets stuck into a wall in which case it stops moving forward. The agent can further issue a control signal between 0 and 1 (which may convolve with additional actuator noise), that determines the change in steering angle, or the car's global orientation  $\theta$ . Note that if the car crashes into a wall, it may still escape after some time, because it is allowed to correct its steering. A control output of 1 changes the global orientation  $\theta$  with  $-2^\circ$  per time tick (steers to the right) and a control output of 0 changes this orientation with  $+2^\circ$  per tick (steers to the left). The agent further has access to a set of eleven rangefinder sensors. These rangefinder sensors are line segments extending from the center of the car

<sup>1</sup>The true forward speed depends on its global orientation, and is equivalent to  $\sqrt{15^2 \cos^2 \theta + 12^2 \sin^2 \theta}$ .



**Figure 2-3:** Depiction of the car’s rangefinder sensors. A sensor’s activation denotes a measure of wall intersection.

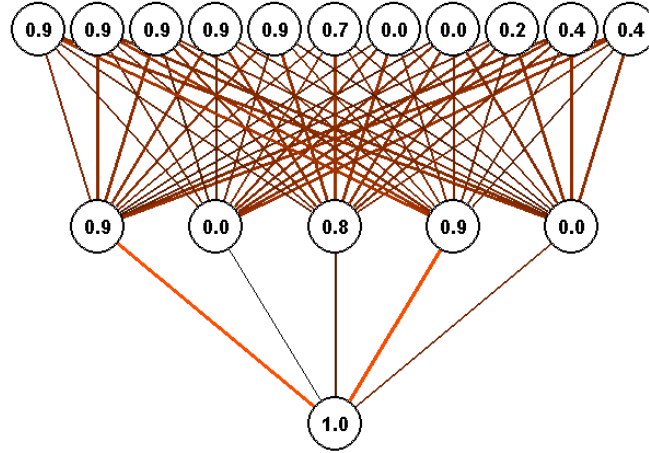
outward with each a length of 350 units, which sense the presence of walls (Figure 2-3). The first and last rangefinder sensors are separated with  $180^\circ$  from each other, and are oriented perpendicular to the car’s facing direction. The other sensors have a spacing of  $18^\circ$  between them. Each rangefinder receives an activation between 0 and 1 at each time tick. A sensor value of 0 indicates no intersection with any wall, and an activation of 1 indicates a full intersection with a wall. Intermediate activations are linear interpolations of where a wall is sensed.

The agent uses the collaborate autoencoder and Hebbian layer as artificial neural network components, in order to learn the optimal control policy in real-time (figure 2-4). It first senses an experience  $y_k \in [0, 1]^{11}$  at discrete time index  $k$ , from which it computes the *perceived* state  $\hat{x}_k = \sigma(W_k y_k)$ , with  $W \in \mathbb{R}^{5 \times 11}$  a matrix containing the connections between the input and the perceived state  $\hat{x}_k$ . This perceived state, or high-level features, are in turn used to compute final control output  $u_k = \sigma(w_k \hat{x}_k)$  that determines the change in turn angle of the car, where  $w_k \in \mathbb{R}^{1 \times 5}$  denotes the Hebbian connections. The nonlinear activation function  $\sigma(\cdot)$  is the logistic sigmoid, and can be viewed as a mechanism that determines the activation of neurons in a more biologically plausible way (close to 0 is a non-firing or inactive neuron, whereas close to 1 is a firing or active neuron).

### 2-2-2 The autoencoder

Since the performance of Hebbian plasticity is given to be very sensitive to how each sample is treated during training, it is naturally important to feed the correct experiences representative of the input domain to the Hebbian plasticity. The functionality of the autoencoder in this sense is twofold. It learns to describe the underlying high-level features of the input space into its perceptive memory, which guarantees that these features are fruitful, given that it can produce a good reconstruction of the original input space. At the same time, this *encoded* representation makes the agent more exploratory, since it considers a more meaningful sense of inputs for its decision making. This part is important, because it solves the main obstacle of Hebbian plasticity, which is its high sensitivity to the choice of inputs.

In order to exploit this input space in an autonomous setting, the novelty-Raahn in [1] uses a *novelty buffer* to select only the most novel experiences for training. This novelty buffer prevents overfilling the perceptive memory with only the same kind of experiences, and hence prevents the agent from settling rather quickly in a bad state. For example, if the agent



**Figure 2-4:** Artificial neural network as used in [1]. The network receives an activation  $y_k$  of input neurons (top part), from which the autoencoder (first layer) computes the high-level encoded features  $\hat{x}_k$  in the hidden neurons (middle part). Finally, these features are used by the Hebbian network (second layer) to compute the control value  $u_k$  in the final output neuron (bottom part).

initially gets stuck into a wall, it may erroneously learn as if its state at that moment is representative of the entire input domain, and consequently remain in that state (it has reached a bad local optimum from which it cannot escape).

The novelty buffer contains at most 500 past experiences, where each experience is assigned a *novelty score*. To determine the novelty score of an arbitrary experience  $i$ , the Euclidean distance  $d_i$  is first computed between it and all the other experiences in the buffer,

$$\forall i \in \mathcal{N} : \quad d_i = \|y_i - y_l\|_2 \quad \forall l \in \mathcal{N} / i, \quad (2-5)$$

with  $\mathcal{N}$  the set of all experiences in the novelty buffer,  $y_i$  the activation vector corresponding to experience  $i$ . The novelty score  $n_i$  of experience  $i$  is then computed as the 20 smallest such distances. If the currently sensed experience  $y_k$  has a novelty score greater than that of the least novel experience in the novelty buffer  $\mathcal{N}$ , it replaces that experience in the buffer.

At each training tick, the autoencoder selects 20 random experiences from the novelty buffer  $\mathcal{N}$  to train on. When training on experience  $y_i$ , the forward activation  $\hat{x}_i$  is first computed.

$$\hat{x}_i = \sigma(W y_i), \quad (2-6)$$

where  $\sigma(\cdot)$  is the logistic sigmoid activation function, and  $W$  a matrix containing the autoencoder weights. This forward activation is used to compute the backward activation (the *reconstruction* of the experience  $y_i$ ).

$$\hat{y}_i = \sigma(W^T \hat{x}_i), \quad (2-7)$$

where the reconstruction error is computed as

$$e_i = y_i - \hat{y}_i. \quad (2-8)$$

Then, the *delta* of the inputs can be computed from this error, which will be used to compute backpropagated error.

$$\delta_y = e_i \odot \sigma'(\hat{y}_i), \quad (2-9)$$

where the  $\odot$  is the component-wise multiplication, and  $\sigma'(\cdot)$  is the logistic derivative. Now, the delta for the hidden features can be computed as

$$\delta_{\hat{x}} = (W\delta_y) \odot \sigma'(\hat{x}_i). \quad (2-10)$$

Now, using both the original error deltas  $\delta_y$  and the backpropagated error deltas  $\delta_{\hat{x}}$ , we can compute the update of the tied autoencoder weights.

$$\Delta w = \alpha(\hat{x}_i\delta_y^T + \delta_{\hat{x}}y_i^T), \quad (2-11)$$

where  $\alpha$  is a learning step size, set to 0.1, and all weights in  $W$  are clipped to the interval  $[-w_{\max}, w_{\max}]$ , with  $w_{\max}$  set to 10, in order to prevent weight saturation.

### 2-2-3 The Hebbian layer

After having inferred the features of the environment through the autoencoder (perception), the agent decides which actions to take (control). The Hebbian layer in this sense can be viewed as an adaptive controller, as it adapts its plastic connections in real-time.

This component uses a modulated form of learning in order to achieve a desired behavior (control policy). This entails that learning is done through the modulated Hebbian learning rule as described in section 2-1. The presynaptic activity, in the case of the autoencoder, is equal to  $\hat{x}_k$  at time tick  $k$ . And the postsynaptic activity is then just the control output  $u_k$ , which is first normalized through a linear transformation  $u_{n,k} = 2u_k - 1$ , in order to be in the range  $[-1, 1]$ . This removes any bias while learning, such that the synapses can learn on their own if the modulation were removed.

$$\Delta w_k = \eta m_k u_{n,k} \hat{x}_k^T + \xi, \quad \xi \sim \mathcal{U}(-0.1, 0, 1)^{1 \times 5}, \quad (2-12)$$

where  $\eta = 1.0$  is the learning step size,  $\xi$  is a vector of random samples drawn from a uniform distribution between -0.1 and 0.1 and  $m_k$  is the modulation term determining when and at what rate to use learning, unlearning or no learning. The modulation scheme in [1] always compares the the next state  $k + 1$  with the current state  $k$ , in order to determine how *good* the corresponding training sample  $(y_k, u_k)$  turned out to be, and influences all Hebbian connections using the same modulation.

#### The modulation scheme

The modulation scheme in [1] is defined as whether the car has turned toward a wall ( $m_k < 0$ , unlearning), turned away from a wall ( $m_k > 0$ , learning) or detecting no walls ( $m = 0$ , no learning). In order to measure the modulation, the authors have first defined a *modulation feeler*, an invisible line extending from the center of the car straight outward in the direction it is facing, i.e. with same global orientation  $\theta$ , and with a maximum range of 400 units. This modulation feeler first senses the presence of walls, just like the rangefinder sensors, but

rather than measuring the distance, it computes the angle  $\beta$  it makes with a certain wall. If the modulation feeler detects a different wall at the next time tick  $k + 1$  than it detected at the current tick  $k$ , then it computes the angle  $\beta$  with respect to the wall it detected at time  $k$ . If no wall was detected at time  $k$  then the corresponding modulation  $m_k = 0$ .

A modulation value  $m_k$  is thus only computed, when a particular wall has been detected at time  $k$ . The measure  $\Delta$  of how much it turned away from this wall is thus computed at time tick  $k + 1$  as

$$\Delta = \begin{cases} \beta_{k+1} - \beta_k & , \text{ if } \beta_k > 90^\circ \\ \beta_k - \beta_{k+1} & , \text{ if } \beta_k \leq 90^\circ \end{cases} \quad (2-13)$$

Finally, the modulation  $m_k$  is computed as the amount in degrees it turned away from the wall normalized with the maximum turning speed per tick, such that the modulation yields a value in the domain of  $[-1, 1]$ .

$$m_k = \Delta / \Theta_{\max}. \quad (2-14)$$

Here,  $\Theta_{\max}$  denotes the maximum turning speed that can occur between two given time ticks, i.e.  $\max(\theta_{k+1} - \theta_k)$ .

Finally, the general algorithm is summarized in pseudocode 1.

---

**Algorithm 1** the novelty-RAAHN algorithm. This methods assumes:

1. multi-dimensional input space
  2. modulation scheme needs be well-defined
  3. unknown system dynamics
- 

```

1: procedure INITIALIZATION
2:   Initialize Autoencoder weights  $W_0 \leftarrow \mathcal{N}(0, 1) \in \mathbb{R}^{m \times n}$ 
3:   Initialize Hebbian weights  $w_0 \leftarrow \mathcal{N}(0, 1) \in \mathbb{R}^m$ 
4:   Initialize novelty-buffer size  $n_{novelty}$  arbitrarily
5:   Initialize number of training samples per tick  $n_{samples}$ 
6: end procedure
7: procedure TRAINING(for every time tick  $k$  do)
8:   take action  $u_k$  using policy derived from  $W_k, w_k$ 
9:   observe measurements  $y_{k+1}$  and reward  $m_{k+1}$ 
10:  update novelty buffer and weights  $W_{k+1}, w_{k+1}$ 
11: end procedure
12: Output: learned agent weights  $W, w$ 

```

---

## 2-2-4 Experiments and results

In the experiments in [1], the simulation is run 200 times for 10,000 ticks each time where an agent with novelty-Raahn is compared against an agent with only a Hebbian network (i.e. without autoencoder). Completing a circle around the center point as shown with the red dot in Figure 2-2, denotes completing one lap.

The results suggest that novelty-Raahn completes 8.8 laps on average, whereas Hebbian alone completes 9.2 laps. Novelty-Raahn performs thus slightly below the lone Hebbian controller,



while requiring extra time in the beginning to acquire a novel set of experiences. Also, considering that novelty-Raahn learns both features and the control policy at the same time, it is likely close to performing as well as possible for such a method.

The authors have furthermore shown that novelty-Raahn performs better when degrading the quality of the sensors to a certain extent. They show that if the length of the rangefinders increases, which makes distinguishing different situations more difficult, novelty-Raahn experiences fewer failures in completing laps than Hebbian, because the Hebbian agent is forced to learn from these degraded inputs, whereas novelty-Raahn learns a new representation.

Additionally, I have performed an experiment with Hebbian and novelty-Raahn myself. The claim is that when training a larger network, the agent needs a significantly longer time to learn, as it learns both perception and action at the same time. This reasoning prompted me to do a different experiment, in which both type of algorithms are run first in a training phase for 12,000 time ticks, with the idea that this is enough time for the agents to reach their near-optimal potential. After this phase, the state of the agents is reset, and they are released to follow their learned controller for 10,000 time ticks without training (evaluation phase). The performance in this scenario is only counted in how many laps each agent makes in this evaluation phase. This process is repeated for both algorithms 10 times, which results in a low enough variance to conclude from it their true performances. Namely, the Hebbian agent performed in the range [9.82, 9.92] with an average performance of 9.86 completed laps. Novelty-Raahn on the other hand performed in the range [9.81-10.04] with an average of 9.91 completed laps. This suggest that after sufficient training, novelty-Raahn even overtakes pure Hebbian slightly in its performance.

### 2-2-5 Current drawbacks

The main issue of the novelty-Raahn algorithm is the **modulation scheme** as currently defined. Firstly, it is only defined for the environment as described in [1], i.e. as in subsection 2-2-1. A closer look at this system, which is relatively simple and static in nature (the entire network only generates the change in steering angle), suggests that achieving the desired behavior requires complex modulation logic. For this system, this means that designing this logic is almost as easy (if not harder) as guessing the weights of the 11-by-1 Hebbian component, since these weights approach rather intuitive steady-state values: using only the Hebbian component without autoencoder results in weights that simply compensate for sensing walls in a particular direction. If such 'simple' agent behavior requires such detailed modulation logic, then we can ask ourselves whether even the slightest changes in the environment will not render the algorithm useless. The algorithm as is, thus **lacks the desired freedom** for the engineer to determine the desired tasks that the agent needs to achieve.

Moreover, the modulation scheme as defined is a **time-local** reward, and therefore the agent develops a short term memory. Rather than just anticipating one step ahead at most, complex autonomous agents may require long-term connections, that enable a more optimal control policy.

Most importantly, the odd observation of this algorithm is that the modulation logic is not explicitly **related** to the **global objective** which is sought to be optimized by the algorithm. In this way, the algorithm suffers from requiring hard-coded logic, which is **not directly** available. Whereas the **directly** available global objective of the system is not used, i.e. the

engineers know what they want to make the agent ultimately achieve, but need an intermediate step to reach this achievement.

### 2-2-6 General proposal

In order to eradicate these difficulties, and thus make the algorithm more useful overall, I will propose in this thesis that the main contribution lies in changing the way in which this modulation scheme is defined. More specifically, **I propose that the modulation scheme should be defined by relating the global objective function explicitly to how each Hebbian input-output training sample should be modulated.**

Assuming that we can successfully achieve the goal in such proposal, we can further recognize some additional advantages: Namely, if we find such a relationship, we can use this global objective to train **any** type of output we wish to control. For example, if we now are to add a control output that sets the agent's acceleration, we can use this more general relationship between how the global objective function should be treated with respect to the modulation of each individual sample, rather than requiring some different detailed modulation logic yet again. Since the change in velocity physically represents something different than the change in steering angle, and suddenly we need not contemplate how this new control variable precisely affects the current modulation behavior.

Furthermore, achieving such a relationship allows for the agent to develop a **long term memory** as well. The idea is that if the algorithm allows for global objective functions, then such an objective function can be one which incorporates the average performance of the agent over multiple time steps in the future. Coincidentally, this means that the agent will implicitly assume **prediction capabilities**. Therefore, we are practically also merging the existing tasks of perception and action, together with prediction, as we have initially considered to be required for prediction error minimization.

# Proposal of the Thesis

The current proposal as put forth in the past section, is to find a relationship between the global objective function given a certain task, and the modulation applied to each training sample in the Hebbian network. In order to find the best conclusion, we must first recall the fundamental function of Hebbian connections, and how modulation affects its plasticity. In this chapter, we will consider the goal of a global reward function and its effect on the behavior of an agent. Subsequently, we will relate one with another based on how they affect the system to achieve its goal, and finally we will conclude the chapter with the obtained answer materialized from the given proposal.

### 3-1 Defining the general global reward

Since we are working in a reinforcement learning setting, we will define the global objective function, which is to be maximized, the *global reward* function  $R^1$ . Specifically, the *particular* reward function, or otherwise the immediate reward  $r_k$ , obtained at discrete time tick  $k$  from the environment, is a measure of how 'good' the current *time-local* input-output or measurement-action sample  $(y_k, u_k)$  is. The global *long-term* reward function  $R_k^N$ , which anticipates a horizon of  $N$  discrete time steps ahead, is simply the average of each particular reward over this horizon.

$$R_k^N = \frac{1}{N} \sum_{i=k}^{k+N-1} r_i. \quad (3-1)$$

The general goal of reinforcement learning is to optimize, specifically maximize, this function  $R_k^N$ . The parameters on which  $R_k^N$  (indirectly) depends, i.e. the neural network weights  $w_k$ , are adjusted in such a manner, that  $R_k^N$  increases overall.

---

<sup>1</sup>Rather than a *fitness* function, which is directly related to the parameters (weights) on which said function depends, we call it the *reward* function because it is directly related to the effect of a state-action pair.

### 3-2 Relating the global reward to modulated Hebbian learning

The hypothesis of Hebbian learning, is that it learns the behavior the agent experiences. That is, the presynaptic activity (sensor input)  $y$  results in postsynaptic activity (control output)  $u$  through the synaptic connection  $w$ , and updates  $w$  by Hebbian learning  $\Delta w = uy^T$ , in such a manner that input  $y$  will more often produce output  $u$ . Using this hypothesis, we find that we can modulate said synaptic plasticity by  $\Delta w = muy^T$ , based on the desired result: If we want that  $y$  produces result  $u$ , then modulate the training sample positively  $m > 0$  such that behavior  $u$ , given  $y$  gets learned. If we want that  $y$  produces result  $u$  no longer, then modulate the training sample negatively  $m < 0$ , such that behavior  $u$  given  $y$  gets unlearned. If we are satisfied with current behavior, do not learn at all, and modulate zero  $m = 0$ , such that there will be no change in behavior.

This is a pretty powerful hypothesis, in that it can effectively reinforce desired behavior. If we now consider the global reward function  $R_k^N$  based on a given sample<sup>2</sup>  $(y_k, u_k)$  then we identify three possible scenarios at each time step. That is, with respect to the previous time step  $k - 1$ , the reward  $R_k^N$  has

- increased (positive modulation),
- decreased (negative modulation),
- stayed the same (zero modulation).

From this we can infer that modulation can be related as

$$m_k^N = g(R_{k+1}^N - R_k^N), \quad (3-2)$$

with  $g(\cdot)$  a nonlinear weight regularization function. Since  $R_k^N$  is defined as the future average of particular rewards  $r$  (equation 3-1), we can make the simplification

$$\begin{aligned} m_k^N &= g\left(\frac{1}{N}\left(\sum_{i=k+1}^{k+N} r_i - \sum_{i=k}^{k+N-1} r_i\right)\right), \\ m_k^N &= g\left(\frac{1}{N}(r_{k+N} - r_k)\right). \end{aligned} \quad (3-3)$$

Note that, since we use an arbitrary time horizon  $N$ , we have to wait  $N$  discrete time steps more with respect to the inference at  $k$ , before we update the weight corresponding to that time sample, i.e.

$$\begin{aligned} w_{k+1} &= w_k + m_k y_k u_k^T \\ w_{k+1} &= w_k + g\left(\frac{1}{N}[r_{k+N} - r_k]\right) y_k u_k^T. \end{aligned} \quad (3-4)$$

The idea is that with this generalized form of modulation, the global reward function subject to the environmental dynamics will be maximized, assuming that the global reward is strictly increasing when a desired change in behavior occurs, and decreasing when an undesired change

---

<sup>2</sup>The input here is denoted as  $y_k$ , but may be an inferred state  $\hat{x}_k$  depending on any intermediate filters. The input in this sense is just what gets fed into the Hebbian network.

in behavior occurs between each time step. This assumption needs to be satisfied to learn properly a desired control policy, because of the sensitivity of Hebbian plasticity. But it does not imply that defining a global reward with this algorithm is more prone to failure than any other algorithm that incorporates objective functions. Rather, I would say that any other control method that utilizes explicit objective functions needs to be properly defined too.

### 3-3 Chapter conclusion: proposal of the thesis

Based upon the previous chapters we have arrived at the following thesis proposal:

#### Proposal 3-3.1: Thesis proposal

Given a *long term global* reward  $R_k^N = \sum_{i=k}^{k+N-1} r_i$ , with  $r_k$  the *particular* reward, the modulation scheme can appropriately be defined as a function of increment of  $R_{k+1}^N$  in the next time step  $k + 1$  with respect to the current reward  $R_k^N$ . Or equivalently, the modulation  $m_k^N$  that is assigned to input-output sample  $(y_k, u_k)$  is equal to  $g(\frac{1}{N}(r_{k+N} - r_k))$ .

Proposal 3-3.1 goes hand in hand with the following assumption:

- A.4** The global reward  $R_k^N$  increases if the environmental state changes desirably and decreases if the environmental state changes undesirably between time ticks  $k$  and  $k + 1$ .

In the next section we will set up the exact modulation scheme as used for the the experimental environment in 2-2-1.

#### 3-3-1 application and use-case

If we consider the environment as in 2-2-1, we remark that the true performance of the autonomous car is that in an allotted amount of time the car must create as many laps as possible around the center point. However, we cannot directly use this performance as our global reward  $R$ , since the map is shaped in such a way that at many places the agent must first decrease the performance, in order to increase the overall performance later on, like driving around a turn that goes a bit outward. So we must first make an adjustment in choosing our appropriate global reward. Given that the first experiment concerns only to correctly learn steering, we can initially use a particular reward  $r_k$  not dissimilar to [1]. If we define  $d_k = 1 - y_k$  as distances<sup>3</sup> to objects, then we can say that desirable behavior occurs when the car's front-pointing sensor (6th out of 11 sensors), detects an increase in distance. That is, we define the particular reward  $r_k$  as

$$r_k = 1 - y_k(6), \quad (3-5)$$

<sup>3</sup>The 1 in this definition is not strictly necessary, since the algorithm cares about increments of rewards only, but is more intuitive as the input domain now physically represents normalized distances.

where  $y_k(6)$  denotes the sixth element out of the 11 rangefinder sensor activations in  $y_k$ . A good regularization function  $g$  as presented in equation 3-3 is the tangent hyperbolic function

$$g(x) = \tanh(\gamma x), \quad (3-6)$$

where  $\gamma$  is just a positive gain constant. The reason for regularization is that if rewards change too radically, the weights do not get too large updates, as the tangent hyperbolic function maps to a value between -1 and 1 (just like the original modulation scheme). Moreover, the tangent hyperbolic function is nearly linear when its argument is close to zero. The gain function  $\gamma$  on the other hand is used to stretch out the weight updates if these are too small. A good value of  $\gamma$  generally is when the modulation becomes approximately 1 with the maximum possible improvement, and -1 with the maximum possible deterioration in behavior.

---

## Chapter 4

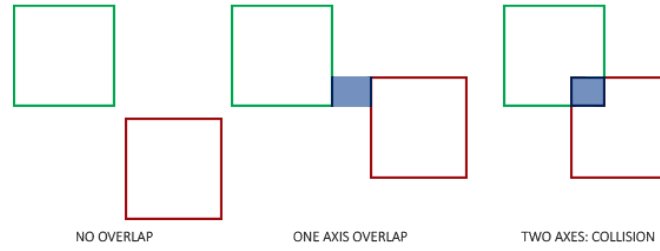
---

# Implementation

Before showing the final results in the following parts of this thesis, we will first cover the details of the entire implementation as necessary to generate the subsequent experiments. The original work of [1] was first recreated in Python. This was done mainly because of my own higher competence in Python rather than the original language, which was written in C#. This allowed me to very easily make any additions to the program without further investigation. Also, rewriting the code was crucial in understanding the exact mechanics of the entire simulation, which in turn allowed for further improvements, and helped me achieve the new additions to the existing algorithm.

### 4-1 Implementing the environment for the autonomous car

The first step of implementation was to create exactly the discrete environment as developed by [1] in Python3, in order to enable stepping through it by giving an control action  $u$  at each particular state. This environment consists of the classes `Car.py`, `Wall.py`, `rangefinder.py` and `configs.py`. See appendix -1-1 for the source code. Both the `Wall` and `rangefinder` class inherit from a lower-level `LineSegment` class (`lines.py`). Any line segment object (instance from `LineSegment`) has the ability to firstly test if it intersects any other line segments present using the AABB (Axis Aligned Bounding Box) method. Secondly, if two such line objects do intersect, then intersections are computed using simple linear line intersection math. An AABB collision is a collision detection method in computer graphics, in which two rectangular shapes collide if they overlap in all dimension axes (axis-aligned). See Figure 4-1 for an example. The same logic can be applied for line segments, in which the upper and lower bounds of said segment is simply seen as the rectangle boundaries in AABB. Of course when using line segments, their AABB collision does not guarantee an actual intersection, so intersection math is still necessary not only to check whether intersections are present between line segments, but also at what point they occur. Skipping the AABB detection part is theoretically possible but not optimal in simulation, because of the considerable difference in computational load, especially when testing multiple line segments against many other line segments.



**Figure 4-1:** Example showcasing AABB collision in 2D. Collision between two rectangular objects only occurs if they overlap in all dimension axes. Adopted from [2].

Furthermore, a line segment object is able to compute the exact intersection coordinates between itself and any other arbitrary line segment. This concept is extended in the rangefinder class, where its objects have additional attributes like a default length, a current activation value  $y_i$ , and a car instance to which it is attached. Each instance of rangefinder updates its positioning dynamically, along with its host car, of which its center point corresponds to the rangefinder's starting position. The end position is then related to a fixed relative angle between the car and itself, and its default length. A rangefinder line segment  $i$  has the additional ability to sense the value of its activation  $y_i$  by first computing the exact intersection point, and simply calculate the distance from its end position to this intersection, divided by its default length.

Whenever the simulation is initialized, first the car and all the map walls are spawned. This is done from the file `configs.py`. The map configuration is stored in an xml file `XMap.xml`, which contains both the car starting position (1561.93, 505.79), and all the wall coordinates (an (X,Y) coordinate and a relative distance in x,y-direction). Now the simulation can step through the environment by receiving an control output  $u$  in the value  $[-0.1, 1.1]$ . The car object is programmed initially to move with constant speed, but only change its global orientation  $\theta$  by adding  $4u - 2$  degrees to it, so that the increment yields a value in the range  $[-2.2, 2.2]$  degrees. The velocity vector is computed as

$$v_k = \left( s_x \cos \theta, \frac{4}{3} s_x \sin \theta \right), \quad (4-1)$$

where  $s_x = 15.0$  is the constant speed component in the initial experiments. Later on, I added the dynamic speed and acceleration control variable of the car, such that an environmental step accepts the value  $u = (u_\theta, u_a)$ . Here the control values change the global orientation  $\theta$  and car's speed  $s_x$  respectively. The control values in  $u$  are both expected to be given in the range  $[-0.1, 1.1]$ . The control value  $u_a$  changes the car's speed attribute by incrementing it with  $1.16u_a - 0.08$  so that the increment in speed is in the range  $[-0.08, 0.08]$  units. The car's speed is bounded to the upper limits of (15.0, 12.0) (in x and y direction resp.) and the lower limits of (10.0, 8.0). This is done to prevent indefinite increase or decrease of speed, and so that the car cannot outperform its previous version by simply driving faster (the upper speed limit is also the starting speed in case of no acceleration).

The environment as described here is used in the following experiments, and in case of constant speed experiments the acceleration control variable  $u_a$  is simply set to 0.5 (constant speed).



## 4-2 Implementing the agent for the autonomous car

The initial version of the autonomous agent is simply a recreation of the existing agent by [1]. The main components are `net3.py`, `modulation.py`, and `agent_functions.py`. A network object (from class `net3`) defines the nonlinear relationship between an arbitrary observation it receives as input, and the control output mapped from this observation. I.e. the network is built to describe exactly the equations as described in the neuronal dynamics in chapter 2. It consists of two layers (an autoencoder layer, and a Hebbian layer) and produces an optimal control policy by interacting with the environment. The network interacts with the environment in two ways. At every time tick  $k$  it samples the observed experience  $y_k$  (the rangefinder activations) to output the control  $u_k$  and send it back to the car, and secondly it adapts its network parameters through sensing a modulation according to the `modulation.py` class. This is done as described in section 2-2-3. Later on, I added the `performance.py` class with the addition of global rewards. The modulation in that case is derived from these global rewards as proposed in chapter 3. Training of the neural network is done through the `agent_functions.py` module, which contains both the autoencoder training process, and the Hebbian learning process as training options. One ambiguity I had found in [1] is the fact that Hebbian learning was actually done by first normalizing the postsynaptic activation. I.e.  $u_k$ , the real postsynaptic value was one between 0 and 1, whereas the postsynaptic activation as used during training is first normalized to a range between -1 and 1 to represent the measure of increment (or decrement) in the car's orientation. A detail not precisely mentioned in [1], but actually crucial for proper learning.

After having implemented the original works of [1] in Python, I proceeded with the additional experiments as described in section 2-2-4. These are experiments where the original algorithm is compared with and without autoencoder. Subsequently I tested the new algorithm with my own contribution where global rewards have been introduced. The corresponding particular reward  $r_k$  in this experiment was defined somewhat intuitively, where we gave relatively high importance to averting walls, and relatively low importance to increasing speed. That is, this particular reward (particular to each world state) is just a weighted average of several objectives that the agent needs to learn simultaneously (both steering and accelerating). The exact relative weights for these objectives are then tunable, based on what the engineer finds important. From  $r_k$  the global reward  $R_k^N$  is then constructed according equation 3-1 in `performance.py`, and the corresponding modulation is the difference of two consecutive global rewards as per equation 3-3. This modulation measurement is defined in `modulation.py` and is fed in an  $N$ -step delayed fashion to the Hebbian layer's training method in `training.py`. Finally, fundamental experiments are done as shown in chapter 5 to prove the new method works. The final algorithm is shown in pseudocode 2, where the main difference is the generalizing addition of global rewards.

## 4-3 Implementing the cart-pole environment

The cart-pole environment as used in the later experiments, is a gym-environment from python OpenAI [11]. The graphical representation of the environment is found in Figure 5-4. It concerns a pole attached by an un-actuated joint to the cart, which moves along a horizontal plane. At each moment in the environment the full four-dimensional state  $x$  is available to

---

**Algorithm 2** the Hebbian learning with global rewards method. This methods assumes:

1. multi-dimensional input space
  2. a global reward function based on the agent’s task
  3. unknown system dynamics
- 

```

1: procedure INITIALIZATION
2:   Initialize Autoencoder weights  $W_0 \leftarrow \mathcal{N}(0, 1) \in \mathbb{R}^{m \times n}$ 
3:   Initialize Hebbian weights  $w_0 \leftarrow \mathcal{N}(0, 1) \in \mathbb{R}^m$ 
4:   Initialize novelty-buffer size  $n_{novelty}$  arbitrarily
5:   Initialize number of training samples per tick  $n_{samples}$ 
6: end procedure
7: procedure TRAINING(for every time tick  $k$  do)
8:   take action  $u_k$  using policy derived from  $W_k, w_k$ 
9:   observe measurements  $y_{k+1}$  and global reward  $R_{k+1}^N$ 
10:  compute generalized modulation  $m_k \leftarrow R_{k+1}^N - R_k^N$ 
11:  update novelty buffer and weights  $W_{k+1}, w_{k+1}$ 
12: end procedure
13: Output: learned agent weights  $W, w$ 

```

---

**Table 4-1:** Summary of the observation and action spaces of the cart-pole environment. The real-valued fully observable state is shown in the left part, and the discrete-valued one-dimensional action space is shown in the right part.

Observations		Actions			
Nr	State	Min	Max	Control Action	Discrete Value
1	Cart Position	-4.8	4.8	Right Force	0
2	Cart Velocity	$-\infty$	$\infty$	Left Force	1
3	Pole Angle	-24	24		
4	Angular Velocity	$-\infty$	$\infty$		

the agent. This state is a sensor observation of the cart position, cart velocity, pole angle and pole velocity at the tip respectively. And a 'zero' state corresponds to an upright position of the pendulum. When an episode starts, the pendulum starts upright with a minor random variation for each state variable, sampled uniformly from a distribution in  $[-0.05, 0.05]$ . The control action applied onto the environment is a force to the right (1) or an equivalent force to the left (0). The environment is forced to receive either of the control values and may not apply zero force. A neat overview of the full observation space and action space is summarized in table 4-1.

The environmental time evolution occurs in episodes. The state initializes at the start of each episode as described above, and simulates a sequence of states based on the control outputs it receives. The episode is terminated when the pole angle deviates more than 15 degrees from its equilibrium state (the upright position), or when the cart moves more than 2.4 units from its center position. If the pole remains alive (balanced) for at least 200 time ticks, the episode terminates due to a successful episode. According to [11], cart-pole is considered solved when

it remains alive for 195 or more time ticks on average for a duration of 100 successive episodes. That same goal will be used for our purposes.

## 4-4 Implementing the cart-pole agent

The code to implement the agent can be reused from the autonomous car environment. The only differences will be in hyper parameters (structure, objective function, etc.). The main noticeable difference with the cart-pole compared to the autonomous car, is the fact that the sensor inputs measure values in the range of real numbers. So the autoencoder cannot properly learn if the state is not pre-processed (assuming that the nonlinearities are sigmoids still). Another obvious difference is the type of sensor input. Rather than high-level environmental features, the agent receives the physical state that is part of describing its own dynamics. For this reason, an autoencoder is generally not necessary, and might only pollute the network with needless additional information.

In experimentation, two different structures of the agent network are used. First, the exact same network structure as in the autonomous car example, to show that it does not work because of violation of the initial assumptions. Secondly, a different structure with multiple layers and different nonlinear activation functions will be used, such that the controller consists of more complexity, which is required for the proper policy.

All the implementation code is found in `cartpole.py` in appendix 7.



---

# Chapter 5

---

## Experiments

In this chapter, we will demonstrate fundamental experiments using the new generalized modulation rule, and report its results. I have reproduced the original work of [1] in a Python simulation, in which I have validated their method (Figure 2-2). Additionally I have added experiments of my own which will be discussed in this chapter (see source code in appendix 7). In the first section I will demonstrate experiments using the autonomous car environment, where we have made changes that make the system more complex. Showing that the new method can indeed learn to control this environment is the most important contribution of this paper. In the second section we will attempt to demonstrate the generality of this method, by applying it on a cart-pole system with inverted pendulum.

### 5-1 Experiments with the autonomous car environment

Firstly, let us introduce into the existing dynamics an additional control variable that sets the acceleration of the car. Now  $u_k \in \mathbb{R}^2$ , where  $u_k(2)$  is the second element of the vector  $u_k$  and which corresponds to the change of speed of the car. The control output  $u_k(2)$  is first linearly mapped to the range  $[-0.08, 0.08]$  to denote the change in car speed. The constant bounds of 0.08 were chosen such that the dynamics of accelerating look realistic relative to the car's speed. The precise choice of this value matters not too much in showing that the new method can regulate multiple control variables. The speed however, is capped to a maximum value of its original forward speed in order to prevent indefinite increase of speed. In this way, the new method is not allowed to outperform the old algorithm by simply driving at a faster speed. On the other hand, the lower limit of speed is set to approximately 12 units per tick, to prevent the car from a complete standstill. Also, the car's speed is set to its minimum value once it crashes into a wall, which additionally causes a sudden decrease in reward. Since the car is allowed to drive slower than its maximum speed, it acts as a disturbance to the overall training process as now it needs to learn both control variables simultaneously. Hence, we will demonstrate that even though both training and the dynamics are now more complex, the method is still able to learn reliably the optimal control policy.

Given the new dynamics of the system, we will choose the particular reward  $r_k$  to be equivalent to a weighted sum of the front-pointing distance  $d_k$  and the absolute speed  $s_k$  of the car.

$$r_k = 1.3 \cdot d_k + 0.3 \cdot s_k. \quad (5-1)$$

The precise relative weights are designed by tuning, but are set initially to match an intuitive understanding of the tasks that the agent needs to achieve. We can argue that in order to maximize the amount of laps the car completes, averting walls is indeed the most important aspect. So the larger relative importance of 1.3 is given to the sensor distances, whereas acceleration is only required if the agent is certain that the car will not crash. Hence the smaller value of 0.3 corresponds to the reward contribution of the speed, which was found to work well.

We furthermore chose a time horizon of  $N = 5$  steps ahead, and  $\gamma = 2$ . Taking a too large value for the anticipation horizon  $N$  may lead to loss of causality. That is, the agent may seek to reward an action too far in history that had no significant effect on the current world state. Of course, a value as low as 5 time steps anticipation reduces the quality of conclusion about the actual prediction capabilities the car has. As the sampling rate of an arbitrary system gets faster, the less impact a delay will have on that system, since the measurements will be closer to each other in value. Nevertheless the tuned value  $N = 5$  worked well in the case of the autonomous car, and it does not take away from the quality of conclusion about the generality of the method.

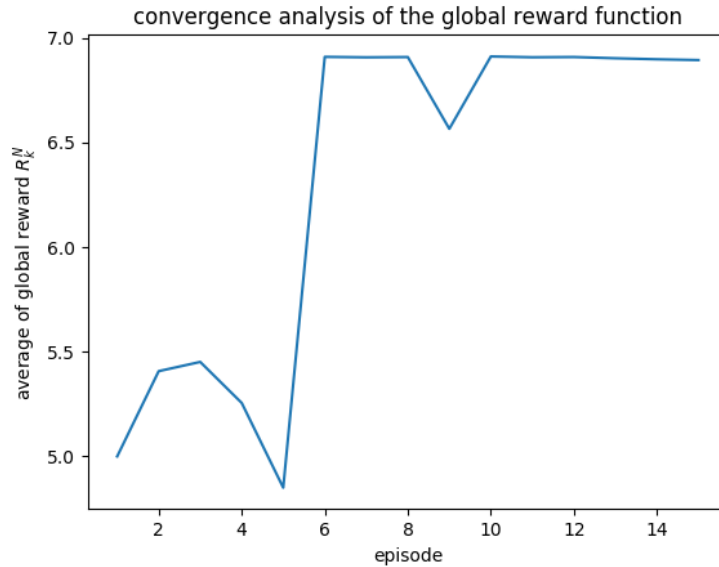
Next, to encourage exploitation of the input space more, we will introduce episodic learning, in which the environment's state gets reset after a certain amount of time ticks. At each reset, the agent's network keeps in memory, such that it keeps improving its policy. Specifically, we will set the agent to train for 15 episodes, where each episode lasts for 1200 simulation time ticks.

### 5-1-1 Convergence test

If we consider the true global optimum, as permitted by the environmental dynamics, where the maximum absolute speed  $s_k$  equals 19.21 (Euclidean distance of 15 and 12 units per tick), and the maximum distance (1.0 in case of no wall intersection), then we find  $\max r_k \approx 6.96$ . If we now collect the average values of the global rewards per episode, we find a convergence behavior as shown in Figure 5-1. From this figure we can deduce that indeed the global reward function  $R_k^N$  is optimized by the new modulation scheme, as the reward value converges to its optimum of 6.9.

### 5-1-2 Performance test

In the following experiment, we will demonstrate that the new method can achieve a similar performance as the existing algorithms. First, let us determine the performance range, in which the algorithms can be considered to perform similarly, i.e. near optimal. Past simulations have revealed that a *good* agent is able to complete anywhere between 0.9 up to approximately 1.0 laps in a span of 1000 time ticks. For example, both pure Hebbian and the novelty-Raahn algorithm performed close to 10 laps (at least 9.8 in all cases) in the evaluated



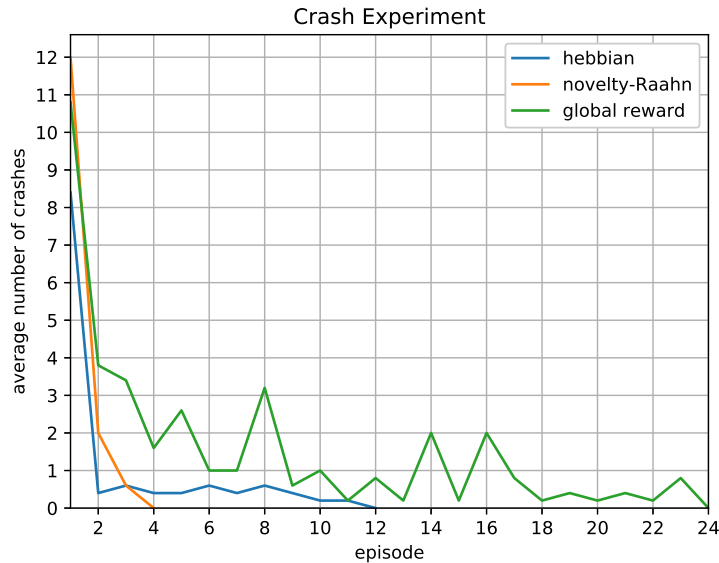
**Figure 5-1:** Graph of the average global reward occurring at each episode.

10,000 time ticks as tested in section 2-2-4. To define a more strict minimum boundary of when an agent can be considered near-optimal (i.e. satisfactory performance), we will determine an additional metric that measures this. That is, the agent behavior is considered unsatisfactory if it ever crashes. Between crashing and no crashing at all is clear difference in overall performance, because a crash may occupy the car at the same location for a significantly longer time, before it escapes. To test this metric, and hence find the minimum satisfactory performance, we will first run an experiment where the number of crashes is tracked for each algorithm.

The set-up of the experiment is as follows. We will again use episodic learning, where the goal is to achieve zero crashes in the final episode. Each episode takes at least 2500 time ticks, to ensure enough training time. If the car has crashed within the initial period of 2500 ticks, the episode will terminate once the 2500 ticks have passed. If the car crashes later than 2500 ticks, the episode terminates immediately. If no crashes occur at all during said episode for 10,000 time ticks, then learning will terminate for the current algorithm, and the final performance of the episode (corresponding to the final 10,000 time ticks) will be recorded. The number of crashes for each episode is recorded as well, which will be shown in a graph. The entire training experiment is repeated 5 times for all three algorithms, i.e. for pure Hebbian, novelty-Raahn and the novelty-Raahn structure in combination with the global rewards as defined in equation 5-1 (called *global reward*). Afterward, the average number of crashes along the 5 runs is computed for each episode and plotted in Figure 5-2.

We can conclude from the crash experiment that indeed the global reward method learns eventually to not crash along with the other algorithms. Expectedly, the global reward method takes longer to learn, since we train it on two control variables simultaneously, and so the network becomes more complex as well.

Furthermore, we find that indeed we have a similar performance when using global rewards. Namely, the minimum performance of novelty-Raahn and pure Hebbian without crashing is



**Figure 5-2:** Average number of crashes for each algorithm during training. From this figure we can conclude firstly that the global reward method learns the optimal behavior of not crashing along with the other algorithms shown. Secondly, global reward has indeed learned a similar performance, because its final performance is in the range of optimal performances as derived from the other algorithms when they do not crash (not shown in this figure).

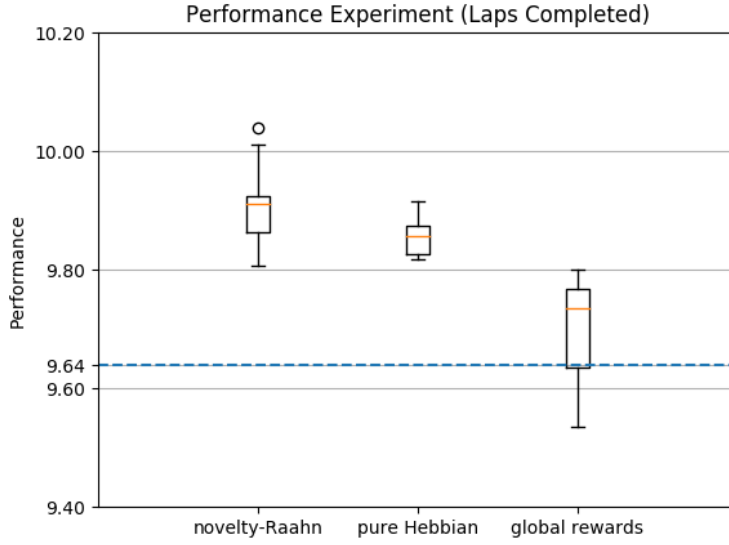
reached with a final performance of 9.64. With these methods, the car is forced to always drive at maximum speed, so we can use this value as the minimum satisfactory threshold. Now we can see that similarly, the global reward method achieves a performance of 9.73 on average, which is well above the satisfactory threshold, and from this we conclude that the global reward method can indeed achieve similar performance as the state-of-the-art algorithms while training multiple control variables simultaneously. All the final performances are compared in box plots as shown in Figure 5-3. If we observe the learning process visually, then we see that the autonomous car agent learns to accelerate immediately after it crashes, as it discovers that this is more optimal. It learns that driving at maximum speed is the best in combination with proper steering, and it does not need to slow down to achieve the optimal control policy. A well trained agent will never have to crash, and simply drive at maximum speed, as if it had no influence over the acceleration, just like the initial algorithm assumed.

## 5-2 Experiments with the cart-pole

In this section, we will do a secondary experiment with the cart-pole environment. The following experiments will show the generality of using Hebbian learning with global rewards in complex control systems.

In the first attempt, we will consider as if the components can be used in the same way as in





**Figure 5-3:** Depiction of the amount of laps completed over a certain amount of experiments. The figure shows that the global reward method achieves similar performances as the state-of-the-art algorithms after sufficient training, since the satisfactory threshold of 9.64 laps is achieved on average.

the autonomous car case. We have a cart-pole with four states

$$x_k = \begin{pmatrix} s_k \\ \dot{s}_k \\ \phi_k \\ \dot{\phi}_k \end{pmatrix}, \quad (5-2)$$

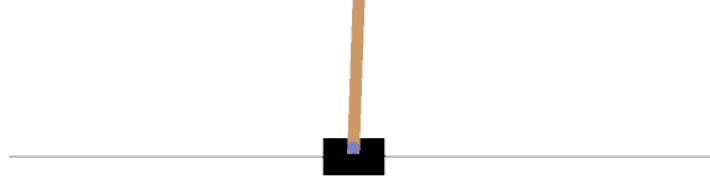
where  $s_k$  denotes the horizontal positioning,  $\phi_k$  the angle with respect to the upper equilibrium point, and the dot notation represent their respective time-derivative. The graphical representation of the corresponding environment is shown in Figure 5-4.

Note that  $x_k \in \mathbb{R}^4$  may have arbitrary real numbers, whereas the assumptions of novelty-Raahn consider an input range of  $[0, 1]$ . We can solve this easily by pre-processing the state to a suitable range of values between 0 and 1. After some tests with the environment, the pre-processed state

$$x_{p,k} = \begin{pmatrix} 1 \\ 1/4 \\ 2 \\ 1/6 \end{pmatrix} \odot x_k + \begin{pmatrix} 0.5 \\ 0.5 \\ 0.25 \\ 0.5 \end{pmatrix} \quad (5-3)$$

turns out to produce suitable values close to the range  $[0,1]$ , such that the autoencoder error could theoretically reach close to zero.

The remaining layers of the agent network remain the same as before, along with the training process. The only obvious difference now is the choice of global reward. If we consider the state  $x_k$ , then we can consider the main objective to bring this state to its desired equilibrium



**Figure 5-4:** Graphical representation of the cart-pole environment, with four known states. Objective is to stabilize the pole in its upper equilibrium (unstable point).

state  $x_d$ .

$$x_d = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}. \quad (5-4)$$

The measure of negative discrepancy between the actual state and the desired state is considered the particular reward. Namely, the Euclidean distance of the weighted state:

$$r_k = -\sqrt{w_r^T(x_d - x_k)} = -\sqrt{w_r^T x_k}. \quad (5-5)$$

After some trial and error with tuning, the relative weights  $w_r$  consists of the values

$$w_r = \begin{pmatrix} 0.1 \\ 0.1 \\ 1.3 \\ 0.1 \end{pmatrix} \quad (5-6)$$

because most importance is given to stabilizing the angle  $\phi_k$ . Various anticipation steps ahead  $N$  have also been tried, but unfortunately the algorithm did learn to stabilize the environment with these settings. After reflecting back upon the precise learning method in combination with the given environment, I have concluded the following: Firstly, one significant difference now is that the state  $x_k$  is already representative of the input space. The four state values represent the physical features of all important aspects, so we could argue that an autoencoder is not necessary, just like with the pure Hebbian controller in section 2-2-4.

Another main problem is still the uncontrollability or instability of the environment. If we observe visually at what happens during training, we see that because the system is unstable without a working controller initially, it keeps receiving negative modulation and hence keeps unlearning all it observes. This is now a major obstacle to proper learning, because the methods of raahn and global rewards assume that the environment has enough time for the modulation signal to directly drive it to the next state in which it receives a higher global reward. Whereas the unstable cart-pole environment becomes quickly unstable ones the pole deviates too much from the equilibrium state.

Finally, the network structure as proposed in [1] also does not contain enough complexity for the purposes of cart-pole. With the network structure as proposed in [1], i.e. as described in chapter 2, only a single layer is effectively used for learning control. The other part (autoencoder), if present, learns the high level perception of the same input space that the sensors supply. We have to take note that a single-layered feedback controller is not sufficient in complexity to control the cart-pole as used in this experiment. Contrary to controlling classic cart-pole systems, which are controllable via a single-layered linear mapping, where the main difference is that our controller is forced to be non-linear. The reason being is that our used cart-pole example has a discrete action space of two values ( $u=1$  or  $u=0$ ), so deriving a controller through linearization is not possible.

With this in mind, we will make some major changes in the agent's architecture. Firstly, we have established above that an autoencoder is not really necessary. The states however, are still in the neurally unsuitable real domain (i.e. it does not corresponds to biologically plausible neuronal activations, like a voltage spike that determines firing or non-firing). For this, we will change the neuronal dynamics of the hidden layers to be either -1 (non-firing) or +1 (firing), where specifically the activation function at each layer is the sign function. Such a structure is similar as in literature in [12], which is used as a supervised learning network. Apart from handling the input neurons correctly the main reasons for such a structure choice are the following: Firstly, this structure allows for a multi-layered network, which is required for more control complexity as we have established before. The version of Hebbian learning in the raahn algorithm normalizes the control output first for proper learning, which is not defined for a multilayered structure. The Hebbian weight updates in [12] are now simply  $\eta m$  or  $-\eta m$ , based on the pre- and postsynaptic activations. Here,  $\eta$  denotes the learning rate and  $m$  the global modulation signal. The network structure as such is strictly defined to learn in a supervised fashion, from which we can infer a different understanding in the sense of reinforcement learning.

In order to make the agent react to global rewards, we will now introduce what the global reward  $R_k^N$  will be. Since an episode of cart-pole terminates ones it deviates too much from the equilibrium, we can assign a particular reward  $r_k$  of 1, at every time tick that it is alive. The anticipation step value  $N$  is now dynamic, and based on the duration of each episode. We can then construct the global reward  $R_k^N$  to be equal to

$$R_k^N = \frac{1}{200} \sum_{i=k}^N r_i. \quad (5-7)$$

The 200 in the denominator corresponds to a perfect episode of 200 steps alive, giving a value of  $R_k^N = 1$ . If we consider the performance of cart-pole, then we anticipate that a lot of learning is necessary initially when the performance is bad, and little to no learning is required once the performance is good. In other words, this means that we already now the optimal global reward that is achievable ( $R_k^N = 1$ ), and so rather than defining the modulation as

$$m_k = g(R_{k+1}^N - R_k^N), \quad (5-8)$$

we can directly write the modulation as

$$m_k = 1 - R_k^N. \quad (5-9)$$

Note that no negative modulation can occur ( $R_k^N$  is between 0 and 1), and so the network will simply learn based on how well it performs. Note also, that we now use true episodic optimization where we wait each episode before learning, and update the weights based on the rolled out policy in that episode, whereas in the previous experiments with autonomous car we used a strictly real-time adaptive learner.

For this experiment, we will also change the amount of layers to three, with first layer a mapping from the four states  $x_k$  to 12 hidden features  $h_{1,k}$ . The second layer a mapping from the previous features to 24 hidden features  $h_{2,k}$ . And finally the last layer sigmoidal logistic mapping (like in raahn) from the previous 24 features to a single scalar value  $u_k$  in the range  $[0, 1]$ , where 1 denotes a right-force to the cart, and 0 denotes a left-force to the cart. The actual control value is rounded, and so the cart is not allowed to stand still, i.e. it must get either a left- or right-force.

With these new agent settings, training works quite well. The cart-pole is considered *solved* if it remains alive<sup>1</sup> for at least 195 time steps, for 100 consecutive episodes. The algorithm was tested 10 times, and cart-pole was solved in an average of 343.3 episodes, using the global reward method as described above.

---

<sup>1</sup>The cart-pole environment stays alive as long as the pole does not deviate more than 15 degrees from its upper equilibrium, and the position not more than 2.4 meters away from the center.

# Discussion and Future Work

We have shown in this thesis that Hebbian learning can effectively be used to control autonomous systems. The autonomous car has learned to control both its steering angle and its velocity simultaneously, through the maximization of a single global reward function. The sensitivity of the performance is in this way shifted from the specific modulation scheme, to a global reward function. Even though care still needs to be taken to correctly set up a global reward function, I would argue that it must be easier, since it is designed in a general context, rather than a specific one for a certain control variable. For future work, we can look at potentially even better ways to come up with this general global reward. We could for example look to use a different idea of the long term effect on the agent with some form of eligibility trace, like in different methods in literature. More importantly, I would say that establishing reward-based Hebbian learning theory using convergence proofs would have a high priority, since this would allow for an easier approach given a set of assumptions about the problem. Furthermore, we could potentially think about different methods of storing sample specific rewards, like a value function approximation as in other reinforcement learning approaches, and use our general modulation scheme along with these reward values to learn more robustly. Overall, the idea of using global rewards opens up a new door to novel algorithms with Hebbian learning, as the definition space of objective functions is very broad.



---

## Chapter 7

---

# Conclusion

I proposed in this thesis that global reward functions can be used to train autonomous systems, which solves the main drawbacks of the existing novelty-Raahn algorithm. In this way, it makes sure that the agent can learn all control agent simultaneously. This is supported by the results in chapter 5. These experiments have shown a couple of interesting things. Firstly, we have shown that autonomous control with Hebbian learning now actually becomes an optimization problem, in which the reward converges. Secondly, the new algorithm enables a similar performance as the existing algorithms, while having to control a more complex system. Finally, we have shown that the method is indeed more general, since it is able to learn to control a completely different system, with an unstable equilibrium. The type of design effort required is now different in that it focuses on a more general concept of optimization, rather than how the desired behavior is related to the exact dynamics. From this we can conclude that the algorithm has become more useful, which is also promising for future research.





---

## Bibliography

- [1] K. Stanley, J. Pugh, and J. Bowren, “Fully autonomous real-time autoencoder-augmented hebbian learning through the collection of novel experiences,” *The 2018 Conference on Artificial Life: A Hybrid of the European Conference on Artificial Life (ECAL) and the International Conference on the Synthesis and Simulation of Living Systems (ALIFE)*, pp. 382–389, 2016.
- [2] J. de Vries, “Collision detection,” in <https://learnopengl.com/In-Practice/2D-Game/Collisions/Collision-detection>, 2018.
- [3] J. Hohwy, *The Predictive Mind*. OUP Oxford, 2013.
- [4] K. J. Friston, “The free-energy principle: a rough guide to the brain?,” *Trends in cognitive sciences*, vol. 13 7, pp. 293–301, 2009.
- [5] D. O. Hebb, *The organization of behavior : a neuropsychological theory*. New York [etc.]: Wiley,, 1949.
- [6] S. Lowel and W. Singer, “Selection of intrinsic horizontal connections in the visual cortex by correlated neuronal activity,” *Science*, vol. 255, no. 5041, pp. 209–212, 1992.
- [7] S. Haykin, *Neural networks : a comprehensive foundation*. Upper Saddle River, N.J.: Prentice Hall,, 1999.
- [8] B. Baruque and E. Corchado, *Fusion methods for unsupervised learning ensembles*. Studies in computational intelligence ; v. 322; Studies in computational intelligence ; v. 322., Berlin ;; Springer,, 2011.
- [9] J. Burms, K. Caluwaerts, and J. Dambre, “Reward-modulated hebbian plasticity as leverage for partially embodied control in compliant robotics,” *Frontiers in Neurorobotics*, vol. 9, p. 9, 2015.
- [10] Q. Sheng, Z. Xianyi, W. Changhong, X. Z. Gao, and L. Zilong, “Design and implementation of an adaptive pid controller using single neuron learning algorithm,” in *Proceedings*

---

*of the 4th World Congress on Intelligent Control and Automation (Cat. No.02EX527)*, vol. 3, pp. 2279–2283 vol.3, 2002.

- [11] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai gym,” *CoRR*, vol. abs/1606.01540, 2016.
- [12] R. W. Strom, “Hebbian learning in multilayer neural networks,” Master’s thesis, California State University, Los Angeles, 2007.

---

# Glossary

## List of Acronyms



---

# Index

autoencoder, 2

biological plausibility, 2

learning, 1

local, 2

neuromodulation, 2

real-time, 2

scalable, 2

synapse, 2



---

# Source Code

Appendices are found in the back.

## -1 Python code for the autonomous car agent

### -1-1 The environment

The following code is part of the autonomous car environment.

#### raahn.py

```
1
2
3 # -*- coding: utf-8 -*-
4 """
5 Created on Tue Jun  5 15:12:41 2018
6
7
8 @author: ajdin
9 """
10 import pygame
11
12 from CarMazeEnv import CarMazeEnv
13
14 display_width = 1600
15 display_height = 1440
16
17 res = (display_width, display_height)
18
19 simul = None
20
21
22 def main():
23     global simul
24     simul = CarMazeEnv()
```

```
25     render = False
26     RUNNING, PAUSE = 0, 1
27     state = RUNNING
28
29     while not simul.game_ext:
30         for event in pygame.event.get():
31             if event.type == pygame.QUIT or simul.game_ext:
32                 simul.close()
33                 return
34
35             if event.type == pygame.KEYDOWN:
36                 if event.key == pygame.K_p and state == RUNNING: state = PAUSE
37                 elif event.key == pygame.K_p and state == PAUSE: state = RUNNING
38
39             if state == RUNNING:
40                 simul.crashcounting_test()
41             try:
42                 if render:
43                     simul.render(res)
44             except Exception as E:
45                 simul.close()
46                 raise E
47
48
49
50 main()
51 pygame.quit()
52 #quit()
53
54 #%%
55 #import envmanager
56 #import matplotlib.pyplot as plt
57 ###
58 ###
59 #import numpy as np
60 #
61 # plt.scatter(envmanager.performance_vector, envmanager.modulation_vector
62             )
63 #performance_vector = envmanager.performance_vector
64 #modulation_vector = envmanager.modulation_vector
65 #p = performance_vector.argsort()
66 #plt.plot(np.tanh(12*performance_vector[p]), modulation_vector[p])
67 #sensor_vector = envmanager.sensor_vector
68 #
69 #
70 #delta_performance_vector = np.zeros(len(performance_vector))
71 #delta_modulation_vector = np.zeros(len(modulation_vector))
72 #delta_sensor_vector = np.zeros(len(sensor_vector))
73 #batchsize = 2
74 #for x in range(len(performance_vector)-batchsize+1):
75 #     delta_performance_vector[x+1] = performance_vector[x+batchsize-1] -
76         performance_vector[x]
```



```

76 #     delta_modulation_vector[x+1] = np.mean(modulation_vector[x:x+
       batchsize])
77 #     delta_sensor_vector[x+1] = -sensor_vector[x+batchsize-1] +
       sensor_vector[x]
78 #
79 ##sort the indices
80 #p = delta_sensor_vector.argsort()
81 #
82 #
83 #plt.plot(delta_sensor_vector[p], modulation_vector[p])

```

### Car.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Wed Jun 13 21:11:27 2018
4
5  @author: ajdin
6  """
7  #from Line import Line
8
9  import numpy as np
10
11 import xml.etree.ElementTree as ET
12 import pygame
13
14 from lines import LineSegment
15
16 redcolor = (0xff,0,0)
17
18 class CarConfig:
19     """
20     Configures initial parameters or constants, used by the Car class.
21     """
22     root = ET.parse('Maps/XMap.xml').getroot()
23     ROOT = root.find('Robot')
24     def __init__(self, root=ROOT):
25         self.x = float(root.find('X').text)
26         self.y = float(root.find('Y').text)
27         self.angle = float(root.find('Angle').text)
28
29
30 class Car(object):
31     """
32     Agent that moves through a 2D environment, which can be visualized
33         using a pygame surface.
34
35     Parameters
36     -----
37     center: tuple, optional.
38         Sets the initial position in a 2D plane using Cartesian coordinates.
39         default is (0,0).

```

```

39     angle: float, optional.
40         Sets the initial angle in degrees of the agent relative to a 2D frame
41         . default is 0.
42
43     Attributes
44     -----
45     center:
46
47     velocity:
48
49     speed:
50
51     acceleration:
52
53     image:
54
55     angle:
56
57     """
58     CONTROL_THRESHOLD = 0.5
59     MIN_SPEED_X = 10.0
60     MIN_SPEED_Y = 8.0
61     MAX_SPEED_X = 15.0
62     MAX_SPEED_Y = 12.0
63     MAX_ROTATE = 1.0
64     MIN_ROTATE = 0.0
65     ROTATE_SPEED = 2.0
66     ROTATE_RANGE = 2.0*ROTATE_SPEED
67     ACCELERATION = 0.08
68     ACCELERATION_RANGE = 2.0*ACCELERATION
69     RADIUS = 3.0
70
71     def __init__(self, center=(0,0), angle=0):
72         self.center = center
73         self.lastpos = center
74         self.velocity = (0, 0)
75         self.speed = (self.MAX_SPEED_X, self.MAX_SPEED_Y)
76         self.speed_x = self.speed[0]
77         self.acceleration = 0
78         self.image = pygame.image.load('Textures/CarResized.png')
79         self.angle = angle % 360
80         self.last_can_move = True
81         self.can_move = True
82         self.reset_speed = False
83         self.num_crashes = 0
84         self.upper_bounds = (self.MAX_SPEED_X, self.MAX_SPEED_Y)
85         self.lower_bounds = (self.MIN_SPEED_X, self.MIN_SPEED_Y)
86
87     def update(self, walls):
88         self.lastpos = self.center
89         radians = np.radians(self.angle)
90         yx_ratio = self.MAX_SPEED_Y / self.MAX_SPEED_X

```

```

91     resulting_speed = (self.speed_x, yx_ratio*self.speed_x)
92     self.speed = np.clip(resulting_speed, self.lower_bounds, self.
93         upper_bounds)
94     self.speed_x = np.clip(self.speed_x, self.MIN_SPEED_X, self.
95         MAX_SPEED_X)
96     self.velocity = (np.cos(radians)*self.speed[0], np.sin(radians)*self.
97         speed[1])
98     original = self.center
99     projected = tuple(np.add(original, self.velocity))
100    collision_line = LineSegment(original, projected)
101    walls_in_bounds = collision_line.entities_in_bounds(walls)
102    self.can_move = True
103    for wall in walls_in_bounds:
104        intersections = collision_line.intersects(wall)
105        if len(intersections) > 0:
106            if self.last_can_move:
107                self.num_crashes += 1
108                self.can_move = False
109                if self.reset_speed:
110                    self.speed = (self.MIN_SPEED_X, self.MIN_SPEED_Y)
111                break
112    if self.can_move:
113        self.center = tuple(np.add(self.center, self.velocity))
114        self.last_can_move = self.can_move
115
116    def draw(self, display, options):
117        x_scale = options['x_scale']
118        x_translate = options['x_translate']
119        x = self.center[0]*x_scale + x_translate
120        y = self.center[1]*x_scale + x_translate
121        angle = self.angle
122        angle %= 360
123        rotated_img = pygame.transform.rotate(self.image, -angle)
124        original_rect = self.image.get_rect(center=(x,y))
125        rotated_rect = rotated_img.get_rect(center=original_rect.center)
126        display.blit(rotated_img, rotated_rect)
127        pygame.draw.circle(display, redcolor, tuple(int(i) for i in (x, y)),
128            int(self.RADIUS), 2)

```

### configs.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Tue Oct 16 08:22:29 2018
4
5  @author: ajdin
6  """
7
8  import json
9
10 file = 'Networks/HebbianNet.json'
11
12 class Net3Config:

```

```

13
14 def __init__(self, filename=file):
15     self.layers = []
16     with open(filename) as rfile:
17         container = json.load(rfile)
18         for layer in container['NeuralNetwork3']['NetworkLayers']:
19             try:
20                 neuron_count = layer['neuron_count']
21                 learning_rate = layer['learning_rate']
22                 training_method = layer['training_method']
23                 modulation_scheme = layer['modulation_scheme']
24                 layerconfig = LayerConfig(neuron_count, learning_rate,
25                                           training_method, modulation_scheme)
26                 self.layers.append(layerconfig)
27             except KeyError:
28                 break
29         self.input_count = container['NeuralNetwork3']['input_count']
30         self.output_noise_mag = container['NeuralNetwork3']['output_noise_mag']
31         self.weight_noise_mag = container['NeuralNetwork3']['weight_noise_mag']
32         self.weight_cap = container['NeuralNetwork3']['weight_cap']
33
34 class LayerConfig:
35
36     def __init__(self, neuron_count, learning_rate, training_method,
37                 modulation_scheme):
38         self.neuron_count = neuron_count
39         self.learning_rate = learning_rate
40         self.training_method = training_method
41         self.modulation_scheme = modulation_scheme

```

### lines.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Aug 25 16:27:41 2018
4
5  @author: ajdin
6  """
7  import numpy as np
8
9  def get_dist(p1 , p2):
10     return ((p2[1]-p1[1])**2 + (p2[0]-p1[0])**2)**(1/2)
11
12 class LineSegment:
13
14     def __init__(self, startPoint, endPoint):
15         self.startPoint = startPoint
16         self.endPoint = endPoint
17         self.setUp()
18

```

```
19 def setUp(self):
20     deltaX = self.endPoint[0] - self.startPoint[0]
21
22     #Set up the bounds
23     self.upperBoundY = max(self.endPoint[1], self.startPoint[1])
24     self.lowerBoundY = min(self.endPoint[1], self.startPoint[1])
25
26     self.upperBoundX = max(self.endPoint[0], self.startPoint[0])
27     self.lowerBoundX = min(self.endPoint[0], self.startPoint[0])
28
29     #If the change in x is 0, the slope is undefined
30     if abs(deltaX) <= 1e-5:
31         self.vertical = True
32         self.slope = None
33         #If the line segment is just a point, it's x coordinate is stored
           in lowerBounds
34         self.lowerBoundX = self.startPoint[0]
35     #End if
36     else:
37         self.vertical = False
38         self.slope = (self.endPoint[1] - self.startPoint[1]) / deltaX
39         self.yIntercept = self.startPoint[1] - (self.slope * self.
           startPoint[0])
40
41 def getY(self, x):
42     if x >= self.lowerBoundX and x <= self.upperBoundX and not self.
           vertical:
43         return self.slope*x + self.yIntercept
44     else:
45         return np.inf
46
47 def intersects(self, line):
48     intersection = []
49     bothVertical = self.vertical and line.vertical
50     parallel = self.slope == line.slope and not self.vertical and not
           line.vertical
51
52     if bothVertical or parallel:
53         if (bothVertical and self.lowerBoundX == line.lowerBoundX) or (
           parallel and self.yIntercept == line.yIntercept):
54             lowerInBounds = self.valueInBounds(self.lowerBoundY, line.
           lowerBoundY, line.upperBoundY)
55             upperInBounds = self.valueInBounds(self.upperBoundY, line.
           lowerBoundY, line.upperBoundY)
56
57             if lowerInBounds and upperInBounds:
58                 intersection.append((self.lowerBoundX, self.lowerBoundY))
59                 intersection.append((self.lowerBoundX, self.upperBoundY))
60             elif lowerInBounds:
61                 intersection.append((self.lowerBoundX, self.lowerBoundY))
62                 intersection.append((self.lowerBoundX, line.upperBoundY))
63             elif upperInBounds:
64                 intersection.append((self.lowerBoundX, line.lowerBoundY))
```

```

65         intersection.append((self.lowerBoundX, self.upperBoundY))
66     elif self.valueInBounds(line.lowerBoundY, self.lowerBoundY, self.
        upperBoundY) and self.valueInBounds(line.upperBoundY, self.
        lowerBoundY, self.upperBoundY):
67         intersection.append((self.lowerBoundX, line.lowerBoundY))
68         intersection.append((self.lowerBoundX, line.upperBoundY))
69     return intersection
70 elif self.vertical:
71     y = line.getY(self.lowerBoundX)
72     #Maker sure the returned y is valid
73     #GetY not returning infinity makes sure that the point is in bounds
        of this line
74     if y == np.inf or not line.valueInBounds(y, self.lowerBoundY, self.
        upperBoundY):
75         return intersection
76
77     intersection.append((self.lowerBoundX, y))
78     return intersection
79 elif line.vertical:
80     y = self.getY(self.lowerBoundX)
81     #Make sure the returned y is valid
82     #GetY not returning infinity makes sure that the point is in bounds
        of line
83     if y == np.inf or not self.valueInBounds(y, line.lowerBoundY, line.
        upperBoundY):
84         return intersection
85
86     intersection.append((line.lowerBoundX, y))
87     return intersection
88 else:
89     intersectionX = (line.yIntercept - self.yIntercept) / (self.slope -
        line.slope)
90     if self.valueInBounds(intersectionX, self.lowerBoundX, self.
        upperBoundX) and self.valueInBounds(intersectionX, line.
        lowerBoundX, line.upperBoundX):
91         intersection.append(((intersectionX), self.getY(intersectionX)))
92
93     return intersection
94
95
96 def valueInBounds(self, value, lowerBound, upperBound):
97     if value >= lowerBound and value <= upperBound:
98         return True
99     else:
100        return False
101
102 def angle_between(self, line2):
103     vector1 = (self.endPoint[0] - self.startPoint[0], self.endPoint[1] -
        self.startPoint[1])
104     vector2 = (line2.endPoint[0] - line2.startPoint[0], line2.endPoint[1]
        - line2.startPoint[1])
105     inner_product = vector1[0]*vector2[0] + vector1[1]*vector2[1]
106     len1 = np.linalg.norm(vector1)

```

```

107     len2 = np.linalg.norm(vector2)
108     return np.degrees(np.arccos(inner_product/(len1*len2)))
109
110     @staticmethod
111     def overlapsInDimension(pointsA, pointsB):
112         Amin = min(pointsA)
113         Amax = max(pointsA)
114         Bmin = min(pointsB)
115         Bmax = max(pointsB)
116         if Amin > Bmax:
117             return False
118         if Amax < Bmin:
119             return False
120         return True
121
122     # def entities_in_bounds(self, entities):
123     #     entities_in_bounds = []
124     #     for entity in entities:
125     #         lowerInBoundsY = self.valueInBounds(self.lowerBoundY, entity.
126     #             lowerBoundY, entity.upperBoundY)
127     #         upperInBoundsY = self.valueInBounds(self.upperBoundY, entity.
128     #             lowerBoundY, entity.upperBoundY)
129     #         lowerInBoundsX = self.valueInBounds(self.lowerBoundX, entity.
130     #             lowerBoundX, entity.upperBoundX)
131     #         upperInBoundsX = self.valueInBounds(self.upperBoundX, entity.
132     #             lowerBoundX, entity.upperBoundX)
133     #         entityInBounds = lowerInBoundsY or upperInBoundsY or
134     #             lowerInBoundsX or upperInBoundsX
135     #         if entityInBounds:
136     #             entities_in_bounds.append(entity)
137     #     return entities_in_bounds
138
139     def entities_in_bounds(self, entities):
140         entities_in_bounds = []
141         pointsAx = (self.lowerBoundX, self.upperBoundX)
142         pointsAy = (self.lowerBoundY, self.upperBoundY)
143         for entity in entities:
144             pointsBx = (entity.lowerBoundX, entity.upperBoundX)
145             pointsBy = (entity.lowerBoundY, entity.upperBoundY)
146             overlapsInX = LineSegment.overlapsInDimension(pointsAx, pointsBx)
147             overlapsInY = LineSegment.overlapsInDimension(pointsAy, pointsBy)
148             entityInBounds = overlapsInX and overlapsInY
149             if entityInBounds:
150                 entities_in_bounds.append(entity)
151         return entities_in_bounds

```

### rangefinder.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Sun Aug 26 20:08:34 2018
4
5  @author: ajdin

```

```

6 """
7 from lines import LineSegment
8 from lines import get_dist
9
10 import numpy as np
11 import pygame
12
13 class RangeFinder(LineSegment):
14
15     def __init__(self, car, default_length=350, relativeAngle=0):
16         self.car = car
17         self.default_length = default_length
18         self.length = self.default_length
19         self.relativeAngle = relativeAngle
20         super().__init__(self.getStartPoint(), self.getEndPoint())
21         self.activation = 0.0
22
23     def getStartPoint(self):
24         return self.car.center
25
26     def getEndPoint(self):
27         rads = np.radians(self.car.angle + self.relativeAngle)
28         change = (np.cos(rads), np.sin(rads))
29         return tuple(np.add(self.car.center, tuple(x*(self.default_length)
30             for x in change)))
31
32     def update_position(self):
33         self.startPoint = self.getStartPoint()
34         self.endPoint = self.getEndPoint()
35         self.setUp()
36
37     def update(self, walls):
38         self.update_position()
39         walls_in_bounds = self.entities_in_bounds(walls)
40         nearestWallDistance = self.default_length
41         for wall in walls_in_bounds:
42             intersections = self.intersects(wall)
43             if len(intersections) > 0:
44                 distance = get_dist(self.car.center, intersections[0])
45                 nearestWallDistance = min(distance, nearestWallDistance)
46         self.length = nearestWallDistance
47         self.activation = (self.default_length - self.length) / self.
48             default_length
49
50     def draw(self, display, options):
51         color = tuple([self.activation*x for x in (255, 80, 0)])
52         x_scale = options['x_scale']
53         x_translate = options['x_translate']
54         startPoint = tuple(x_scale*x+x_translate for x in self.startPoint)
55         endPoint= tuple(x_scale*x+x_translate for x in self.endPoint)
56         pygame.draw.aaline(display, color, startPoint, endPoint)

```



```

57 class RangeFinderGroup(list):
58
59     def __init__(self, car, size=11):
60         self.car = car
61         self.count = size
62         self.activations = []
63         self.configure()
64
65     def configure(self, default_length=350, angleOffset=-90, angle_between
        =18):
66         self.default_length = default_length
67         self.startAngle = angleOffset
68         self.angleSpacing = angle_between
69
70         for i in range(self.count):
71             relativeAngle = self.startAngle + (self.angleSpacing*i)
72             self.append(RangeFinder(self.car, self.default_length,
                relativeAngle))
73             self.activations.append(self[i].activation)
74
75     def update(self, walls):
76         for i in range(self.count):
77             self[i].update(walls)
78             self.activations[i] = self[i].activation
79
80     def draw(self, display, options):
81         for i in range(self.count):
82             self[i].draw(display, options)

```

## Wall.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu Sep 13 13:26:34 2018
4
5  @author: ajdin
6  """
7  from lines import LineSegment
8  import pygame
9  import xml.etree.ElementTree as ET
10
11 root = ET.parse('Maps/XMap.xml').getroot()
12 ROOT = root.find('Entity')
13
14 class EntityConfig:
15     def __init__(self, root=ROOT):
16         self.x = float(root.find('X').text)
17         self.y = float(root.find('Y').text)
18         self.relX = float(root.find('RelX').text)
19         self.rely = float(root.find('RelY').text)
20         self.angle = float(root.find('Angle').text)
21         self.type_ = None
22         if root.get('Type'):

```

```

23     self.type_ = root.get('Type')
24
25 black = (0, 0, 0)
26
27 class Wall(LineSegment):
28     def __init__(self, startPoint, endPoint):
29         super().__init__(startPoint, endPoint)
30
31     def draw(self, display, options):
32         x_scale = options["x_scale"]
33         x_translate = options["x_translate"]
34         startPoint = tuple(x_scale*x+x_translate for x in self.startPoint)
35         endPoint = tuple(x_scale*x+x_translate for x in self.endPoint)
36         pygame.draw.aaline(display, black, startPoint, endPoint)

```

## -1-2 Xmap.xml

```

1 <?xml version="1.0" encoding="utf-8"?>
2 <Map xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:xsd="
   http://www.w3.org/2001/XMLSchema">
3   <Robot>
4     <X>1561.9331742243444</X>
5     <Y>505.79080279525812</Y>
6     <Angle>0</Angle>
7   </Robot>
8   <Entity Type="Wall">
9     <X>1301.3842482100238</X>
10    <Y>675.13128038897889</Y>
11    <RelX>769.83293556080639</RelX>
12    <RelY>1.4438228390645236E-11</RelY>
13    <Angle>0</Angle>
14  </Entity>
15  <Entity Type="Wall">
16    <X>2071.21718377083</X>
17    <Y>675.13128038899333</Y>
18    <RelX>210.78758949880694</RelX>
19    <RelY>-52.512155591572082</RelY>
20    <Angle>0</Angle>
21  </Entity>
22  <Entity Type="Wall">
23    <X>2282.0047732696371</X>
24    <Y>622.61912479742125</Y>
25    <RelX>155.79952267303088</RelX>
26    <RelY>-171.53970826580246</RelY>
27    <Angle>0</Angle>
28  </Entity>
29  <Entity Type="Wall">
30    <X>2437.804295942668</X>
31    <Y>451.07941653161879</Y>
32    <RelX>128.30548926014308</RelX>
33    <RelY>-297.56888168557509</RelY>
34    <Angle>0</Angle>
35  </Entity>

```

```
36 <Entity Type="Wall">
37   <X>2566.1097852028111</X>
38   <Y>153.5105348460437</Y>
39   <RelX>142.05250596658698</RelX>
40   <RelY>-150.53484602917376</RelY>
41   <Angle>0</Angle>
42 </Entity>
43 <Entity Type="Wall">
44   <X>2708.162291169398</X>
45   <Y>2.9756888168699334</Y>
46   <RelX>164.96420047732727</RelX>
47   <RelY>-38.508914100486209</RelY>
48   <Angle>0</Angle>
49 </Entity>
50 <Entity Type="Wall">
51   <X>2873.1264916467253</X>
52   <Y>-35.533225283616275</Y>
53   <RelX>155.79952267303088</RelX>
54   <RelY>45.510534846029145</RelY>
55   <Angle>0</Angle>
56 </Entity>
57 <Entity Type="Wall">
58   <X>3028.9260143197562</X>
59   <Y>9.97730956241287</Y>
60   <RelX>105.39379474940324</RelX>
61   <RelY>94.521880064829816</RelY>
62   <Angle>0</Angle>
63 </Entity>
64 <Entity Type="Wall">
65   <X>3134.3198090691594</X>
66   <Y>104.49918962724269</Y>
67   <RelX>-13.747016706443901</RelX>
68   <RelY>143.53322528363049</RelY>
69   <Angle>0</Angle>
70 </Entity>
71 <Entity Type="Wall">
72   <X>3120.5727923627155</X>
73   <Y>248.03241491087317</Y>
74   <RelX>-128.30548926014308</RelX>
75   <RelY>196.04538087520177</RelY>
76   <Angle>0</Angle>
77 </Entity>
78 <Entity Type="Wall">
79   <X>2992.2673031025724</X>
80   <Y>444.07779578607494</Y>
81   <RelX>-160.38186157517885</RelX>
82   <RelY>227.55267423014607</RelY>
83   <Angle>0</Angle>
84 </Entity>
85 <Entity Type="Wall">
86   <X>2831.8854415273936</X>
87   <Y>671.630470016221</Y>
88   <RelX>-201.622911694511</RelX>
```

```
89     <RelY>161.03727714748766</RelY>
90     <Angle>0</ Angle>
91 </ Entity>
92 <Entity Type="Wall ">
93     <X>2630.2625298328826</X>
94     <Y>832.66774716370867</Y>
95     <RelX>-114.55847255369918</RelX>
96     <RelY>154.03565640194506</RelY>
97     <Angle>0</ Angle>
98 </ Entity>
99 <Entity Type="Wall ">
100    <X>2515.7040572791834</X>
101    <Y>986.70340356565373</Y>
102    <RelX>-68.7350835322195</RelX>
103    <RelY>147.03403565640167</RelY>
104    <Angle>0</ Angle>
105 </ Entity>
106 <Entity Type="Wall ">
107    <X>2446.9689737469639</X>
108    <Y>1133.7374392220554</Y>
109    <RelX>-9.164677804295934</RelX>
110    <RelY>129.52998379254495</RelY>
111    <Angle>0</ Angle>
112 </ Entity>
113 <Entity Type="Wall ">
114    <X>2437.804295942668</X>
115    <Y>1263.2674230146004</Y>
116    <RelX>36.658711217183736</RelX>
117    <RelY>154.03565640194483</RelY>
118    <Angle>0</ Angle>
119 </ Entity>
120 <Entity Type="Wall ">
121    <X>2474.4630071598517</X>
122    <Y>1417.3030794165452</Y>
123    <RelX>206.20525059665852</RelX>
124    <RelY>105.02431118314416</RelY>
125    <Angle>0</ Angle>
126 </ Entity>
127 <Entity Type="Wall ">
128    <X>2680.6682577565102</X>
129    <Y>1522.3273905996894</Y>
130    <RelX>531.55131264916554</RelX>
131    <RelY>367.58508914100548</RelY>
132    <Angle>0</ Angle>
133 </ Entity>
134 <Entity Type="Wall ">
135    <X>3212.2195704056758</X>
136    <Y>1889.9124797406948</Y>
137    <RelX>206.20525059665852</RelX>
138    <RelY>234.55429497568866</RelY>
139    <Angle>0</ Angle>
140 </ Entity>
141 <Entity Type="Wall ">
```

```
142     <X>3418.4248210023343</X>
143     <Y>2124.4667747163835</Y>
144     <RelX>-9.164677804295934</RelX>
145     <RelY>133.03079416531591</RelY>
146     <Angle>0</ Angle>
147 </ Entity>
148 <Entity Type="Wall ">
149     <X>3409.2601431980383</X>
150     <Y>2257.4975688816994</Y>
151     <RelX>-100.81145584725527</RelX>
152     <RelY>133.03079416531546</RelY>
153     <Angle>0</ Angle>
154 </ Entity>
155 <Entity Type="Wall ">
156     <X>3308.4486873507831</X>
157     <Y>2390.5283630470149</Y>
158     <RelX>-160.38186157517976</RelX>
159     <RelY>38.508914100485526</RelY>
160     <Angle>0</ Angle>
161 </ Entity>
162 <Entity Type="Wall ">
163     <X>3148.0668257756033</X>
164     <Y>2429.0372771475004</Y>
165     <RelX>-288.68735083532192</RelX>
166     <RelY>-35.008103727714115</RelY>
167     <Angle>0</ Angle>
168 </ Entity>
169 <Entity Type="Wall ">
170     <X>2859.3794749402814</X>
171     <Y>2394.0291734197863</Y>
172     <RelX>-192.45823389021461</RelX>
173     <RelY>-101.5235008103723</RelY>
174     <Angle>0</ Angle>
175 </ Entity>
176 <Entity Type="Wall ">
177     <X>2666.9212410500668</X>
178     <Y>2292.505672609414</Y>
179     <RelX>-375.75178997613375</RelX>
180     <RelY>-287.06645056726074</RelY>
181     <Angle>0</ Angle>
182 </ Entity>
183 <Entity Type="Wall ">
184     <X>2291.169451073933</X>
185     <Y>2005.4392220421532</Y>
186     <RelX>-261.19331742243435</RelX>
187     <RelY>-147.034035656402</RelY>
188     <Angle>0</ Angle>
189 </ Entity>
190 <Entity Type="Wall ">
191     <X>2029.9761336514987</X>
192     <Y>1858.4051863857512</Y>
193     <RelX>-169.546539379475</RelX>
194     <RelY>-24.505672609400335</RelY>
```

```
195     <Angle>0</ Angle>
196 </ Entity>
197 <Entity Type="Wall ">
198     <X>1860.4295942720237</X>
199     <Y>1833.8995137763509</Y>
200     <RelX>-490.31026252983293</RelX>
201     <RelY>-3.5008103727714115</RelY>
202     <Angle>0</ Angle>
203 </ Entity>
204 <Entity Type="Wall ">
205     <X>1370.1193317421908</X>
206     <Y>1830.3987034035795</Y>
207     <RelX>-174.1288782816232</RelX>
208     <RelY>70.016207455429139</RelY>
209     <Angle>0</ Angle>
210 </ Entity>
211 <Entity Type="Wall ">
212     <X>1195.9904534605676</X>
213     <Y>1900.4149108590086</Y>
214     <RelX>-229.11694510739824</RelX>
215     <RelY>178.54132901134517</RelY>
216     <Angle>0</ Angle>
217 </ Entity>
218 <Entity Type="Wall ">
219     <X>966.87350835316931</X>
220     <Y>2078.9562398703538</Y>
221     <RelX>-155.79952267303111</RelX>
222     <RelY>203.04700162074505</RelY>
223     <Angle>0</ Angle>
224 </ Entity>
225 <Entity Type="Wall ">
226     <X>811.07398568013821</X>
227     <Y>2282.0032414910988</Y>
228     <RelX>-219.95226730310253</RelX>
229     <RelY>147.034035656402</RelY>
230     <Angle>0</ Angle>
231 </ Entity>
232 <Entity Type="Wall ">
233     <X>591.12171837703568</X>
234     <Y>2429.0372771475008</Y>
235     <RelX>-229.11694510739858</RelX>
236     <RelY>45.510534846029259</RelY>
237     <Angle>0</ Angle>
238 </ Entity>
239 <Entity Type="Wall ">
240     <X>362.0047732696371</X>
241     <Y>2474.54781199353</Y>
242     <RelX>-256.61097852028638</RelX>
243     <RelY>-28.006482982171747</RelY>
244     <Angle>0</ Angle>
245 </ Entity>
246 <Entity Type="Wall ">
247     <X>105.39379474935072</X>
```

```
248     <Y>2446.5413290113584</Y>
249     <RelX>-151.21718377088303</RelX>
250     <RelY>-105.02431118314416</RelY>
251     <Angle>0</Angle>
252 </Entity>
253 <Entity Type="Wall">
254     <X>-45.823389021532307</X>
255     <Y>2341.5170178282142</Y>
256     <RelX>2</RelX>
257     <RelY>-154.03565640194483</RelY>
258     <Angle>0</Angle>
259 </Entity>
260 <Entity Type="Wall">
261     <X>-43.823389021532307</X>
262     <Y>2187.4813614262694</Y>
263     <RelX>192.45823389021473</RelX>
264     <RelY>-238.05510534846053</RelY>
265     <Angle>0</Angle>
266 </Entity>
267 <Entity Type="Wall">
268     <X>146.63484486868242</X>
269     <Y>1949.4262560778088</Y>
270     <RelX>279.52267303102622</RelX>
271     <RelY>-238.05510534845939</RelY>
272     <Angle>0</Angle>
273 </Entity>
274 <Entity Type="Wall">
275     <X>426.15751789970864</X>
276     <Y>1711.3711507293494</Y>
277     <RelX>215.36992840095468</RelX>
278     <RelY>-185.54294975688777</RelY>
279     <Angle>0</Angle>
280 </Entity>
281 <Entity Type="Wall">
282     <X>641.52744630066331</X>
283     <Y>1525.8282009724617</Y>
284     <RelX>119.14081145584737</RelX>
285     <RelY>-168.03889789303093</RelY>
286     <Angle>0</Angle>
287 </Entity>
288 <Entity Type="Wall">
289     <X>760.66825775651068</X>
290     <Y>1357.7893030794307</Y>
291     <RelX>-4.582338902147967</RelX>
292     <RelY>-185.54294975688754</RelY>
293     <Angle>0</Angle>
294 </Entity>
295 <Entity Type="Wall">
296     <X>1301.3842482100238</X>
297     <Y>675.13128038897889</Y>
298     <RelX>-183.29355608597143</RelX>
299     <RelY>-115.52674230144339</RelY>
300     <Angle>0</Angle>
```

```
301 </Entity>
302 <Entity Type="Wall">
303   <X>1118.0906921240523</X>
304   <Y>559.6045380875355</Y>
305   <RelX>-187.87589498806688</RelX>
306   <RelY>-171.53970826580229</RelY>
307   <Angle>0</Angle>
308 </Entity>
309 <Entity Type="Wall">
310   <X>930.21479713598546</X>
311   <Y>388.06482982173321</Y>
312   <RelX>-183.29355608591868</RelX>
313   <RelY>-94.521880064829816</RelY>
314   <Angle>0</Angle>
315 </Entity>
316 <Entity Type="Wall">
317   <X>746.92124105006678</X>
318   <Y>293.5429497569034</Y>
319   <RelX>-270.35799522673051</RelX>
320   <RelY>-94.521880064829588</RelY>
321   <Angle>0</Angle>
322 </Entity>
323 <Entity Type="Wall">
324   <X>476.56324582333627</X>
325   <Y>199.02106969207381</Y>
326   <RelX>-155.79952267303088</RelX>
327   <RelY>14.00324149108576</RelY>
328   <Angle>0</Angle>
329 </Entity>
330 <Entity Type="Wall">
331   <X>320.7637231503054</X>
332   <Y>213.02431118315957</Y>
333   <RelX>-128.30548926014319</RelX>
334   <RelY>84.019448946515354</RelY>
335   <Angle>0</Angle>
336 </Entity>
337 <Entity Type="Wall">
338   <X>192.45823389016221</X>
339   <Y>297.04376012967492</Y>
340   <RelX>-64.152744630071652</RelX>
341   <RelY>126.02917341977309</RelY>
342   <Angle>0</Angle>
343 </Entity>
344 <Entity Type="Wall">
345   <X>128.30548926009055</X>
346   <Y>423.072933549448</Y>
347   <RelX>36.658711217183736</RelX>
348   <RelY>133.03079416531608</RelY>
349   <Angle>0</Angle>
350 </Entity>
351 <Entity Type="Wall">
352   <X>164.96420047727429</X>
353   <Y>556.10372771476409</Y>
```



```
354     <RelX>91.646778042959454</RelX>
355     <RelY>101.52350081037207</RelY>
356     <Angle>0</Angle>
357 </Entity>
358 <Entity Type="Wall">
359     <X>256.61097852023374</X>
360     <Y>657.62722852513616</Y>
361     <RelX>192.45823389021473</RelX>
362     <RelY>126.02917341977172</RelY>
363     <Angle>0</Angle>
364 </Entity>
365 <Entity Type="Wall">
366     <X>449.06921241044847</X>
367     <Y>783.65640194490788</Y>
368     <RelX>197.04057279236281</RelX>
369     <RelY>126.02917341977309</RelY>
370     <Angle>0</Angle>
371 </Entity>
372 <Entity Type="Wall">
373     <X>646.10978520281128</X>
374     <Y>909.685575364681</Y>
375     <RelX>91.646778042959568</RelX>
376     <RelY>94.5218800648297</RelY>
377     <Angle>0</Angle>
378 </Entity>
379 <Entity Type="Wall">
380     <X>737.75656324577085</X>
381     <Y>1004.2074554295107</Y>
382     <RelX>18.329355608591868</RelX>
383     <RelY>168.03889789303253</RelY>
384     <Angle>0</Angle>
385 </Entity>
386 <Entity Type="Wall">
387     <X>1397.8854415274473</X>
388     <Y>462.13495863941552</Y>
389     <RelX>535.21718377084085</RelX>
390     <RelY>-7.5740259740153988</RelY>
391     <Angle>0</Angle>
392 </Entity>
393 <Entity Type="Wall">
394     <X>1933.1026252982881</X>
395     <Y>454.56093266540012</Y>
396     <RelX>109.97613365155121</RelX>
397     <RelY>-53.2987012987013</RelY>
398     <Angle>0</Angle>
399 </Entity>
400 <Entity Type="Wall">
401     <X>2043.0787589498393</X>
402     <Y>401.26223136669881</Y>
403     <RelX>91.646778042959568</RelX>
404     <RelY>-129.03896103896102</RelY>
405     <Angle>0</Angle>
406 </Entity>
```

```
407 <Entity Type="Wall">
408   <X>2134.7255369927989</X>
409   <Y>272.22327032773779</Y>
410   <RelX>91.64677804295934</RelX>
411   <RelY>-162.7012987012987</RelY>
412   <Angle>0</Angle>
413 </Entity>
414 <Entity Type="Wall">
415   <X>1397.8854415274473</X>
416   <Y>462.13495863941552</Y>
417   <RelX>-150.30071599049575</RelX>
418   <RelY>-114.171428571418</RelY>
419   <Angle>0</Angle>
420 </Entity>
421 <Entity Type="Wall">
422   <X>1247.5847255369515</X>
423   <Y>347.96353006799751</Y>
424   <RelX>-128.3054892601433</RelX>
425   <RelY>-143.06493506493496</RelY>
426   <Angle>0</Angle>
427 </Entity>
428 <Entity Type="Wall">
429   <X>1119.2792362768082</X>
430   <Y>204.89859500306255</Y>
431   <RelX>-190.62529832935547</RelX>
432   <RelY>-145.87012987012997</RelY>
433   <Angle>0</Angle>
434 </Entity>
435 <Entity Type="Wall">
436   <X>928.65393794745273</X>
437   <Y>59.028465132932581</Y>
438   <RelX>-234.615751789976</RelX>
439   <RelY>-84.155844155844363</RelY>
440   <Angle>0</Angle>
441 </Entity>
442 <Entity Type="Wall">
443   <X>694.03818615747673</X>
444   <Y>-25.127379022911782</Y>
445   <RelX>-340.92601431980904</RelX>
446   <RelY>-64.519480519480737</RelY>
447   <Angle>0</Angle>
448 </Entity>
449 <Entity Type="Wall">
450   <X>353.11217183766769</X>
451   <Y>-89.646859542392519</Y>
452   <RelX>-157.63245823389019</RelX>
453   <RelY>16.831168831168725</RelY>
454   <Angle>0</Angle>
455 </Entity>
456 <Entity Type="Wall">
457   <X>195.4797136037775</X>
458   <Y>-72.8156907112238</Y>
459   <RelX>-260.27684964200449</RelX>
```

```
460     <RelY>53.298701298701076</RelY>
461     <Angle>0</Angle>
462 </Entity>
463 <Entity Type="Wall">
464     <X>-64.797136038227</X>
465     <Y>-19.516989412522719</Y>
466     <RelX>-164.96420047732698</RelX>
467     <RelY>129.03896103896096</RelY>
468     <Angle>0</Angle>
469 </Entity>
470 <Entity Type="Wall">
471     <X>-229.76133651555398</X>
472     <Y>109.52197162643824</Y>
473     <RelX>-80.649164677804322</RelX>
474     <RelY>165.50649350649354</RelY>
475     <Angle>0</Angle>
476 </Entity>
477 <Entity Type="Wall">
478     <X>-310.4105011933583</X>
479     <Y>275.02846513293179</Y>
480     <RelX>7.3317422434367927</RelX>
481     <RelY>232.83116883116918</RelY>
482     <Angle>0</Angle>
483 </Entity>
484 <Entity Type="Wall">
485     <X>-303.07875894992151</X>
486     <Y>507.85963396410096</Y>
487     <RelX>120.97374701670645</RelX>
488     <RelY>207.584415584415</RelY>
489     <Angle>0</Angle>
490 </Entity>
491 <Entity Type="Wall">
492     <X>-182.10501193321505</X>
493     <Y>715.444049548516</Y>
494     <RelX>278.60620525059659</RelX>
495     <RelY>193.55844155844147</RelY>
496     <Angle>0</Angle>
497 </Entity>
498 <Entity Type="Wall">
499     <X>96.501193317381535</X>
500     <Y>909.00249110695745</Y>
501     <RelX>208.95465393794734</RelX>
502     <RelY>165.5064935064936</RelY>
503     <Angle>0</Angle>
504 </Entity>
505 <Entity Type="Wall">
506     <X>305.45584725532888</X>
507     <Y>1074.508984613451</Y>
508     <RelX>109.97613365155144</RelX>
509     <RelY>145.87012987012963</RelY>
510     <Angle>0</Angle>
511 </Entity>
512 <Entity Type="Wall">
```

```
513     <X>415.43198090688031</X>
514     <Y>1220.3791144835807</Y>
515     <RelX>-10.997613365154905</RelX>
516     <RelY>120.623376623377</RelY>
517     <Angle>0</Angle>
518 </Entity>
519 <Entity Type="Wall">
520     <X>404.43436754172541</X>
521     <Y>1341.0024911069577</Y>
522     <RelX>-153.96658711217174</RelX>
523     <RelY>151.48051948051943</RelY>
524     <Angle>0</Angle>
525 </Entity>
526 <Entity Type="Wall">
527     <X>250.46778042955367</X>
528     <Y>1492.4830105874771</Y>
529     <RelX>-245.61336515513125</RelX>
530     <RelY>193.55844155844147</RelY>
531     <Angle>0</Angle>
532 </Entity>
533 <Entity Type="Wall">
534     <X>4.854415274422422</X>
535     <Y>1686.0414521459186</Y>
536     <RelX>-245.61336515513119</RelX>
537     <RelY>235.63636363636351</RelY>
538     <Angle>0</Angle>
539 </Entity>
540 <Entity Type="Wall">
541     <X>-240.75894988070877</X>
542     <Y>1921.6778157822821</Y>
543     <RelX>-142.96897374701678</RelX>
544     <RelY>238.44155844155807</RelY>
545     <Angle>0</Angle>
546 </Entity>
547 <Entity Type="Wall">
548     <X>-383.72792362772554</X>
549     <Y>2160.11937422384</Y>
550     <RelX>1</RelX>
551     <RelY>266.49350649350663</RelY>
552     <Angle>0</Angle>
553 </Entity>
554 <Entity Type="Wall">
555     <X>-382.72792362772554</X>
556     <Y>2426.6128807173468</Y>
557     <RelX>113.64200477326949</RelX>
558     <RelY>143.06493506493507</RelY>
559     <Angle>0</Angle>
560 </Entity>
561 <Entity Type="Wall">
562     <X>-270.08591885445605</X>
563     <Y>2569.6778157822819</Y>
564     <RelX>263.94272076372317</RelX>
565     <RelY>103.79220779220759</RelY>
```

```
566     <Angle>0</ Angle>
567 </ Entity>
568 <Entity Type="Wall ">
569     <X> -6.1431980907328807</X>
570     <Y>2673.4700235744895</Y>
571     <RelX>421.57517899761331</RelX>
572     <RelY>14.025974025973483</RelY>
573     <Angle>0</ Angle>
574 </ Entity>
575 <Entity Type="Wall ">
576     <X>415.43198090688043</X>
577     <Y>2687.4959976004629</Y>
578     <RelX>344.59188544152744</RelX>
579     <RelY> -109.40259740259808</RelY>
580     <Angle>0</ Angle>
581 </ Entity>
582 <Entity Type="Wall ">
583     <X>760.02386634840786</X>
584     <Y>2578.0934001978649</Y>
585     <RelX>260.27684964200466</RelX>
586     <RelY> -154.28571428571377</RelY>
587     <Angle>0</ Angle>
588 </ Entity>
589 <Entity Type="Wall ">
590     <X>1020.3007159904125</X>
591     <Y>2423.8076859121511</Y>
592     <RelX>304.26730310262519</RelX>
593     <RelY> -291.74025974025972</RelY>
594     <Angle>0</ Angle>
595 </ Entity>
596 <Entity Type="Wall ">
597     <X>1324.5680190930377</X>
598     <Y>2132.0674261718914</Y>
599     <RelX>201.62291169450987</RelX>
600     <RelY> -100.9870129870128</RelY>
601     <Angle>0</ Angle>
602 </ Entity>
603 <Entity Type="Wall ">
604     <X>1526.1909307875476</X>
605     <Y>2031.0804131848786</Y>
606     <RelX>263.94272076372363</RelX>
607     <RelY>5.610389610389575</RelY>
608     <Angle>0</ Angle>
609 </ Entity>
610 <Entity Type="Wall ">
611     <X>1790.1336515512712</X>
612     <Y>2036.6908027952682</Y>
613     <RelX>546.2147971360373</RelX>
614     <RelY>356.25974025973892</RelY>
615     <Angle>0</ Angle>
616 </ Entity>
617 <Entity Type="Wall ">
618     <X>2336.3484486873085</X>
```

```
619     <Y>2392.9505430550071</Y>
620     <RelX>318.93078758949878</RelX>
621     <RelY>148.6753246753251</RelY>
622     <Angle>0</ Angle>
623 </ Entity>
624 <Entity Type="Wall ">
625     <X>2655.2792362768073</X>
626     <Y>2541.6258677303322</Y>
627     <RelX>421.57517899761342</RelX>
628     <RelY>84.155844155843624</RelY>
629     <Angle>0</ Angle>
630 </ Entity>
631 <Entity Type="Wall ">
632     <X>3076.8544152744207</X>
633     <Y>2625.7817118861758</Y>
634     <RelX>377.58472553699266</RelX>
635     <RelY>-22.4415584415583</RelY>
636     <Angle>0</ Angle>
637 </ Entity>
638 <Entity Type="Wall ">
639     <X>3454.4391408114134</X>
640     <Y>2603.3401534446175</Y>
641     <RelX>194.29116945107398</RelX>
642     <RelY>-58.909090909090082</RelY>
643     <Angle>0</ Angle>
644 </ Entity>
645 <Entity Type="Wall ">
646     <X>3648.7303102624874</X>
647     <Y>2544.4310625355274</Y>
648     <RelX>87.9809069212406</RelX>
649     <RelY>-218.80519480519524</RelY>
650     <Angle>0</ Angle>
651 </ Entity>
652 <Entity Type="Wall ">
653     <X>3736.711217183728</X>
654     <Y>2325.6258677303322</Y>
655     <RelX>-40.324582338901564</RelX>
656     <RelY>-277.71428571428532</RelY>
657     <Angle>0</ Angle>
658 </ Entity>
659 <Entity Type="Wall ">
660     <X>3696.3866348448264</X>
661     <Y>2047.9115820160469</Y>
662     <RelX>-109.97613365155121</RelX>
663     <RelY>-246.85714285714266</RelY>
664     <Angle>0</ Angle>
665 </ Entity>
666 <Entity Type="Wall ">
667     <X>3586.4105011932752</X>
668     <Y>1801.0544391589042</Y>
669     <RelX>-373.91885441527484</RelX>
670     <RelY>-336.62337662337677</RelY>
671     <Angle>0</ Angle>
```

```
672 </Entity>
673 <Entity Type="Wall">
674   <X>3212.4916467780004</X>
675   <Y>1464.4310625355274</Y>
676   <RelX>-362.92124105011862</RelX>
677   <RelY>-210.38961038961043</RelY>
678   <Angle>0</Angle>
679 </Entity>
680 <Entity Type="Wall">
681   <X>2849.5704057278817</X>
682   <Y>1254.041452145917</Y>
683   <RelX>-40.324582338902474</RelX>
684   <RelY>-117.81818181818198</RelY>
685   <Angle>0</Angle>
686 </Entity>
687 <Entity Type="Wall">
688   <X>2809.2458233889793</X>
689   <Y>1136.223270327735</Y>
690   <RelX>120.97374701670651</RelX>
691   <RelY>-224.41558441558402</RelY>
692   <Angle>0</Angle>
693 </Entity>
694 <Entity Type="Wall">
695   <X>2930.2195704056858</X>
696   <Y>911.807685912151</Y>
697   <RelX>359.2553699284008</RelX>
698   <RelY>-336.62337662337654</RelY>
699   <Angle>0</Angle>
700 </Entity>
701 <Entity Type="Wall">
702   <X>3289.4749403340866</X>
703   <Y>575.18430928877444</Y>
704   <RelX>168.63007159904555</RelX>
705   <RelY>-361.87012987012986</RelY>
706   <Angle>0</Angle>
707 </Entity>
708 <Entity Type="Wall">
709   <X>3458.1050119331321</X>
710   <Y>213.31417941864458</Y>
711   <RelX>7.3317422434374748</RelX>
712   <RelY>-272.10389610389575</RelY>
713   <Angle>0</Angle>
714 </Entity>
715 <Entity Type="Wall">
716   <X>3465.4367541765696</X>
717   <Y>-58.789716685251165</Y>
718   <RelX>-131.97136038186181</RelX>
719   <RelY>-204.77922077922074</RelY>
720   <Angle>0</Angle>
721 </Entity>
722 <Entity Type="Wall">
723   <X>3333.4653937947078</X>
724   <Y>-263.5689374644719</Y>
```

```
725     <RelX> -366.58711217183782</RelX>
726     <RelY> -123.4285714285715</RelY>
727     <Angle>0</ Angle>
728 </ Entity>
729 <Entity Type="Wall ">
730     <X>2966.87828162287</X>
731     <Y> -386.9975088930434</Y>
732     <RelX> -315.26491646778049</RelX>
733     <RelY>25.246753246753258</RelY>
734     <Angle>0</ Angle>
735 </ Entity>
736 <Entity Type="Wall ">
737     <X>2651.6133651550895</X>
738     <Y> -361.75075564629014</Y>
739     <RelX> -205.28878281622883</RelX>
740     <RelY>100.98701298701297</RelY>
741     <Angle>0</ Angle>
742 </ Entity>
743 <Entity Type="Wall ">
744     <X>2446.3245823388606</X>
745     <Y> -260.76374265927717</Y>
746     <RelX> -131.97136038186181</RelX>
747     <RelY>109.40259740259688</RelY>
748     <Angle>0</ Angle>
749 </ Entity>
750 <Entity Type="Wall ">
751     <X>2314.3532219569988</X>
752     <Y> -151.36114525668029</Y>
753     <RelX>0</RelX>
754     <RelY>0</RelY>
755     <Angle>0</ Angle>
756 </ Entity>
757 <Entity Type="Wall ">
758     <X>2314.3532219569988</X>
759     <Y> -151.36114525668029</Y>
760     <RelX> -87.9809069212406</RelX>
761     <RelY>260.88311688311938</RelY>
762     <Angle>0</ Angle>
763 </ Entity>
764 <Entity Type="Point">
765     <X>1619.212410501194</X>
766     <Y>1275.6414521459071</Y>
767     <RelX>0</RelX>
768     <RelY>0</RelY>
769     <Angle>0</ Angle>
770 </ Entity>
771 </Map>
```

### -1-3 Agent

#### CarMazeEnv.py



```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu Jun 14 00:28:13 2018
4
5  @author: ajdin
6  """
7  import numpy as np
8
9  import pygame
10 import visualnet
11 import agent_functions as funks
12 import novelty
13
14 from configs import MapBuilder
15 from envmanager import EnvManager
16 from controllers import ControlScheme as control
17 from modulation import ModulationScheme, ModulationSignal
18 from rangefinder import RangeFinderGroup as sensors
19 from net3 import NeuralNetwork3 as nn3
20 from net3 import Agent
21 from performance import Performance
22
23 #####
24
25 black = (0, 0, 0)
26 white = (255, 255, 255)
27 red = (255, 0, 0)
28 redcolor = (0xff,0,0)
29
30 def drawText(font, text, pos, display):
31     textSurf = font.render(text, False, black)
32     display.blit(textSurf, pos)
33
34 def drawCenter(pos, display, options):
35     x_scale = options['x_scale']
36     x_translate = options['x_translate']
37     position = tuple(int(x_scale*x+x_translate) for x in pos)
38     pygame.draw.circle(display, redcolor, position, 5, 2)
39
40
41 mapbuilder = MapBuilder()
42
43 class CarEnv:
44     OPTIONS = {'x_scale': 0.2184233207295375,
45               'yScale': 0.22767977831846448,
46               'x_translate': 133.81512733541817,
47               'yTranslate': 138.1115070345661,
48               'rotation': 180}
49
50     def __init__(self):
51         MAP_BUILDER = MapBuilder()
52         self.car = MAP_BUILDER.car
53         self.walls = MAP_BUILDER.walls
```

```

54     self.center_point = MAP_BUILDER.point
55     self.rangefinder_group = sensors(self.car)
56     self.performer = Performance(self.rangefinder_group, self.
        center_point)
57     # initialize environment variables
58     self.rangefinder_group.update(self.walls)
59     self.performance = 0
60     self.ticks = 0
61     # render variables
62     self.options = self.OPTIONS
63     self.game_display = None
64     self.font = None
65     self.clock = None
66
67     def step(self, action):
68         if len(action) == 1:
69             u_theta, u_speed = action[0], 0.5
70             self.car.reset_speed = False
71         elif len(action) == 2:
72             self.car.reset_speed = True
73             u_theta, u_speed = action
74         else:
75             self.car.reset_speed = False
76             u_theta, u_speed = 0.5, 0.5
77         self.car.angle += float(u_theta)*self.car.ROTATE_RANGE - self.car.
            ROTATE_SPEED
78         self.car.speed_x += float(u_speed)*self.car.ACCELERATION_RANGE - self
            .car.ACCELERATION
79         self.car.update(self.walls)
80         self.rangefinder_group.update(self.walls)
81         self.performance += self.performer.normlaps() / 360
82         self.ticks += 1
83         info = {'position': self.car.center,
84                'orientation': self.car.angle,
85                'can move': self.car.can_move,
86                'time tick': self.ticks}
87         return self.rangefinder_group.activations, np.abs(self.performance),
            False, info
88
89     def reset(self):
90         MAP_BUILDER = MapBuilder()
91         self.car = MAP_BUILDER.car
92         self.rangefinder_group = sensors(self.car)
93         self.rangefinder_group.update(self.walls)
94         self.ticks = 0
95         return self.rangefinder_group.activations
96
97     def draw_environment(self):
98         for wall in self.walls:
99             wall.draw(self.game_display, self.options)
100        self.rangefinder_group.draw(self.game_display, self.options)
101        self.car.draw(self.game_display, self.options)
102

```

```

103 def render(self, size=(1600, 1440), fps=45, drawparams={'functions':
    [], 'arguments': []}):
104     # Draw the display and objects
105     if self.game_display is None:
106         pygame.init()
107         pygame.font.init()
108         self.game_display = pygame.display.set_mode(size)
109         pygame.display.set_caption('2D Autonomous Car Driving Simulation')
110     if self.font is None:
111         self.font = pygame.font.SysFont('Arial', 30)
112     if self.clock is None:
113         self.clock = pygame.time.Clock()
114     self.game_display.fill(white)
115     self.draw_environment()
116     # Draw the texts
117     drawText(self.font, 'tick: '+str(self.ticks), (10, 10), self.
        game_display)
118     drawText(self.font, 'FPS: '+ "%.1f" % self.clock.get_fps(), (200, 10),
        self.game_display)
119     drawText(self.font, 'Performance: '+ "%.2f" % np.abs(self.performance
        ), (400, 10), self.game_display)
120     drawText(self.font, 'Forward velocity: ' + "%.1f" % np.linalg.norm(
        list(self.car.speed)),
        (700, 10), self.game_display)
121     drawCenter(self.center_point, self.game_display, self.options)
122     if self.game_display is not None:
123         for function, args in zip(drawparams['functions'], drawparams['
            arguments']):
124             function(*args)
125     # Render the next frame
126     pygame.display.flip()
127     self.clock.tick(fps)
128
129
130 def close(self):
131     pygame.quit()
132
133
134 class CarMazeEnv:
135     # MAP_BUILDER = MapBuilder()
136     DISPLAY_WIDTH = 1000
137     DISPLAY_HEIGHT = 800
138     hebbian_hyper_parameters = {"num_nodes": [11, 1],
139                                "activation_functions": [funs.linu, funs.logistic],
140                                "training_functions": [funs.hebbian1],
141                                "learning_rates": [1.0],
142                                "output_noise": [0, 0.1],
143                                "weight_noise": [0.1],
144                                "buffers": [None],
145                                "tag": 'Pure Hebbian'}
146
147     raahn_hyper_parameters = {"num_nodes": [11, 5, 1],
148                              "activation_functions": [funs.linu, funs.logistic,
149                                                       funs.logistic],

```

```

149         "training_functions": [funs.autoencoder ,funs.
           hebbian1],
150         "learning_rates": [0.1, 1.0],
151         "output_noise": [0, 0, 0.1],
152         "weight_noise": [0, 0.1],
153         "buffers": [novelty.NoveltyBuffer(500), None],
154         "tag": 'novelty-Raahn'}
155
156 grewards_hyper_params = {"num_nodes": [11, 5, 2],
157                          "activation_functions": [funs.linu, funs.logistic,
           funs.logistic],
158                          "training_functions": [funs.autoencoder ,funs.
           hebbian1],
159                          "learning_rates": [0.1, 1.0],
160                          "output_noise": [0, 0, 0.1],
161                          "weight_noise": [0, 0.1],
162                          "buffers": [novelty.NoveltyBuffer(500), None],
163                          "tag": 'Hebbian with global rewards'}
164
165 def __init__(self):
166     # initialize map objects
167     self.environment = CarEnv() # --> new!
168     # initialize game variables
169     self.resolution = (self.DISPLAY_WIDTH, self.DISPLAY_HEIGHT)
170     self.game_ext = False
171     self.episode = 1
172     # initialize algorithm settings
173     self.agent = Agent(self.raahn_hyper_parameters)
174     self.activations = None
175     self.modulation_scheme = ModulationScheme(self.environment.
           rangefinder_group, self.environment.walls)
176     # Initialize visual settings
177     self.options = self.environment.OPTIONS
178     self.bounds = None
179     self.calculate_bounds()
180     self.network_visualizer = visualnet.NetworkVisualizer(self.
           network_parameters,
181                                                         [self.bounds
           [0], 0,
           1600, self.
           bounds [1]])
182
183     self.crashes = []
184     self.performances = []
185     self.avg_objectives = []
186     pygame.init()
187     pygame.font.init()
188
189 def step(self):
190     # num_crashes = self.env.car.num_crashes
191     # control.range_finder_control(self.network, self.rangefinder_group,
           angle_control=True)
192     observation = self.environment.rangefinder_group.activations

```

```

193     # self.network.add_experience(observation)
194     self.activations = self.agent.propagate(observation)
195     # action = tuple(self.network.activations[-1])
196     action = tuple(self.activations[-1])
197     self.environment.step(action)
198     self.modulation_scheme.wall_avoidance()
199     modulation = self.modulation_scheme.modulations[0]
200     # self.network.modulation_signal.modulations[0] = modulation
201     # self.network.train()
202     self.agent.train(modulation, self.activations)
203
204 def convergence_test(self):
205     if self.episode == 15:
206         self.game_ext = True
207         observation = self.environment.rangefinder_group.activations
208         self.activations = self.agent.propagate(observation)
209         action = tuple(self.activations[-1])
210         self.environment.step(action)
211         self.modulation_scheme.wall_avoidance2()
212         global_reward = self.environment.performer.objective_func()
213         self.avg_objectives.append(global_reward)
214         modulation = self.modulation_scheme.modulations[0]
215         if self.environment.ticks == 1200:
216             print(np.mean(self.avg_objectives))
217             self.avg_objectives = []
218             del self.environment
219             self.environment = CarEnv()
220             self.modulation_scheme = ModulationScheme(self.environment.
                rangefinder_group, self.environment.walls)
221             self.episode += 1
222         else:
223             self.agent.train(modulation, self.activations)
224
225
226 def crashcounting_test(self):
227     if self.episode >= 50 or self.environment.ticks == 10000:
228         self.game_ext = True
229         self.performances.append(self.environment.performance)
230         self.crashes.append(self.environment.car.num_crashes)
231         print('Episode: ', self.episode, ' Num crashes: ', self.crashes
            [-1],
            ' Perf.: ', abs(self.performances[-1]), '\n')
232     # num_crashes = self.environment.car.num_crashes
233     observation = self.environment.rangefinder_group.activations
234     self.activations = self.agent.propagate(observation)
235     action = tuple(self.activations[-1])
236     self.environment.step(action)
237     self.modulation_scheme.wall_avoidance()
238     global_reward = self.environment.performer.objective_func()
239     self.avg_objectives.append(global_reward)
240     modulation = self.modulation_scheme.modulations[0]
241     if self.environment.car.num_crashes > 0 and self.environment.ticks >=
        2500:

```

```

243     self.crashes.append(self.environment.car.num_crashes)
244     self.performances.append(self.environment.performance)
245     print('Episode: ', self.episode, ' Num crashes: ', self.crashes
246           [-1],
247           ' Perf.: ', abs(self.performances[-1]), '\n')
248     del self.environment
249     self.environment = CarEnv()
250     self.modulation_scheme = ModulationScheme(self.environment.
251         rangefinder_group, self.environment.walls)
252     self.episode += 1
253 else:
254     self.agent.train(modulation, self.activations)
255
256 # self.Reward = self.envmanager.Reward
257 # self.avg_objectives.append(self.Reward)
258 # if self.car.num_crashes > num_crashes and self.ticks >= 1000:
259 #     num_crashes = self.car.num_crashes
260 #     self.performances.append(self.performance)
261 #     self.crashes.append(num_crashes)
262 #     self.reset()
263 #     self.performance = 0
264 #     print(self.crashes)
265 # elif self.episode >= 100 or self.ticks == 11000:
266 #     self.game_ext = True
267 #     self.performances.append(self.performance)
268 #     self.crashes.append(self.car.num_crashes)
269 # else:
270 #     self.envmanager.handle_env(training=True)
271 #     if self.ticks == 1000:
272 #         self.performances.append(self.performance)
273 #         self.crashes.append((self.episode, num_crashes))
274
275 def reset(self, keep_performance=True):
276     MAP_BUILDER = MapBuilder()
277     if keep_performance:
278         performances = self.performances
279     self.car = MAP_BUILDER.car
280     self.rangefinder_group = sensors(self.car)
281     self.rangefinder_group.update(self.walls)
282     self.envmanager = EnvManager(self)
283     self.performances = performances
284     self.episode += 1
285     self.ticks = 0
286     return self.rangefinder_group.activations
287
288 def render(self, size=(None)):
289     drawparams = {'functions': [], 'arguments': []}
290     if self.environment.game_display is not None:
291         text_args1 = (self.environment.font, 'Modulation: ' + '%.3f' % self
292             .modulation_scheme.modulations[0],
293             (400, 40), self.environment.game_display)
294         text_args2 = (self.environment.font, 'Episode: ' + str(self.episode
295             ),

```

```

292         (10,40), self.environment.game_display)
293     text_args3 = (self.environment.font, 'Global Reward: ' + '%.2f' %
                self.avg_objectives[-1],
294                 (400, 70), self.environment.game_display)
295     visual_args = (self.network_parameters, self.environment.
                game_display)
296     drawparams = {'functions': [drawText, drawText, drawText, self.
                network_visualizer.visualize],
                'arguments': [text_args1, text_args2, text_args3,
                visual_args]}
297     # self.network_visualizer.visualize(self.network_parameters, self.
                environment.game_display, self.environment.font)
298     self.environment.render(drawparams=drawparams)
299
300
301     def close(self):
302         self.game_ext = True
303         self.environment.close()
304
305     def calculate_bounds(self):
306         if not self.bounds:
307             xbound = 0
308             ybound = 0
309             for wall in self.environment.walls:
310                 maxboundWallx = max(wall.startPoint[0], wall.endPoint[0])
311                 maxboundWally = max(wall.startPoint[1], wall.endPoint[1])
312                 xbound = max(xbound, maxboundWallx)
313                 ybound = max(ybound, maxboundWally)
314             xbound = self.options['x_scale']*xbound+self.options['x_translate']
315             ybound = self.options['x_scale']*ybound+self.options['x_translate']
316             self.bounds = (xbound, ybound)
317
318     @property
319     def network_parameters(self):
320         network_parameters = {'num_nodes': self.agent.num_neurons,
                'weights': self.agent.weights,
321                             'activations': self.activations,
322                             'weight_cap': self.agent.weight_cap}
323
324         return network_parameters
325
326
327     def visualize_network(self, display, options):
328         self.network_visualizer.visualize(display)
329
330     """ Experiments with supervised Hebbian learning
331     #from collections import deque
332     #
333     #if __name__ == "__main__":
334     #    supervised_params = {"num_nodes": [11, 15, 15, 5],
335     #                        "activation_functions": [funs.linu, np.sign, np.sign,
336     #                                                np.sign],
337     #                        "training_functions": [funs.hebbian3, funs.hebbian3,
338     #                                                funs.hebbian3],
339     #                        "learning_rates": [1.0, 1.0, 1.0],

```

```

338 #             "output_noise": [0, 0, 0, 0],
339 #             "weight_noise": [0, 0, 0],
340 #             "buffers": [None, None, None]}
341 # agent = Agent(supervised_params)
342 # shape = (supervised_params["num_nodes"][-1], supervised_params["
num_nodes"][0])
343 # weight_des = np.random.uniform(-1,1, shape)
344 # errors = deque(maxlen=100)
345 # for i in range(10000):
346 #     activations = agent.propagate(np.random.choice([-1,1], (11,)))
347 #     output_des = weight_des @ activations[0]
348 #     output = activations[-1]
349 #     distance1 = np.abs(output-output_des)
350 #     modulation1 = np.multiply(distance1, output)
351 #     distance2 = np.linalg.norm(output-output_des)
352 #     modulation2 = distance2*output
353 #     agent.train(modulation2
354 #                 , activations)
355 #     errors.append(distance1)
356 #     if i%20==0:
357 #         print('error:', np.mean(errors))

```

### agent\_functions.py

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Tue Mar  6 15:22:56 2018
4
5 @author: ajdin
6 """
7 import numpy as np
8
9 # Activation functions
10 def linu(x):
11     #return arg
12     return x
13
14 def step(x):
15     #step function
16     return 1 * (x > 0)
17
18 def relu(x):
19     #rectified linear unit
20     return np.maximum(x, 0)
21
22 def logistic(x):
23     return 1.0 / (1.0 + np.exp(-x))
24
25 def logistic_derivative(x):
26     return x * (1.0 - x)
27
28 #tangent hyperbolic
29 tanh = np.tanh

```



```

30
31 #Exponentially Varied Weight Adjustment
32 def evwa(self, learning_rate, reward):
33     rho = 0.02
34     adjustment = []
35     for i in range(self.num_layers):
36         adjustment.append(learning_rate*(rho**reward))
37     return adjustment
38
39 #Linearly Varied Weight Adjustment
40 def lvwa(self, learning_rate, reward):
41     adjustment = []
42     for i in range(self.num_layers):
43         adjustment.append(learning_rate*(1-reward))
44     return adjustment
45
46 # Compute node values in layers
47 def layer_output(activation_fun, parameters, input_nodes, noise_magtd
    =0.0):
48     output_nodes = activation_fun(np.dot(parameters, input_nodes))
49     if noise_magtd:
50         output_nodes += np.random.uniform(-noise_magtd, noise_magtd, np.shape
            (output_nodes))
51     return output_nodes
52
53 # train weights
54 def train(mod_vector, input_samples, output_samples, learning_rate):
55     for i in range(len(mod_vector)):
56         modulation = mod_vector[i]
57         input_sample = np.reshape(input_samples[i], (-1,1))
58         output_sample = np.reshape(output_samples[i], (-1,1))
59         plasticity = output_sample @ input_sample.T
60         if i == 0:
61             weight_delta = learning_rate*modulation*plasticity
62         else:
63             weight_delta += learning_rate*modulation*plasticity
64     return weight_delta
65
66 def hebbian1(layer_parameters, presynaptic, postsynaptic):
67     lr = layer_parameters["learning_rate"]
68     modulation = layer_parameters["modulation"]
69     noise_magtd = layer_parameters["weight_noise"]
70     presynaptic = np.array(presynaptic)
71     postsynaptic = np.array(postsynaptic)*2-1
72     if presynaptic.ndim == 1:
73         presynaptic = np.reshape(presynaptic, (-1,1))
74     if postsynaptic.ndim == 1:
75         postsynaptic = np.reshape(postsynaptic, (-1,1))
76     update = lr*modulation*(postsynaptic @ presynaptic.T)
77     if noise_magtd:
78         update += np.random.uniform(-noise_magtd, noise_magtd, np.shape(
            update))
79     return update

```

```
80
81 def hebbian2(layer_parameters, presynaptic, postsynaptic):
82     lr = layer_parameters["learning_rate"]
83     modulation = layer_parameters["modulation"]
84     noise_magtd = layer_parameters["weight_noise"]
85     presynaptic = np.array(presynaptic)
86     postsynaptic = np.array(postsynaptic)
87     if presynaptic.ndim == 1:
88         presynaptic = np.reshape(presynaptic, (-1,1))
89     if postsynaptic.ndim == 1:
90         postsynaptic = np.reshape(postsynaptic, (-1,1))
91     update = lr*modulation*(postsynaptic @ presynaptic.T)
92     if noise_magtd:
93         update += np.random.uniform(-noise_magtd, noise_magtd, np.shape(
94             update))
95     return update
96
97 def hebbian3(layer_parameters, presynaptic, postsynaptic):
98     """
99     Supervised-type hebbian algorithm.
100    Modulation equals the 2 norm of the discrepancy between actual and
101    desired output,
102    times the desired output.
103    """
104    lr = layer_parameters["learning_rate"]
105    modulation = layer_parameters["modulation"]
106    noise_magtd = layer_parameters["weight_noise"]
107    presynaptic = np.array(presynaptic)
108    postsynaptic = np.array(postsynaptic)
109    if presynaptic.ndim == 1:
110        presynaptic = np.reshape(presynaptic, (-1,1))
111    if postsynaptic.ndim == 1:
112        postsynaptic = np.reshape(postsynaptic, (-1,1))
113    modulation.shape = postsynaptic.shape
114    update = lr*(modulation @ presynaptic.T)
115    if noise_magtd:
116        update += np.random.uniform(-noise_magtd, noise_magtd, np.shape(
117            update))
118    return update
119
120 def autoencoder(layer_parameters, presynaptic, postsynaptic):
121     lr = layer_parameters["learning_rate"]
122     weights = layer_parameters["weights"]
123     noise_magtd = layer_parameters["weight_noise"]
124     presynaptic = np.array(presynaptic)
125     postsynaptic = np.array(postsynaptic)
126     if presynaptic.ndim == 1:
127         presynaptic = np.reshape(presynaptic, (-1,1))
128     if postsynaptic.ndim == 1:
129         postsynaptic = np.reshape(postsynaptic, (-1,1))
130     reconstruction = logistic(weights.T @ postsynaptic)
131     error = presynaptic - reconstruction
```

```

130 deltas = np.multiply(error, logistic_derivative(reconstruction))
131 back_prop_deltas = np.multiply(logistic_derivative(postsynaptic), (
    weights @ deltas))
132 error_weight_delta = lr*(postsynaptic @ deltas.T)
133 backprop_weight_delta = lr*(back_prop_deltas @ presynaptic.T)
134 weight_delta = error_weight_delta + backprop_weight_delta
135 if noise_magtd:
136     weight_delta += np.random.uniform(-noise_magtd, noise_magtd, np.shape
        (weight_delta))
137 return weight_delta
138
139 def select_several(layer_parameters):
140     novelty_buffer = layer_parameters['buffers']
141     # error = 0.0
142     samples = [] #linkedList
143     experiences = []
144     for occupant in iter(novelty_buffer):
145         samples.append(occupant.experience)
146     sample_count = min(20, len(novelty_buffer))
147     for i in range(sample_count):
148         index = np.random.randint(len(samples))
149         sample = samples.pop(index)
150         experiences.append(np.array(sample))
151     return experiences

```

### modulation.py

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Aug 19 18:40:16 2018
4
5 @author: ajdin
6 """
7 import numpy as np
8
9 import buffer
10
11 from rangefinder import RangeFinder
12 from lines import get_dist
13 from collections import deque
14
15
16 class ModulationSignal:
17     NO_MODULATION = 0.0
18
19     def __init__(self):
20         self.modulations = []
21
22     def add_signal(self, *args):
23         if len(args) == 0:
24             self.modulations.append(self.NO_MODULATION)
25         elif len(args) == 1:
26             self.modulations.append(args[0])

```

```

27     return len(self.modulations) - 1
28
29
30 class ModulationScheme:
31     MODULATION_STRENGTH = 1.0
32     MODULATION_RESET = 0.0
33     MODULATION_NOT_RESET = -1.0
34     PERPENDICULAR = 90.0
35     SCHEME_STRINGS = ['WallAvoidance', 'Acceleration', 'NaiveWallAvoidance',
36                       , 'NaiveAcceleration']
37
38 def __init__(self, rangefinder_group, walls):
39     self.rangefinder_group = rangefinder_group
40     self.previous_fitness = None
41     self.N_size = 5
42     self.reward_history = deque(maxlen=self.N_size)
43     self.Reward = 0
44     self.acceleration
45     self.walls = walls
46     self.view_distance = 400.0
47     self.viewline = RangeFinder(self.rangefinder_group.car, self.
48                                 view_distance)
49     self.modulations = [0.0]*len(self.SCHEME_STRINGS)
50     # Initializes for no former experience
51     self.last_angle_between = None
52     self.last_wall_in_range = None
53     self.compare_wall = None
54     self.last_nearest_dist = None
55
56 def wall_avoidance2(self):
57     """
58     Minimize the average sensor values. Total reward is Negative mean of
59     sensor outputs.
60     Change in gradient (objective increase) equals the negative mean
61     sensor values of current time
62     tick, minus the negative mean sensor values of previous time tick.
63     """
64     params = (self.rangefinder_group.activations, self.rangefinder_group.
65              car.speed)
66     self.reward_history.append(self.local_reward(*params))
67     self.modulations[0] = 0.0
68
69     if len(self.reward_history)==self.reward_history.maxlen:
70         change_in_objective = self.reward_history[-1] - self.Reward
71         self.modulations[0] = np.tanh(10*change_in_objective)
72         self.modulations[0] = change_in_objective
73         self.Reward = self.reward_history[0]
74
75 def local_reward(self, *params):
76     activations, speed = params
77     avg_distance = np.mean(1 - activations[5])
78     norm_speed = np.linalg.norm(speed)
79     local_reward = (1.3*avg_distance + 0.3*norm_speed)/self.N_size

```

```
75     return local_reward
76
77 def modulate(self):
78     """
79     Modulates all schemes ever defined, so that a combination of schemes
80     may be integrated.
81     """
82     self.viewline.update_position()
83     walls_in_bounds = self.viewline.entities_in_bounds(self.walls)
84     last_angle = None
85     nearest_wall = None
86     nearest_dist = self.view_distance
87     # Get nearest wall or None
88     for wall in walls_in_bounds:
89         intersections = self.viewline.intersects(wall)
90         if len(intersections) > 0:
91             dist = get_dist(self.rangefinder_group.car.center, intersections
92                             [0])
93             if dist < nearest_dist:
94                 nearest_dist = dist
95                 nearest_wall = wall
96
97     # No previous wall yields nothing to modulate
98     if not self.last_wall_in_range:
99         self.modulations[0] = 0.0
100        self.modulations[1] = 0.0
101    else:
102        angle_between = self.viewline.angle_between(self.last_wall_in_range
103            )
104        last_intersection = self.viewline.intersects(self.
105            last_wall_in_range)[0]
106        distance = get_dist(self.rangefinder_group.car.center,
107            last_intersection)
108        delta = angle_between - self.last_angle_between
109        gamma = distance - self.last_nearest_dist
110        if angle_between > self.PERPENDICULAR:
111            modulation1 = self.MODULATION_STRENGTH*delta / self.
112                rangefinder_group.car.ROTATE_SPEED
113        else:
114            modulation1 = -self.MODULATION_STRENGTH*delta / self.
115                rangefinder_group.car.ROTATE_SPEED
116        modulation2 = self.MODULATION_STRENGTH*gamma / self.
117            rangefinder_group.car.ACCELERATION
118        self.modulations[0] = modulation1
119        self.modulations[1] = modulation2
120
121    # No current wall yields no angle to store
122    if nearest_wall:
123        last_angle = self.viewline.angle_between(nearest_wall)
124    # Store last wall and angle in any case (may be None)
125    self.last_angle_between = last_angle
126    self.last_nearest_dist = nearest_dist
127    self.last_wall_in_range = nearest_wall
```

```
120
121 def wall_avoidance(self):
122     # print('computing wall avoidance modulation: ', self.modulations[0])
123     self.viewline.update_position()
124     walls_in_bounds = self.viewline.entities_in_bounds(self.walls)
125     compare_wall = None
126     nearest_dist = self.viewline.default_length
127     # Get the nearest wall in the view_distance if any
128     for wall in walls_in_bounds:
129         intersections = self.viewline.intersects(wall)
130         if len(intersections) > 0:
131             dist = get_dist(self.rangefinder_group.car.center, intersections
132                             [0])
133             if dist < nearest_dist:
134                 nearest_dist = dist
135                 compare_wall = wall
136     self.compare_wall = compare_wall
137     # The angle to use for modulation. Should never be zero when the
138     # angle delta is calculated.
139     # If it is, then there must be a bug.
140     angle_between = 0.0
141     new_last_angle = self.MODULATION_NOT_RESET
142
143     # If there is no nearest wall
144     if not compare_wall:
145         # If there is no previous wall, set the modulation to zero and
146         # reset the last angle
147         if not self.last_wall_in_range:
148             if self.last_angle_between != self.MODULATION_RESET:
149                 self.last_angle_between = self.MODULATION_RESET
150                 self.modulations[0] = 0.0
151             # nothing to modulate, and nothing to save, Don't Continue
152             return
153         # Just left a wall
154     else:
155         angle_between = self.viewline.angle_between(self.
156             last_wall_in_range)
157
158     # If the wall has changed
159     elif compare_wall != self.last_wall_in_range:
160         # There was a last wall that is different from the current wall
161         if self.last_wall_in_range:
162             angle_between = self.viewline.angle_between(self.
163                 last_wall_in_range)
164             new_last_angle = self.viewline.angle_between(compare_wall)
165         # It is the first time any wall was hit, don't continue
166         # Save the angle between last and current wall
167         else:
168             angle_between = self.viewline.angle_between(compare_wall)
169             self.last_angle_between = angle_between
170             self.last_wall_in_range = compare_wall
171         return
```

```

168     # The usual case, the last wall is equal to the current wall
169     else:
170         angle_between = self.viewline.angle_between(compare_wall)
171         new_last_angle = angle_between
172         delta = angle_between - self.last_angle_between
173         modulation = self.MODULATION_STRENGTH
174
175         if angle_between > self.PERPENDICULAR:
176             modulation *= delta / self.rangefinder_group.car.ROTATE_SPEED
177         else:
178             modulation *= -delta / self.rangefinder_group.car.ROTATE_SPEED
179         self.modulations[0] = modulation
180         self.last_angle_between = new_last_angle
181         self.last_wall_in_range = compare_wall
182
183     def performance_increase(self):
184         pass
185
186     def naive_wall_avoidance(self):
187         print(self.modulations)
188         self.modulations[0] = self.last_forward_activation - self.
189             rangefinder_group.activations[5]
190         self.last_forward_activation = self.rangefinder_group.activations[5]
191         pass
192
193     def naive_acceleration(self):
194         # if self.last_wall_in_range:
195         # intersections = self.viewline.intersects(self.last_wall_in_range)
196         self.modulations[1] = 1-2*abs(self.modulations[0])
197         # print('computing acceleration modulation: ', self.modulations[1])
198         return
199
200     def acceleration(self):
201         # If there is no nearest wall in range
202         if not self.last_wall_in_range:
203             if self.last_nearest_dist:
204                 self.last_nearest_dist = None
205                 self.modulations[1] = 0.0
206             return
207         else:
208             intersection = self.viewline.intersects(self.last_wall_in_range)[0]
209             nearest_dist = get_dist(self.rangefinder_group.car.center,
210                 intersection)
211
212         pass
213
214     SCHEMES = [wall_avoidance, acceleration, naive_wall_avoidance,
215         naive_acceleration]
216
217     @staticmethod
218     def get_scheme_from_string(scheme_string):
219         for i in range(len(ModulationScheme.SCHEME_STRINGS)):
220             if scheme_string == ModulationScheme.SCHEME_STRINGS[i]:

```

```

218         return i
219     return -1
220
221     @staticmethod
222     def getSchemeFunction(scheme):
223         if scheme >= 0 and scheme < len(ModulationScheme.SCHEMES):
224             return ModulationScheme.SCHEMES[scheme]
225         else:
226             return None
227
228     def reset(self):
229         self.last_angle_between = self.MODULATION_RESET
230         self.last_wall_in_range = None

```

### net3.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Fri Oct 12 12:52:38 2018
4
5  @author: ajdin
6  """
7
8  import numpy as np
9  import configs
10 import novelty
11
12 # from RaahnXmlConfig import NeuralNetworkConfig, LayerConfig
13 from collections import deque
14
15 from functions import Activation
16 from training import TrainingMethod
17 from modulation import ModulationScheme, ModulationSignal
18 from buffer import Buffer
19 from eligibility import Eligibility
20
21 class NeuralNetwork3:
22     """
23     Creates a NeuralNetwork that strictly uses three layers at max:
24     inputlayer, hiddenlayer,
25     and outputlayer. hiddenlayer is optional.
26     parameters (optional keyword arguments):
27     inputcount: number of neurons in the input layer
28     hiddencount: number of neurons in the hidden layer.
29     outputcount: number of neurons in the output layer
30
31     attributes:
32     activations: list of neuron values for each layer
33     layers: list of NetworkLayers containing the neuron connections
34     output_noise_mag: magnitude of random noise to be applied after
35     computing activation
36     weight_noise_mag: magnitude of random noise to be applied after
37     updating weights

```



```

35     weight_cap: weights will not grow beyond upper limit weight_cap and
        lower limit -weight_cap.
36     # use_novelty: boolean determining whether a novelty buffer will be
        used.
37
38     methods:
39         set_up: initializes the NeuralNetwork3 layers based on parameters
        arguments.
40         configure: configures hyper-parameters of the network, based on
        RaahnXmlConfig.
41         init_neurons: returns numpy array of neurons with values 0.0
42         add_experience: adds a sample experience to the network's input
        layer
43         propagate_signal: maps the input activation to the output
        activation
44         train: trains all connections that are present within the network
45     """
46     WEIGHT_SCALE = 6.0
47     ACTIVATION = Activation.logistic
48     ACTIVATION_DERIVATIVE = Activation.logistic_derivative
49     TRAINING_METHODS = {'Hebbian': TrainingMethod.hebbian_learning,
50                         'HebbianHistory': TrainingMethod.
51                             hebbian_history_learning,
52                         'HebbianEligibility': TrainingMethod.
53                             hebbian_eligibility,
54                         'HebbianLongterm': TrainingMethod.
55                             hebbian_longterm,
56                         'HebbianEpisodic': TrainingMethod.
57                             hebbian_episodic,
58                         'NoTraining': TrainingMethod.no_training,
59                         'Autoencoder': TrainingMethod.sparse_autoencoder,
60                         'LinAutoencoder': TrainingMethod.
61                             linear_autoencoder}
62     MOD_SCHEMES = {'WallAvoidance': ModulationScheme.wall_avoidance,
63                  'WallAvoidance2': ModulationScheme.wall_avoidance2}
64
65     def __init__(self, **kwargs):
66         self.activations = []
67         self.layers = []
68         # set up the layers if specified
69         self.set_up(**kwargs)
70         # configure hyper-parameters of the network.
71         self.configure()
72         # ready up for training
73         self.compile_network()
74
75     @classmethod
76     def default(cls, filename=None):
77         if filename:
78             config = configs.Net3Config(filename)
79         else:
80             config = configs.Net3Config()
81         net3 = NeuralNetwork3()

```

```

77     net3.activations.append(net3.init_neurons(config.input_count))
78
79     return config
80
81
82     def set_up(self, **kwargs):
83         inputcount = kwargs.get('inputcount')
84         hiddencount = kwargs.get('hiddencount')
85         outputcount = kwargs.get('outputcount')
86         if self.valid_neuron_count(inputcount) and self.valid_neuron_count(
            outputcount):
87             self.activations.append(self.init_neurons(inputcount))
88             if self.valid_neuron_count(hiddencount):
89                 self.activations.append(self.init_neurons(hiddencount))
90                 add_layer = NetworkLayer(self, hiddencount)
91                 self.layers.append(add_layer)
92             self.activations.append(self.init_neurons(outputcount))
93             add_layer = NetworkLayer(self, outputcount)
94             self.layers.append(add_layer)
95
96     def valid_neuron_count(self, count):
97         if isinstance(count, int):
98             if count > 0:
99                 return True
100         return False
101
102     def configure(self, output_noise_mag=0.1, weight_noise_mag=0.1):
103         self.output_noise_mag = output_noise_mag
104         self.weight_noise_mag = weight_noise_mag
105         self.weight_cap = 10.0
106         self.modulation_signal = ModulationSignal()
107         self.modulation_signal.add_signal()
108         #activation functions
109         self.activation_func = NeuralNetwork3.ACTIVATION
110         self.activation_deriv = NeuralNetwork3.ACTIVATION_DERIVATIVE
111         # specify buffer settings if any
112         self.initlen = 1
113         self.maxlen = 200
114         self.growth = 1.0
115
116     def init_neurons(self, count):
117         return np.ones(count)*0.0
118
119     def add_experience(self, experience):
120         # add experience to the very first layer, and propagate it
121         self.activations[0] = np.array(experience)
122         self.propagate_signal()
123         # add experience to any buffers if specified
124         for layer in self.layers:
125             input_sample = np.array(layer.network.activations[layer.
                current_layer])
126             if not layer.usenovelty:

```

```

127         output_sample = np.array(layer.network.activations[layer.
128             current_layer + 1])
129         sample = (input_sample, output_sample)
130         layer.history_buffer.append(sample)
131         # layer.estim_window.add_sample(sample)
132         if layer.usenovelty:
133             new_occupant = novelty.NoveltyOccupant(experience=input_sample)
134             new_distances = layer.novelty_buffer.compute_new_distances(
135                 new_occupant)
136             if len(layer.novelty_buffer) == layer.history_size:
137                 least_novel = layer.novelty_buffer[0]
138                 if new_occupant.novelty_score > least_novel.novelty_score:
139                     layer.novelty_buffer.remove_novelty(least_novel)
140                     layer.novelty_buffer.add_novelty(new_occupant,
141                         new_distances)
142             else:
143                 layer.novelty_buffer.add_novelty(new_occupant, new_distances)
144
145     def propagate_signal(self):
146         for layer in self.layers:
147             layer.propagate_signal()
148
149     def compile_network(self, training_method='HebbianLongterm',
150         modulation_scheme='WallAvoidance2'):
151         for layer in self.layers:
152             layer.training_method = NeuralNetwork3.TRAINING_METHODS[
153                 training_method]
154             layer.modulation_scheme = NeuralNetwork3.MOD_SCHEMES[
155                 modulation_scheme]
156             # layer.history_buffer = Buffer(self.initlen, maxlen=self.maxlen,
157                 growth_fact=self.growth)
158             layer.learning_rate = 1.0
159             if len(self.layers) > 1:
160                 self.layers[0].usenovelty = True
161                 self.layers[0].history_size = 500
162                 self.layers[0].samples_per_tick = 20
163                 self.layers[0].training_method = NeuralNetwork3.TRAINING_METHODS[
164                     'Autoencoder']
165                 self.layers[0].learning_rate = 0.1
166                 self.layers[0].novelty_buffer = novelty.NoveltyBuffer(self.layers
167                     [0].history_size)
168                 self.layers[0].errorbuffer = deque(maxlen=self.layers[0].
169                     history_size)
170                 self.layers[1].history_buffer = deque(maxlen=self.initlen)
171
172     def train(self):
173         for layer in self.layers:
174             layer.train()
175
176 class NetworkLayer:
177     """
178     Creates a fully connected NetworkLayer in a given NeuralNetwork3

```

```

170 parameters:
171     network: NeuralNetwork3 in which the NetworkLayer is added
172     neuron_count: number of neurons in current layer
173     learning_rate: sets the learning rate for training (optional),
174                   default is 1.0
175     training_method: set the training method for current layer, None
176                   value averts training
177
178 attributes:
179     current_layer: denotes the index of the current layer
180     shape: denotes the dimension that the weights Matrix will assume
181     weights: weight matrix, used to map the preceding neurons linearly to
182             the current neurons
183
184 """
185 def __init__(self, network, neuron_count, learning_rate=1.0,
186             training_method=None):
187     self.network = network
188     self.current_layer = len(self.network.layers)
189     self.neuron_count = neuron_count
190     self.usenovelty = False
191     self.learning_rate = learning_rate
192     self.training_method = training_method
193     self.history_buffer = deque(maxlen=1)
194     self.history_size = 1
195     # self.estim_window = Buffer(140, maxlen=140)
196     self.averages = np.zeros((neuron_count, 1))
197     # initialize the weights
198     self.init_weights()
199     # self.elig_trace = Eligibility(self)
200
201     self.average_error = 0.0
202
203 def init_weights(self):
204     input_count = len(self.network.activations[self.current_layer])
205     self.shape = (self.neuron_count, input_count)
206     total_neurons = sum(self.shape)
207     range_ = np.sqrt(self.network.WEIGHT_SCALE / total_neurons)
208     self.weights = np.random.uniform(-range_, range_, self.shape)
209
210 def propagate_signal(self):
211     forwarded_activation = self.weights @ self.network.activations[self.
212         current_layer]
213     magnitude = self.network.output_noise_mag
214     noise = np.random.uniform(-magnitude, magnitude, forwarded_activation
215         .shape) if not self.usenovelty else 0
216     forwarded_activation = NeuralNetwork3.ACTIVATION(forwarded_activation
217         ) + noise
218     self.network.activations[self.current_layer + 1] =
219         forwarded_activation
220
221 def train(self):
222     if self.usenovelty:

```

```

215     error = self.train_several()
216     self.errorbuffer.append(error)
217     self.average_error = np.mean(self.errorbuffer)
218     else:
219         return self.training_method(self, self.network.modulation_signal.
                modulations[0])
220
221 def train_several(self):
222     error = 0.0
223     samples = [] #linkedList
224     for occupant in iter(self.novelty_buffer):
225         samples.append(occupant.experience)
226     sample_count = min(self.samples_per_tick, len(self.novelty_buffer))
227     for i in range(sample_count):
228         index = np.random.randint(len(samples))
229         sample = samples.pop(index)
230         self.network.activations[0] = np.array(sample)
231         self.network.propagate_signal()
232         self.update_averages()
233         error += self.training_method(self)
234     error /= sample_count
235     return error
236
237 def train_recent(self):
238     return
239
240 def update_averages(self):
241     exponent = 1.0 / self.history_size
242     decay = 0.01**exponent
243     neurons = np.reshape(self.network.activations[self.current_layer +
                1], (-1,1))
244     self.averages = (decay*self.averages) + (1.0-decay)*neurons
245
246 """
247 import agent_functions as funs
248 # import novelty
249
250 hyper_parameters = {"num_nodes": [4, 12, 24, 1],
251                     "activation_functions": [funs.linu, np.tanh, np.tanh,
                np.tanh],
252                     "training_functions": [funs.hebbian2, funs.hebbian2,
                funs.hebbian2],
253                     "learning_rates": [1.0, 1.0, 1.0],
254                     "output_noise": [0, 0, 0, 0.1],
255                     "weight_noise": [0, 0, 0.1],
256                     "buffers": [novelty.NoveltyBuffer(500),None,None]}
257
258 class Agent:
259     def __init__(self, hyper_params):
260         self.tag = hyper_params.get('tag', 'Not Specified')
261         self.num_layers = len(hyper_params["num_nodes"]) - 1
262         self.num_neurons = hyper_params["num_nodes"]
263         self.activation_funs = hyper_params["activation_functions"]

```

```

264     self.training_funs = hyper_params["training_functions"]
265     self.learning_rates = hyper_params["learning_rates"]
266     self.buffers = hyper_params["buffers"]
267     self.output_noise = hyper_params["output_noise"]
268     self.weight_noise = hyper_params["weight_noise"]
269     self.weight_cap = 10
270     self.weights = []
271     for i in range(self.num_layers):
272         randomw = np.random.randn(self.num_neurons[i+1], self.num_neurons[i
273             ])
274         self.weights.append(randomw)
275     print('Initialized Agent: ' + self.tag + '\n')
276     def propagate(self, experience):
277         input_activation = np.array(experience)
278         activations = []
279         activations.append(self.activation_funs[0](input_activation))
280         for i in range(self.num_layers):
281             activations.append(funs.layer_output(self.activation_funs[i+1],
282                 self.weights[i],
283                     activations[i], noise_magtd=
284                         self.output_noise[i+1]))
285         if self.buffers[i] is not None:
286             self.update_novelty(activations[i], self.buffers[i])
287         return activations
288     def update_novelty(self, experience, novelty_buffer):
289         input_sample = np.array(experience)
290         new_occupant = novelty.NoveltyOccupant(experience=input_sample)
291         new_distances = novelty_buffer.compute_new_distances(new_occupant)
292         if len(novelty_buffer) == novelty_buffer.history_size:
293             least_novel = novelty_buffer[0]
294             if new_occupant.novelty_score > least_novel.novelty_score:
295                 novelty_buffer.remove_novelty(least_novel)
296                 novelty_buffer.add_novelty(new_occupant, new_distances)
297             else:
298                 novelty_buffer.add_novelty(new_occupant, new_distances)
299     def train(self, modulation, activations):
300         for i in range(self.num_layers):
301             layer_parameters = {'weights': self.weights[i],
302                 'weight_noise': self.weight_noise[i],
303                 'learning_rate': self.learning_rates[i],
304                 'buffers': self.buffers[i],
305                 'modulation': modulation}
306             if self.buffers[i] is not None:
307                 experiences = funs.select_several(layer_parameters)
308                 for experience in experiences:
309                     postsynaptic = funs.layer_output(self.activation_funs[i+1],
310                         self.weights[i],
311                             experience, noise_magtd=self.
312                                 output_noise[i+1])

```

```

311         delta_weights = self.training_funs[i](layer_parameters,
312         experience, postsynaptic)
313         self.weights[i] = np.clip(self.weights[i] + delta_weights, -
314         self.weight_cap, self.weight_cap)
315         layer_parameters['weights'] = self.weights[i]
316     else:
317         delta_weights = self.training_funs[i](layer_parameters,
318         activations[i], activations[i+1])
319         self.weights[i] = np.clip(self.weights[i] + delta_weights, -self.
320         weight_cap, self.weight_cap)
321
322 # -----
323 if __name__ == "__main__":
324     agent = Agent(hyper_parameters)
325     rand_exp = np.random.uniform(0,1, (4,))
326     activations = agent.propagate(rand_exp)
327     for i in range(agent.num_layers):
328         print("weight shape: ", agent.weights[i].shape)
329         print("activation shape: ", activations[i+1].shape, "\n")

```

### performance.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Sun Oct 14 21:05:14 2018
4
5  @author: ajdin
6  """
7  import numpy as np
8
9  from collections import deque
10
11 class Performance:
12
13     def __init__(self, rangefinder_group, center):
14         self.rangefinder_group = rangefinder_group
15         self.center_point = center
16         self.N_size = 5
17         self.gamma = 1.0
18         self.reward_history = deque(maxlen=self.N_size)
19         self.Reward = 0
20
21     def normlaps(self):
22         center2_last = tuple(np.subtract(self.rangefinder_group.car.lastpos,
23         self.center_point))
24         center2_current = tuple(np.subtract(self.rangefinder_group.car.center
25         , self.center_point))
26         cur_angle = np.arctan2(*center2_current[::-1])
27         last_angle = np.arctan2(*center2_last[::-1])
28         # anti-clockwise gives positive performance
29         angle_diff = np.degrees(cur_angle - last_angle)
30         if np.abs(angle_diff) > 90:

```

```

29     change = 360 - np.abs(angle_diff)
30     if angle_diff < 0.0:
31         return change
32     else:
33         return -change
34     return angle_diff
35
36 def objective_func(self):
37     params = (self.rangefinder_group.activations, self.rangefinder_group.
38               car.speed)
39     self.reward_history.append(self.local_reward(*params))
40     weighted_avg = np.multiply(np.array(self.reward_history), self.gamma
41                               **np.arange(len(self.reward_history)))
42     return np.mean(weighted_avg)
43
44 def local_reward(self, *params):
45     activations, speed = params
46     avg_distance = np.mean(1 - activations[5])
47     norm_speed = np.linalg.norm(speed)
48     local_reward = (1.3*avg_distance + 0.3*norm_speed)
49     return local_reward

```

### visualnet.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Sep 15 12:51:39 2018
4
5  @author: ajdin
6  """
7  import pygame
8  import numpy as np
9
10 RADIUS = 25
11 THICKNESS = 2
12 BLACK = (0, 0, 0)
13 WHITE = (255, 255, 255)
14
15 class NetworkVisualizer:
16     FONT = 'Arial'
17     FONT_SIZE = 18
18
19     def __init__(self, network, box):
20         self.network = network
21         self.box = box #(lower_x, lower_y, upper_x, upper_y)
22         self.neuron_descriptions = []
23         self.connection_descriptions = []
24         self.layer_count = len(network.layers) + 1
25
26         self.set_up()
27         self.set_up_grid()
28         self.set_up_connections()
29

```



```

30     pygame.font.init()
31     self.font = pygame.font.SysFont(self.FONT, self.FONT_SIZE, bold=True)
32
33     def set_up(self):
34         for i in range(self.layer_count):
35             neuron_description = NeuronGroupDescription(self.network, i)
36             self.neuron_descriptions.append(neuron_description)
37
38     def set_up_grid(self):
39         y_values = np.linspace(self.box[1], self.box[3], self.layer_count+1,
40                               endpoint=False)[1:]
41         for i in range(self.layer_count):
42             self.neuron_descriptions[i].y_value = y_values[i]
43             neuron_count = len(self.neuron_descriptions[i].values)
44             self.neuron_descriptions[i].x_values = np.linspace(self.box[0],
45                                                                self.box[2], \
46                                                                neuron_count+1,
47                                                                endpoint=
48                                                                False)[1:]
49
50     def set_up_connections(self):
51         for i in range(self.layer_count - 1):
52             connection_description = ConnectionDescription(self.network, i)
53             connection_description.input_neuron_description = self.
54                 neuron_descriptions[i]
55             connection_description.output_neuron_description = self.
56                 neuron_descriptions[i+1]
57             self.connection_descriptions.append(connection_description)
58
59     def visualize(self, display):
60         self.draw_connections(display)
61         self.draw_neurons(display)
62         self.draw_legend(display)
63
64     def draw_neurons(self, display):
65         for i in range(self.layer_count):
66             self.neuron_descriptions[i].draw(display, self.font)
67
68     def draw_connections(self, display):
69         for i in range(len(self.connection_descriptions)):
70             self.connection_descriptions[i].draw(display, self.font)
71
72     def draw_legend(self, display):
73         y_offset = 50
74         bar_length = 400
75         box_center = (self.box[0]+self.box[2]) / 2
76         x_start = box_center - 1/2*bar_length
77         x_end = box_center + 1/2*bar_length
78         color_positions = np.arange(x_start, x_end, 1)
79         cap = self.network.weight_cap
80         color_bar = np.linspace(-cap, cap, len(color_positions)-1 )
81         for i in range(len(color_bar)):
82             weight = color_bar[i]

```

```

77     color = [(cap+weight)/(2*cap)*x for x in (255, 80, 0)]
78     thickness = 2*weight/cap + 3
79     start_point = (color_positions[i], y_offset)
80     end_point = (color_positions[i+1], y_offset)
81     pygame.draw.line(display, color, start_point, end_point, int(
82         thickness))
83     if i==0 or (i+1)% 50 == 0 or i==398:
84         text = '%.1f' % weight
85         text_surface = self.font.render(text, False, BLACK)
86         display.blit(text_surface, start_point)
87     pass
88
89 class NeuronGroupDescription:
90
91     def __init__(self, network, values_index):
92         self.network = network
93         self.idx = values_index
94         self.y_value = None
95         self.x_values = None
96
97     @property
98     def values(self):
99         return self.network.activations[self.idx]
100
101     def draw(self, display, font):
102         color = BLACK
103         radius = RADIUS
104         thickness = THICKNESS
105         values = self.values
106         for i in range(len(values)):
107             position = (int(self.x_values[i]), int(self.y_value))
108             pygame.draw.circle(display, WHITE, position, radius, 0)
109             pygame.draw.circle(display, color, position, radius, thickness)
110             text = '%.1f' % values[i]
111             text_surface = font.render(text, False, BLACK)
112             rect = text_surface.get_rect(center=position)
113             display.blit(text_surface, rect)
114
115
116 class ConnectionDescription:
117
118     def __init__(self, network, layer_index):
119         self.network = network
120         self.idx = layer_index
121         self.input_neuron_description = None
122         self.output_neuron_description = None
123         self.cap = self.network.weight_cap
124
125     @property
126     def connections(self):
127         return self.network.layers[self.idx].weights
128

```

```

129 def draw(self, display, font):
130     for i in range(self.connections.shape[0]):
131         for j in range(self.connections.shape[1]):
132             in_index = j
133             out_index = i
134             weight = self.connections[i][j]
135             start_point = (self.input_neuron_description.x_values[in_index],
136                           self.input_neuron_description.y_value)
137             end_point = (self.output_neuron_description.x_values[out_index],
138                          self.output_neuron_description.y_value)
139             color = [(self.cap+weight)/(2*self.cap)*x for x in (255, 80, 0)]
140             thick = 2*weight/self.cap + 3
141             pygame.draw.line(display, color, start_point, end_point, int(
                thick))

```

## -2 Python Code for the cartpole implementation

### -2-1 Cartpole.py

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Nov 24 13:20:24 2018
4
5  @author: ajdin
6  """
7
8  import numpy as np
9  import gym
10 import time
11 import functions
12 import agent_functions as funs
13
14 from net3 import NeuralNetwork3 as nn3
15 from collections import deque
16
17 # make the environment
18 cartpole = gym.make('CartPole-v0')
19
20 # agent settings
21 max_episodes = 100
22 network = nn3(inputcount=4, hiddencount=7, outputcount=1)
23 N_size = 1
24 reward_history = deque(maxlen=N_size)
25 R = 0
26 old_R = 0
27 Reward = lambda queue : np.mean(queue)
28 p_reward = lambda state : -np.linalg.norm(np.multiply(state, np.array
    ([0.1, 0.1, 1.3, 0.1])))
29 gamma = 0.9
30 raw_rewards = []
31
32 # options

```

```

33 rendering = False
34
35 # data sampling
36 s1 = []
37 s2 = []
38 s3 = []
39 s4 = []
40
41
42
43
44 # discount rewards
45 def discount_rewards(r, gamma=0.9):
46     # Init discount reward matrix
47     discounted_reward= np.zeros_like(r)
48     # Running_add: store sum of reward
49     running_add = 0
50     # Foreach rewards
51     for t in reversed(range(0, len(r))):
52         running_add = running_add * gamma + r[t] # sum * y (gamma) + reward
53         discounted_reward[t] = running_add
54     return discounted_reward
55
56 #linearly varied weight adjustment
57 def lvwa(learning_rate, reward):
58     return learning_rate*(1-reward)
59
60 last_num_steps = None
61 for episode in range(1, max_episodes+1):
62     old_R = None
63     modulations = []
64     R = 0
65     state = cartpole.reset()
66     network.modulation_signal.modulations[0] = 0.0
67     network.layers[1].history_buffer.clear()
68     network.layers[1].history_size = 200
69     features = []
70     outputs = []
71     for n in range(200):
72         if rendering:
73             start = time.time()
74             cartpole.render()
75             preprocess_state = np.multiply(state, np.array([1, 1/4, 2, 1/6])) +
76                 np.array([0.5, 0.5, 0.25, 0.5])
77             reward_history.append(p_reward(state))
78             modulation = 0.0
79
80             if len(reward_history) == reward_history.maxlen:
81                 R = Reward(reward_history)
82                 if old_R is not None:
83                     modulation = np.tanh(2*(R - old_R))**(1-last_num_steps/200)
84                     modulation = modulation*1 if modulation > 0 else modulation
85                 else:

```

```

85     modulation = 0
86     # print(modulation)
87     old_R = R
88 #     print('modul: ', modulation)
89
90     network.add_experience(preprocess_state)
91     features.append(network.activations[1])
92     outputs.append(network.activations[-1])
93     modulations.append(modulation)
94
95     state, reward, flag, info = cartpole.step(int(round(network.
        activations[-1][0])))
96     network.modulation_signal.modulations[0] = modulation
97     network.train()
98     # print(np.mean(network.layers[1].weights))
99
100    network.layers[1].weights += funs.train(modulations, features,
        outputs, 1.0)
101    modulations.clear()
102    features.clear()
103    outputs.clear()
104
105
106
107    if rendering:
108        time.sleep(max(1./40 - (time.time() - start), 0))
109    if flag:
110        last_num_steps = (n+1)
111        if last_num_steps is not None:
112            mean_mod = 1 - last_num_steps/200
113        else:
114            mean_mod = 0.0
115
116        reward = last_num_steps/200
117        # mean_mod = np.mean(modulations)
118        modulations = [mean_mod for x in modulations]
119
120 #     modulation = lvwa(1.0, reward)
121 #     network.modulation_signal.modulations[0] = modulation
122 #     network.layers[1].history_size = last_num_steps
123
124 #     particular_reward = (n+1) - old_R
125 #     raw_rewards = [0.02*particular_reward]*(n + 1)
126 #     old_R = n+1
127 #     discounted = discount_rewards(raw_rewards, gamma=0.99)
128 #     network.modulation_signal.modulations[0] = discounted
129
130 #     print('episode: ', ' steps: ', n+1, ' history: ', network.layers
        [1].history_buffer)
131 #     print('final training: ', episode ,'\n\n')
132 #     print('num samples: ', len(network.layers[1].history_buffer),
        discounted)
133

```

```
134     print('episode: ', episode, '. Steps completed: ', n+1, ' avg error
      : ', '%.4f' % network.layers[0].average_error,
135         ' mean mod: %.3f' % mean_mod, 'max mod: %.2f' %0.0)
136 #     print('\n', ' discounted: ', discounted)
137 #     print(network.layers[1].history_buffer, '\n\n')
138
139 #     network.train()
140     reward_history.clear()
141     break
142 # network.train()
143
144 cartpole.close()
145
146
147 ### render the learned policy
148
149
150 cartpole = gym.make('CartPole-v0')
151 state = cartpole.reset()
152 for i in range(200):
153     start = time.time()
154     cartpole.render()
155     preprocess_state = np.multiply(state, np.array([1, 1/4, 2, 1/6])) + np.
        array([0.5, 0.5, 0.25, 0.5])
156     network.add_experience(preprocess_state)
157     state, reward, done, info = cartpole.step(int(round(network.activations
        [-1][0])))
158     time.sleep(max(1./40 - (time.time() - start), 0))
159     if done:
160         break
161 print('done after ', i+1, ' steps.')
162 cartpole.close()
```